

- 10) Clock :- It is used to synchronize operations.
- 11) Reset :- Initializes all modules.
- * Word is collection of bit $8086 \rightarrow 8$ bit. $8086 \rightarrow 16/32/64$
No. of bits at a time can transferred from one

CHAPTER 4 - Cache Memory

4.1 Computer Memory system Overview

A memory is characterised by following parameters

Location :- It refers to whether memory is internal to or external to the computer.

* Internal memory include main memory, registers, cache memory whereas external memory refers to peripheral devices, disc, magnetic tape etc

2) Capacity :- For internal memory, the capacity is defined in terms of bytes / words whereas for external memory the capacity is always expressed in terms of bytes only.

3) Unit of transfer :- For internal memory, the unit of transfer is equal to the no. of electrical wire into and out of the memory module i.e. it deals with the data bus width.

Word :- It is a natural unit of organisation. A size of a word is typically equal to the no. of bits used to represent an integer and instruction length.

Addressable unit - It is the smallest location which can be addressed uniquely. If 'A' is the length of an address in ~~size~~ bits, then the no. of addressable unit is given as $N = 2^A$ ①

Ex - 8086 \rightarrow 20 no. of address line

Addressable unit $= 2^{20} \approx 1$ MB.

- 4) Access Method - There are four types:-
- Sequential access \rightarrow It starts at the beginning and proceeds in sequential order. Here the access time depends on the location of data and previous location. ex - magnetic tape.
 - Direct access \rightarrow Individual blocks have unique address hence here access is given by jumping into the respective block plus the sequential search. Here also, access time depends on the location & previous location. ex - 3000 2000 2020
 - Random access \rightarrow Individual access can exactly identify different location. Here as the access time is independent of the location & previous access
ex - RAM \rightarrow MOV Ax, [3000H]
 - Associative access \rightarrow Here data is located as a portion of the main memory. The access time is independent of the location and the previous location. ex - cache memory

5) Performance -

i) Access Time / Latency \rightarrow It is the time of time spent in presenting the address and getting the valid data.
ex - MOV Ax, [5000H] time required to access data.

Imp ii) Memory cycle time \rightarrow It refers to the access time + any additional time required before the next access

Imp iii) Transfer rate \rightarrow It is the rate at which data can be moved into and out of a memory unit.

For random access, the transfer rate is $\frac{1}{\text{cycle time}}$ ②

$$T_A = T_C + \frac{n}{R}$$

where $T_n \leftarrow$ Average time to read or write n bits
 $T_A \leftarrow$ Is the average access time
 $n \leftarrow$ No. of bits
 $R \leftarrow$ Is the transfer rate in bits per second (bps)

6) Physical types of memory -

- (i) Semiconductor memory (RAM, ROM, ...)
- (ii) Optical memory (CD, DVD) Compact disc
↳ faster, write on both side - more storage
- (iii) Magnetic memory (Magnetic disk, Magnetic tape)

7) Physical characteristics -

- (i) Volatile -
- (ii) Non-volatile -

8) Organisation - It refers to the physical arrangement of bits in words. It is not frequently used.

Memory Hierarchy

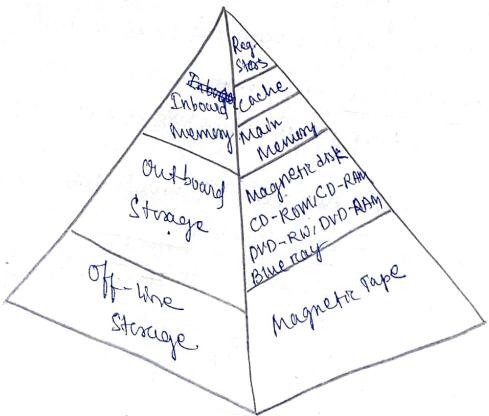


fig:- Memory Hierarchy

The following relationship holds w.r.t capacity, time and cost.

↳ Faster access time \rightarrow Greater cost per bit

↳ Greater capacity \rightarrow smaller cost per bit

↳ Greater capacity \rightarrow slower access time

To remove this dilemma we are using a memory hierarchy as we move down the hierarchy the following occurs.

↳ Decreasing cost per bit

↳ Increasing capacity

↳ Increasing access time

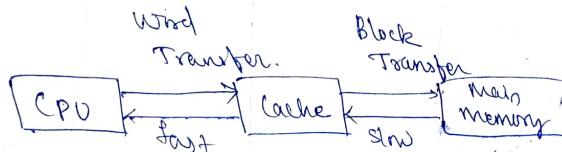
↳ Decreasing frequency of access of the memory by the processor
(this principle is known as locality of reference) nearby to the processor.

Main memory is the principle internal memory system of the computer. Each location in main memory has one unique address. Main memory is usually extended with a higher speed and smaller cache.

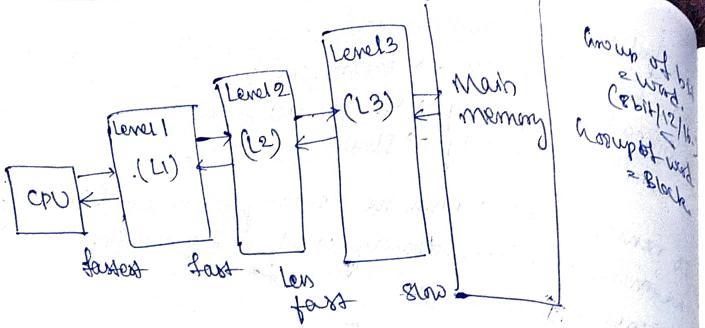
9.2 Cache Memory principle.

Cache memory is designed to combine the memory access time of expensive, high speed memory, combined with large memory size of less expensive and lower speed memory.

Cache is small amount of fast memory located on CPU chip or module. It usually lies between main memory and CPU.



↳ Single cache

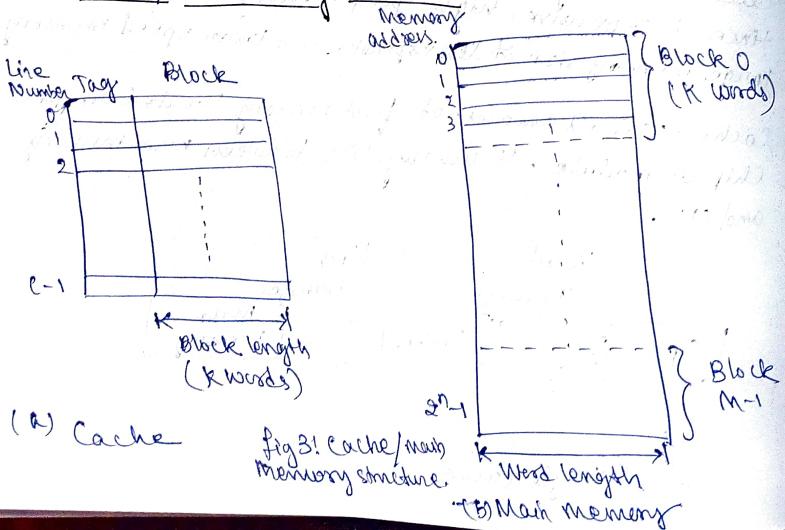


(b) 3-level cache Organization.

Fig 2:- Cache and Main memory

When the processor attempts to read a word of memory, a check is made to determine if the word is present in the cache. If so, the word is delivered to the processor. If not, a block of main memory consisting of some fixed number of words is read into the cache and then delivered to the processor.

Cache / Main memory Structure



Main memory consists of upto 2^n addressable words with each word having one unique n-bit address. The memory consists of a number of fixed length blocks of K words each where we have $\frac{2^n}{K}$ blocks in main memory.

The cache consists of small m blocks called lines ($0, 1, 2, \dots, l-1$). Each lines contains K words along with few tag bits and control bits. The length of a line not including the tag and control bits is known as the line size.

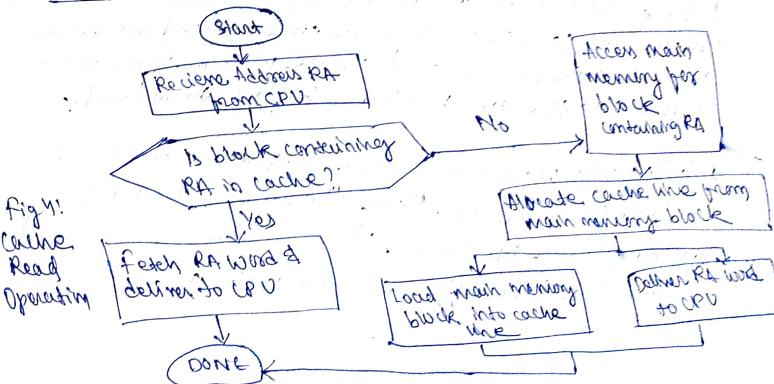
Line = Tag + Control + Block

Line size = Block

The number of lines is very less than the number of main memory blocks i.e. $m \ll M$.

At any time, some subset of the block of memory resides in lines in the cache. If a word in a block of memory is needed, that block is transferred to one of the lines of the cache since there are more blocks than lines and individual blocks cannot be uniquely and permanently dedicated to a particular block hence each line includes a tag that identifies which particular block is currently being stored.

Cache read Operation



The processor generates Read address for a word if the word is contained in the cache then it is delivered to the processor otherwise the block containing the word is loaded into the cache and the same word is delivered to the processor simultaneously.

Cache organization

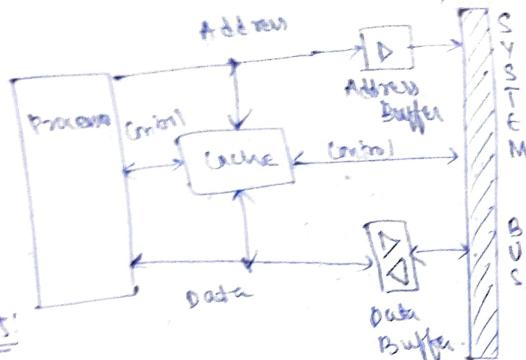


fig:

The cache is connected to the processor by means of data, control and address lines. The data and address lines are attached to the data and address buffers through the separators to reach the main memory.

i) 'Cache hit' occurs when the required data is available in the cache. In this case communication is direct between the processor and cache with no system bus traffic. Also the data and address buffers are disabled.

ii) 'Cache miss' occurs when the requested data is not available in the cache. Here, the desired address is loaded into the system bus. Then the data are returned through the data buffer to both cache and the processor.

4.3 Elements of Cache design

a) Cache Address - logical cache + Physical cache

b) Cache size *

c) Mapping function - Direct mapping, Associative, Set associative

d) Replacement Algorithm - LRU (least recently used)

- FIFO (first in first out)

- LFU (least frequently used)

- Random

e) Write Policy - write through, write back

f) Line size

g) Number of caches

Main memory
Management Unit

Cache Address

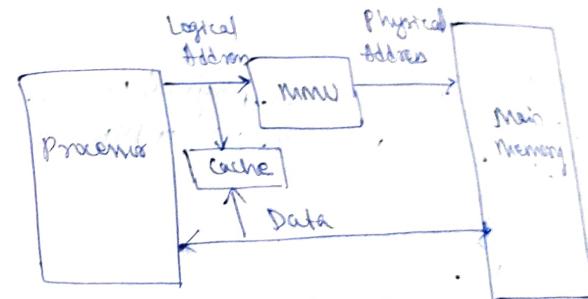
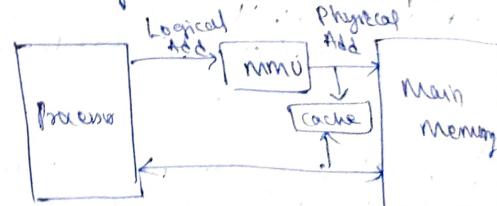


fig 6: (a) Logical Cache / Virtual cache



(b) Physical Cache / Virtual Cache

Virtual memory is a facility that allows the programs to address memory from a logical point of view irrespective of the actual amount of main memory physically available. When virtual memory is used, the address field of machine instructions contains the virtual addresses.

For reads or writes from main memory, a hardware memory management unit (MMU) translates each virtual address into a physical address of main memory.

A logical cache also known as virtual cache stores data using virtual addresses. The processor then access the cache directly without going through the MMU.

A physical cache stores data using main memory physical addresses.

Advantage:
Logical cache is faster than the physical cache because the cache can respond before the MMU performs an address translation.

Virt. to Phys. Mapping function

Since there are few number of cache lines than the available main memory blocks hence an algorithm is needed for mapping main memory blocks into cache lines.

1) Direct Mapping:

In this mapping each block of main memory maps into only one possible cache line. The mapping is expressed as $b \equiv \text{Modulo } m$

where

$i \leftarrow$ Cache line no.

$j \leftarrow$ Main memory block number

$m \leftarrow$ Number of lines in the cache

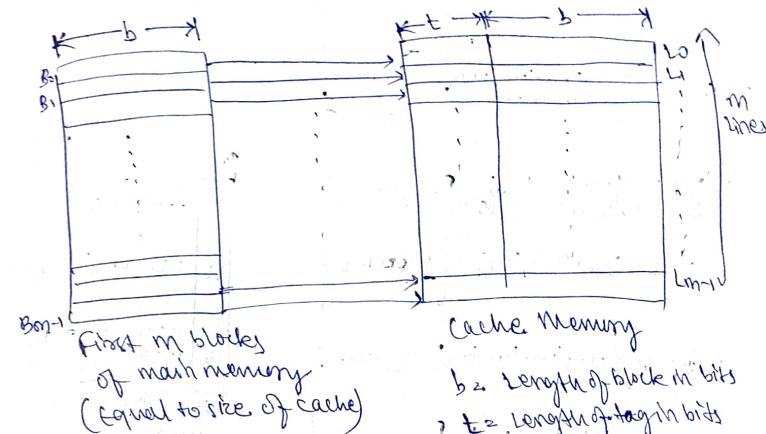
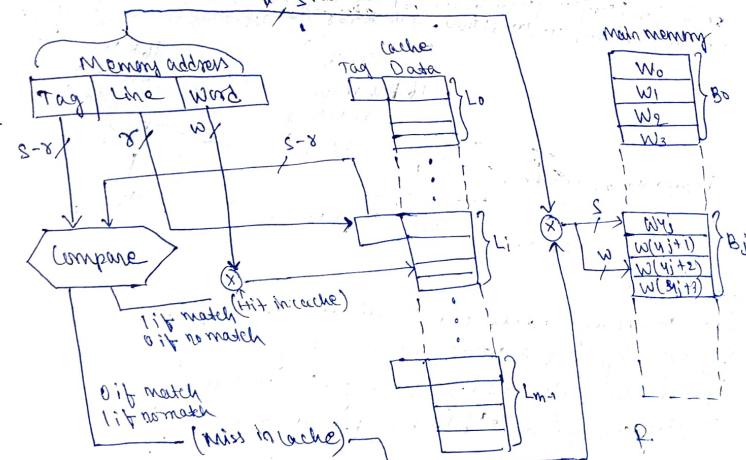
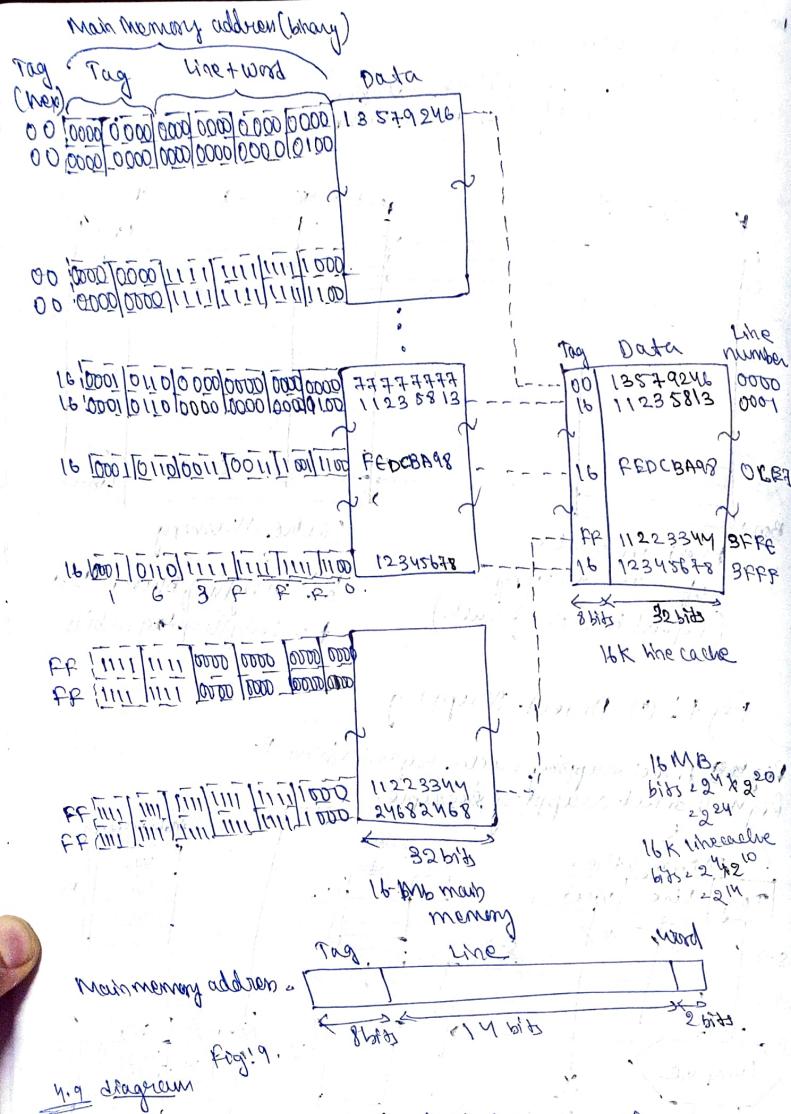


Fig 7: (1) Direct Mapping

Fig 7.8 Direct mapping cache organization





Each memory address can be divided into 3 fields.
The least significant 8 bit can uniquely identify a word or a byte within a block of main memory. The remaining 5 bits specify one of the 2^5 blocks of main memory.

The s bits are divided into two fields i.e. ($s-w$) bits of tag field and w bits of the field.

- \Rightarrow Address length = $(s+w)$ bits.
- \Rightarrow No. of addressable units = 2^{s+w} words / bytes
- \Rightarrow Block size = LRU size = 2^w words
- \Rightarrow Number of blocks in main memory = $\frac{2^{s+w}}{2^w} = 2^s$
- \Rightarrow No. of lines in cache = $2^s = m$ lines.
- \Rightarrow Size of cache = 2^{s+w} words / bytes [excluding the tag field]
- \Rightarrow Size of tag = $(s-w)$ bits

Cache Memory

~~each sentence~~ each sentence is 16 words

Q Let cache size = 16 words, block size = 4 words

Determine the ~~bias~~ no. of lines in the cache.

1990-1991

No. of bits used = 16 - No. of bits

No. of bits used to represent number size = 32

Time is finite. $\therefore \frac{16 \text{ words}}{\text{min}} = y$

No. of letters in each word

Q) Let us consider a 64 words of main memory is mapped into 16 word of cache memory using direct mapping function assume a block size of 4 words

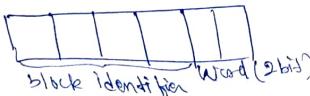
Mem memory = 64 words.

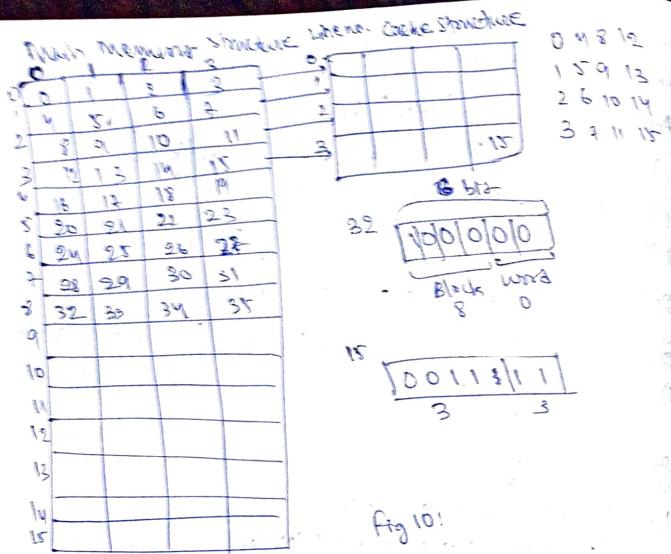
$$\text{No. of addressable lines} = \log_2 64 = \log_2 2^6 = 6 \log_2 2 = 6.$$

Physical address of data = 

given block size of M words

No. of bits used to represent/identify each word = 2 bit.





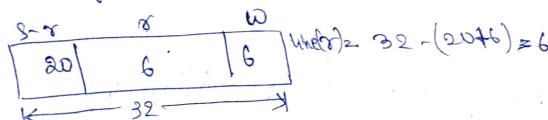
Q) Consider a memory system that uses a 32 bit address to address at the byte level, & a cache that uses 8 byte line size. Assume a direct mapped cache with a tag field in the address of 20 bits. Show the address format and determine the no. of addressable units, no. of blocks in main memory, no. of lines in cache and size of the tag.

Given - The address length ≥ 32 bits \Rightarrow SW ≥ 32 bits

$$\text{Cache size} = \text{Block size} = 64 \text{ byte}$$

$w = 6$

Given tags $\{ \pm 20 \text{ bits} \}$



Q) Consider a direct mapped cache of size 32 KB, with block size of 32 bytes. The CPU generates 32 bit addresses. Identify the no. of bits needed for cache indexing and the no. of tag bits required.

Given Cache size = 32 kB

Block size = 32 Bytes

$$\text{No. of Blocks} = \frac{32KB}{1KB} = 32$$

bits for cache indexing / No. of lines (r) = 10



4.10 Explanations

2/4/24

Disadvantage of direct Mapping

Disadvantage of direct mapping
In case of direct mapping multiple blocks of main memory are assigned to lines of cache as follows.

		Main memory block assigned			Line no	
Address	Block	Block Address	Block Length	Block Address	Block Length	Line no
0	0	M1	3M+1	0 → 0	4 8 12	0
1	1	M2	3M+1	1 →		1
2	2	M3	3M+1			2
3	3					3
⋮	⋮	⋮	⋮			
M-1	M	3M+1	3M+1			

Accent more
C. as he hit

In direct mapping there is a fixed cache location for ~~any~~ any given block. Thus if a program repeatedly refers to two different blocks that map into the same line the blocks will be continuously swapped in the cache and the hit ratio will be low.

This phenomenon is also known as thrashing also the miss penalty will be high.

To reduce the miss penalty we need to store the discarded data 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. For this we are using the associative mapping.

Associative Mapping

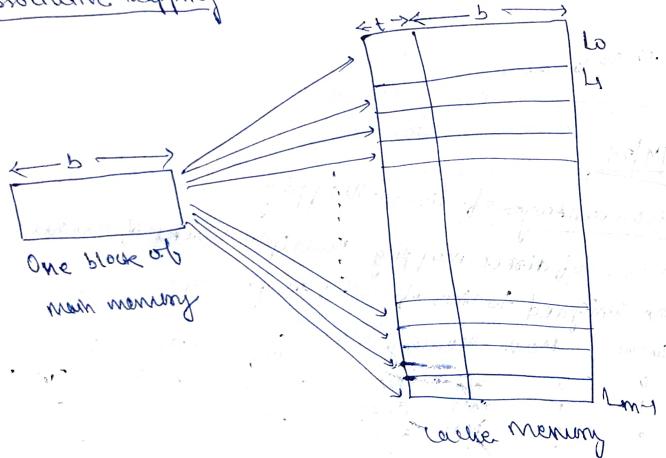


Fig 11: Associative Mapping.

Fig 12: Associative mapping maps each main memory block to any line of the cache as shown in fig 11. In this case a memory address consists only a tag and a word.

field. The tag field uniquely identifying a block of main memory. To determine whether a block is in the cache, the cache control logic must examine the every line tag for a match as shown in fig 12.

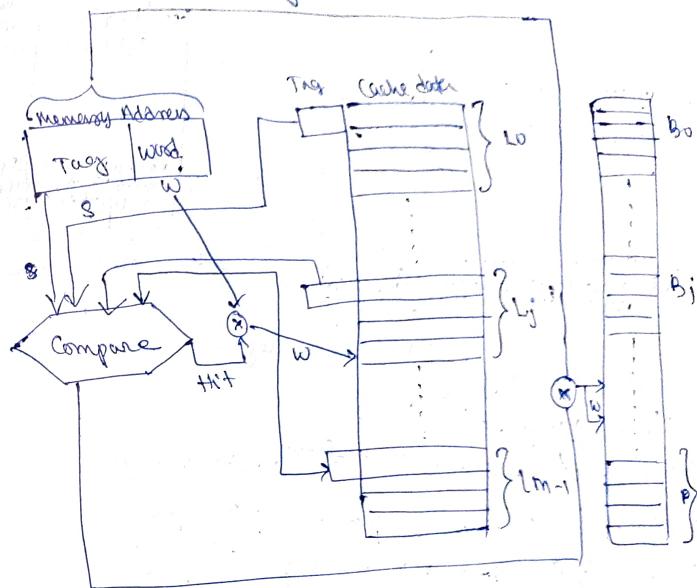


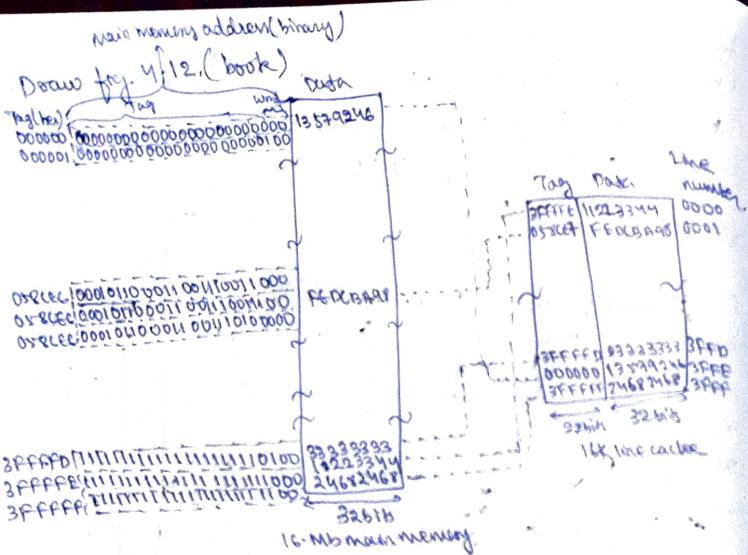
Fig 12: Associative Mapping

Let us consider the example for associative mapping. A main memory address consist of 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.

Given a 32-bit hexadecimal address 16339C, identify the tag to be stored.

0001 0110 0011 0011 1001 1100
Tag word

Tag 0001011000110011100111 058 EET 0 BC E 7 0



for associative mapping

- 7) Address length = $(S + W)$ bits
 - 8) No of addressable unit = $2^{(S+W)}$ words or bytes
 - 9) Block size = Line size = 2^W words or bytes
 - 10) No of blocks in main memory = $\frac{2^{S+W}}{2^W} = 2^S$
 - 11) No of lines in cache = Not Defined
 - 12) Size of tag is bits

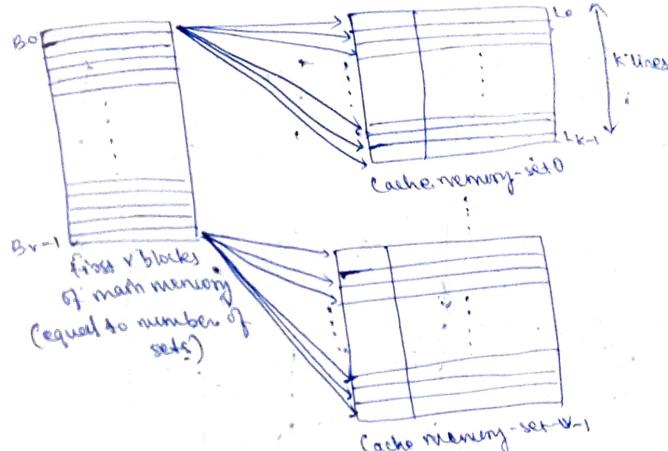
8/11/24
3) Set associative mapping

fig 4.13, 4.14 and 4.15 to be drawn

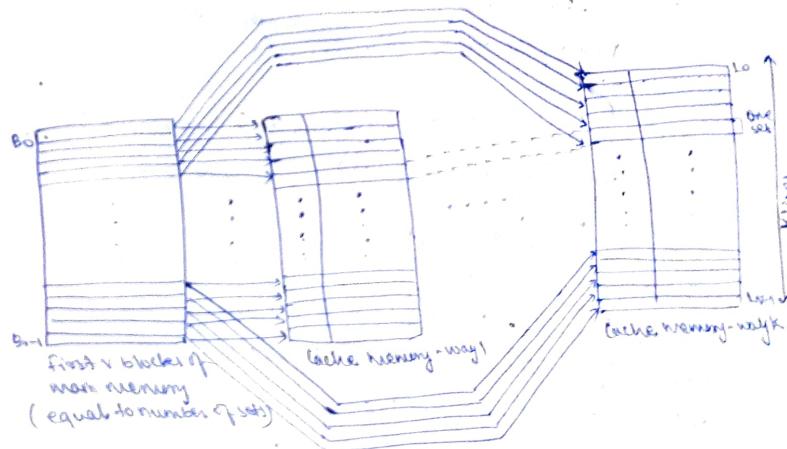
It is a combination of both the direct and associative mapping. In this case the value consist of no. of sets each of which consists of no. of lines.

The relationship is given by $m \equiv k \pmod{v}$ and $i \equiv j \pmod{v}$

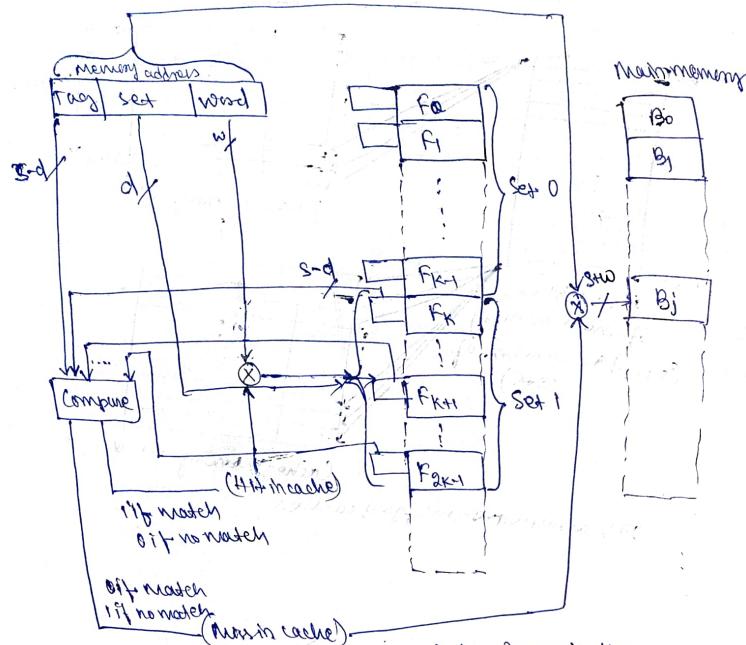
Where, i_2 = cache set number
 i_1 = main memory block no.



(a) V architecture-mapped cache



(*) K direct-mapped cache.



K-way set-associative Cache Organization

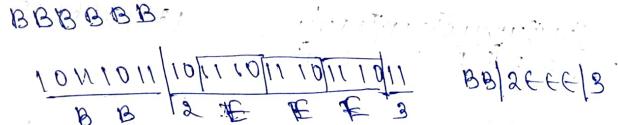
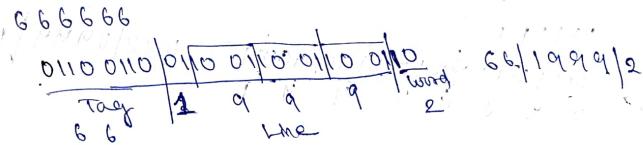
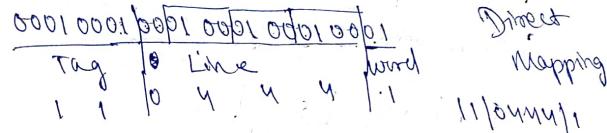
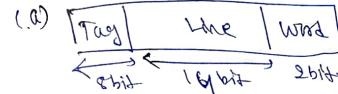
Here the memory address is referred as a combination of 3-fields i.e. Tag, set and word.

Tag	Set	Word
$s-d$	d	w

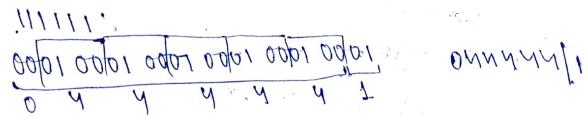
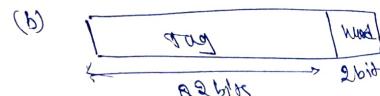
- 1) Address length = $s-d + d + w = (s+w)$ bits
- 2) No. of addressable units = $2^{(s+w)}$ words or bytes
- 3) Block size = Line size = 2^w words or bytes
- 4) No. of blocks in main memory = $\frac{2^{(s+w)}}{2^w} = 2^s$
- 5) No. of lines in cache set = K
- 6) Size of tag: Number of set = $V \cdot 2^{d-2}$
- 7) No. of lines in cache = $M = K \cdot V = K \cdot 2^d$

8) Size of cache = $K \cdot V \cdot 2^{d-w}$ words or bytes

Q. 4/2
Soln:- Address: 111111, 666666, BBBB BBBB



Associative cache mapping.



666666

0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9

199999|2

BBBBBBB

1	0	1	1	0	1	1	0	1	1	0	1	1	0	1	1	0	1	1
2	E	E	E	E	E	E	E	E	E	E	E	E	E	E	E	E	E	E

, 2FFFFFF|3

(c) Tag / set / word

Set Associative Mapping

9	13	12
Tag	Set	Word

111111

0	0	0	1	0	0	0	1	0	0	0	1	0	0	0	1	0	0	0
0	2	2	0	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4

022/04444|1

BBBBBBB

1	0	1	1	0	1	1	0	1	1	0	1	1	0	1	1	0	1	1
1	7	7	10	E	E	E	E	E	E	E	E	E	E	E	E	E	E	E

177/0FEE|3

666666

0	1	0	0	1	0	0	1	0	0	1	0	0	1	0	0	1	0	0
0	8	8	1	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9

0CC/1999|2

(d) Given Address length 16 bit

(a)

$\Rightarrow 2^{16-3} = 13$

block size of 8 bytes \Rightarrow block size $= 2^W = 8 \Rightarrow W = 3$

$\Rightarrow S = 16 - 3 = 13 \Rightarrow S = 13$

No. of lines in cache = 32 lines

$$\Rightarrow 2^8 = 32$$

$$\Rightarrow \lceil \sqrt{25} \rceil \text{ line}$$

size of the tag = $s - w = 8$

Tag	Line	Word
8	5	3

Line no
0001 0001 0001 1011 → 3
0000 0001 0011 0100 → 6
1101 0000 0001 1001 → 3
1010 1010 1010 1010 → 21 ← binary of 15 in hexa

(e) 0001 1010 0001 1010

The seven different addresses are

0001 1010 0001 1010

3 bits: 000
001
010
011
100
101
110
111

000
001
011
100
101
110
111

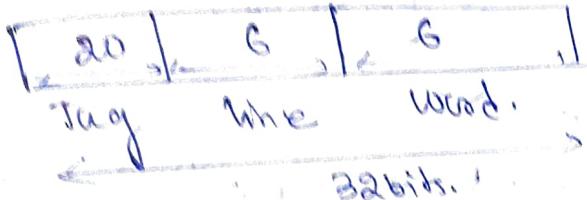
Because the two data or lines with two different main memory locations (addresses) can be stored in the same place for the same line number of the cache hence the tag is used in the cache to identify or distinguish the data

(d) $2^{16} \leftarrow$

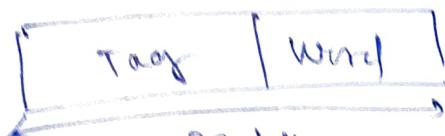
4.1)

a) Given, Tag = 20 bits

Cache = 64 byte

Line size = ~~32~~³² bits, $G = 8$ No. of addressable unit = ~~32~~³² = 2^{32} bytes.

b)



32, b4,

Tag size = $32 - 6 = 26$ No. of addressable units = 2^{26} bytes

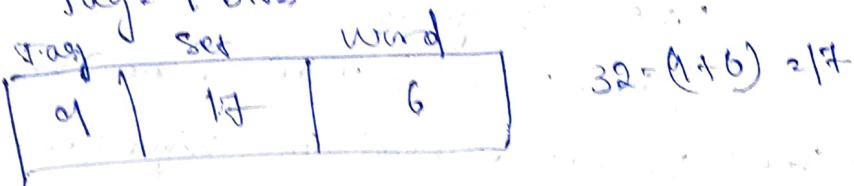
No. of line in the cache = No defined.

No. of Tag bits = 26.

4) Mux-set-associative Mapping

 $K = 4$

Tag = 9 bits.

No. of line in the cache set = $K = 4$ No. of sets in the cache = $2^4 = 2^4 \cdot 2^{10} = 128 K$

No. of tag bits = 9 bits.

S-d = 32
d = set
w = word

$$K = \frac{64}{4} = 16$$

$$64 = K \times V = 16 \times 4$$

$\frac{1}{2}, \frac{2}{2}, \frac{3}{2}, \frac{4}{2}, \frac{5}{2}, \frac{6}{2}$
 $\frac{7}{2}, \frac{8}{2}$

u-1

The cache is divided into 16 set of 4 line each
No. of bits for set = $\log_2 16 = 4$

$$\text{Main memory} = 4 \text{K} \times 8^10 = 2^2 \times 2^{10} = 2^{12}$$

$$\text{tag} + \text{set} = 12 \text{ bits}$$

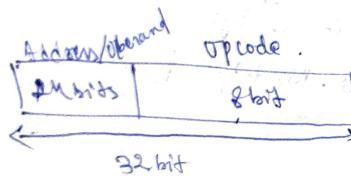
$$\text{tag} = 12 - 4 = 8 \text{ bits}$$

hence, 128 words

$$\begin{aligned} \text{No. of bits to identify word} &= \log_2 128 \\ &\approx \log_2 2^7 = 7 \end{aligned}$$

Chapter-3

3-3



a) Maximum achievable capacity $\approx 2^{24} \times 16 \text{ MB}$.

b) i) 32-bit local address bus + 16-bit local data bus.
 $1+2 = 3$ cycles

Instruction and data transfer would take 3 bus cycles
each 1 for the address and 2 for the data.

If the address bus is 32 bits, the whole address can
be transferred to memory at once however since the
data bus is only 16 bits, it will require 2 bus cycles
to fetch the 32-bit instruction or operand.

ii) 16-bit local address bus + 16-bit data bus
 $8+2 = 10$

instruction and data transfer would take 4 bus cycles each 2
for the address and 2 for the data.

$$\text{d) } \text{PC} = 2^4$$

the PC need 4 bits for 16-bit address and the IR needs 32
bits for 32-bit addresses.

3-4. processor size is 16 bits

$$\text{so addressable bit} = 2^{16}$$

(b) The max memory address $\approx 2^{16} \times 2^6 \times 2^{10} = 64 \text{ KB}$.

$$\text{p) } 2^{16}$$

c) ~~we need~~ Separate I/O instructions are needed because
during its execution it will generate separate I/O signals these
signals will be different from the memory signals which are
generated during the execution of memory instructions.
Therefore one more op. $\#$ pin will be needed to carry
the I/O signals.

$$S = M \mid \overline{D} \mid \boxed{16}$$

9-1. M data

$$S=0 \quad \overline{D}=0=1 \cdot \overline{D} \text{ data.}$$

$$\text{d) No. of I/O ports} \approx 8 \text{ bits}$$

Hence it can support 2^8 I/O ports.

3-5 microprocessor size ≈ 32 bit
data bus size $= 16$ bit