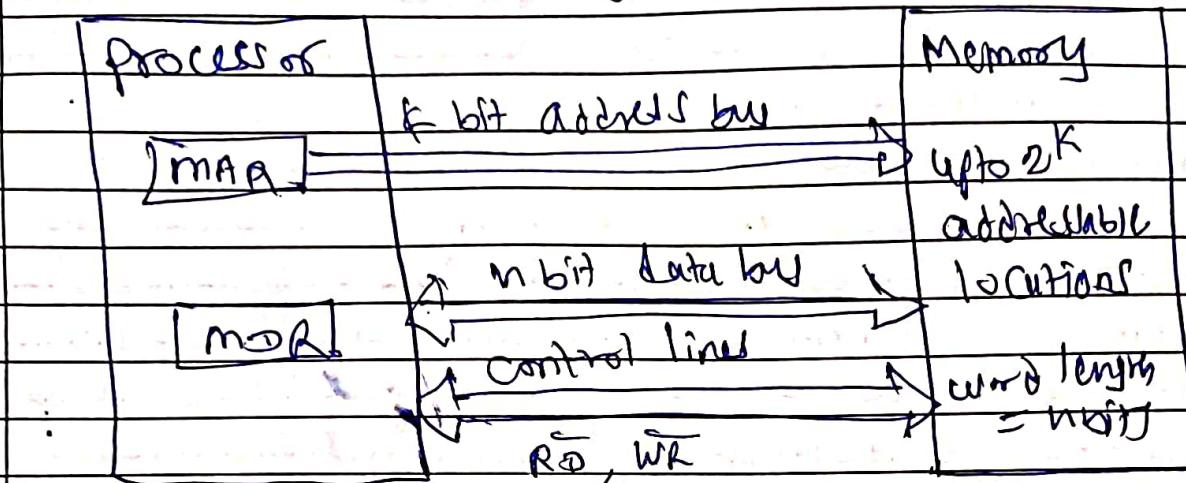


## UNIT IV Memory Organization



connection between memory & processor.

### Characteristics of Memory System

- 1 Location — CPU — CPU register; Cache memory  
internal — main memory  
External — Secondary memory - magnetic disk, I/O controller.

### 2 Capacity —

Word size — 8, 16 & 32 bits

No. of word — memory capacity is

$4K \times 8$  word size is 8. & No. of words are  $4K = 4096$

- 3 Unit of Transfer — It is maximum number of bits that can be read or written into the memory at a time.

- (1) word
- (2) block

#### 4. Access Method -

The order or sequence in which information can be accessed.

① Random Access - If storage locations can be accessed in any order and access time is independent of the location being accessed. The access method is known as random access. (RAM)

② Serial Access - If storage locations can be accessed only in a certain predetermined sequence, the access method is known as serial access.  
e.g. magnetic disk or tapes, CD-Roms.

#### V Performance -

① Access Time - It is time taken by memory to complete read/write operation from the instant that an address is sent to the memory.

② Memory cycle time -

③ Transfer rate - It is defined as the rate at which data can be transferred into or out of a memory unit.

#### VI Physical characteristics -

① Volatile / non volatile -

If memory can hold data even if power is turned off, it is called nonvolatile memory. Otherwise it is called volatile memory.

② Erasable / Non erasable - The memories in which data is once programmed cannot be erased are called as nonerasable memories. If data is the memory is erasable then memory is called erasable memory.

Comparison between serial & Random access memories.

Serial access memories

Random access memories

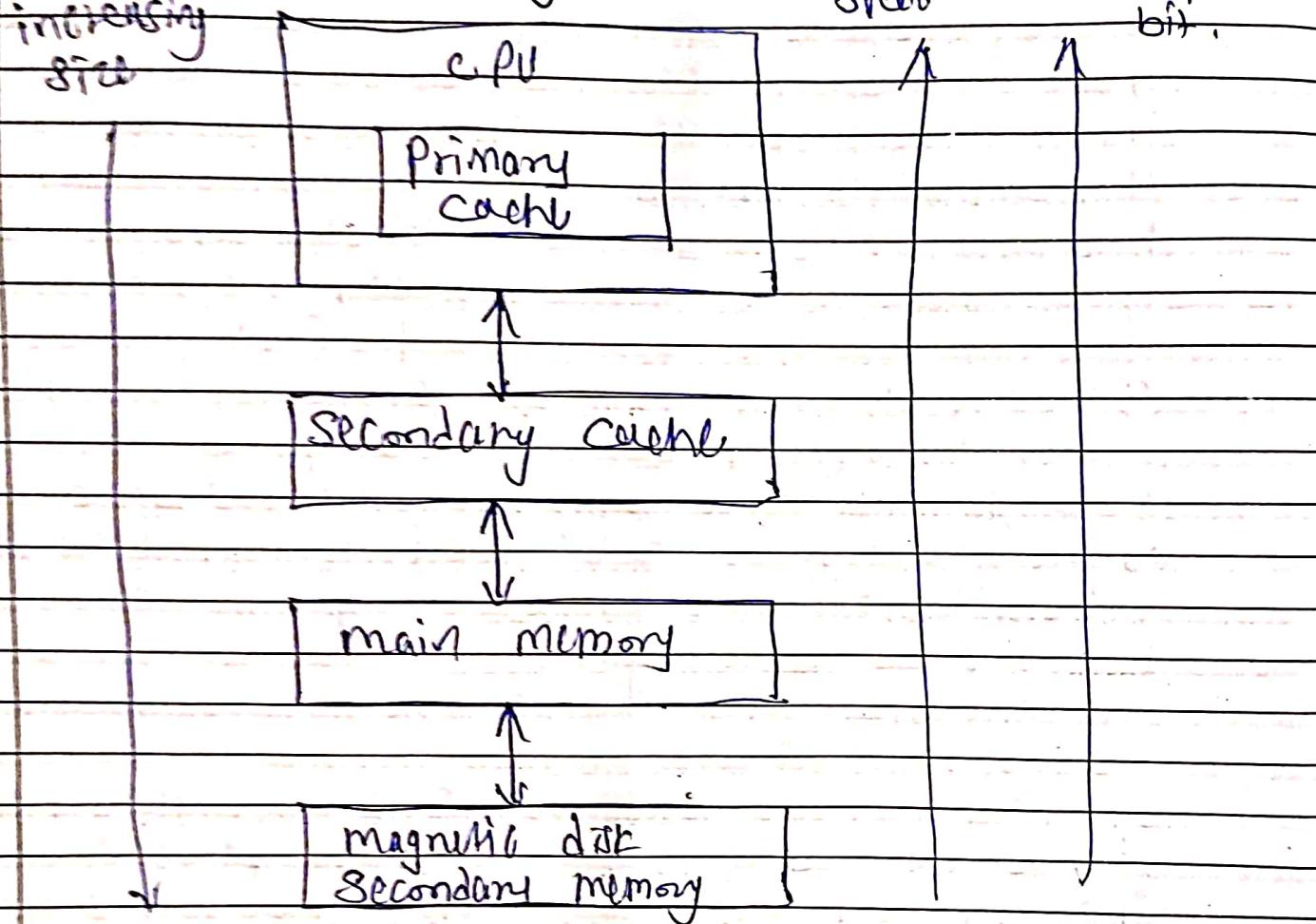
- |   |   |
|---|---|
| (1) use serial access method  | (2) use random access method  |
| (2) memory is organized into units of data, called records.   | (3) each storage location in the memory has an unique address. It can be accessed independently of the other locations. |
| (3) memory access time is dependent on the position of storage location.  | (4) memory access time is independent of storage location being accessed.   |
| (4) The time required to bring the desired location into correspondance with a read-write head increases effective access time. slower. | (5) memory access time is less.   |
| (5) cheaper than random access memories   | (6) Random access memories are comparatively costly.  |
| (6) Nonvolatile memories  | (7) may be volatile or nonvolatile depending on physical char.  |
| (7) magnetic tape is an example of serial access memory.  | (8) Semiconductor memories are random access memories.  |

Memory Hierarchy

Memory  
increasing  
size

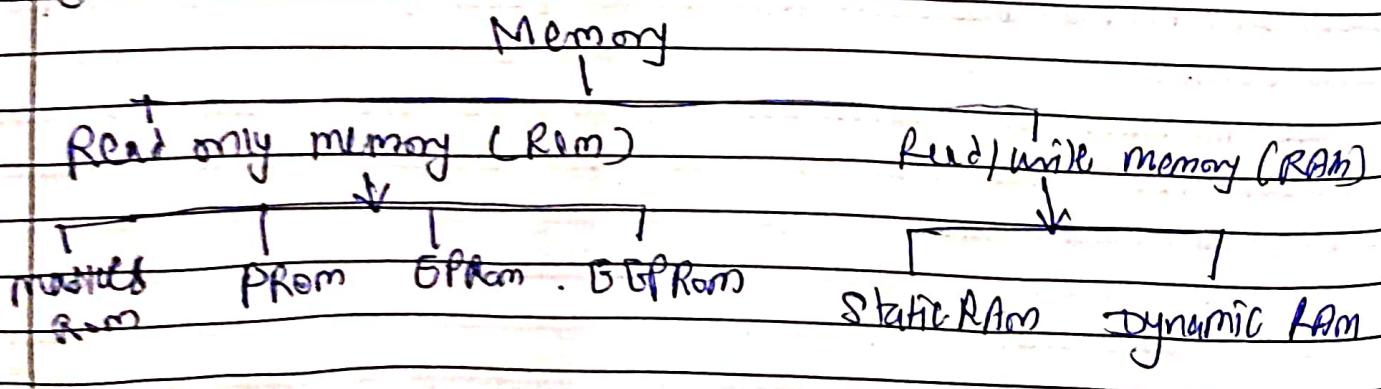
Hierarchy

increasing  
speed  
cost per  
bit.



Main Memory organization

main memory consist of frames supported  
with SRAM cache.



## I. ROM - Read only memory

It is a read only memory. It is non volatile memory i.e. it can hold data even if power is turned off. Rom is used to store the binary codes for the sequence of instructions you want the computer to carry out and data such as look up tables. There are 4 types of Rom.

1. Masked Rom

2. PROM

3. EPROM

4. EEPROM.

## II. PROM - Programmable Read only memory.

III. EPROM - Erasable Programmable Read only memory

IV. EEPROM - Electrical Erasable Programmable Read only memory.

2. RAM - It is called read / write memory. It is volatile memory i.e. it cannot hold data when power is turned off. There are 2 types of RAM:

1. STATIC RAM

2. DYNAMIC RAM

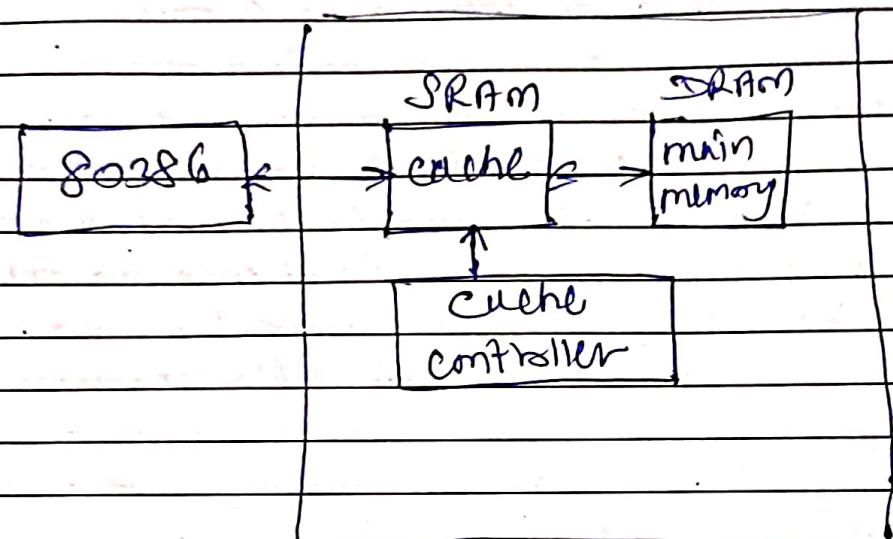
## Comparison between SRAM & DRAM

### Static RAM

### Dynamic RAM

- |   |  |
|---|--|
| <ul style="list-style-type: none"><li>① Static RAM contains less memory cells per unit area.</li><li>② It has less access time hence faster memories.</li><li>③ Static RAM consists of number of FF's - each FF stores one bit.</li><li>④ Refreshing circuitry is not required.</li><li>⑤ Cost is more.</li></ul> | <ul style="list-style-type: none"><li>① Dynamic RAM contains more memory cells as compared to static RAM per unit area.</li><li>② Its access time is greater than static RAM.</li><li>③ Dynamic RAM stores the data at a charge on the capacitor. It consists of MOSFET &amp; capacitor for each cell.</li><li>④ Refreshing circuitry is required to maintain the charge on the capacitor.</li><li>⑤ Cost is less.</li></ul> |
|---|--|

## Cache Memory System.



## Cache memory system.

A cache memory system includes a small amount of fast memory (SRAM) and a large amount of slow memory (DRAM) - This system is configured to simulate a large amount of fast memory.

1. Cache - This block consists of static RAM.
2. Main memory - This block consists of dynamic RAM
3. Cache controller - This block implements the cache logic.

## Elements of cache design

- 1 Cache size
- 2 mapping function
- 3 replacement algorithm write policy
- 4 block size
- 5 No. of cache.

- 1 Cache size - The size of the cache should be small enough so that the overall average cost per bit is close to that of main memory alone.

large enough so that the overall average access time is close to that of the cache alone.

## II Mapping Function —

The cache memory can store a reasonable number of blocks at any given time but this number is small compared to the total number of blocks in the main memory. There are two mapping functions commonly used.

1) Direct mapping

2) Associative mapping

## III Replacement Algorithm —

When new block is brought into the cache one of the existing blocks must be replaced by a new block.

There are 4 replacement algorithms.

1) Least Recently Used (LRU)

2) First in first out (FIFO)

3) Least Frequently Used (LFU)

4) Random

## IV Write Policy — It is also known as cache updating policy. In cache system, two copies of the same data can exist at a time one in cache and one in main memory. If one copy is altered and other is not, two different sets of data become associated with the same address. To prevent this the cache system has updating systems such as: write through systems, buffered write through systems and write back systems. The choice of cache write policy also changes the design of cache.

## I Block Size -

It should be optimum for cache memory system. Its importance is already

II Number of Cache - When on chip cache is insufficient the secondary cache is used. The cache design changes as number of cache used in the system changes.

## Mapping Functions

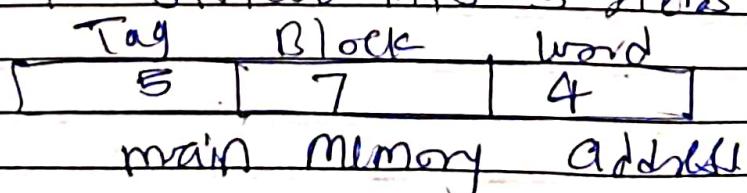
The cache memory can store a reasonable number of blocks at any give time but the no. is small compared to the total number of blocks in the main memory. There are two main mapping techniques

1. direct mapping technique
2. Associative mapping technique.

## I Direct mapping -

In this technique, each block from the main memory has only one possible location in the cache organisation. The block  $i$  of main memory maps to block  $i \bmod 128$  of the cache, whenever one of main memory blocks 0, 128, 256 ... is loaded in the cache, it is stored in cache block 0. blocks 1, 129, 257, are stored in cache block 1, etc.

To implement cache system, the address is divided into 3 fields.



## II Associative Mapping - (Fully Associative mapping)

In this technique, a main memory block can be placed into any cache block position. Although there is no fix block, the memory address has only 2 fields word & tag. The 12 tag bits are required to identify a memory block when it is residing in the Cache. The high order 12 bits of an address received from the CPU are compared to the tag bits of each block of the cache to see if the desired block is present. Once the desired block is present the 4 bit word is used to identify the necessary word from the cache.

Tag	Word
12	4

main memory address.

## II Set associative mapping.

The set associative mapping is a combination of both direct and associative mapping. A block of data from any page in the main memory can go into particular block location of any direct mapping cache.

Tag	Set	Word
6	6	4

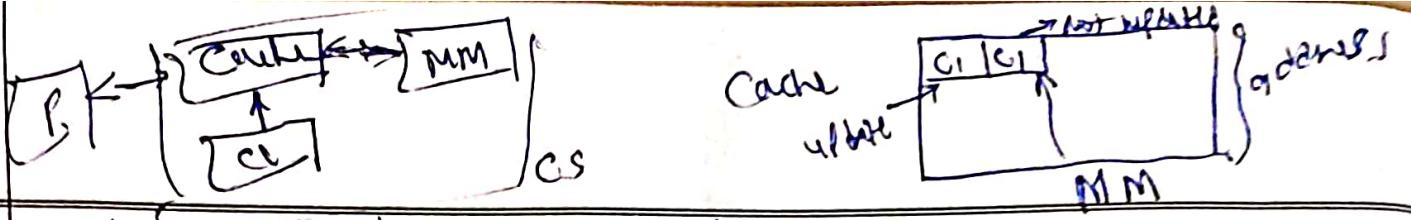
## Comparison between mapping techniques

→ direct mapping

→ associative mapping

→ set associative mapping

- ① each block from the main memory has only one feasible position in the cache.
- ② needs only one comparison.
- ③ cache hit ratio depends on processor needs to access same memory location from two different pages of the main memory frequently.
- ④ main memory address is divided into three fields - tag, block, word.
- ⑤ searching time is less
- ⑥ A block of data from main memory can be placed into any cache block position.
- ⑦ A block of data from main memory can be placed into a particular block position of any direct mapping cache.
- ⑧ needs no. of comparisons equal to no. of blocks per set.
- ⑨ cache hit ratio is affected by processor needs to access same memory location from two different pages of main memory frequently.
- ⑩ 2 fields - tag & word
- ⑪ 3 fields - tag, set, word.
- ⑫ searching time is more
- ⑬ searching time is increased with number of blocks for sets.



## Cache write/ updating

In a Cache system, two copies of same data can exist at a time, one in cache and one in main memory. If one copy is altered and other is not, two different sets of data become associated with same address.

To prevent this the cache system has updating systems such as write through system, buffered write through system and write back system.

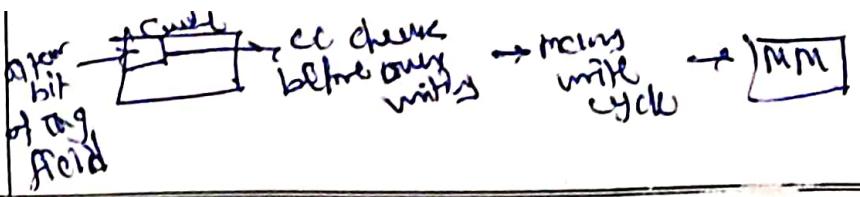
### 1 Write through system -

In write through updating system, the cache controller copies data to the main memory immediately after it is written to the cache. Due to that main memory always contains a valid data and block in the cache can be overwritten immediately without data loss.

The write through is a simple approach. This approach requires time to write data in main memory with increase in bus traffic. This in effect reduces the system performance.

### 2 Buffered write through system

In buffered write through system, the processor can start a new cycle before the write cycle to the main memory is completed. This means that the write accesses to the main memory are buffered. In such systems, read access which is a cache hit can be performed simultaneously when main memory is updated. However two consecutive write operations to the main memory or read operation with cache "miss" require the processor to wait.



### 3. Write Back System. -

In a write back system, the alter bit in tag field is used to keep information of the new data. Cache controller checks this bit before overwriting any block in the cache. if it is set, the controller copied the block to main memory before loading new data into the cache.

due to one time write operation number of write cycles are reduced in write back systems. but this system has following disadvantages.

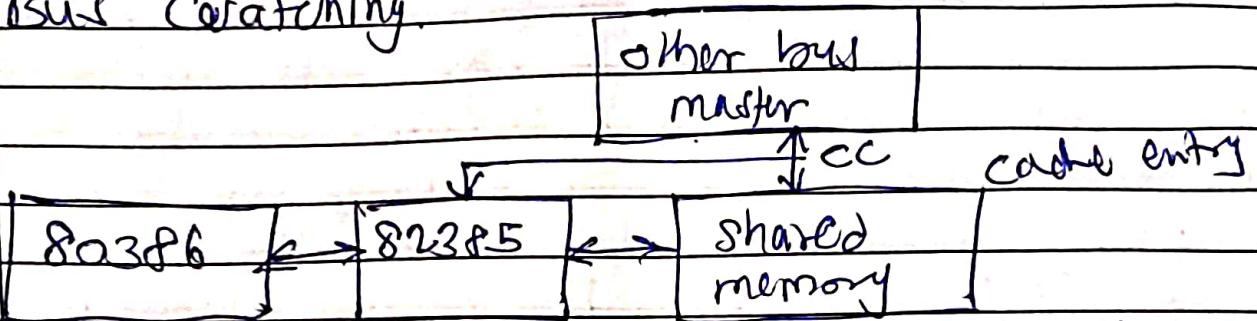
1. write back cache controller logic is more complex.
2. It is necessary that all altered blocks must be written to the main memory before another busie can access these blocks in main memory.
3. In case of power failure, the data in the cache memory is lost.

### Cache Coherency

Cache Updating Systems eliminates data inconsistency in the main memory caused by cache write operations. in multiprocessor systems another bus master can take over control of the system bus. This bus master could write data into a main memory blocks which are already held in the cache of another processor. When this happens, the data in the cache no longer match those held in main memory creating inconsistency. There are different approaches to prevent data inconsistency that is to protect cache coherency.

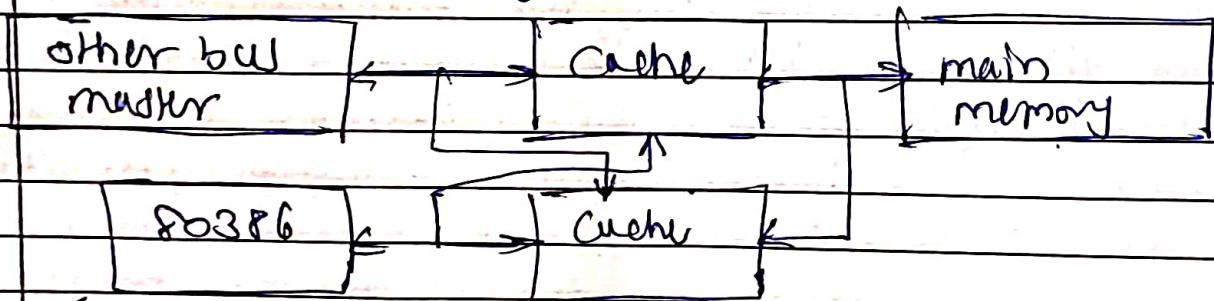
1. bus watching (snooping)
2. bus transparency
3. Non Cacheable memory
4. Cache flushing

## I Bus watching.



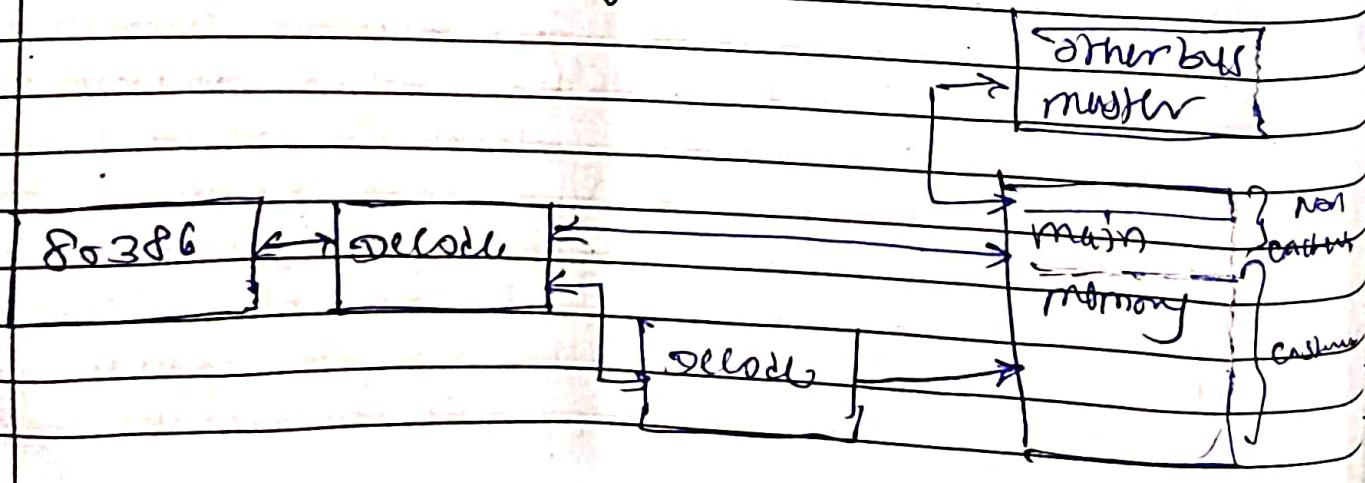
In bus watching, Cache controller invalidates the cache entry, if another master writes to a location in shared memory which also resides in the cache memory.

## II HW transparency -



In HW transparency, accesses of all devices to the main memory are directed through the same cache, or by copying all cache units both to the main memory & to all other caches that share the same memory.

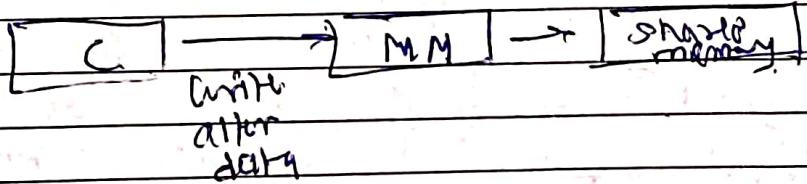
## III Non Cacheable memory -



The 80386 DX can partition its main memory into a cacheable and non cacheable memory. By designating shared memory as non cacheable memory cache coherency can be maintained. Since shared memory is never copied into cache.

W Cache flushing -

To avoid data inconsistency a cache flush writes any altered data to main memory and caches in the system are flushed before a device writes to shared memory.



## Replacement Algorithms-

When a new block is brought into the cache, one of the existing blocks must be replaced by a new block. In case of direct mapping cache, we know that each block from the main memory has only one possible location in the cache. Hence there is no choice. The previous data is replaced by the data from the same memory location from new page of the main memory but for associative and set associative technique there is a choice of replacing existing block. The choice of replacement of the existing block should be such that the probability of accessing same block must be very less. The replacement algorithms do the task of selecting the existing block which must be replaced.

There are 4 replacement algorithm.

- 1 Least Recently Used (LRU)
- 2 First-in First-out (FIFO)
- 3 Least Frequently Used (LFU)
- 4 Random.

### 1 Least Recently Used - LRU

In this technique, the block in the set which has been in the cache longest with no reference to it is selected for the replacement. We assume that most recently used memory locations are more likely to be referenced again.

This technique can be easily implemented in the two way set associative cache organization.

II First - in - first - out - FIFO

→ This technique uses same concept that stack implementation used in the microprocessors. In this technique, the block which is first loaded in the cache amongst the present blocks in the cache is selected for the replacement.

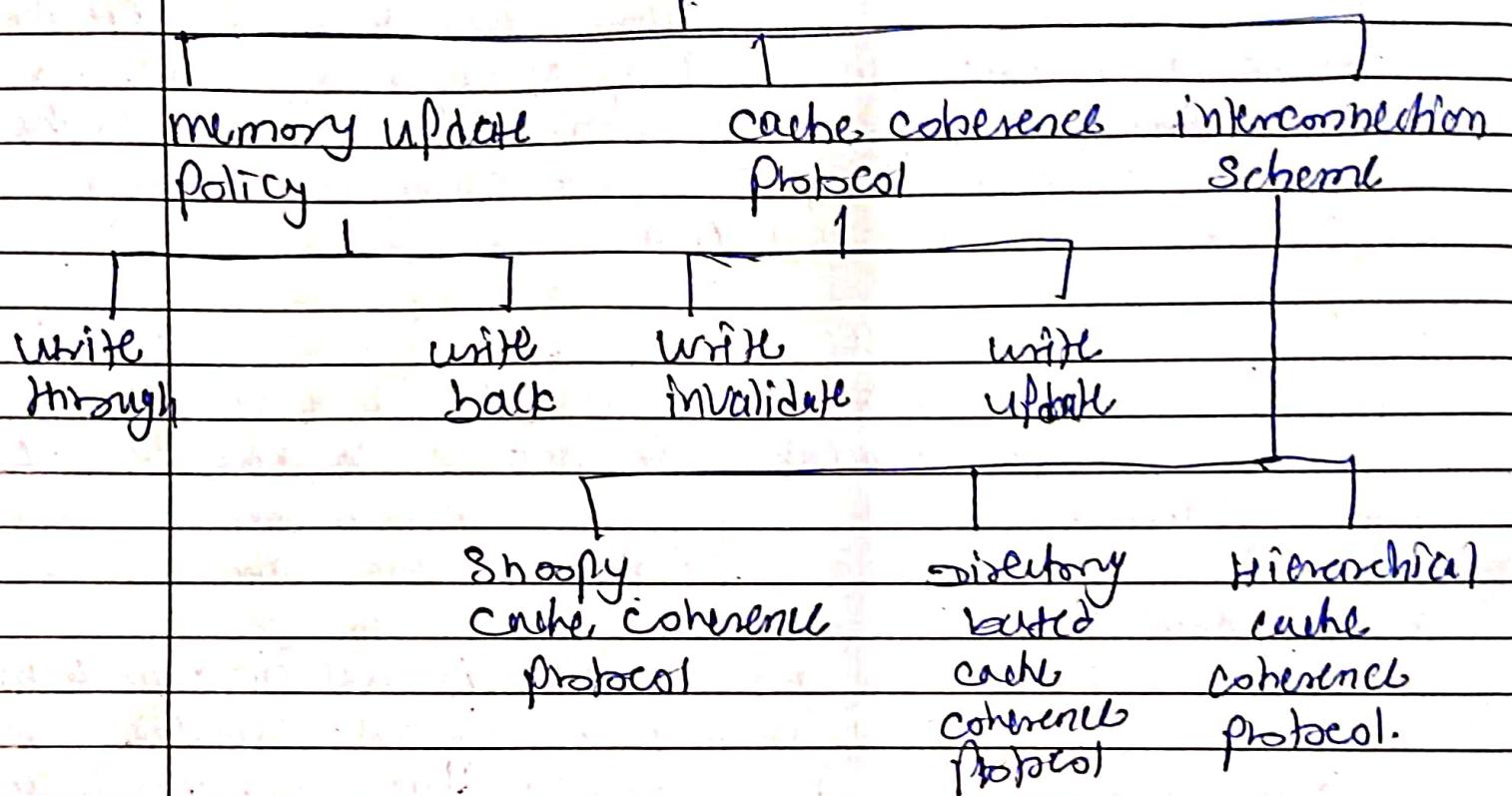
III Least Frequently Used - LFU

→ In this technique, the block in the set which has the fewest references is selected for the replacement.

IV Random -

There is no specific criteria for replacement of any block. The existing blocks are replaced randomly. Simulation studies have proved that random replacement algorithm provides only slightly inferior performance to algorithms like

## Hardware based cache coherence protocols



In a coherent microprocessor system, the caches provide both migration and replication of shared data items. Replication allows that multiple copies of the same data item reside simultaneously in different local caches, in order to reduce both latency of access and contention for a read shared data item. Migration implies a single copy of a data item that must be moved to the accessing site for exclusive use, that reduce the latency to access a shared data item that is allocated remotely.

In multiprocessor systems supporting both migration and replication of shared data items, coherence is maintained by introducing bus based or slob based protocols.

The protocols used to maintain coherence for multiple processors are called cache coherence protocols. Hardwired based protocols can be classified according to their:

- 1 memory update policy
- 2 Cache coherency protocol
- 3 interconnection scheme.

data, after } same address  
data }

### I Memory update Policy —

In cache system two copies of same data may exist at a time, one in cache and one in main memory, if one copy is altered and other is not, different sets of data become associated with same address. To prevent this memory system uses two types of memory update policies:

- 1 write through policy
- 2 write back policy

In write through policy the cache controller copies data to the main memory immediately after it is written to the cache due to the main memory always contains a valid data and any block in the cache can be overwritten immediately without data loss. This write through policy maintains consistency between the main memory and cache.

In write back policy, the alter bit in the tag field is used to keep information of the new data memory contents are not updated immediately. memory contents are updated eventually when the modified block in the cache is replaced or invalidated.

## II Cache Coherence Protocol -

Cache coherence protocols are used to update copies in the multiprocessor system. Two protocols are used for updating the cache copies of a data.

1. write update protocol
2. write invalidate protocol.

### write invalidate protocol

1. multiple writes to the same word with no intervening reads require only one initial invalidation.

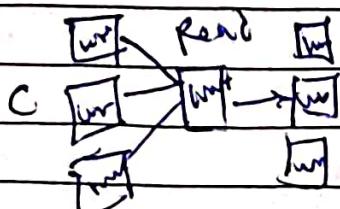
2. multiword cache block. Only the first write to any word in the block needs to generate an invalidate. This protocol works on cache blocks.

3. Any read to invalidated data is forced to fetch a new copy of data. This takes comparatively longer time.

1. write update protocol multiple writes to the same word with no intervening read require multiple write broadcasts.

2. multiword cache block, each word written in a cache block requires a write broadcast. This protocol works on individual word.

3. The written data is immediately updated in other cache. Thus any read to the data in the cache is valid. Thus less time is required to read the data.



### III Interconnection Scheme -

Depending on the nature of the interconnection network how based protocols are classified into 3 basic classes.

#### i) Snoopy cache coherence protocol -

Snoopy cache coherence protocols are very popular in shared bus multiprocessors due to their relative simplicity. Every cache that has a copy of the data from a block of physical memory also has copy of a sharing status of the block. A no centralized state is kept. This protocol is typically used in single bus based shared memory system where consistency commands are broadcast via the bus, & all cache controller monitor or loop on the bus to determine whether or not they have a copy of a block that is requested on the bus.

#### ii) Directory based cache coherence protocol -

Large interconnection networks cannot support broadcasting efficiently and therefore a mechanism is needed that can directly forward consistency command to those caches that contain a copy of the update data. For this purpose the sharing status of a block of physical memory is kept in just one location called the directory.

#### iii) Hierarchical cache coherence protocol.

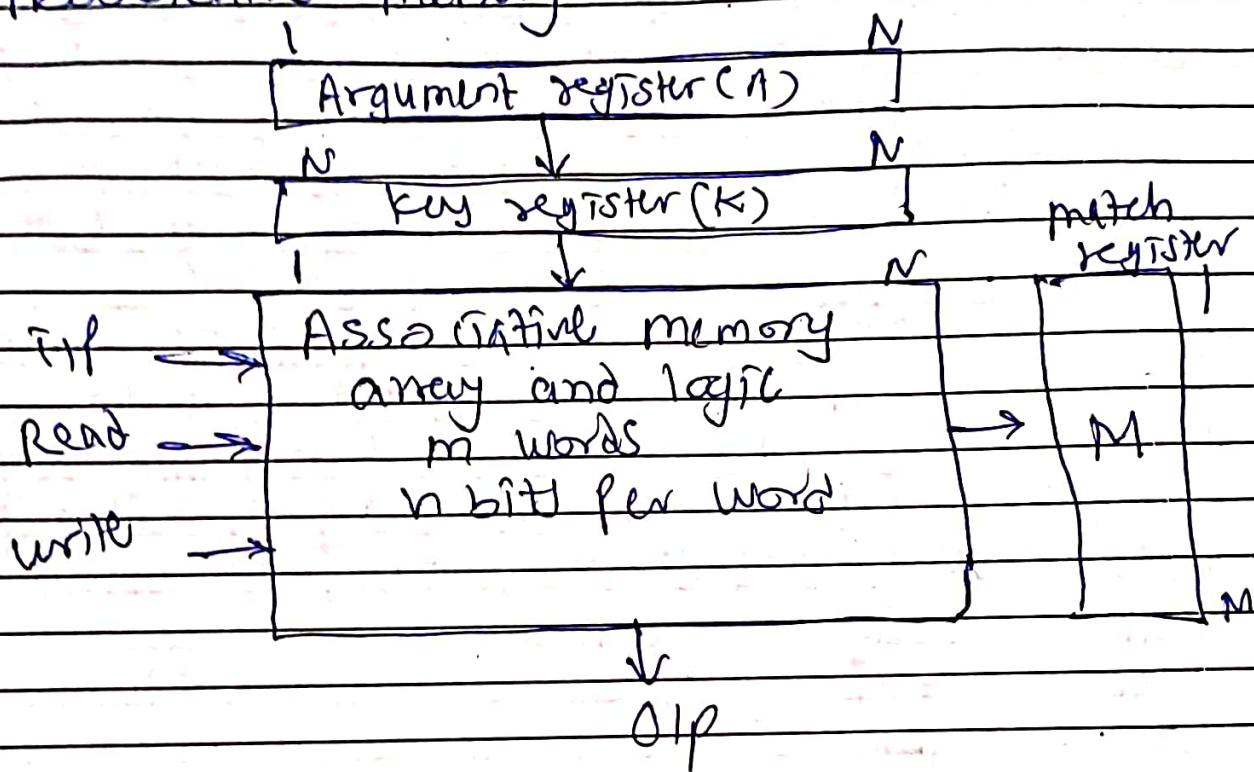
This protocol tries to avoid the application of the costly directory protocol but still provide high scalability. It proposes multiple bus buses with the application of hierarchical cache coherence protocols that are generalized or extended version of the single bus based snoopy cache protocol.

## MESI protocol

Write invalidate protocol is also called as MESI protocol.

1. Modified - The line in the cache, different from main memory is modified & this line is available only in this cache.
2. Exclusive - The line in the cache is same as that in main memory & it is not present in any other cache.
3. Shared - The line in the cache is same as that in main memory & the same line may be present in one or more other caches.
4. Invalid - The line in the cache does not contain valid data.

## Associative memory



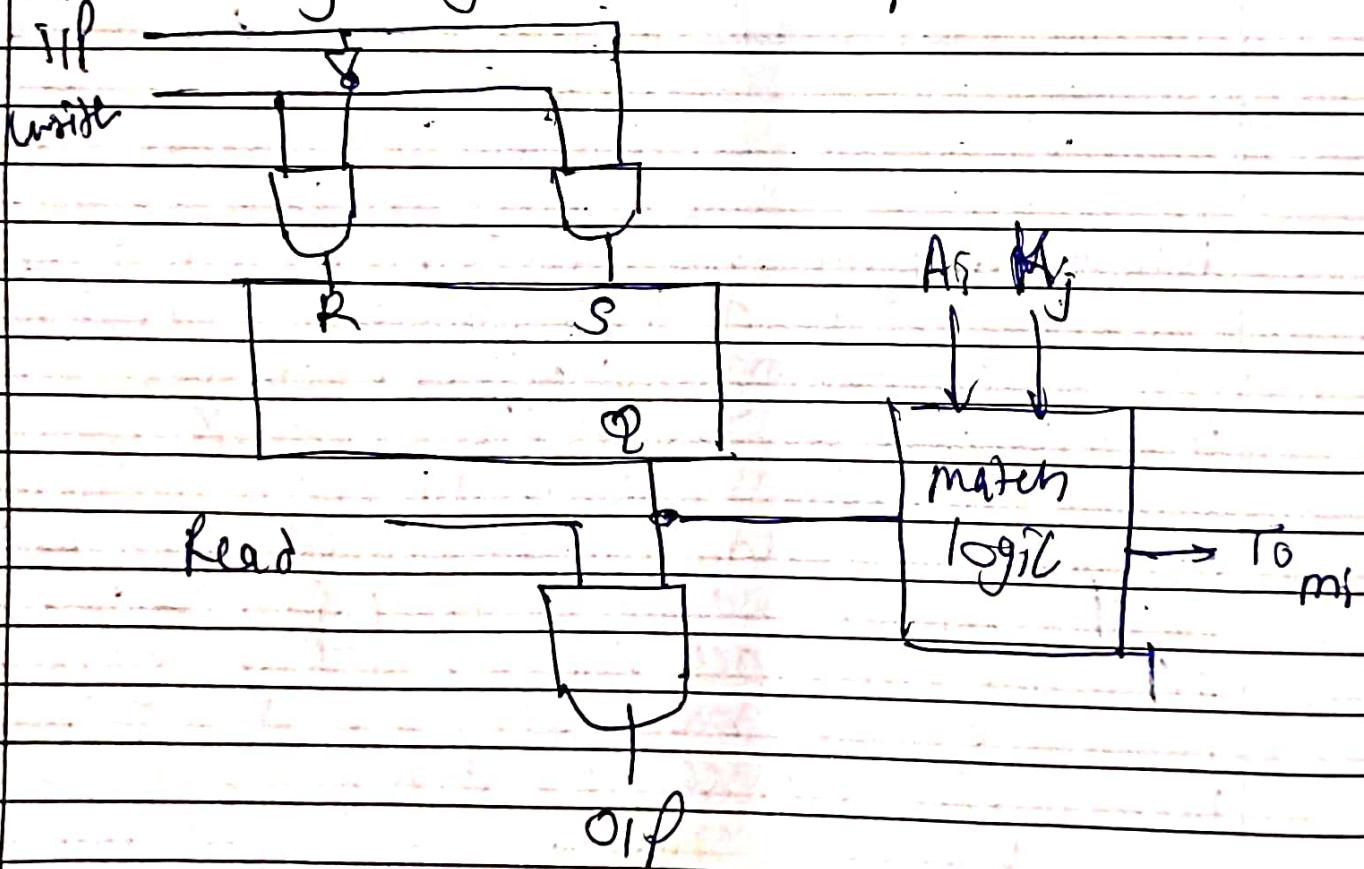
## Block diagram of associative memory.

The time required to find an object stored in memory can be reduced considerably if objects are selected based on their content, not on their locations. A memory unit accessed by the content is called an associative memory or content Addressable memory (CAM). This type of memory is accessed simultaneously and in parallel on the basis of data content determined by specific address or location.

block diagram of associative memory array with match logic for  $m$  bit words & associated registers. The argument register (A) and key register (K) each have  $n$ -bits per word. Each word in memory is compared in parallel with the content of the argument register.

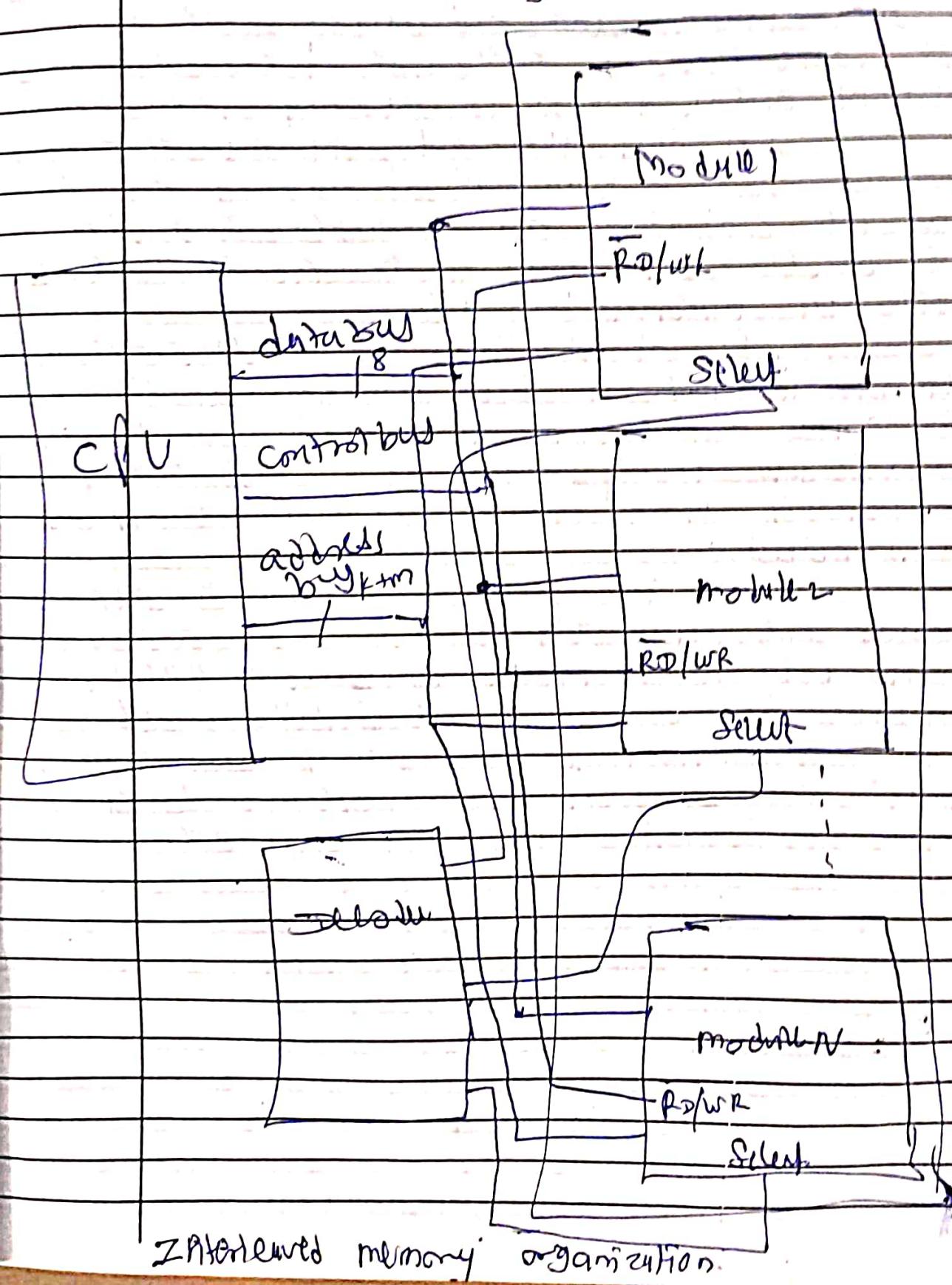
The words that match with the word stored in the argument register set a corresponding bits in the match register. Therefore reading can be accomplished by a sequential access to memory for those words whose corresponding bits in the match register have been set.

The key register provides a mask for choosing a particular field or bits in the argument word. Only those bits in the argument register having 1's in their corresponding position of the key register are compared.



Internal organization of typical off-the-shelf associative memory

## Interleaved memory



Main memory access time is the bottleneck in the system and one way to reduce it is to use cache memory. Alternative technique to reduce memory access time is memory interleaving.

In this technique, the main memory is divided into a number of memory modules and the address are arranged such that the successive words in the address space are placed in different module.

Most of the time CPU access consecutive memory locations. In such situation address will be to the different modules. These modules can be accessed in parallel, the average access time of fetching word from the main memory can be reduced.

The low order bits of the memory address are generally used to select a module & the high order nibbles are used to access a particular location within the selected module.

The effect of interleaving is substantial, it does not speed up memory operation by a factor equal to the number of the module. If