

Unit-V

Topics to cover:

- * Memory Concept and Hierarchy
- * Memory Management
- * Cache Memory : Mapping and Replacement Techniques
- * Virtual Memory ✓
- * DMA ✓
- * PIO
- * Accessing I/O parallel and serial interface ✓
- * Interrupt I/O
- * Inter Connection Standards : USB, SATA

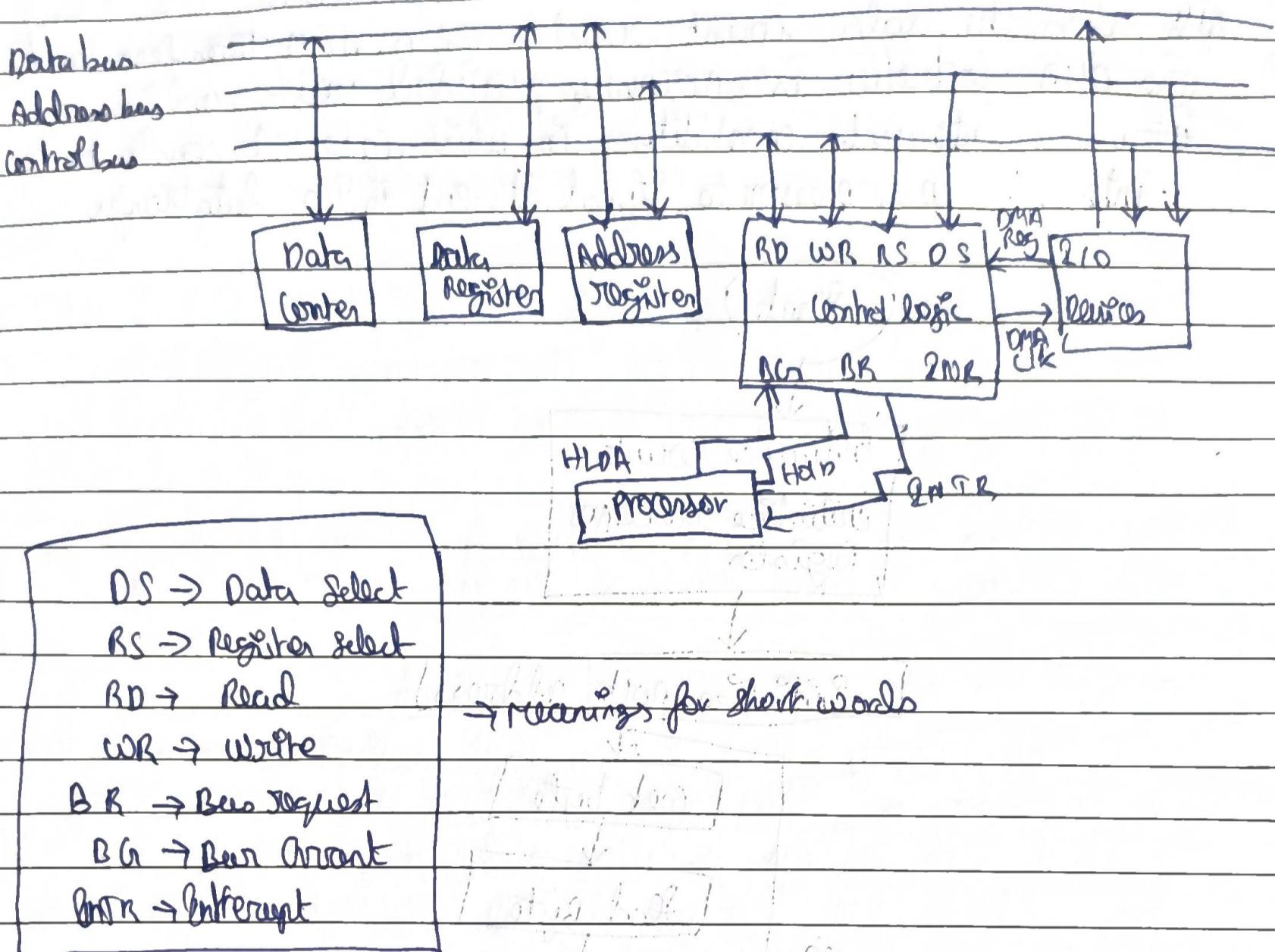
Direct memory access

It allows certain hardware subsystems to access the main system memory independently

DMA:

- * DMA uses hardware for accessing the memory, that hardware is called DMA controller
- * It has the work of transferring the data between Input / Output devices & Main Memory with very less interaction with the processor
- * The direct memory access controller is a control unit, which has the work of transferring data.

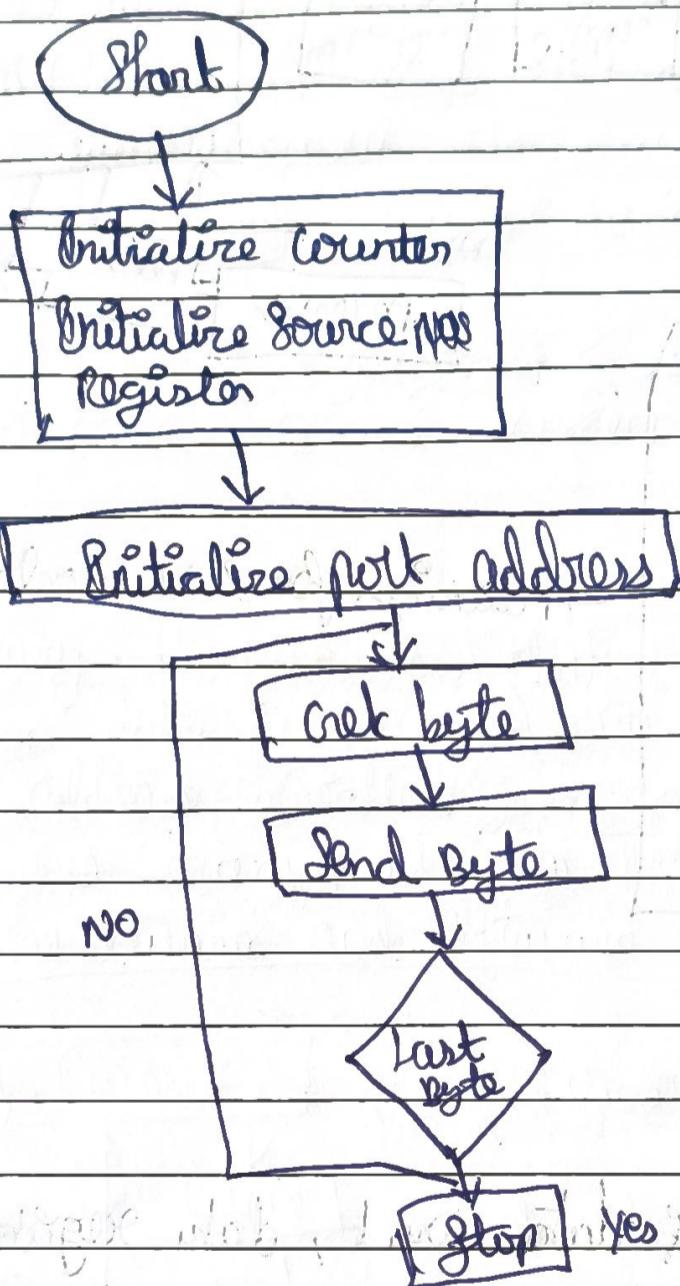
Databases
Addressing
Control bus



Explanation:

- ① DMA consists of data counter, data register, address register and control logic
- ② Data transfer counter register stores the number which gives the number data transfers to be done in one DMA cycle
- ③ It is automatically decremented after each word transfer
- ④ Data Register acts as a buffer where as address register initially holds the starting address of the device
- ⑤ Actually, it stores the address of the next word to be transferred. It is automatically incremented and decremented after each word transfer.

- ⑥ After each transfer, data counter is tested for zero
- ⑦ When the data count reaches zero, DMA transfers halts
- ⑧ The DMA controller is normally provided with an interrupt capability, in which case it sends an interrupt to processor to signal the end of I/O data transfer.



Types of DMA

⇒ Single Ended DMA

⇒ Dual Ended DMA

⇒ Arbitrated Ended DMA

⇒ Interleaved DMA

- i) It is operated by reading and writing from single memory address. They are the simplest DMA (Single Ended).
- ii) It can read & write from two memory address dual Ended DMA is more advanced than single-ended DMA.
- iii) It works by reading and writing to several memory address. It is more advanced than Dual-Ended DMA (Arbitrated-Ended).
- iv) It is DMA that read from one memory address and write from another memory address.

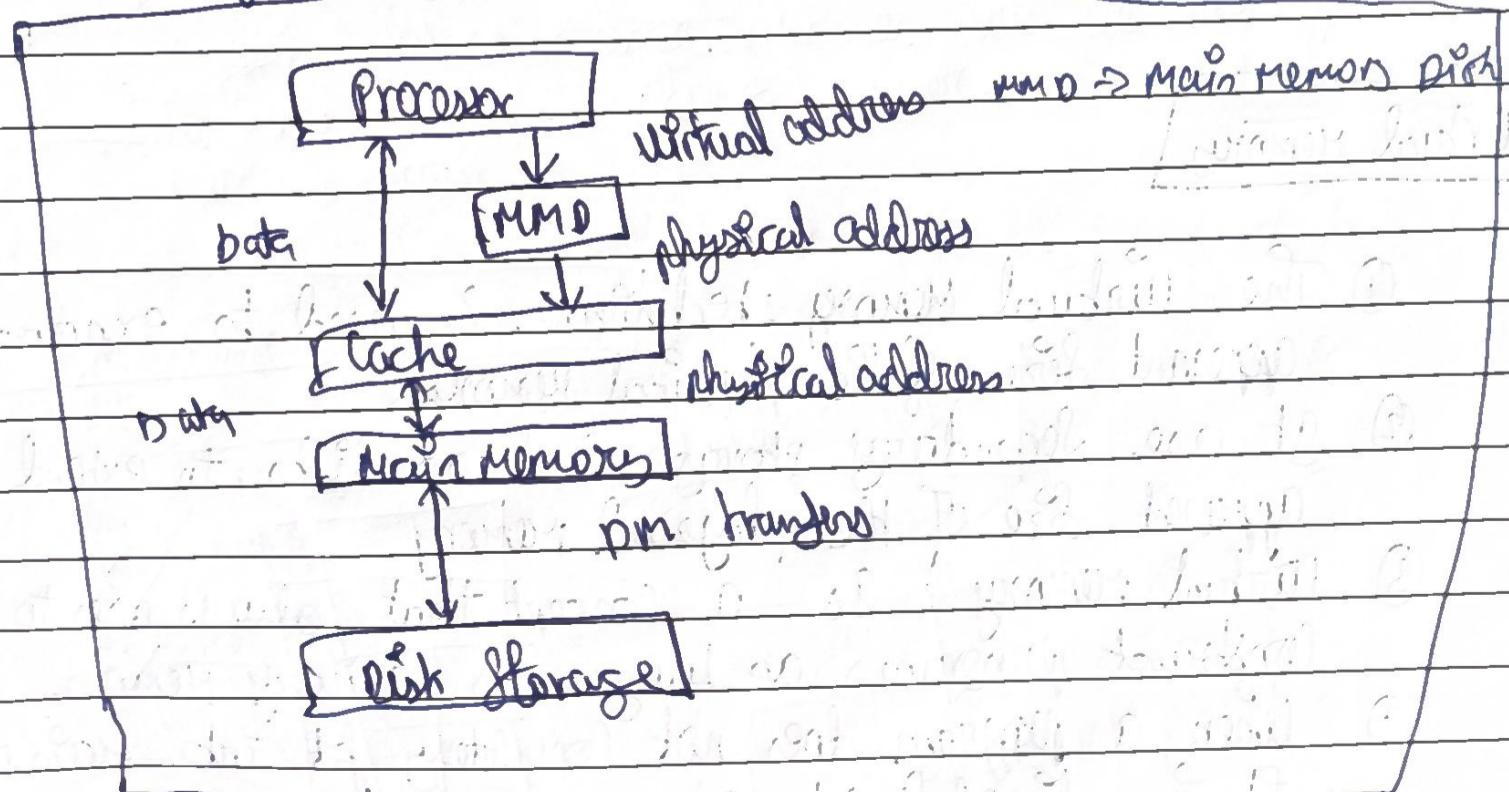
DMA controller have 3 registers:

- i) Address Register : It contains the address to specify the desired location in memory.
- ii) Word count Register : It contains the number of words to be transferred.
- iii) Control Register : It specifies the transfer mode.

Virtual Memory

- ① The Virtual Memory : Technique is used to expand the apparent size of the physical memory.
- ② It uses Secondary storage such as disks, to extend the apparent size of the physical memory.
- ③ Virtual memory : Is a concept that allows user to construct programs as large as auxiliary memory.
- ④ When a program does not completely fit into main memory it is divided into segments.

- (5) The segments which are currently being executed are kept in the main memory & remaining segments are stored in the secondary storage devices such as magnetic disk.
- (6) If an executing program needs a segment which is not currently in the main memory, the required segment is copied from the secondary storage device.
- (7) When a new segment of a program is to be copied into a main memory it must replace another segment in the memory.
- (8) In virtual memory, the address that processor issues to access either instruction or data are called virtual address.
- (9) The set of such address is called address space.
- (10) These address are transferred into physical address by a combination of hardware and software components.



Paging

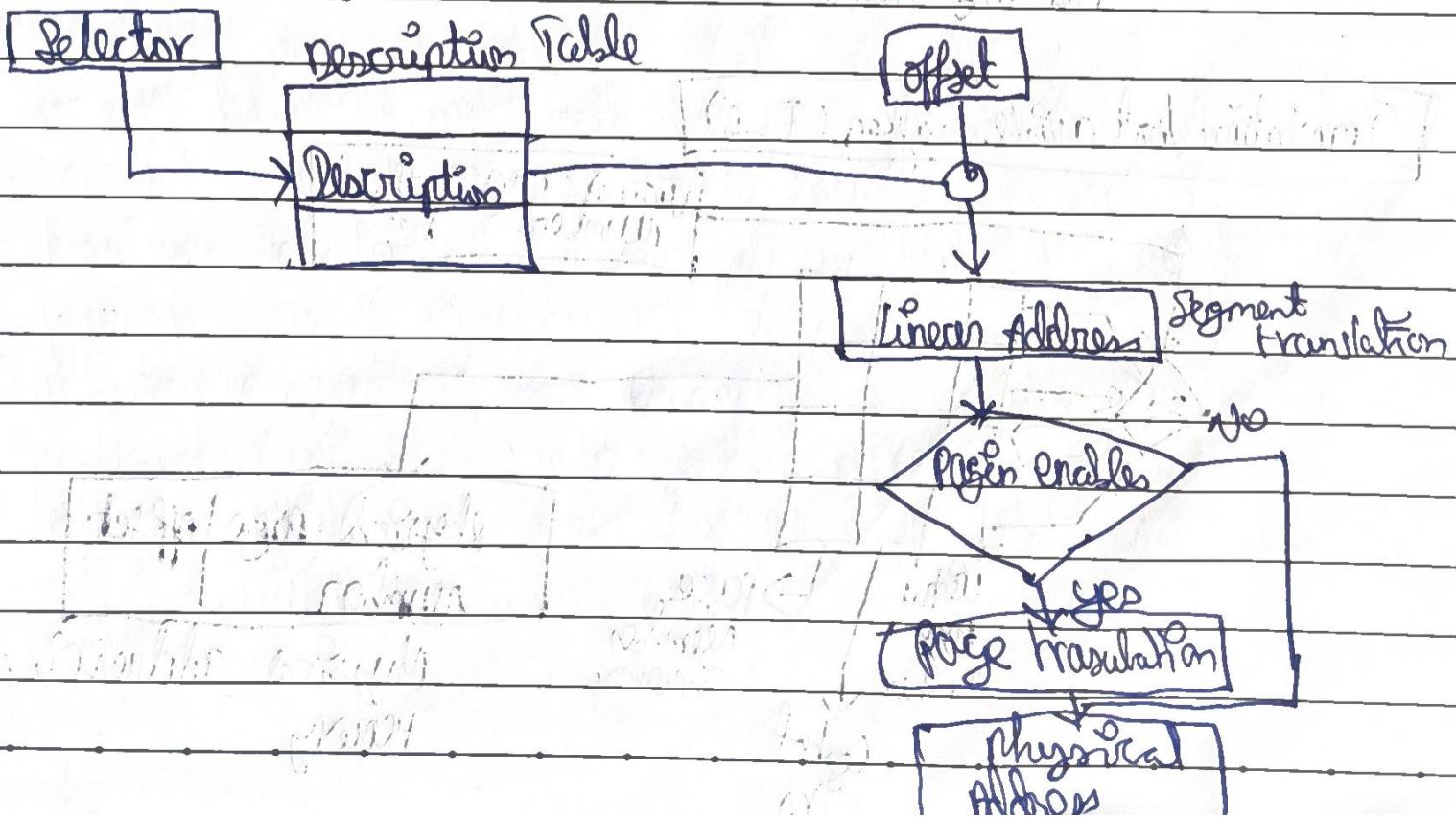
- * It is a simple method for translating virtual address into physical address by assuming that all programs and data are composed of fixed length unit called pages.
- * Pages constitute the basic unit of information that is moved between the main memory and the disk whenever the page translation mechanism determines that a swap is required.

Ex :

Page 100000	4KB	Virtual Address
Page 1000	4KB	
physical address		Segmented Address
Space	4KB	first part of address
Page 1	4KB	second part of address
Page 1	4KB	third part of address
Page 0	4KB	fourth part of address

Virtual to Physical Address Translation

① Segment translation



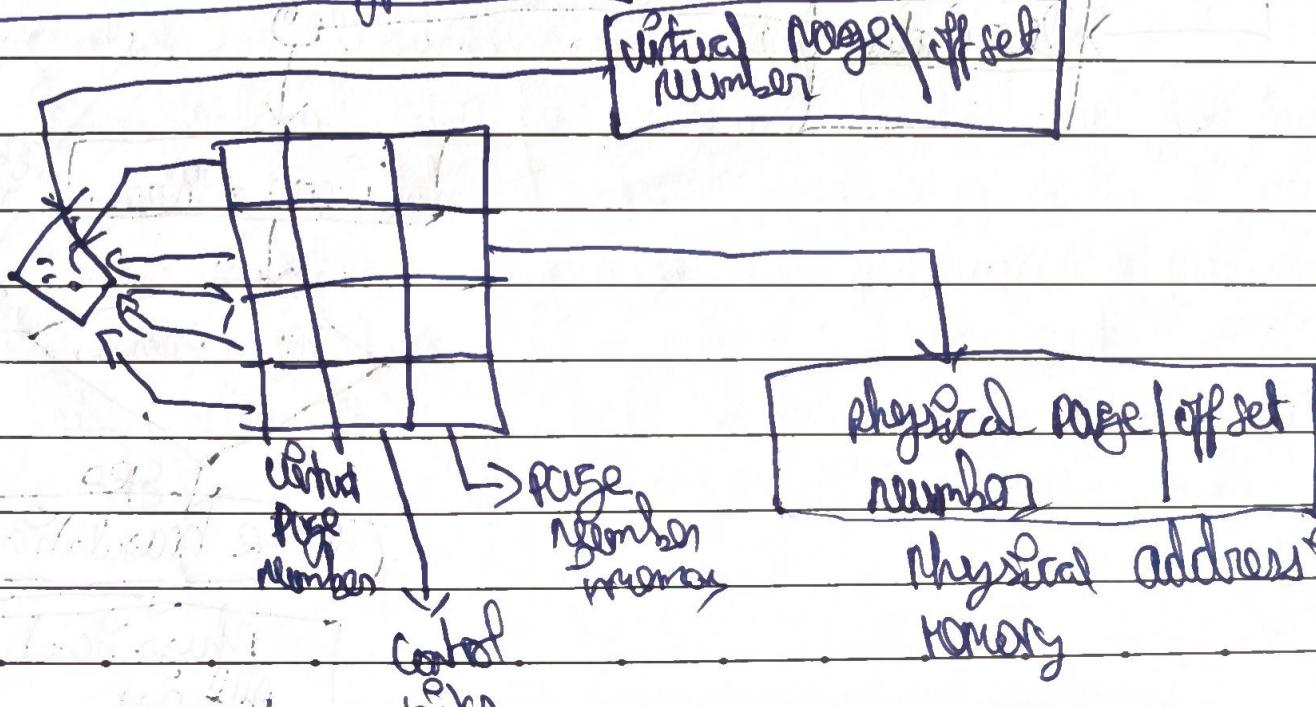
Key Components:

- a) MMU : Memory Management Unit , A hardware device responsible for translating virtual address to physical Address at runtime
- *) Page tables : Data structures maintained by the operating system that Map virtual page numbers to physical frame numbers.
- *) Translation lookaside Buffer (TLB) : A cache used by the MMU to Speed up address translation by storing recent virtual to physical mappings

Process overview :

- 1) The program generates a virtual address, which consists of two parts : The virtual page number and the page offset (location within that page)
- 2) The virtual page number indexes into the page table to find the corresponding physical frame number in RAM
- 3) The physical address is then computed by combining the physical frame frame number with the page offset
- 4) The MMU uses this physical address to access the actual Memory location

Translation lookaside buffer (TLB)



- i) To support demand paging and virtual memory processor has to access page table which is kept in main memory.
- ii) To reduce the access time and degradation of performance, a small portion of the page table is accommodated in the memory management unit.
- iii) This portion is called translation lookaside buffer.

Serial interface and parallel interface

Difference :-

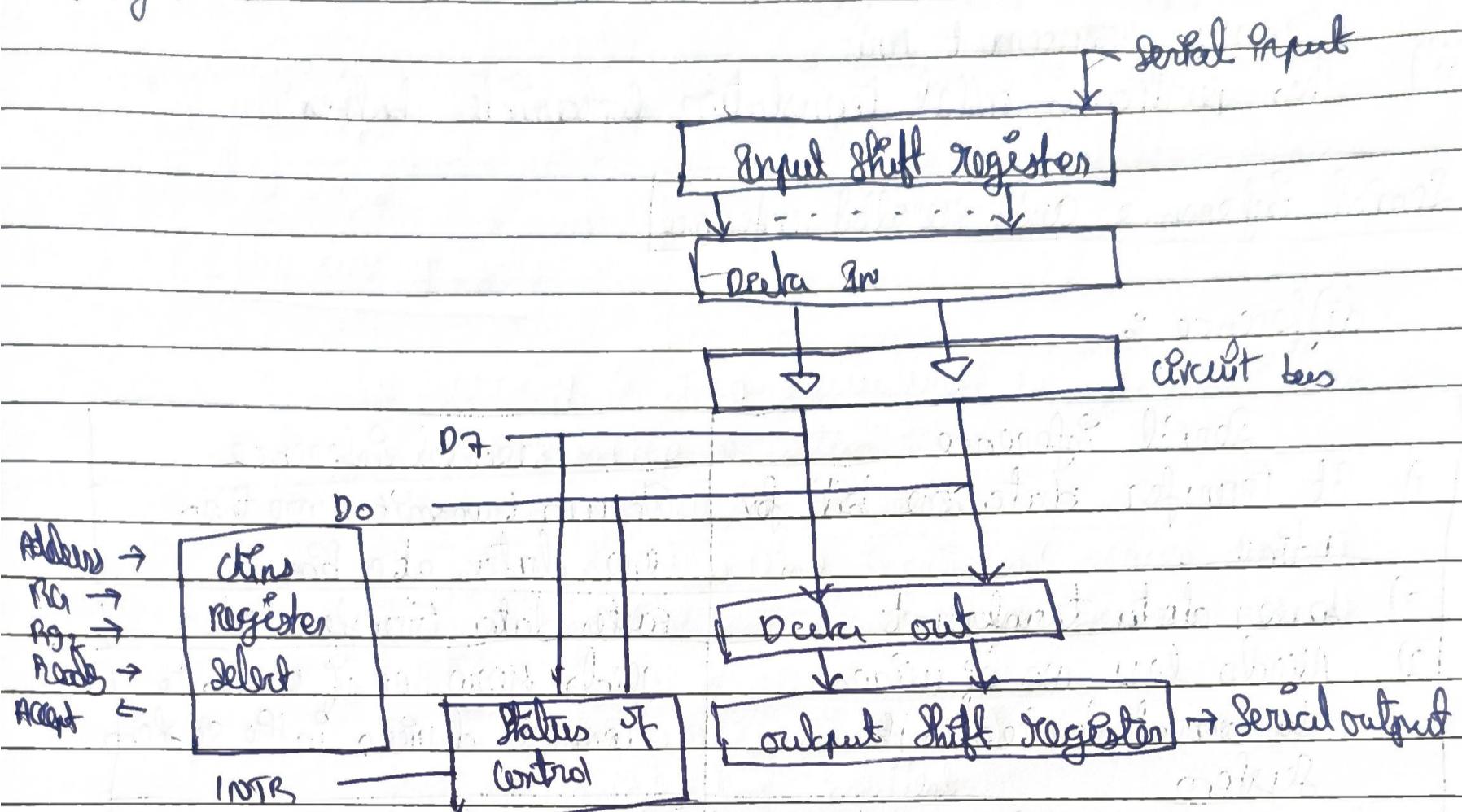
Serial interface	Parallel interface
i) It transfer data one bit at a time	It can transmit more than one data at a time
2) Lower data transfer rate	Faster data transfer rate
3) Needs less no of wires to connect devices in the system	Needs more no of wires to connect devices in the system
4) Well suited for long distances because fewer wires are used	Not suitable for long distance interface

Serial Interface :-

- * A serial interface is used to transmit/receive data serially CIR | one bit at a time
- * A key feature of an interface circuit for a serial port is that it is capable for communicating in a bit serial fashion on processor side.
- * A shift register is used to transform information between the parallel and serial formats.
- * The input shift register accepts serial data bit by bit and converts it into parallel data.

- * The converted parallel data is located in the data register and it is then read by the processor using data bus.

Diagram:



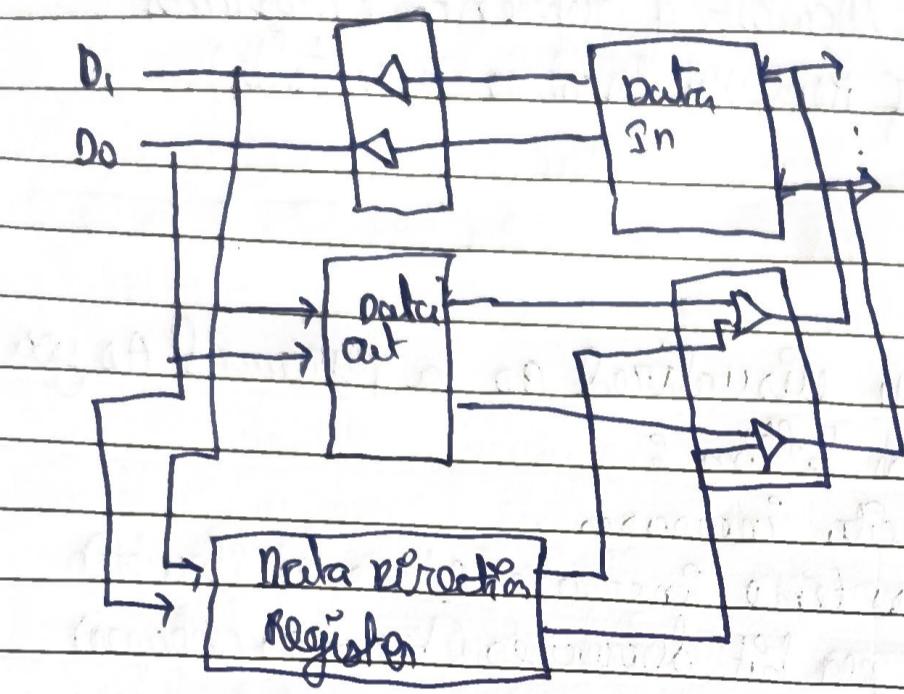
Parallel Interface

- * The input and the output interfaces can be combined into a single interface and the direction of data flow can be controlled by data direction register.

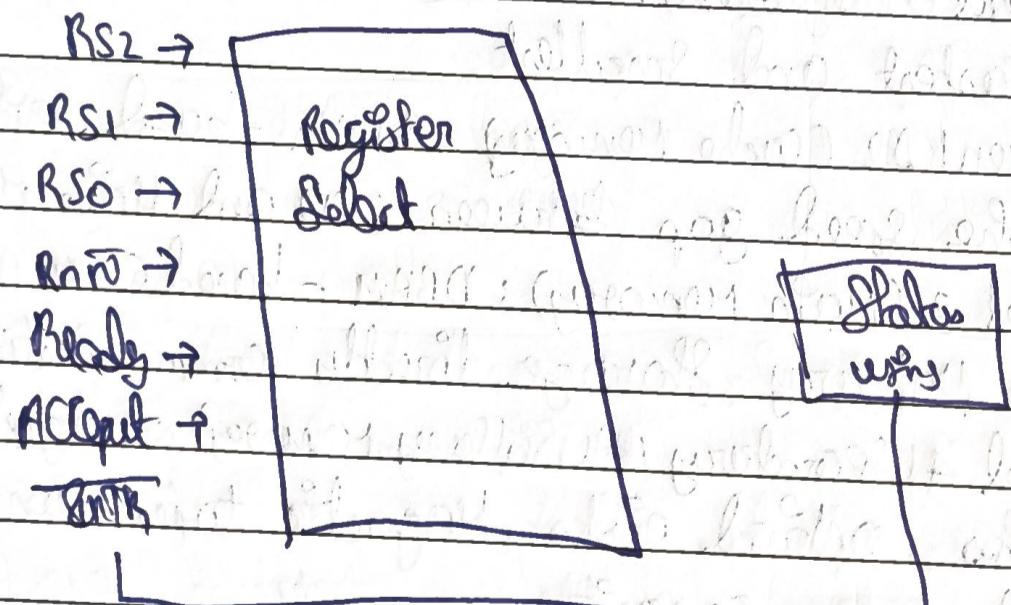
- * A single interface can be programmed to use for input the data or output the data.
- * Such a interface is known as programmable parallel interface.
- * Data and interface lines are bidirectional.
- * Their direction is controlled by data direction register (DDR).

- * M0 & M1 are connected to status control's
- * Ready and accept signals are provided as handshaking signals

Diagram:



Address:



Memory concepts and Hierarchy

The memory system is organized in a hierarchy based on access time, capacity and cost. The goal is to optimize performance by utilizing the principle of locality of reference (programs tend to access the same set of memory location repeatedly).

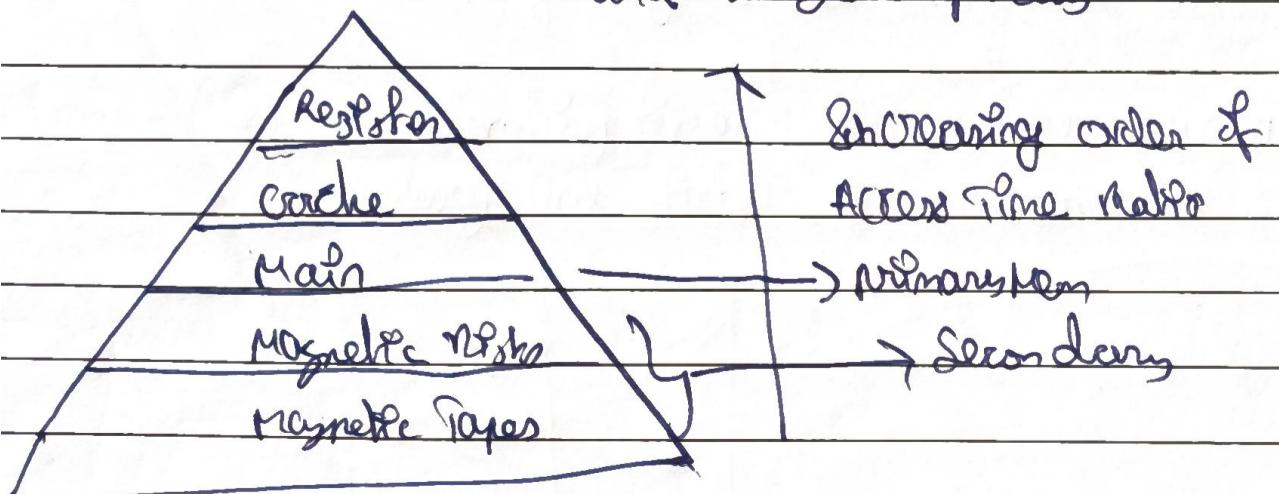
The Hierarchy Levels

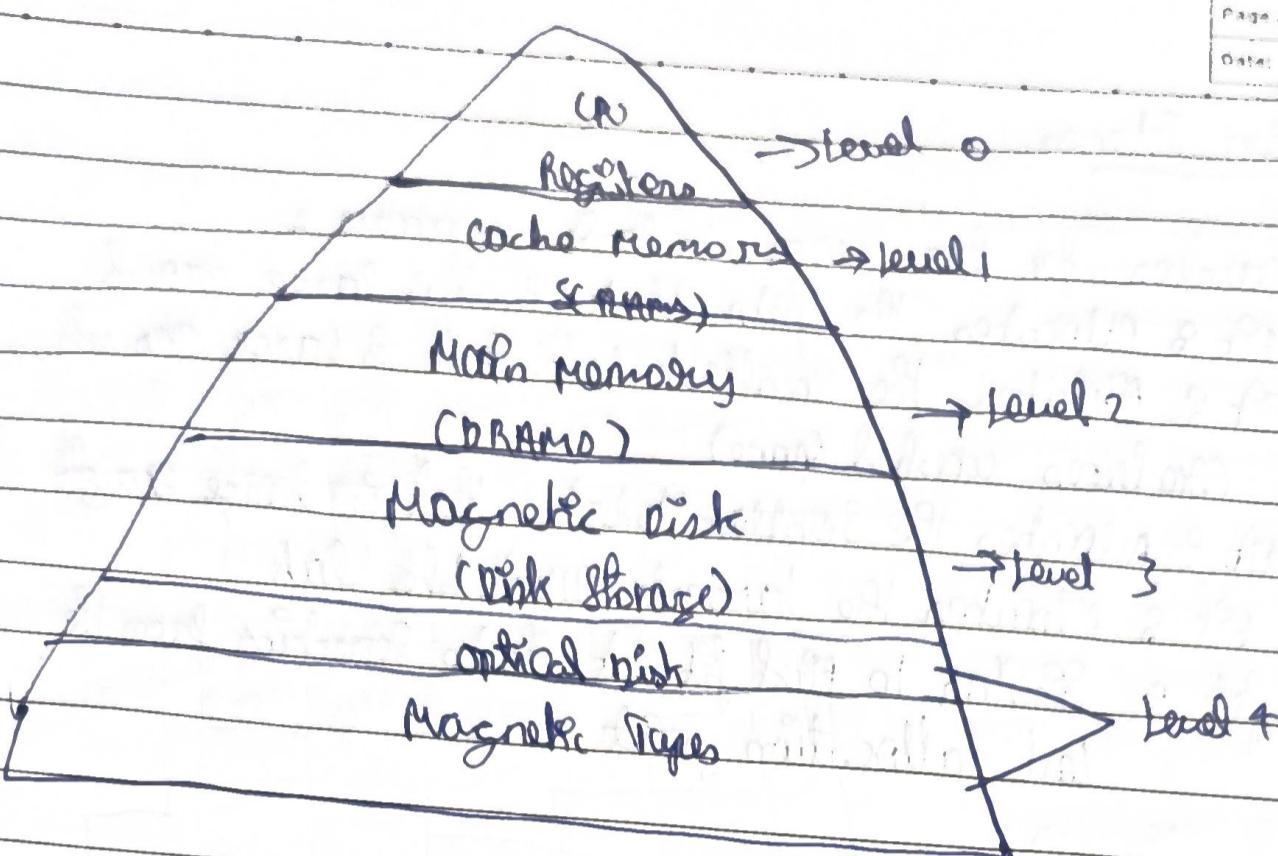
The hierarchy is often visualized as a pyramid. As you move from top to bottom:

- * Capacity increases
- * Access time increases (becomes slower)
- * Cost per bit decreases (becomes cheaper)

Level Breakdown:

- o Level 0 (Registers): Located inside the CPU.
Fastest and Smallest
- o Level 1 & 2 (Cache Memory): SRAM - based Bridges
The speed gap between CPU and Main memory
- o Level 3 (Main memory): DRAM - based (RAM / ROM).
The Primary Storage directly communicating
- o Level 4 (Secondary / Auxiliary memory): Magnetic disks, Optical Disks, Magnetic tapes non-volatile and large capacity





Memory Management

Memory Management is the process of controlling and coordinating Computer Memory, assigning blocks to running programs to optimize System performance.

Key Concepts:

- Scraping: A mechanism where a process is temporarily moved out of Main Memory to Secondary Storage (disk) to make room for other processes, and brought back later.

Fragmentation:

- External Fragmentation: Free memory is separated into small blocks scattered across the system, making them unusable for large processes.

- Internal Fragmentation: Memory is allocated in fixed blocks, but the process is smaller than the block, leaving unused space inside.

Partition Allocation Schemes

How the OS Searches for free space to load a process:

- First fit : Allocates the first block that is large enough
- Best fit : Allocates the smallest block that is large enough (Produces wasted space)
- Worst fit : Allocates the largest block that is large enough
- Next fit : Similar to first fit but starts searching from the last allocation point

Paging Vs Segmentation

- Paging : logical memory is divided into fixed-size blocks called pages . physical memory is divided into frames . eliminates external fragmentation
- Segmentation : memory is divided into variable-length segments based on logical units

Static Loading:

- *) Static loading is used when you want to load your program statically . Then at the time of compilation , the entire program will be linked and compiled without need of any external module or program dependency
- *) At loading time , the entire program is loaded into memory and starts its execution

Dynamic linking

- * In a dynamically linked program, references will be resolved and the linking will be done at the time of execution.
- * routines of the library are loaded into memory only when they are required in the program.

Cache Memory

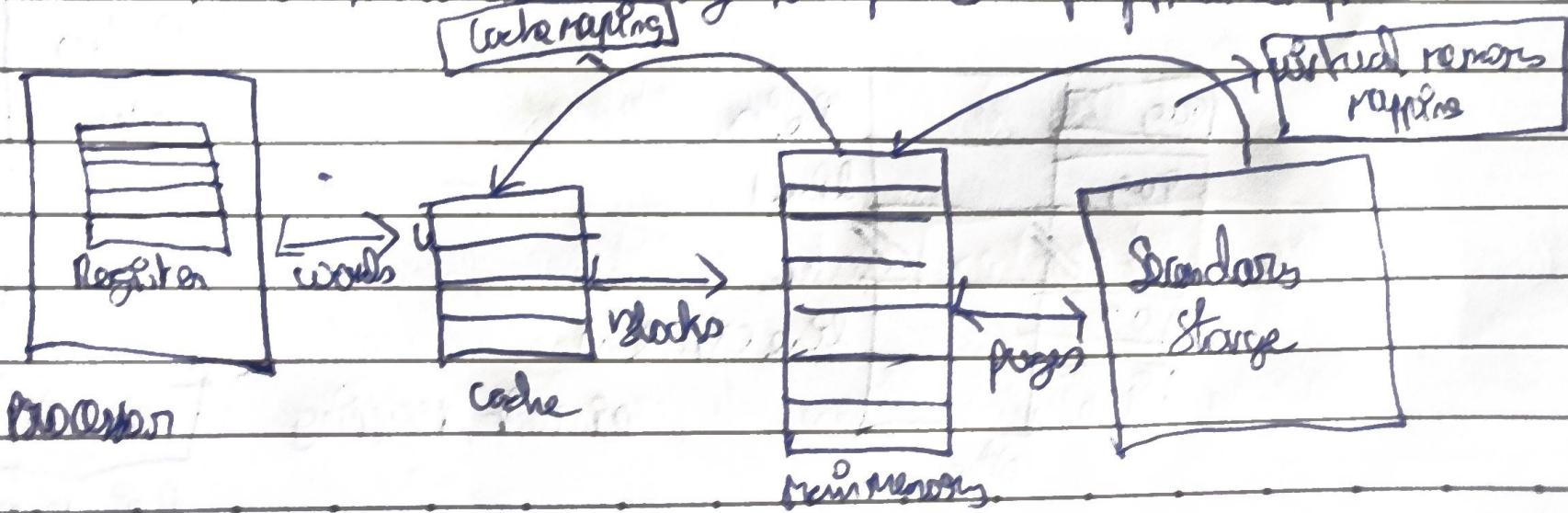
Overview

Cache memory is a small, high-speed memory unit situated between the CPU and Main Memory (MM). Its primary purpose is to bridge the significant speed mismatch between the fast processor and the slower main memory.

- o Function: It holds frequently requested data and instructions so the CPU can access them immediately.
- o Locality of Reference: Cache operates on the principle that programs tend to access the same set of memory locations repeatedly over a short period.

Levels:

- o L1 (Primary) Fastest, smallest, usually embedded in the processor
- o L2 (Secondary) Larger than L1, slightly slower
- o L3 = Specialized memory to improve the performance of L1 and L2



Mapping Techniques

Cache Mapping

Cache mapping defines the specific rules for storing a block of data from the main memory into the cache memory. Since the cache is smaller than main memory, we need a method to determine which main memory block occupies which cache line.

The main memory is divided into lines of the same size.

Direct Mapping:

In this technique, a specific block from main memory maps to only one specific line in the cache.

* Formula:

$$\text{Cache Line} = (\text{Main Memory Block Address})$$

* Mechanism: If the cache has n lines, block j of main memory maps to line $(j \bmod n)$.

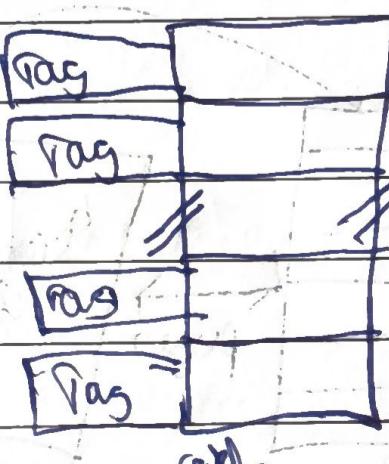
* Replacement: No replacement algorithm needed. If a new block maps to a line already occupied, the old block is simply overwritten.

* Address Structure: The physical address is split into three parts

o Tag | Line number | Block/Line offset

o Line number

o Block/Line offset



line 0

line 1

line (Mod n)

Direct mapping



Main memory

D) Fully Associative Mapping

This is the most flexible technique. A block from main memory can be placed in any freely available line in the cache.

- o Mechanism: The cache control logic checks all lines to find the data.

- o Replacement: A replacement algorithm is mandatory because if the cache is full, the system must decide which existing block to evict to make room for the new one.
- o Address Structure: The physical address is split into two parts (no line number is needed):
 - o Block number (tag).
 - o Block / line offset.

C) k-way Set Associative Mapping

This is a hybrid approach that combines the simplicity of Direct Mapping with the flexibility of Fully Associative Mapping.

- o Mechanism: Cache lines are grouped into sets, where each set contains k lines. A block from Main memory maps to a specific set, but within that set, it can occupy any of the k lines.

- o Formula:

$$\text{Cache Set} = (\text{Block Address}) \text{ mod total sets in cache}$$

- o Replacement: A replacement algorithm is required only for the specific set if all lines in that set are full.

- o Address Structure

- o Tag
- o Set number
- o Block/line offset

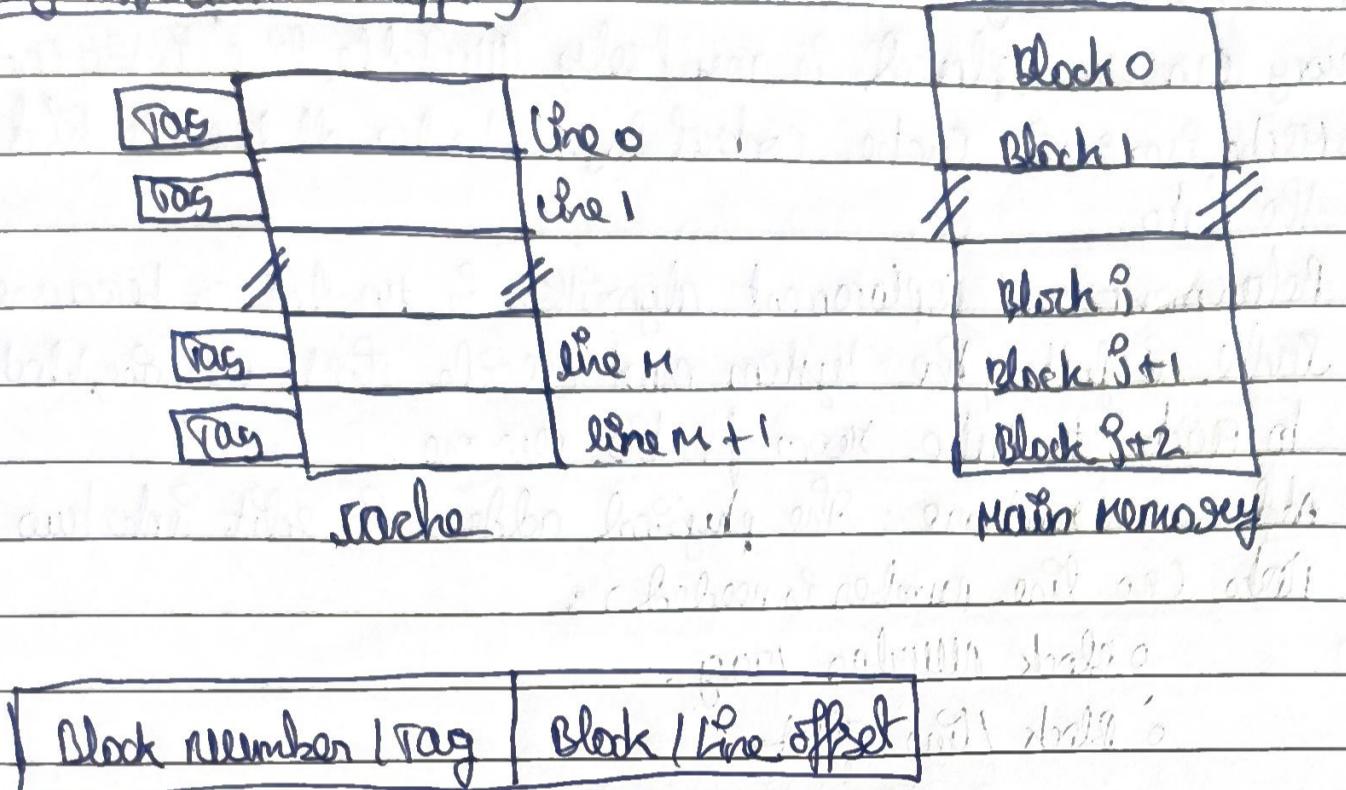
- o Special Cases:

* If $k=1$, it functions exactly like Direct Mapping

* If $k = \text{total lines}$, it functions exactly like Fully Associative Mapping

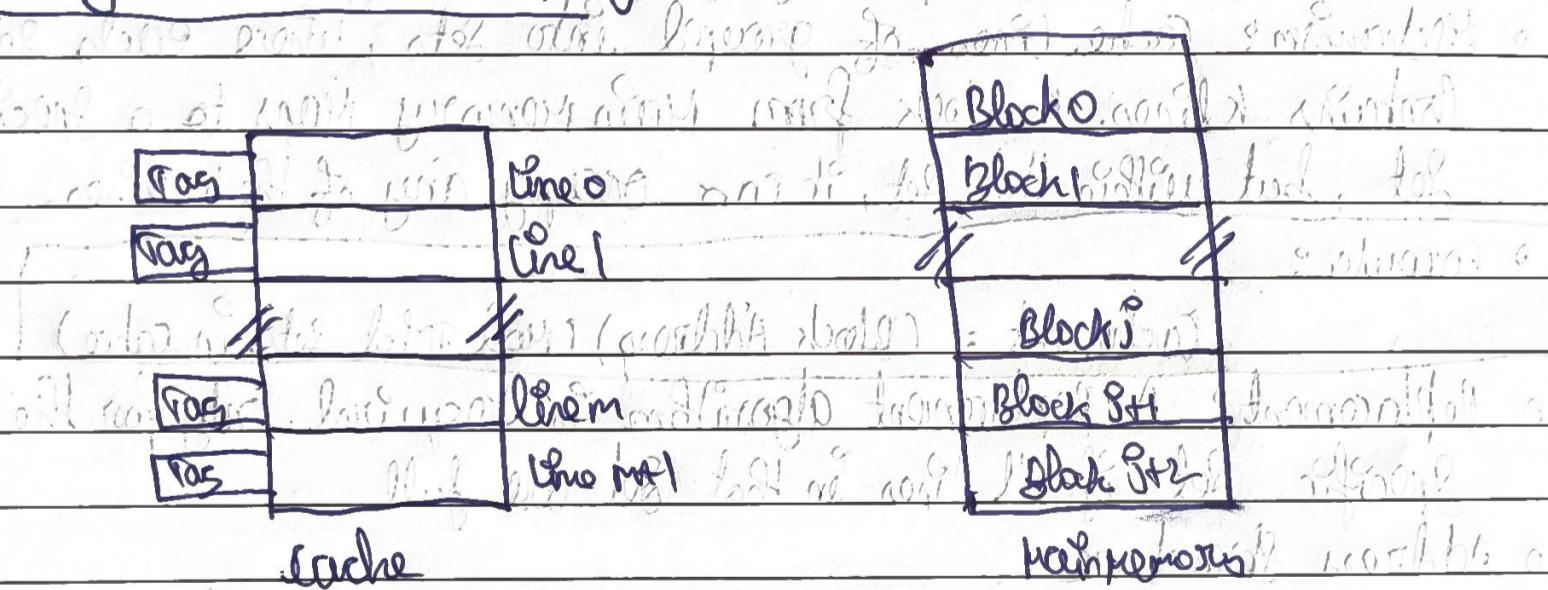
Block Diagram

Fully Associative Mapping



Division of Physical Address in fully associative Mapping

K-way Set Associative Mapping



2-way Set Associative Mapping

Tag	Set number	Block / Line offset
-----	------------	---------------------

Division of physical address in K-way Set Associative Mapping

Replacement Algorithm :

The performance of cache memory is measured by the hit ratio (h)
 - The percentage of memory accesses satisfied by the cache

Average Memory Access Time (T_{avg}) :

$$T_{avg} = h \times T_c + (1-h) \times M$$

o h = Hit rate

o $1-h$ = Miss rate

o T_c = Time to access

o M = Miss penalty

USB and SATA

USB Universal Serial Bus

- * I plug and play \rightarrow starts automatically
- * Hot pluggable \rightarrow can be attached or remove anytime
- * Polling \rightarrow device do not interrupt the host

Ex:

Type A, B, C

SATA (Serial Advanced Technology Attachment)

- * used to connect standard storage device, HDD, SSD, RAM
- * So can not be removed easily, should not be removed while running