



Computer Organization and Architecture

Memory Devices

Introduction

Memory Hierarchy,
Random Access Memory,
Read Only Memory,
Serial Access Memory,
Direct Access Memory,
Cache Memory,
Overview of Virtual Memory and
Auxiliary Memory.

Main Memory

Addresses	Values
0000000000000000	01111001
0000000000000001	10010100
0000000000000010	10000000
•	•
•	•
•	•
1111111111111101	11110000
1111111111111110	11100000
1111111111111111	00000111
Memory	

Introduction

memory is used to store data and instruction. Computer memory is storage space in computer where data is to be processed and instructions required for processing are stored.

The memory is divided into large number of small parts.

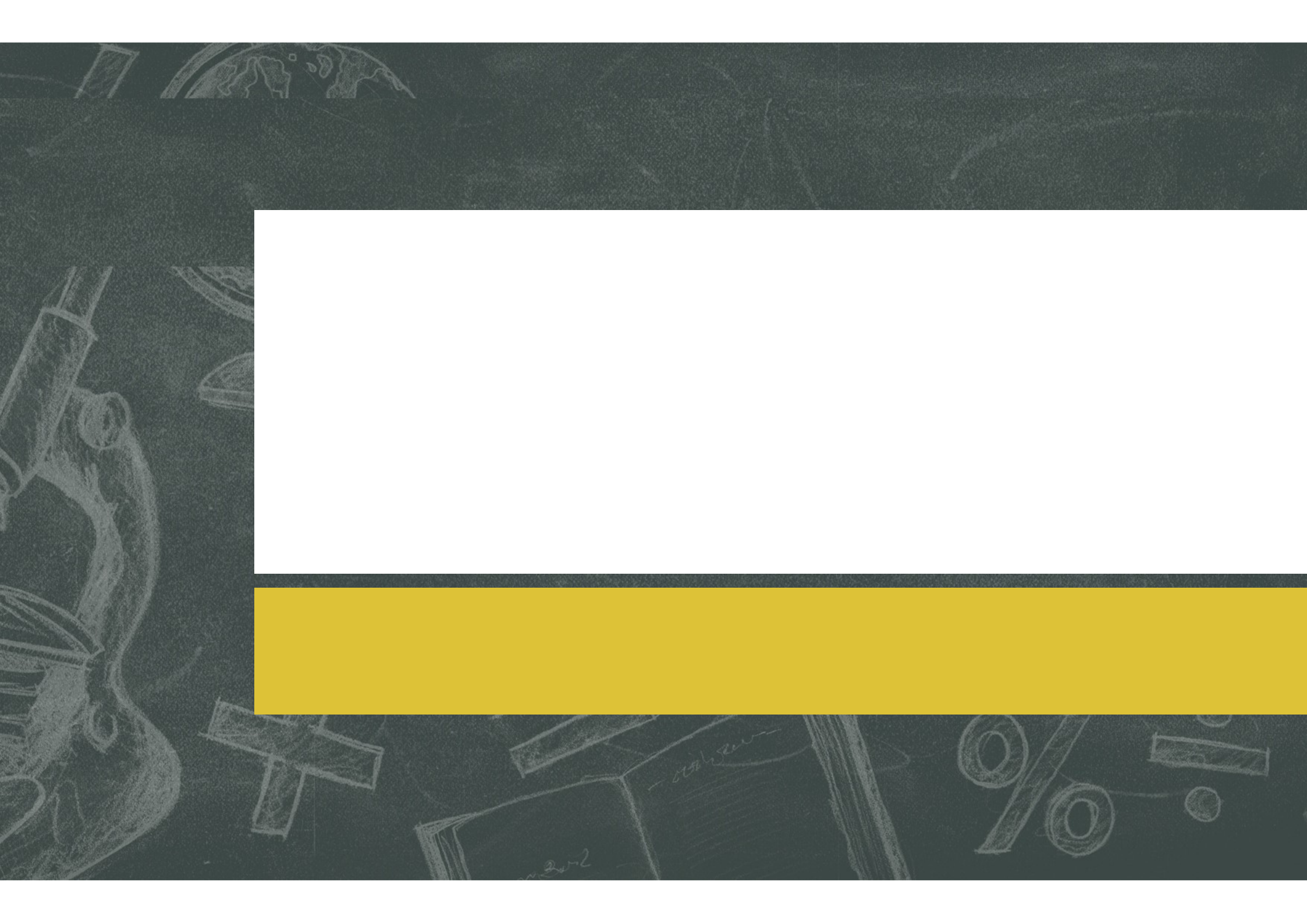
Each part is called a cell. Each location or cell has a unique address which varies from zero to memory size minus one.

For example if computer has 64k words, then this memory unit has $2^{16} = 65536$ memory location. The address of these locations varies from 0 to 65535.

Storage Capacity

Name	Abbreviation	Number of Bytes
Byte	B	1
Kilobyte	KB	1,024 Bytes
Megabyte	MB	1,024 Kilobytes (about 1 million)
Gigabyte	GB	1,024 Megabytes (about 1 billion)
Terabyte	TB	1,024 Gigabytes (about 1 trillion)
Petabyte	PB	1,024 Terabytes (about 1 quadrillion)

<i>Unit</i>	<i>Exact Number of bytes</i>	<i>Approximation</i>
-----	-----	-----
kilobyte	2^{10} bytes	10^3 bytes
megabyte	2^{20} bytes	10^6 bytes
gigabyte	2^{30} bytes	10^9 bytes
terabyte	2^{40} bytes	10^{12} bytes
petabyte	2^{50} bytes	10^{15} bytes
exabyte	2^{60} bytes	10^{18} bytes



Memory Hierarchy

The memory system is a hierarchy of storage devices with different capacities, costs, and access times.

Memory is primarily of two types

Internal Memory – cache memory and primary/main memory

External Memory – magnetic disk / optical disk, etc.

Introduction

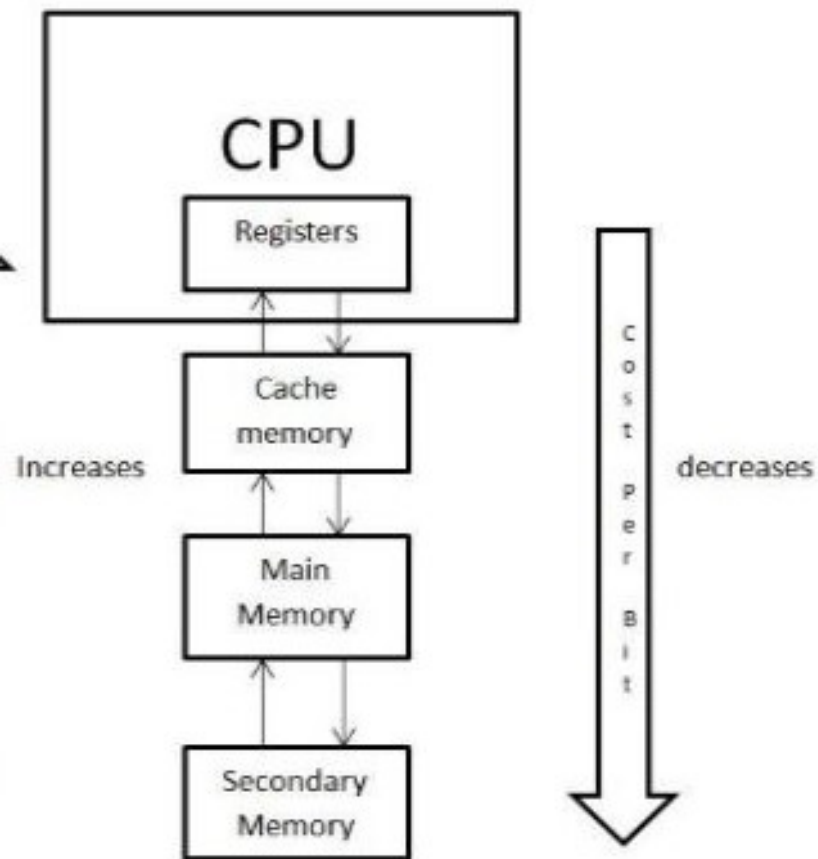
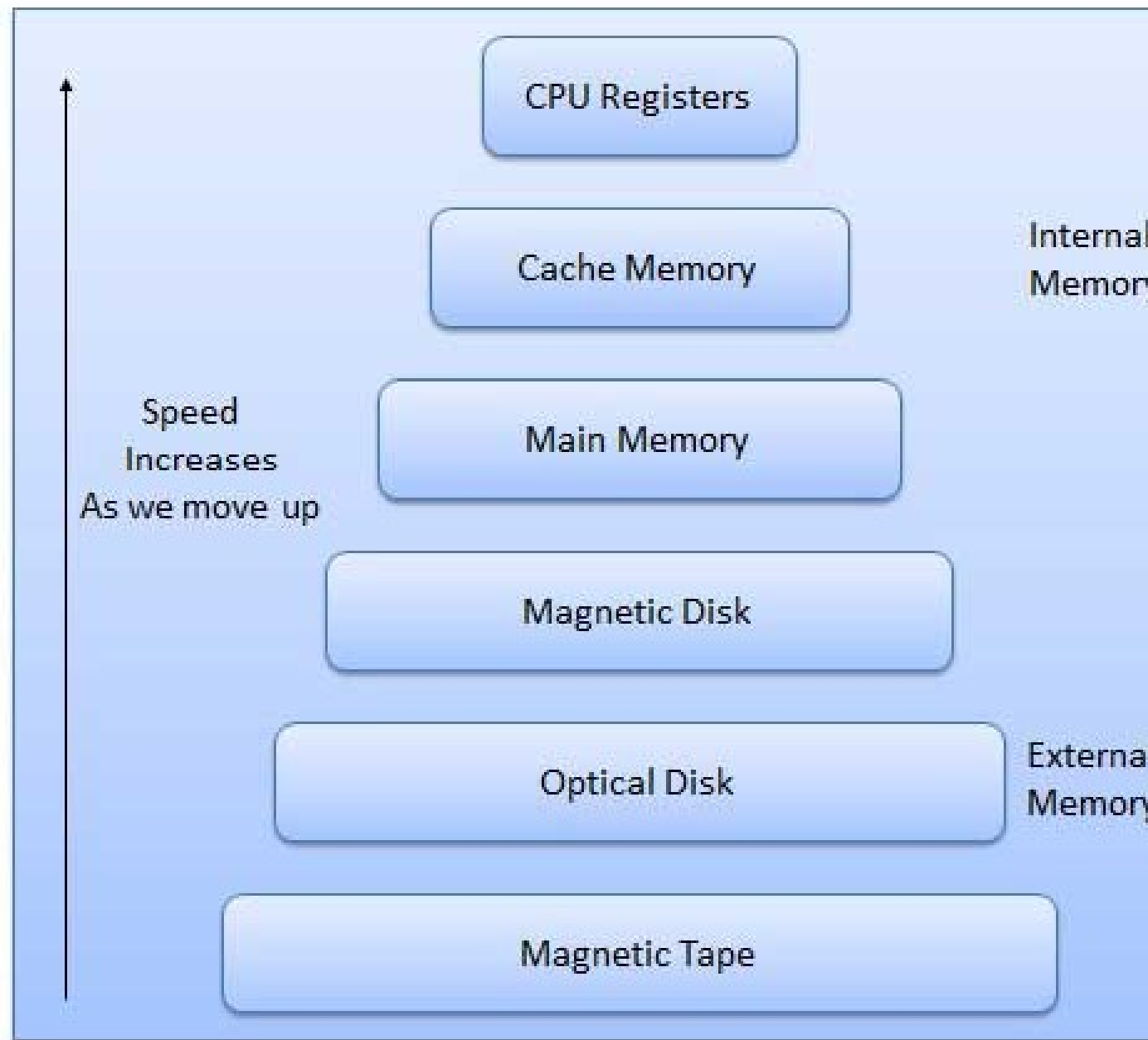


Figure 1



Introduction

Characteristics of Memory Hierarchy are following when we move from top to bottom.

Capacity in terms of storage increases.

Cost per bit of storage decreases.

Frequency of access of the memory by the CPU decreases.

Access time by the CPU increases.

RAM

RAM constitutes the internal memory of the CPU for storing data, program and program result. It is read/write memory. It is called random access memory (RAM).

Since access time in RAM is independent of the address to the word that is, each storage location inside the memory is as easy to reach as other location & takes the same amount of time.

We can reach into the memory at random & extremely fast but it can also be quite expensive.

RAM

RAM is volatile, i.e. data stored in it is lost when we switch off computer or if there is a power failure.

Hence, a backup uninterruptible power system (UPS) is often used with computers. RAM is small, both in terms of its physical size and in the amount of data it can hold.

RAM is of two types

- Static RAM (SRAM)
- Dynamic RAM (DRAM)

RAM

Static RAM (SRAM)

The word **static** indicates that the memory retains its contents as long as power remains applied. However, data is lost when the power gets down due to volatile nature.

SRAM chips use a matrix of 6-transistors and no capacitors. Transistors do not require power to prevent leakage, so SRAM need not have to be refreshed on a regular basis.

Because of the extra space in the matrix, SRAM uses more chips than DRAM for the same amount of storage space, thus making the manufacturing costs higher.

Static RAM is used as cache memory needs to be very fast and small.

RAM

Dynamic RAM (DRAM)

DRAM, unlike SRAM, must be continually **refreshed** in order for it to maintain the data.

This is done by placing the memory on a refresh circuit that rewrites the data several hundred times per second.

DRAM is used for most system memory because it is cheap and small. DRAMs are made up of memory cells. These cells are composed of one capacitor and one transistor.

RAM

Other RAM types

SDRAM (Synchronous DRAM)

EDRAM (Enhanced DRAM)

EDO (Extended Data Out)

FLASH RAM

Ferroelectric RAM

ROM

ROM stands for Read Only Memory. The memory from which we can only read but cannot write on it. This type of memory is non-volatile. The information is stored permanently in such memories during manufacture.

A ROM, stores such instruction as are required to start computer when electricity is first turned on, this operation is referred to as bootstrap.

ROM

PROM (Masked ROM)

The very first ROMs were hard-wired devices that contained a pre-programmed set of data or instructions. These kind of ROMs are known as masked ROMs. It is inexpensive ROM.

PROM (Programmable Read Only Memory)

PROM is read-only memory that can be modified only once by a user. The user buys a blank PROM and enters the desired contents using a PROM programmer.

It can be programmed only once and is not erasable.

ROM

EPROM (Erasable and Programmable Read Only Memory)

The EPROM can be erased by exposing it to ultra-violet light for a duration of upto 40 minutes.

EEPROM (Electrically Erasable and Programmable Read Only Memory)

The EEPROM is programmed and erased electrically. It can be erased and reprogrammed about ten thousand times.

Both erasing and programming take about 4 to 10 ms (millisecond). In EEPROM, any location can be selectively erased and programmed.

Serial Access Memory

Sequential access means the system must search the storage device from the beginning of the memory address until it finds the required piece of data.

Memory device which supports such access is called a Sequential Access Memory or Serial Access Memory. Magnetic tape is an example of serial access memory.

Direct Access Memory

Direct access memory or Random Access Memory, refers to the conditions in which a system can go directly to the information that the user wants.

Memory device which supports such access is called a Direct Access Memory. Magnetic disks, optical disks are examples of direct access memory.

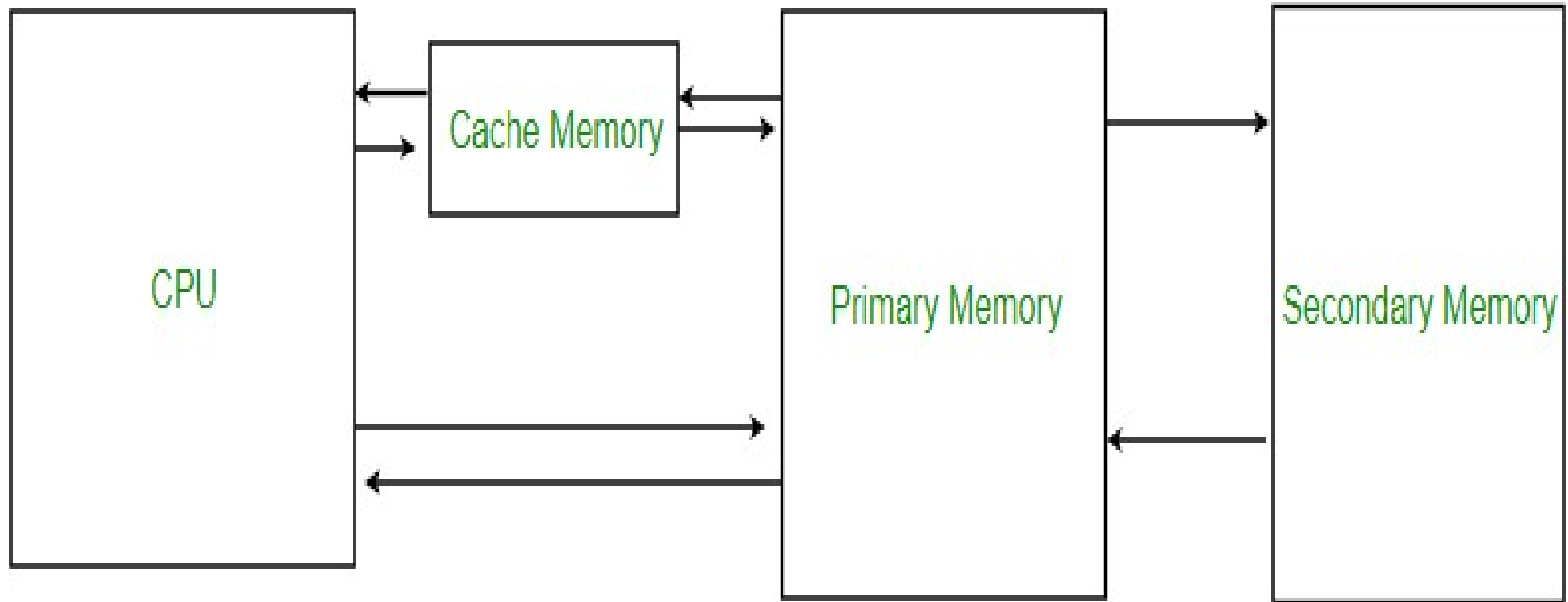
Cache Memory

As CPU has to fetch instruction from main memory speed of CPU depends on fetching speed from main memory. CPU contains registers which have the fastest access but they are limited in number as well as costly. Cache is cheaper so we can access cache. Cache memory is a very high-speed memory that is placed between the CPU and main memory, to operate at the speed of the CPU.

It is used to reduce the average time to access data from the main memory. The cache is a smaller and faster memory which stores copies of the data from frequently used main memory locations.

Most CPUs have different independent caches, including instruction and data.

Cache Memory



Cache Memory

Cache Performance

When the processor needs to read or write a location in memory, it first checks for a corresponding entry in the cache.

If the processor finds that the memory location is in the cache

cache hit has occurred and data is read from cache

If the processor **does not** find the memory location in the cache

a **cache miss** has occurred. For a cache miss, the cache

allocates a new entry and copies in data from main memory

then the request is fulfilled from the contents of the cache.

Cache Memory

The performance of cache memory is frequently measured in terms of a quantity called Hit ratio.

- Hit ratio = $\text{hit} / (\text{hit} + \text{miss}) = \text{no. of hits} / \text{total accesses}$
- Miss ratio = $\text{miss} / (\text{hit} + \text{miss}) = \text{no. of miss} / \text{total accesses}$

Cache Mapping

The three different types of mapping used for the purpose cache memory are as follows:

- Direct mapping
- Associative mapping,
- Set-Associative mapping.

Locality of reference

Since size of cache memory is less as compared to main memory. So to choose which part of main memory should be given priority and loaded in cache is decided based on locality of reference.

. Spatial Locality of reference: If the storage has been accessed then likelihood of accessing the storage nearby that, is high.

. Temporal Locality of reference: It tells us whether memory locations in a program are likely to be accessed again in the near future. A method has temporal locality if it is called repeatedly in a short period of time.

Virtual Memory

Virtual Memory is a storage scheme that provides user an illusion of having a very big main memory. This is done by treating a part of secondary memory as the main memory.

In this scheme, User can load the bigger size processes than the available main memory by having the illusion that the memory is available to load the process.

Instead of loading one big process in the main memory, the Operating System loads the different parts of more than one process in the main memory.

By doing this, the degree of multiprogramming will be increased and therefore, the CPU utilization will also be increased.

Virtual Memory

How Virtual Memory Works?

Whenever some pages need to be loaded in the main memory for the execution and the memory is not available for those many pages, then in that case, instead of stopping the pages from entering in the main memory, the OS searches for the RAM area that are least used in the recent times or that are not referenced and copy that into the secondary memory to make the space for the new pages in the main memory.

Auxiliary Memory

Auxiliary memory is much larger in size than main memory but slower. It normally stores system programs, instruction and data files. It is also known as secondary memory.

Secondary memories cannot be accessed directly by a processor. First, the data/information of auxiliary memory is transferred to the main memory and then that information can be accessed by the CPU.

CDs, DVDs, HDDs are examples of Auxiliary memory.

Auxiliary Memory

Characteristics of Auxiliary Memory are following.

- **Non-volatile memory** – Data is not lost when power is cut off.
- **Reusable** – The data stays in the secondary storage on permanent basis until it is not overwritten or deleted by the user.
- **Reliable** – Data in secondary storage is safe because of high physical stability of secondary storage device.
- **Convenience** – With the help of a computer software, authorized people can locate and access the data quickly.
- **Capacity** – Secondary storage can store large volumes of data in sets of multiple disks.
- **Cost** – It is much lesser expensive to store data on a tape or disk than primary memory.