## Memory Management

Memory is the physical device which is used store program or data on a temporary or permanent basis for use in a computer or other digital electronic device.

There are various types of memories in a computer system and are accessed by various processes for their execution.

→ It is most imp. and more complex task for OS. Memory management involves treating main memory as a resource to be allocated and shared among no. of active processes ie it is the act of managing computer memory.

→ It checks how much memory is to be allocated to processes. It decides which process will get memory at what time. It tracks whenever some memory gets freed or unallocated and correspondingly it updates the status.

### Memory management requirements.

① **Relocation**
→ programmer doesnot know where the program will be placed in memory when it is executed.
→ While the program is executing, it may be swapped to disk and return to main memory at a different location.

② **protection.**
→ processes should not be able to reference memory location in another process without premission.
→ processes should not be able to the crop the OS.

③ **sharing**
→ Allow several processes to access the same portion of memory in a controlled way.

⇒ Better to allow each process accessed to the same copy of the program rather than have their own separate copy

④ **Physical organization.**
⇒ Memory available for a program plus it's data may be insufficient.
⇒ Programmer doesnot know how much space will be available.

⑤ **logical organization.**
⇒ Program return in modules.
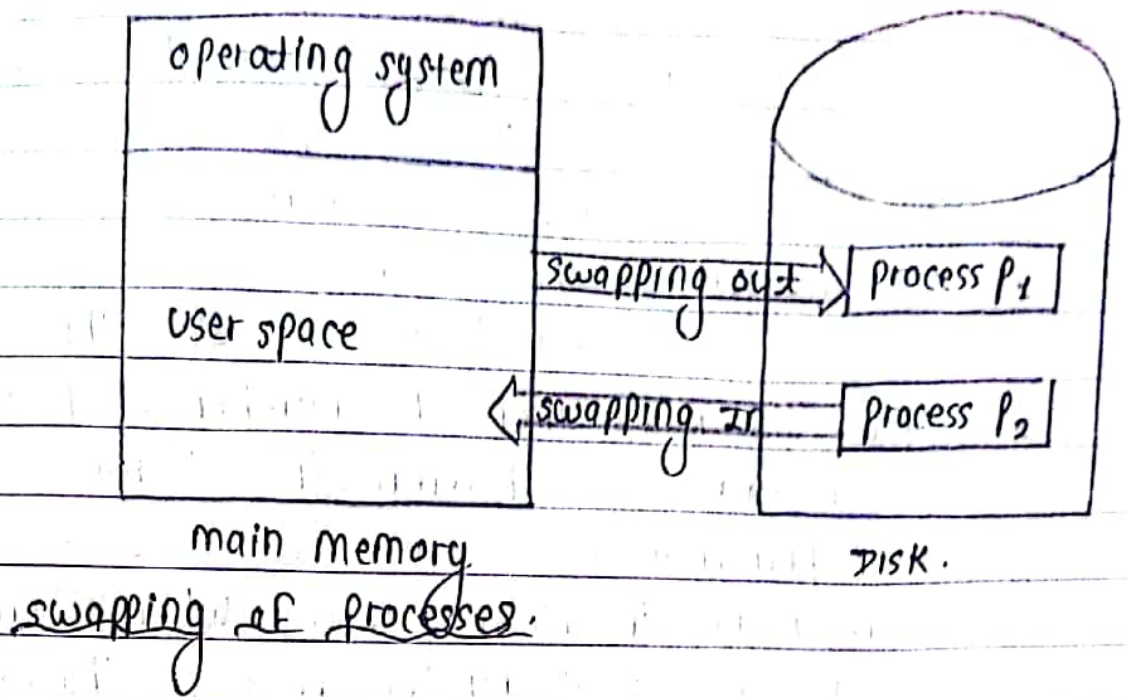⇒ Modules can be return and complied independently.

⑥ **Swapping.**
⇒ Any process which needs to be executed should be kept in main memory.
⇒ However, if the process in the ^main memory is not getting executed for some reason such as waiting for an event to occur or waiting for an I/O, then such process can be moved back to disk (swapout) and any process in the disk which is ready to be executed can be looded into the main memory (swap in)
⇒ This process of swapping out & swapping in of processes beth main memory and disk in order to reduce the CPU ideal time is called as swapping.
⇒ The phenomenon of moving processes from main memory to secoundary memory is called swapping out, while the
⊕ the phenomenon of moving process from secondary memory back to primary memory is called swapping in.

| operating system |
|---|
| user space |

main memory

swapping out → process $P_1$

← swapping in ← process $P_2$

DISK .

<u>swapping of processes.</u>

## Advantages of swapping.

The advantages of swapping are listed below:

→ Swapping helps in acheiving the goal of maximum CPU utilization.

→ Swapping ensures proper memory availability for every process that needs to be executed.

→ Swapping helps to avoid the problem of process starvation means a process should not take much time for execution so that the next process should be executed.

→ CPU performs various tasks simultaneously with the help of swapping so that their processes donot have to wait much longer before execution.

→ Swapping ensures proper RAM (main memory) utilization.

→ Swapping creates a dedicated disk partition in the hard drive for swapped processes which is called swap space.

→ Swapping in OS is an economical process.

→ Swapping method can be explained on priority-based sheduling where a high priority process is swapped in and the low priority process is swapped out which improves the performance.

## Disadvantages of swapping.

Some disadvantages of swapping are listed below:

→ If the system deals with the power-cut during bulky swapping activity then the user may lose all information which is related to the program.

→ If the swapping method uses an algorithms that is not up to the mark then the number of page faults can be increased and therefore this decreases the Complete performance.

→ There may be inefficiency in case when there is some common resources used by the processes that are participating in the swapping process.
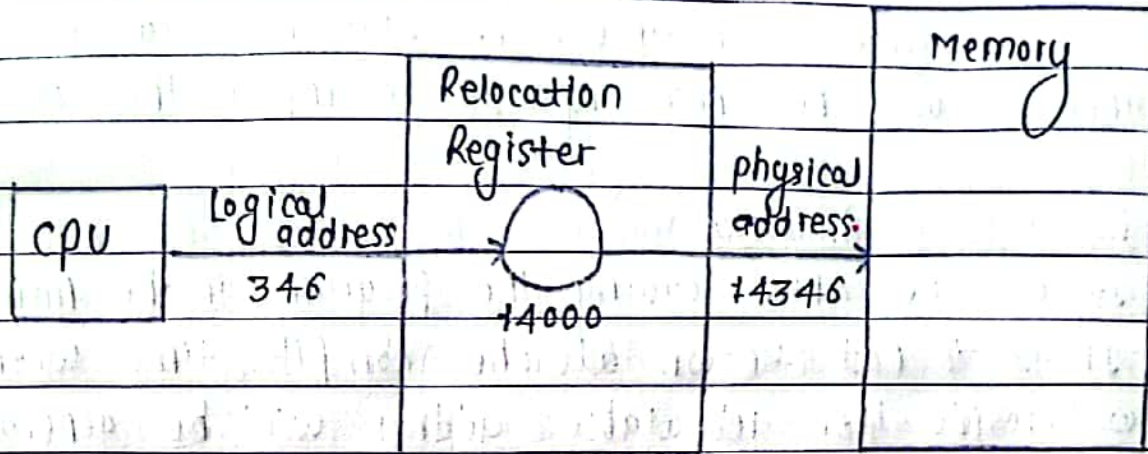
## physical Address and logical address.

The physical address identifies the physical location of required data in memory. The user never directly deals with physical address but can access it by its corresponding logical address.

→ The logical address is an address that is generated by the CPU during program execution.

→ The logical address is a virtual address as it doesn't exist physically and therefore it is also known as virtual address.

## Memory Management Unit (MMU)

→ The user only things that he is accessing the data from the logical address.

→ The transaction from the logical to the physical address is done by special equipment in the cpu that is called memory management unit
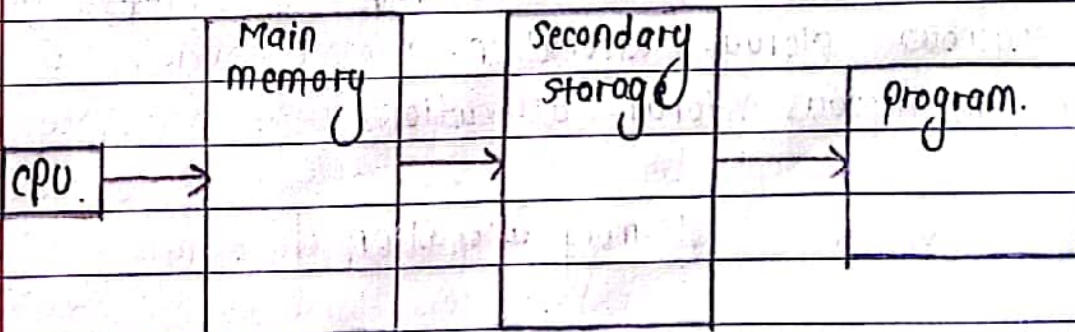
## Mapping virtual address to physical address.

```
                                    Relocation              ┌──────────┐
                                    Register                │ Memory   │
           ┌──────┐  Logical               physical         │          │
           │ CPU  │  address      ( )      address          │          │
           │      │  346                    74346           │          │
           └──────┘            14000                        └──────────┘
                                MMU
```

→ The CPU generate the logical address (346)

→ The MMU will generate the base address (74000) which is stored in the relocation register.

→ The value of relocation register (here 74000) is added to the logical address to get the physical address (ie 74000+74346).

## Address binding.

→ Address binding is the process of mapping from one address space to another address space.

```
              ┌──────────┐      ┌──────────┐
              │ Main     │      │ Secondary │      ┌──────────┐
              │ memory   │      │ storage   │      │ Program. │
  ┌─────┐  →  │          │  →   │          │  →   │          │
  │ CPU.│     │          │      │          │      │          │
  └─────┘     │          │      │          │      └──────────┘
              │          │      └──────────┘
              └──────────┘
```

→ The addresses use in a source code. The variable, names, constants and instruction levels are the basic elements of the symbolic address space.

### (1) Compile Time Binding.

→ If it is known at compile time where Process / Program reside in memory then absolute address is generated.

→ If, at some later time, the starting location change. then it will be necessary to recompile this code.

### (2) Load Time Binding

→ It is done after loading the program in the memory.

→ If it is not known at the compile time where Process will reside then relocatable address will be generated.

→ Loader translate the relocatable address to absolute address.

→ once the process loads, it doesnot move in memory.
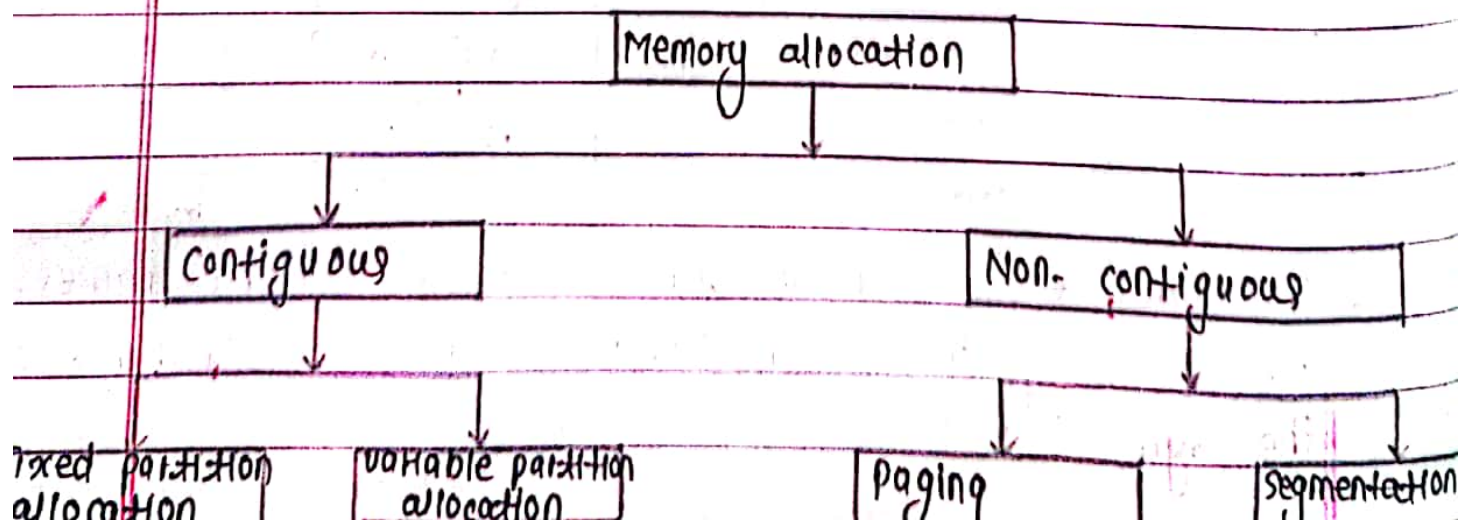
### (3) Run Time Binding.

If process can be moved from one memory to another during execution then binding is done at Run time.

### Memory Allocation Techniques.

Since the main memory must accommodate both the operating system and the various user process, main memory has to be allocated in the most efficient way possible.

Memory allocation is of two types :

① contiguous storage allocation.

② Non-contiguous storage allocation.

```
                        ┌─────────────────────┐
                        │  Memory allocation  │
                        └─────────────────────┘
                                  │
               ┌──────────────────┴──────────────────┐
               ▼                                      ▼
      ┌────────────────┐                    ┌──────────────────┐
      │  Contiguous    │                    │  Non-contiguous  │
      └────────────────┘                    └──────────────────┘
          │                                      │
    ┌─────┴──────┐                        ┌──────┴──────┐
    ▼            ▼                        ▼             ▼
```

| fixed partition allocation | variable partition allocation | Paging | Segmentation |

## Difference bet^n physical address and logical address.

| parameter | Logical address | physical address. |
|---|---|---|
| Basic | Logical address is generated by CPU in perspective of program | Physical address is the location that exists in memory unit |
| Address space | Logical address space is set of all logical addresses generated by CPU in reference to program | Physical address is set of all physical addresses mapped to the corresponding logical addresses. |
| Visibility | User can view the logical address of the program. | User can never view physical address of the program. |
| Generation | Generated by CPU | Computed by MMU. |
| Access | The user can use the logical address to access the physical address. | The user can indirectly access physical address but not directly. |
| Editable | Logical address can be change. | Physical address will not change. |
| Also called | Virtual address | Real address. |

## Contiguous storage allocation.

In contiguous storage allocation, each process occupy a single contiguous section of memory. In a multiprogramming environment, several programs resides in primary memory at a time and the cpu passes it's control rapidly bet^n these programs. To support multiprogramming one idea is to divide the main memory into several partition each of which is allocated to a single process. Depending on how and when partition are created. There are two types of memory partition.

① fixed or static partitioning
② variable or dynamic partitioning

### ① fixed or static partitioning (Multiprogramming with fixed partitioning)

In fixed partitioning, main memory is divided into a no. of static partition at system generation time. There are two alternatives for fixed partitioning is equal size partition and unequal size partition.

| OS |
|----|
| 8 MB |
| 8 MB |
| 8 MB |
| 8 MB |

fig ⓐ Equal size partition

| OS |
|----|
| 2 MB |
| 6 MB |
| 8 MB |
| 16MB |

fig ⓑ unequal size partition.

Fixed Partitioning of 32 MB memory

To understand internal fragmentation in detail, consider an example in which we have a physical memory with the following fixed partition.



When a process or program of size 80K (m) arrives, it is accommodate in partition I but process I is 100K(N) size so, M<N therefore M can be given partition process I the left over unused space is (N-m) = (100-80)K = 20K. This causes internal fragmentation of 20K here.

② **Variable or dynamic partitioning.**
To overcome the problem with fixed partitioning, the concept of dynamic partition was introduced. With dynamic partitioning, the partition are variable length and number. When a brought into main memory, it is allocated exactly as much memory as it require and no more.

**External fragmentation**
External fragmentation exist when there is a enough total memory space to satisfy a requesting process but the available space are non-contiguous; storage is fragmented into a large no. of small holes (free spac

To solve this problem, compaction is employed. Compaction is a technique by which the processes are relocated in such as way that the small chunks of free memory are made contiguous to each other clubbed together into a stg single free portition that may be big enough to accommodate additional process as an example consider a memory that has three holes of size 30K, 20K, 50K that have been compacted into one large hole or block of 100K which is shown in figure.
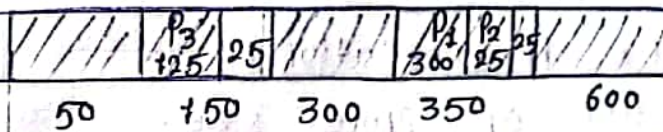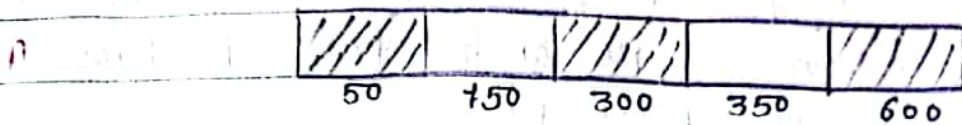


Memory placement algorithms.

In order to load a new process, OS checks for the free memory partition with the help of PDT (partition description table) if the search is found to be successfull then the entry is marked as "allocated". When the process terminates or swapped out, it is updated as "free". Three most common stati ^ to allocate free partition to the new process are

① First fit

② Best fit

③ Worst Fit.



| 50 | 150 | 300 | 350 | 600 |

P2 25 | P3 125 | | P1 300 | P4 50 |

| 50 | 150 | 300 | 350 | 600 |

P3 125 | 25 | | P1 300 | P2 25 |

| 50 | 150 | 300 | 350 | 600 |

P2 25 | P3 125 | | P1 300 | P4 50 |

| 50 | 150 | 300 | 350 | 600 |

- **first fit**

In a first fit approach, the first free partition is allocated large enough to accomodate the process.

- **Best fit**

In best fit approach, the smallest free partition is allocated that meets the requirement of the process

- **worst fit**

In a worst fit approach, the largest available partition is allocated to the newly entered process.

**Advantages of fixed or static partitioning.**

→ Easy implementation.
→ External fragmentation.
→ Internal fragmentation.
→ limiting the process size.
→ Lesser degree of multiprogramming.
→ Lower OS overhead.

## Disadvantages

→ Requires entire program to be stored contiguously.

→ Jobs are allocated space on the basis of first available partition of required size.

→ Work well only if all the Jobs are of some size or if the sizes are known ahead of time.

## Advantages of variable or dynamic partitioning.

→ No internal fragmentation.

→ No limitation on process size.

→ External fragmentation.

→ Complex memory allocation.

→ Dynamic degree of multiprogramming.

## Disadvantages

→ Limited flexibility in accomodating changing memory requirements of processes

→ Wastage of memory resources when a partition is not fully utilized.

→ High fragmentation when processes of different sizes are present in the system.
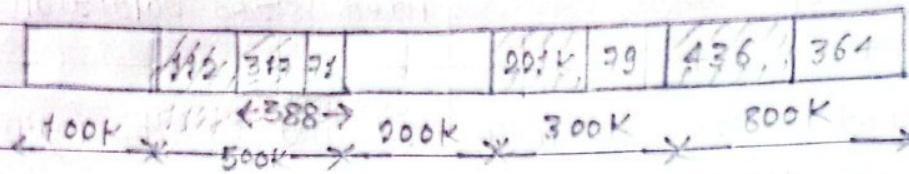
Differences bet<sup>n</sup> fixed partitioning and variable partitioning

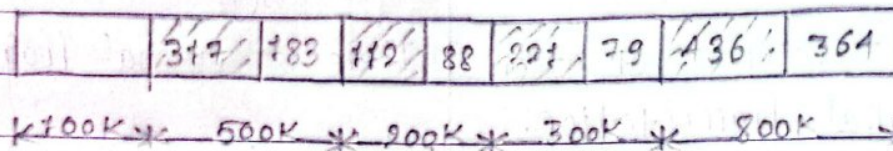| Fixed partitioning | Variable partitioning |
|---|---|
| ① In a multi-programming with fixed partitioning the main memory is divided into fixed sized partitions | In multi-programming with variable partitioning the main memory is not divided into fixed sized partitions. |
| ② It doesnot utilize the main memory effectively. | It utilizes the main memory effectively. |
| ③ There is presence of external and internal fragmentation. | There is external fragmentation. |
| ④ Degree of multi-programming is less. | Degree of multi-programming is higher. |
| ⑤ It is more easier to implement. | It is less easier to implement. |
| ⑥ There is limitation on size of process. | There is no limitation on size of process. |
| ⑦ only one process can be placed in a partition. | In variable partitioning, the process is allocated chunk of free memory. |

Q. Given memory partition is 100K, 500K, 200K, 300K and 800K (in ordered) how would each of first fit, best fit and would worst fit algorithm places. processes of 112K, 317K, 221K and 436K in order.
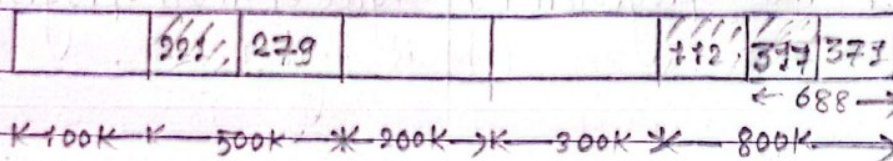
| 100K | 500K | 200K | 300K | 800K |
|------|------|------|------|------|

first

| | 112 | 317 | 21 | | 221 | 79 | 436 | 364 |

100K | ←-588-→ 200K | 300K | 800K
        500K

best

| | 317 | 183 | 112 | 88 | 221 | 79 | 436 | 364 |

100K | 500K | 200K | 300K | 800K

worst

| | 221 | 279 | | | 112 | 317 | 371 |

←688→
100K | 500K | 200K | 300K | 800K

## Non-contiguous storage allocation:

Employing compaction technique to avoid external fragmentation can be expensive. Another possible soln to external fragmentation is to have non-contiguous storage allocation. There are two main methods of non-contiguous storage allocation.

① paging

A memory management skim scheme that premits the physical address space of a process to be non-contiguous is called paging. paging avoids external fragmentation

500
317
183

388
312
71

300
221
79
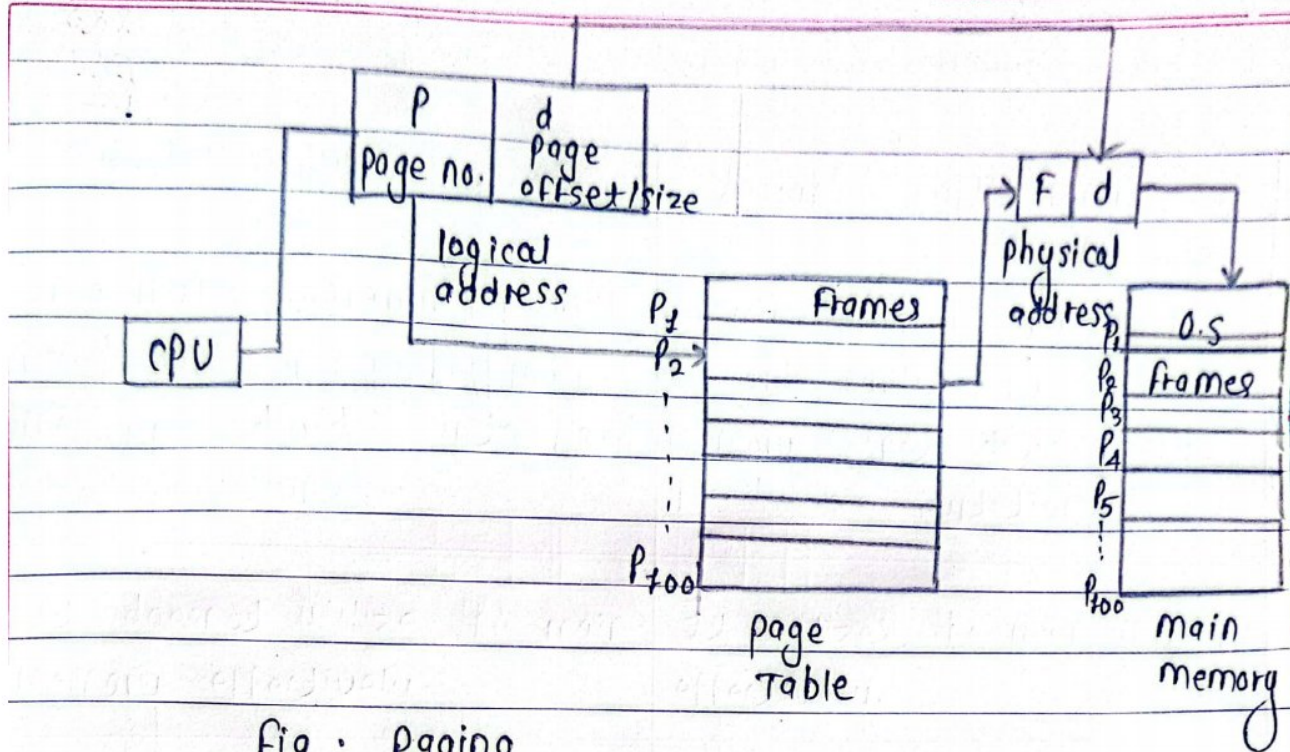
800
112

DATE : __/__/__
PAGE : _____

fig : Paging

as well as the need for compaction.

The basic method for implementing paging involves breaking physical memory into fixed size block called frames and breaking logical memory into block of same size called pages.

The hardware support for paging is illustrate in figure. Every address generated by the CPU is divided into a two parts : a page number (p) and page offset (d). The page no. is used as an index into a page Table. The page Table contains the base address of each page in physical memory. This base address is combined with the page offset to define the physical memory address that is sent to a memory unit.