

UNIT-III

memory management strategies

Introduction:

→ strategies divided into several categories

- * Fetch strategies.
- * Placement strategies
- * Replacement strategies.

*

* Fetch strategies we are taking the data.

- it demand or anticipates data.
- decides which pieces of data to load resulting in memory management

✓ Placement Strategies

it decides where in main memory to place incoming data.

✓ Replacement Strategies

- it decide which data is to remove from main memory to make more space.

Introduction:

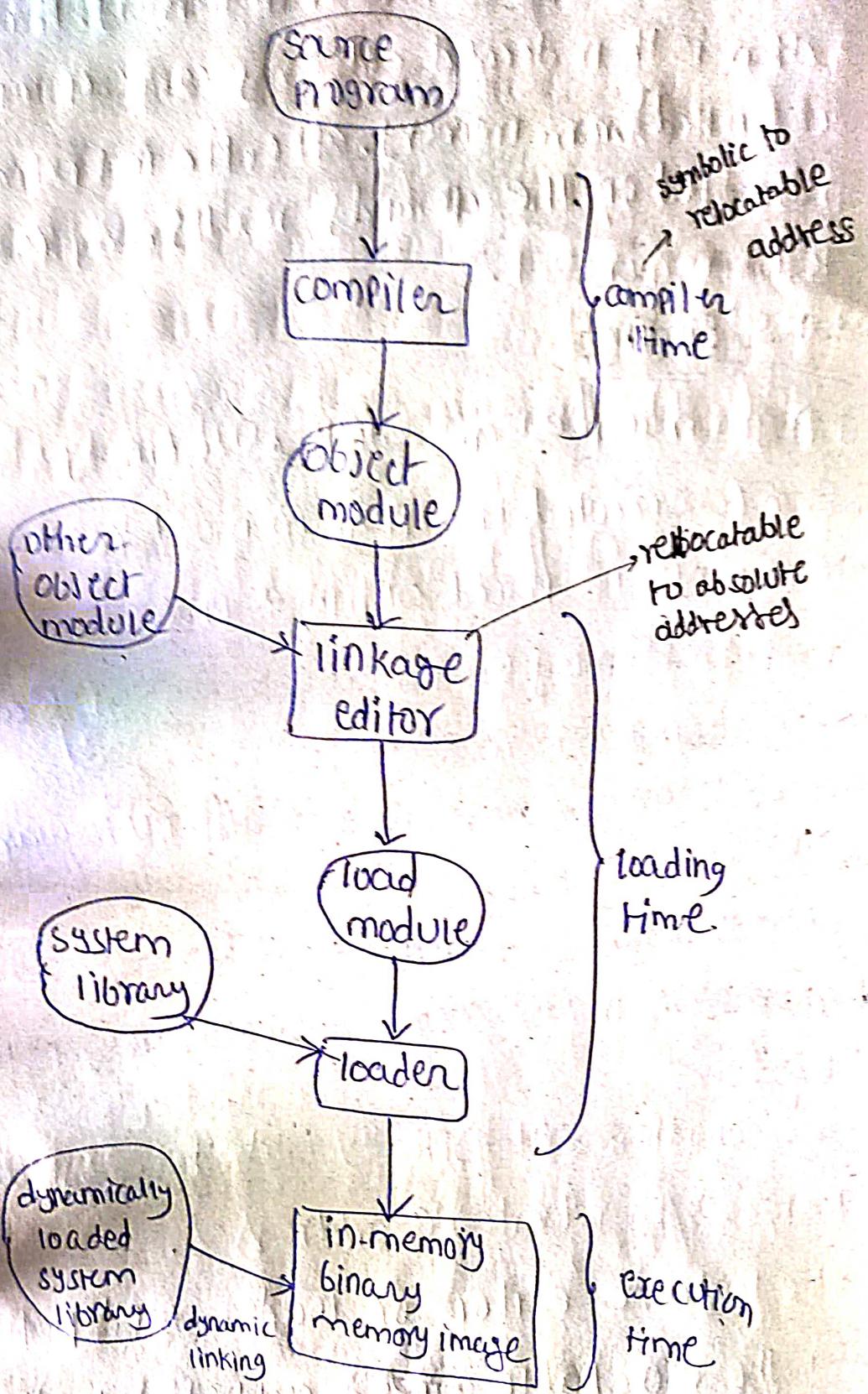
→ memory is a large array of words or bytes each with its own address.

→ the CPU ~~refeferences~~ fetches instructions from memory according to the value of the program counter

→ First fetches an instruction from memory. The instruction is then decoded and may cause operands to be fetched from memory. After the instruction has been executed on the operands, results may be stored back in memory.

Address Binding

- A program resides on a disk as a binary executable file. To be executed the program must be brought into memory and placed within a process.
- The process may be moved between disk and memory during its execution.
- Addresses may be represented in different ways during these steps.
 - Addresses may be represented in different ways
 - Addresses in the source program are generally symbolic.
 - A compiler will typically bind these symbolic addresses to relocatable addresses.
 - The linkage editor or loader will in turn bind the relocatable addresses addressed to absolute addresses.
 - Each binding is mapping from one address space to another.
 - * compile time
 - * load time
 - * execution time.

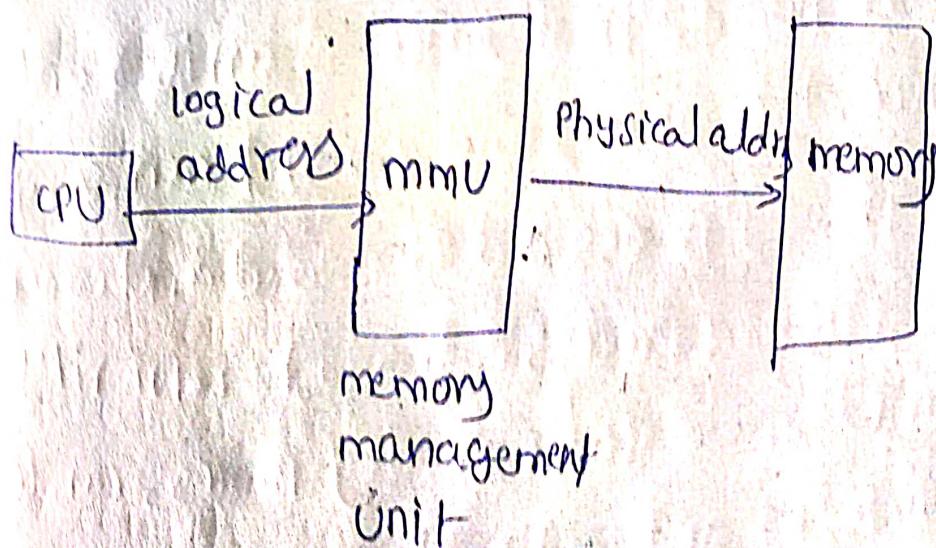


Logfe-
fig 2 multistep processing of a user program.

logical versus physical address space

- the address space is generated by the CPU is called logical address space.
- the address space is generated by the memory unit is called physical address space.

Logical	Physical
→ generated by the CPU; it also referred as virtual address	→ address seen by the memory unit
→ the logical address space is the set of all logical address generated by the program.	→ physical address space is the set of all physical address generated by the program.
⇒ logical & physical address are the "same" in compile time and load time; address binding scheme	
→ the logical and physical address are different in execution time	



Swapping

→ moving data from one device to another device temporarily.

→ it is mechanism in which a process can be swapped temporarily out of main memory to secondary memory. later again swaps back to main memory.

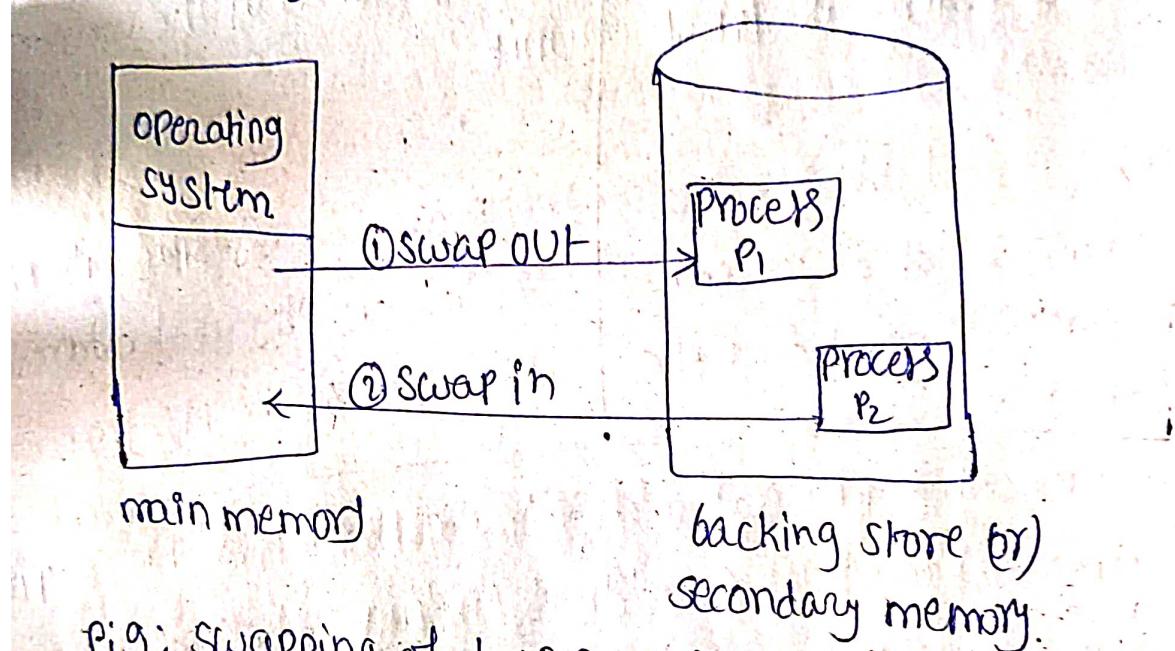
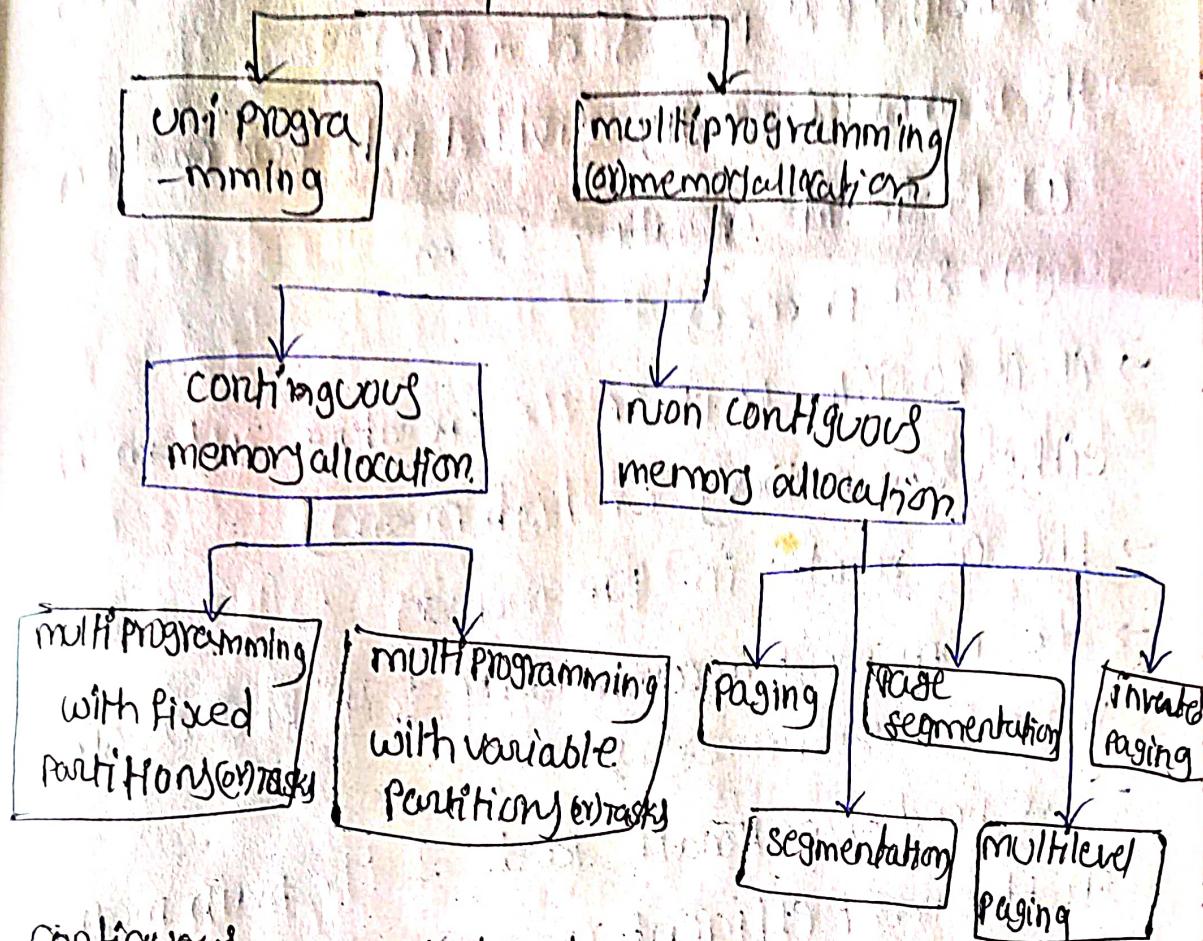


fig: swapping of two processes

Ex: Swapping policy is used for priority based scheduling algorithms. If higher priority process arrives and wants service, the memory management can swap out the lower priority process and then load and execute the higher priority process. When higher priority process finishes, the lower priority process can be swapped back in and continued. This variant of swapping is sometimes called swapout (or) roll out and swapin (or) roll in.

memory allocation

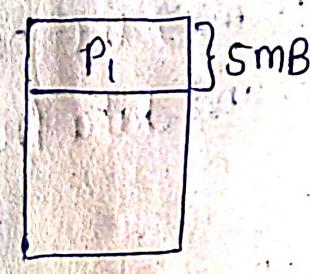
Memory management techniques



contiguous memory allocation

The operating system can allocate a single single contiguous block of memory for the processes

Eg The process is P_1 it having 5mb memory allocated by the OS.



memory

This are two types

① multiprogramming with partitions (or) tasks

② multiprogramming with variable partitions (or) tasks

- contiguous memory allocation is a memory allocation method that allocate a single contiguous section of memory to a process or a file.
- the memory is usually divided into two partitions
 - (1) Resident operating system
 - (2) User processes
- it is possible to place the operating system in either low memory or high memory
- the OS resides in low memory
- the MM must accommodate both the OS and the various user processes.
- we need to allocate parts of the main memory in the most efficient way possible, this section explain one common method, contiguous memory allocation.

* Memory Protection

- if the OS is residing in low memory and the user processes are executing in high memory, we need to protect the operating system code and data from changes by the user processes
- we also need to protect the user processes from one another
- we can provide this protection by using a relocation registers

→ The relocation registers contains the value of the smallest physical address.

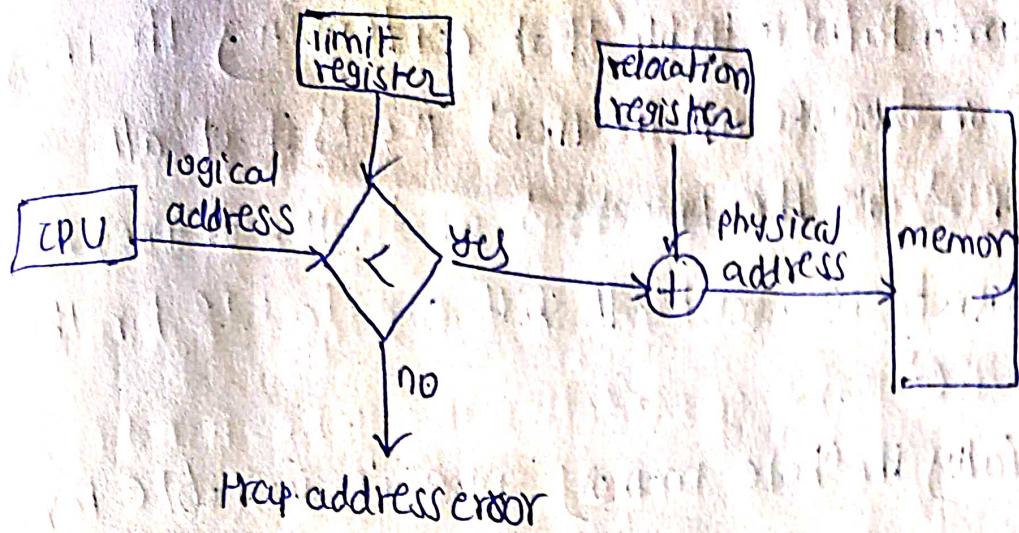


fig: hardware support for relocation and limit registers

* memory allocation

Whenever a process is request to execute that process must be stored into a variable memory, i.e RAM

→ This are two types:

- ① multiprogramming with fixed tasks (MFT).
- ② multiprogramming with a variable no. of tasks (MVT)

① MFT (static)

→ MFT in the fixed size partition, the memory is divided into fixed sized blocks and each block contains exactly one process. But the fixed sized partition will limit the degree of multiprogramming as the no. of the partitions will decide the no. of processes.

④ MVT (Dynamic)

In the variable size partition method, the OS maintains a table that contains information about all memory parts that all occupies by the processes and all ~~memory~~ information about all still available for the processes.

→ Initially the hole memory space is available for user processes as a large block called hole,

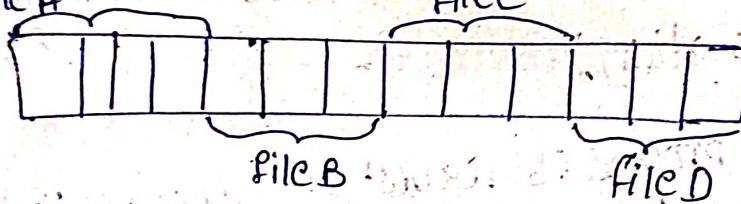


Fig: contiguous memory allocation of 4 files

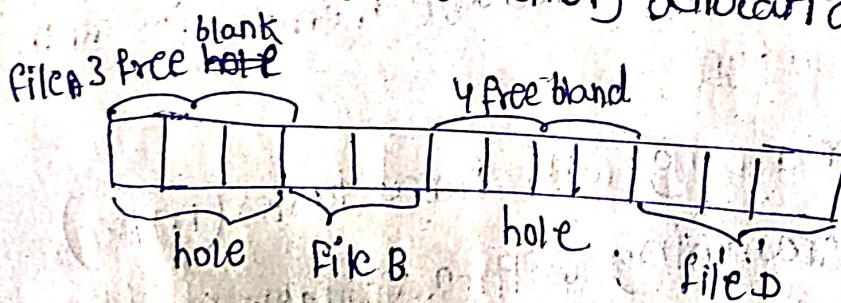


Fig: when the file A and C terminates the release the memory creating hole,
Dynamic storage allocation:

→ The First fit

* best fit

* worst fit

This strategies are the ones most commonly used to select a free hole from the set of available holes

* first fit

- allocate the "first hole" that is big enough.
- searching can start either at the beginning of the set of holes or at the location where the previous first fit search ended.

* Best fit

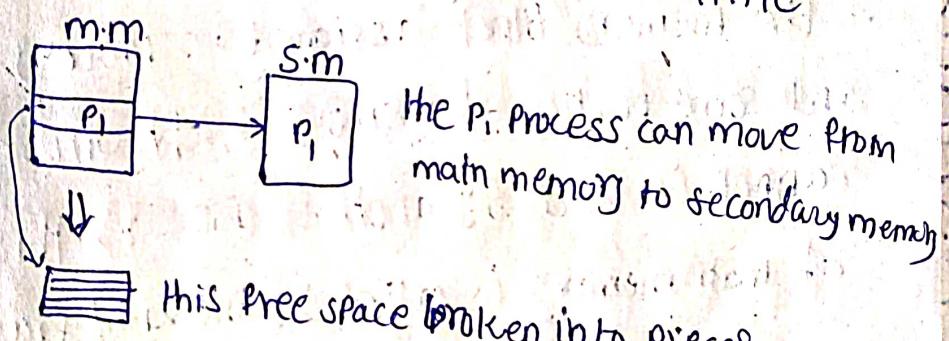
- allocate the "smallest hole" that is big enough.

* worst fit

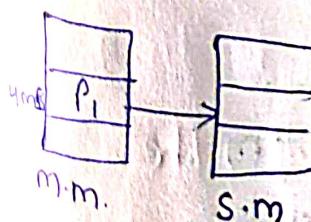
- allocate the "largest hole".

Fragmentation

- a processes are loaded & removed from memory the free memory space is broken into little pieces.

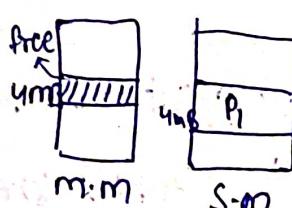


after some times that processes cannot be allocated to memory because of small size and the memory block ~~is~~ remains "unused" this problem is called fragmentation.



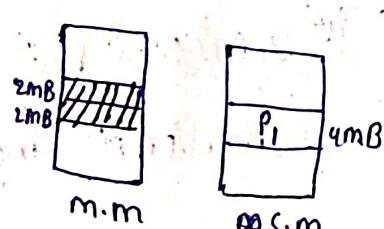
Stage 1

The P_1 is swap
from m.m to s.m



Stage 2

The m.m have a
free space



Stage 3

The P_1 is swap from s.m to m.m
it can't be possible.

→ The fragmentation can be divided into 2

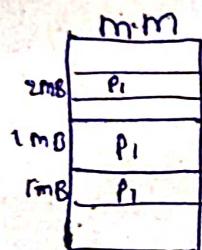
① External Fragmentation

② Internal Fragmentation

① External Fragmentation

Total memory space is enough to satisfy a request or reside a process in it but it is not contiguous so it can not be used. Then it called external fragmentation.

Ex:-

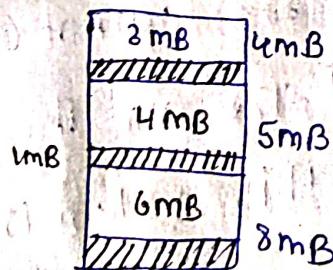


4mB
← P_1 request.

② Internal fragmentation

→ The memory block assigned to process is bigger and some portion of memory is left unused, as it cannot be used by another process. This wastage of main memory is called as internal fragmentation.

Ex:-



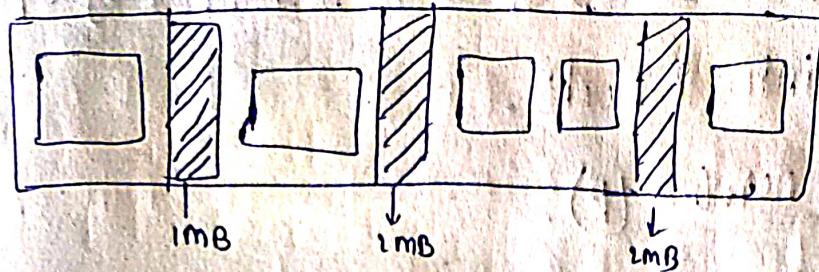
$$P_1 = 4 \text{ mB}$$

$$P_2 = 6 \text{ mB}$$

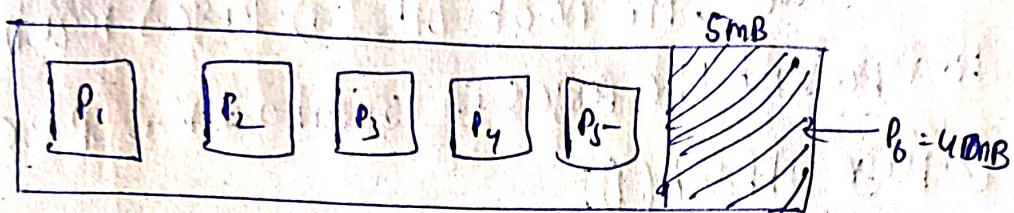
$$P_3 = 3 \text{ mB}$$

→ It overcome this fragmentation problem the fragmentation memory compact (or) shiftful.

Fragmented memory before compaction



Fragmented memory after compaction



- External fragmentation can be reduced by compaction or shuffle free memory together is one large block.
- Internal fragmentation can be reduced by effectively assigning the smallest partitions but large enough for the process.

Paging

- Paging is a memory management technique in logical memory which process address space is broken into blocks of the same size called pages.
- The size of the process is measured in the no. of pages
- Similarly, physical memory is divided into small fixed sized blocks of memory called frames and the size of a frame is same as that of page.
- Optimum utilization of the main memory and to avoid external fragmentation

→ The logical address space is generated by CPU is divided into 2 parts.

① Page number (P):

The page number is used as an index into a page table. The no of bits (or) pages required to represent the pages in logical address space (or)

② page offset (d):

No of bits required to represent particular page size of logical address (or) page offset.

Page number (P)	Page offset (d)
-----------------	-----------------

→ The physical address space is generated or divided by the memory. This are two parts.

① Frame number (F):

The no of bits required to represent the frame of Physical address (or) Frame number

② frame offset (d):

The no of bits required to represent particular frame or frame size of Physical address.

frame number (F)	frame offset (d)
------------------	------------------

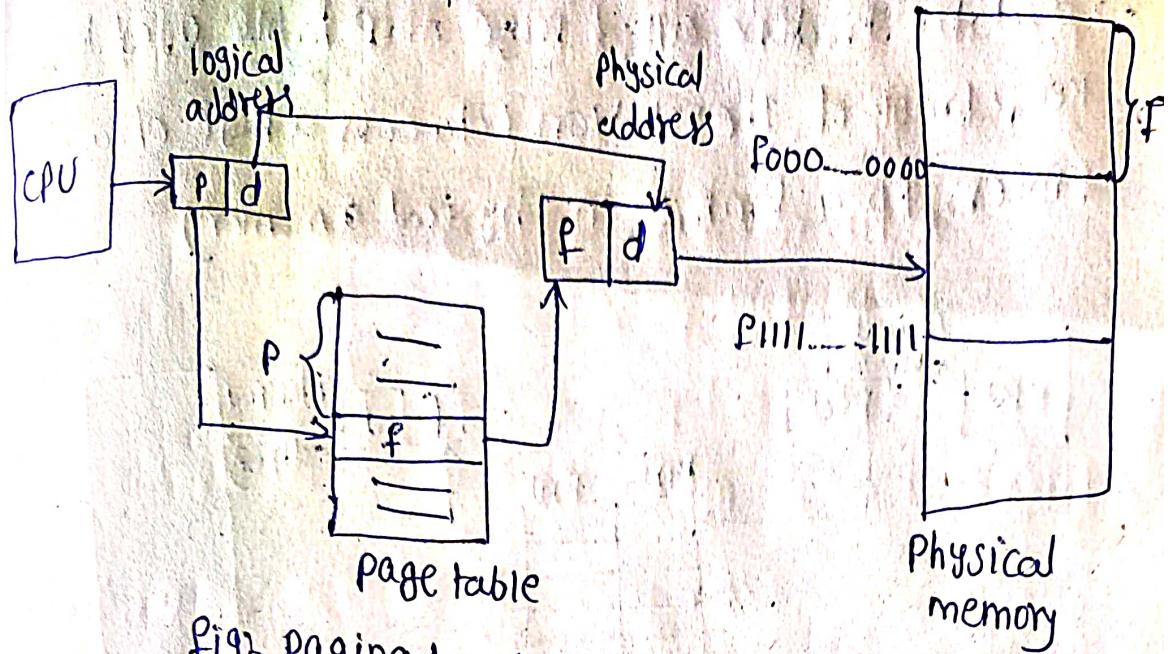


fig7 Paging hardware

- The page number is used as an index into a page table. The page table contains the base address of each page in physical memory.
- This page base address is combined with the page offset to define the physical memory address that is sent to the memory unit.

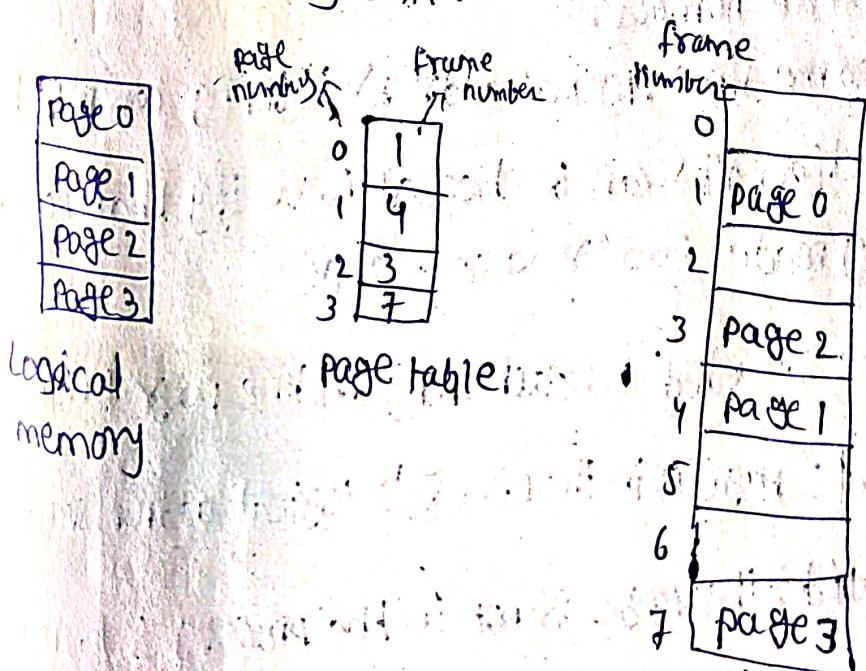
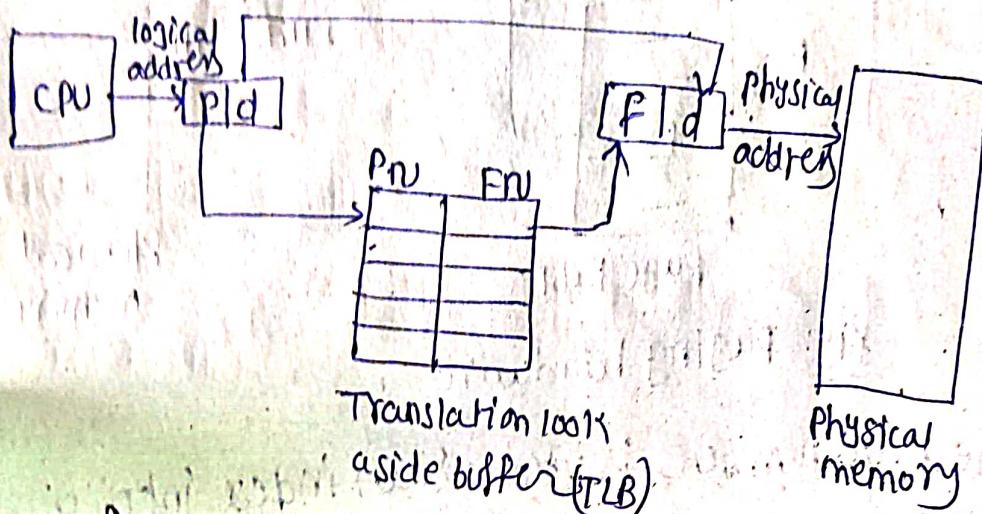


fig8 Paging model of logical and physical memory.

- Hardware support
- If available that particular process will execute
 - a data structure called page map table used to keep track of the relation between a page of a process to a frame in physical memory



For Paging hardware with TLB

Protection

a page to be read-write or read only: every reference to memory goes through the page table to find the correct frame number.

→ at the same time physical address is being computed, the protection bits can be checked to verify that no writes are being made to a read only page.

→ generally attached to each entry in the page table

① Valid: Page is in the process's logical address space

② Invalid: The page is not in the process's logical address space.

Page number	Frame number	Valid - invalid bit	Page
0	2	v	page 0
1	3	v	page 1
2	4	v	page 2
3	7	v	
4	8	v	
5	9	v	
6	0	i	page 3
7	0	i	page 4
			Pages
			Page n

Page Table.

: Fig:- valid(v) or invalid(i) bit in a page table.

Advantages & disadvantages:

- * Paging reduces external fragmentation but still suffer from internal fragmentation.
- * It is simple to implement.
- * Swapping becomes very easy because equal size of pages and frames.
- * It is not suitable for small RAM's.
- * Paging is simple to implement an efficient memory management technique.
- * Page table requires extra memory space, so may not be good for a system have small RAM.

Structure of the Page Table

The most common techniques for structuring the page table are

- * Hierarchical paging
- * Hashed page tables
- * Inverted page tables.

Hierarchical Paging

- another name for Hierarchical Paging is multiple level paging.
- The page table is too big to fit in a contiguous space, so we may have a hierarchy with several levels.
- in this type of paging the logical address space is broken up into multiple pagetables.

→ 16 bit page table is used.

& 32 bit address space

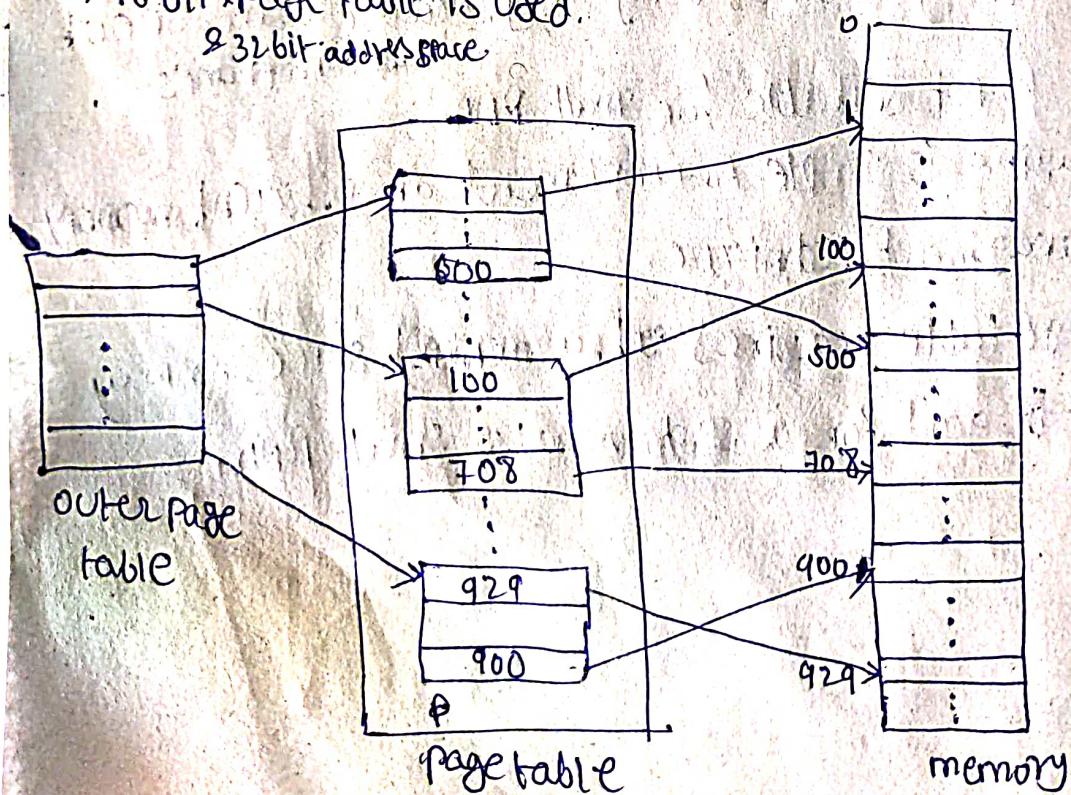


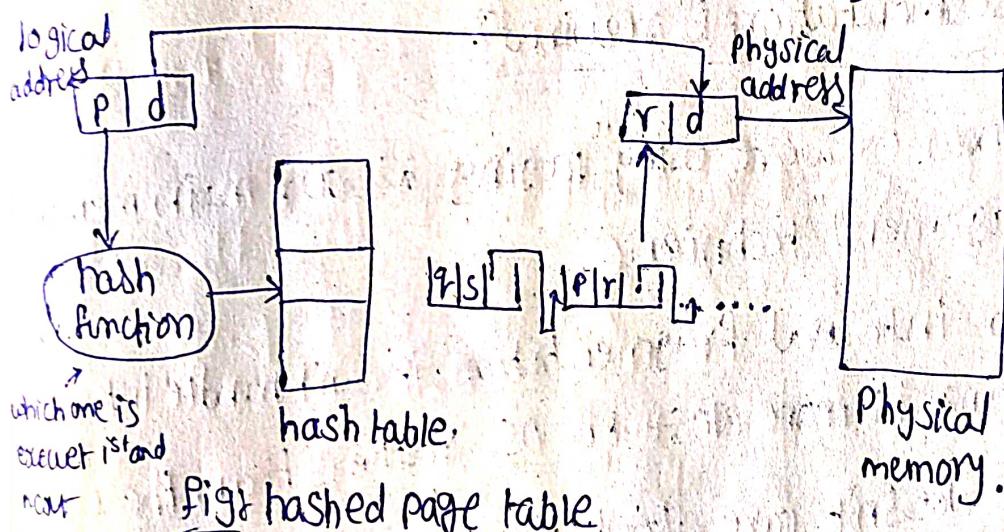
Fig: two level Page table schema

Hashed page tables

This approach is used to handle address spaces that are larger than 32 bits
each element mainly consists of

- * The virtual page number
- * the value of the mapped page frame
- * a pointer to the next element in the linked list.

→ In this virtual page, the no. of hashed into a page table. The page table mainly contains a chain of elements hashing to the same elements.



Inverted page table

The inverted page table basically "combines a page table and a frame table into a single data structure".

→ This technique decreases the memory to store each page table; but it also increases the time to search the table.

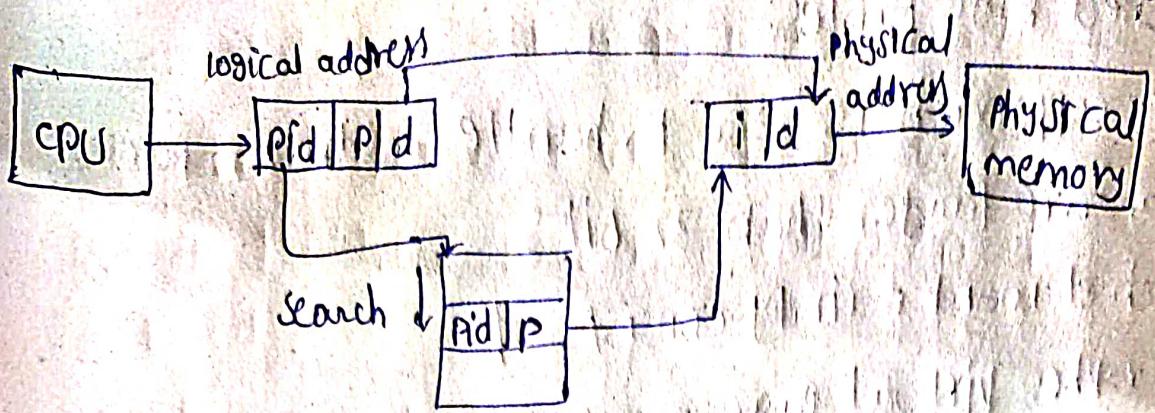


Fig 8 inverted page table

advantages

decreased the memory to store each page table.

disadvantage:

inverted page table create to take more time
so increases the time combine the page and frame table

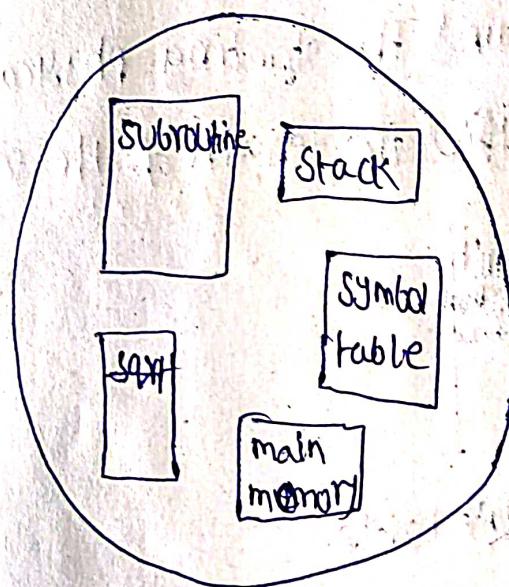
why segmentation is required?

- Till now, we were using "paging" as our main memory management technique.
- No internal fragmentation; it can't handle internal fragmentation.
- Paging is more close to OS rather than the user process.
- OS doesn't care about the user's view of the process.
- It is better to have segmentation which divides the process into the segments.

Segmentation

In operating systems, segmentation is a memory management technique in which, the "memory is divided into the variable size parts". each part is known as "segment".

- the details about each segment are stored in a table called as segment table.
- segment table contains mainly two information about segment:
 - ① Base: it contains starting physical address.
 - ② limit: it is the length of the segment.



logical address

fig 2. User's View of a Program.

- A program is a collection of segments.
- A segment is a collect logical unit such as main program, procedure, function, local variable, global variable, common blocks, stack, symbol table.

- Segmentation is a memory management scheme that supports this later view of memory.
- a logical address space is a collection of segments.
- each segment specifies the 2 quantities
 - (a) segment name (or) number
 - (b) segment offset

→ logical address space consists of a two tuples.

$\langle \text{segment-name, offset} \rangle$

→ Segment table maps two dimensional physical address.

- ① base :- it contains the starting physical address
- ② limit : it specifies the length of the segment

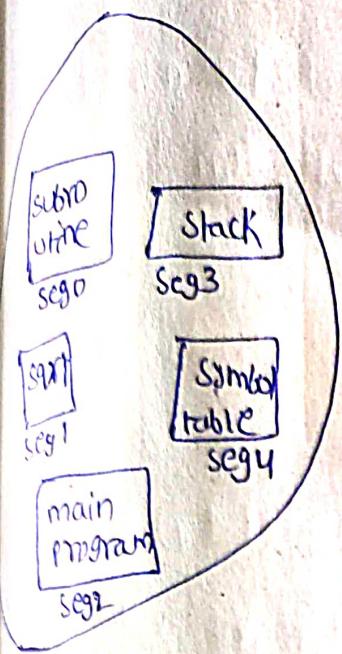
base	limit
address of seg	length of segment

Fig: Segment table

- segment table base register (STBR). it points to the segment table(s) location in memory
- segment table length register (STLR). it indicates no. of segments used by program

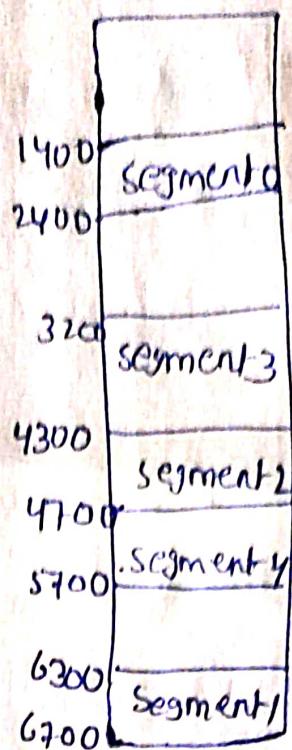
SL STLR

Example:-



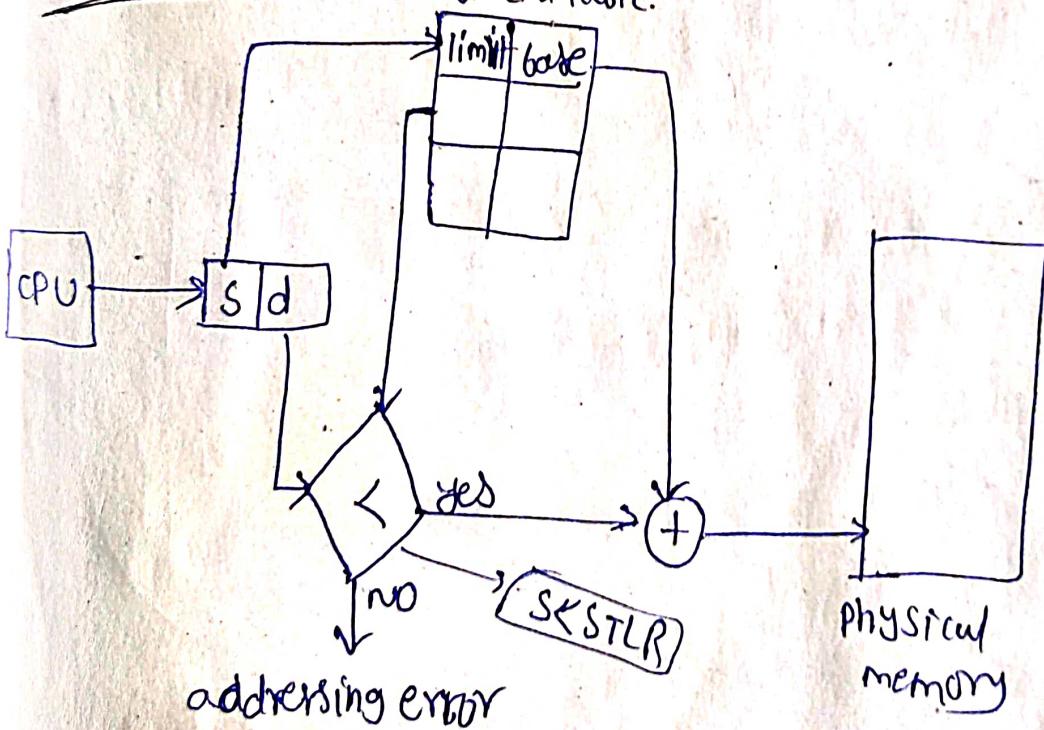
	limit	base
0	1000	1400
1	400	6800
2	400	4300
3	1100	3200
4	1000	4700

Segmentation table



segmentation h/w:

segment table.



Unit - III
Unit - 2 Virtual memory management

virtual means :-

Virtual is same as "real" but "not real". (or) "imaginary"
some part of the secondary storage memory is acting
as main memory.

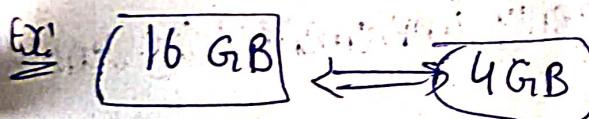
Eg:- mirror it will display your image it is not a
real it is a virtual i.e it just gives illusion to user.
memory means :-

Every task is executed with in the main memory.

Introduction

Virtual memory is a technique in which it is
executed may not be completely in main memory.

→ Virtual memory capability is double of the
Physical memory.



Users use V.M. real memory of
 computer

a programmer can write a program which
requires more than the capacity of main memory
such a program is executed by virtual memory
technique.

- more programs can running concurrently by V.M.
- less I/O needed to load or swap process.

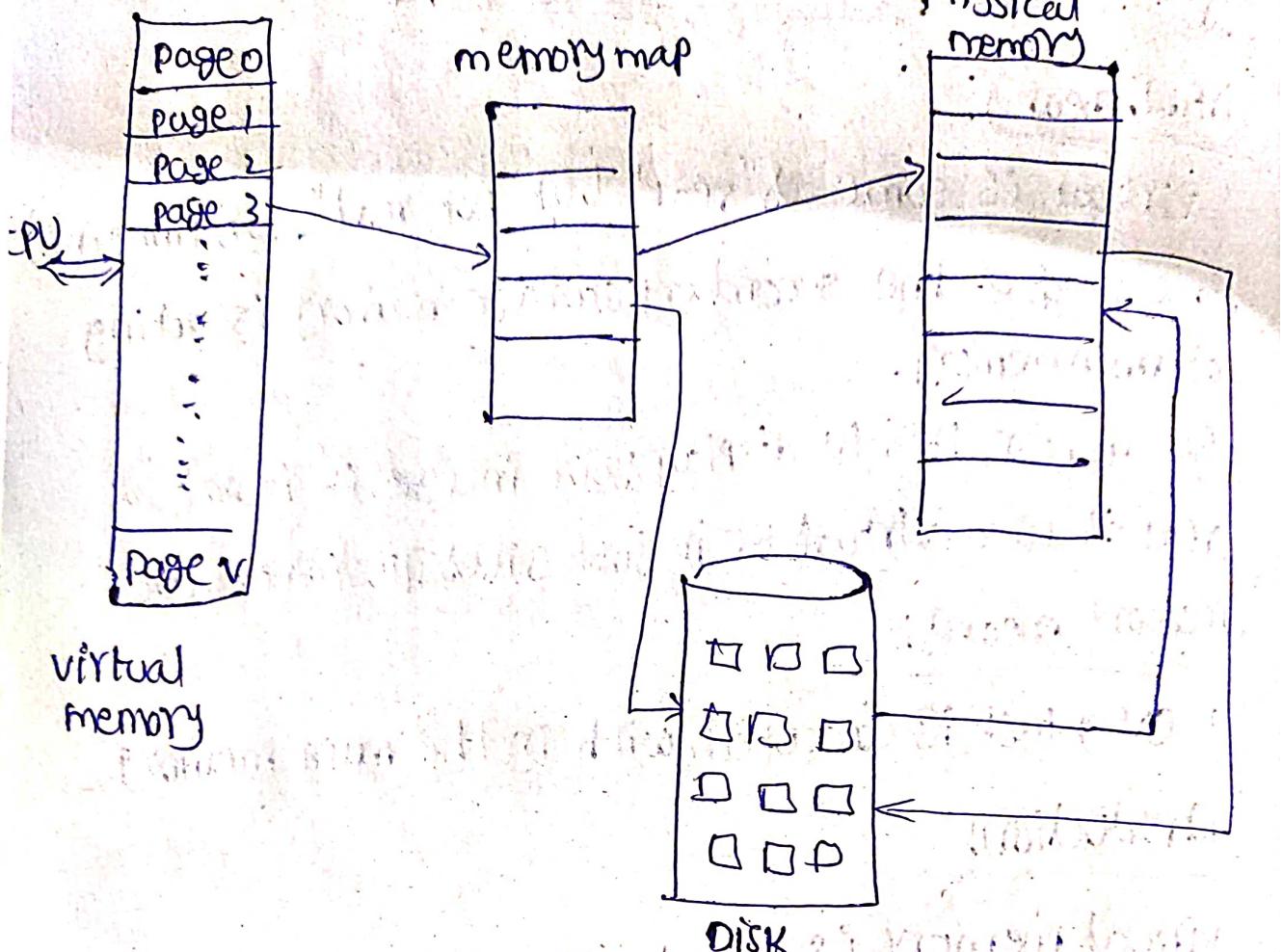


Fig 2 Virtual memory. That is larger than physical memory.

Demand Paging

- virtual memory is commonly implemented by demand paging
- demand paging is a method of virtual memory management
- In a system that uses demand paging the operating system copies a disk page into physical memory.

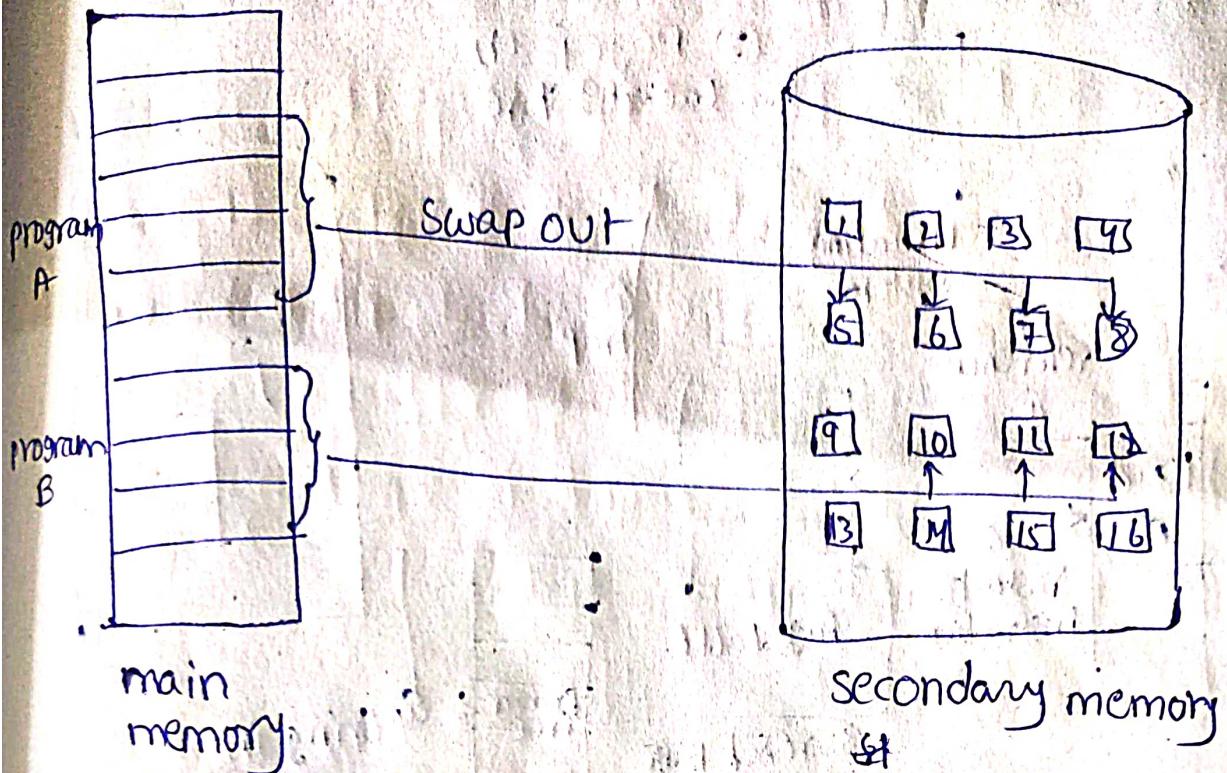


Fig: Transfer of a paged memory to contiguous disk space.

- If the referred page is not present in the main memory then there will be a miss and the concept is called "Page miss (or) Page fault".
- The CPU has to access the missed page from the secondary memory.
- If the no. of page fault is very high then the effective access time of the system will become very high.
- If page fault occurs then our CPU can run effectively.

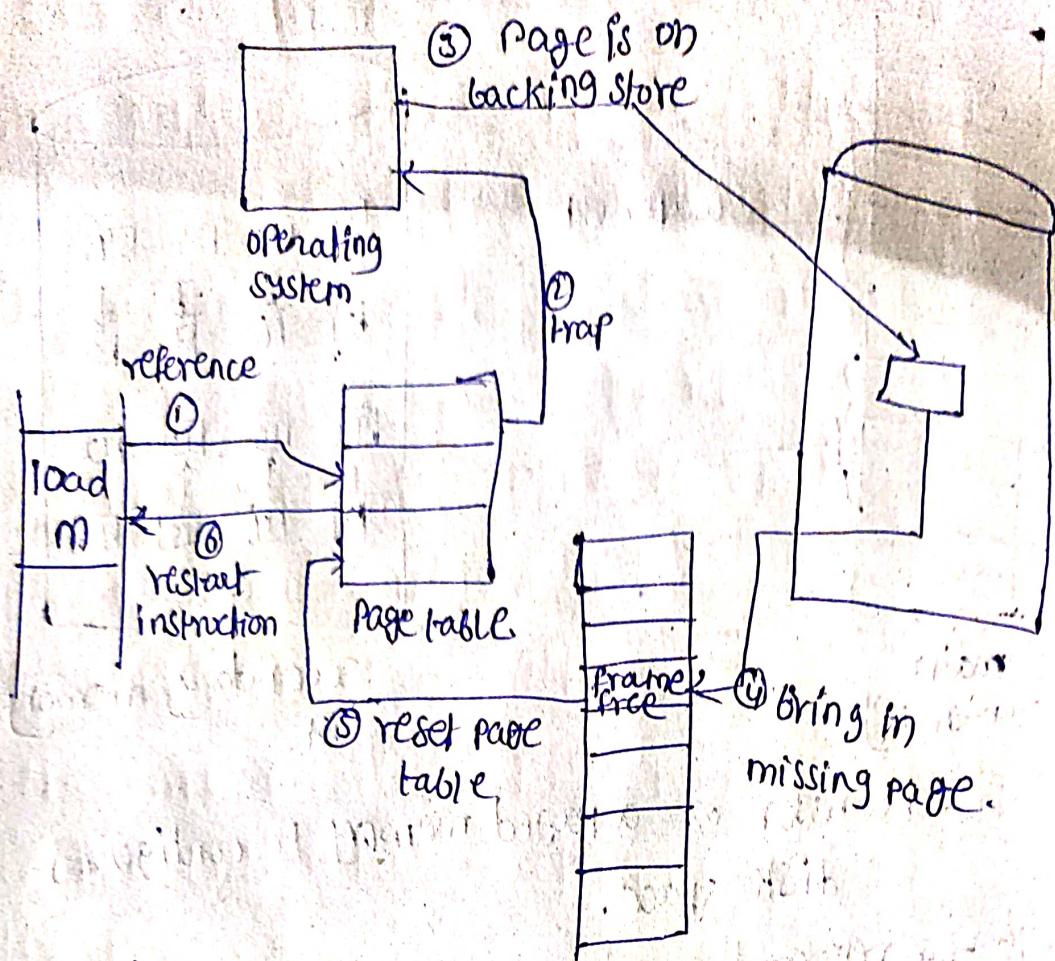


fig: steps in handling a Page fault.

copy on write

- copy on write (or) simply "cow" is a resource management technique,
- one of its main uses is in the implementation of the fork system call in which it shares the virtual memory (page) of the os.
- the idea behind a copy-on-write is that when a parent process creates a child process then both of these processes initially will share the same pages in memory and these shared pages

will be marked as copy-on-write which means that if any of these processes will try to modify the shared pages then only a copy of these pages will be created and the modifications will be done on the copy of pages by that process and thus not affecting the other process.

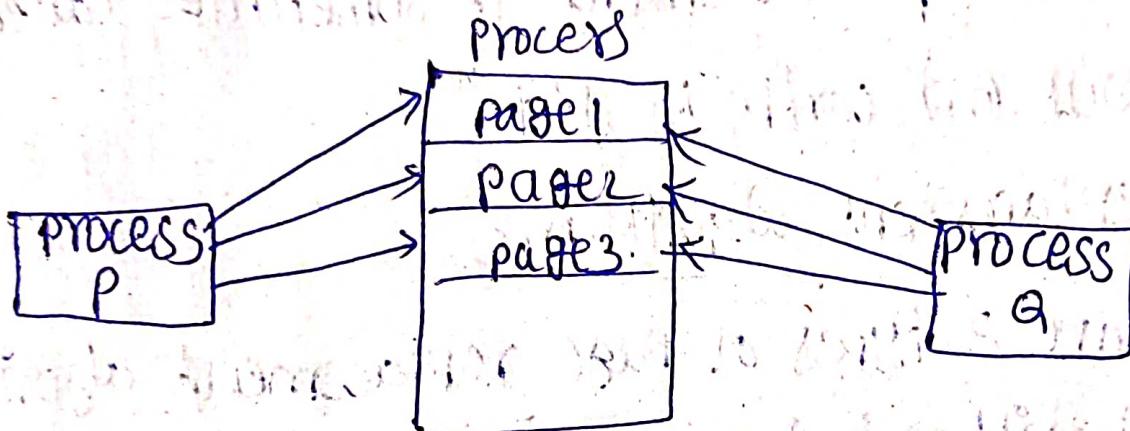


Fig: Before Process P modifies page 3.

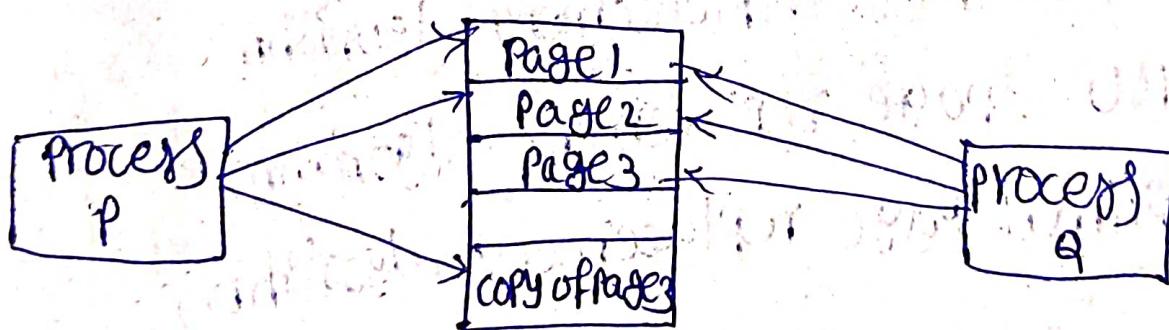


Fig: after process P modifies page 3.

Page Replacement

- page replacement is done when the requested page is not found in the main memory (page fault)
- these are page replacement algorithms decide which memory page is to be replaced
- the process of replacement is sometimes called "swap out (or) write to disk"

Page Replacement algorithm

- there are 3 types of Page replacement algorithm each algorithm has an efficient method to replace the pages.

① FIFO page replacement algorithm.

② LRU page replacement algorithm

③ Optimal page replacement algorithm

① FIFO page replacement algorithm

- FIFO (First come first out) the simplest page replacement algorithm

→ when a page wants to be replaced, the oldest page is chosen. we can create a FIFO queue to hold all pages in memory.

Ex To illustrate the problems that are possible with a FIFO page replacement algorithm.

Consider the 20 reference string and 3 frames in the memory.

7, 0, 1, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2, 1, 2, 0, 1, 7, 0, 1

initially the 3 frames are empty:

page	7	0	1	2	0	3	0	4	2	3	0	3	2	1	2	0	1	7	0	1
fram e 1	7	7	7	2	2	2	4	4	4	0	0	0	0	0	0	0	0	7	7	7
fram e 2	0	0	0	0	3	3	3	2	2	2	2	2	1	1	1	1	1	0	0	0
fram e 3	1	1	1	1	0	0	0	3	3	3	3	3	2	2	2	2	2	1	1	1
m m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m

page fault

$m = \text{miss} \rightarrow$ we are take the page number is not match

$H = \text{Hits} \rightarrow$ the page number is match

$$\text{Hit ratio} = \frac{\text{Total no. of hits}}{\text{Total no. of pages}} = \frac{5}{20} = 0.25$$

$$\text{miss ratio} = \frac{\text{Total no. of miss}}{\text{Total no. of pages}} = \frac{15}{20} = 0.75$$

\rightarrow In FIFO, if the item is not there in memory

then we call it has "miss", and that item will be taken from secondary memory.

\rightarrow if the item is in main memory then we call it has "hit"

"hit"

② LRU page replacement algorithm

→ Least Recently used (LRU) replacement associates with each page the time of their pages last use.

→ When a page must be replaced, LRU chooses that page that has not been used for the longest period of time.

Ex: Given reference strings (pages) are 20 as follow. and no. of frames are 3.

7, 0, 1, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2, 1, 2, 0, 1, 7, 0, 1

so,

page	7	0	1	2	0	3	0	4	2	3	0	3	2	1	2	0	1	1	7	0	1
F ₁	7	7	7	2	2	2	2	4	4	4	0	0	0	1	1	1	1	1	1	1	
F ₂	0	0	0	0	0	0	0	3	3	3	3	3	3	3	0	0	0	0	0	0	
F ₃		1	1	1	3	3	3	2	2	2	2	2	2	2	2	2	2	7	7	7	
m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	m	

$$\text{Hits} = 8, \text{ miss} = 12$$

$$\text{Hit ratio} = \frac{\text{Total no. of hits}}{\text{Total no. of pages}} = \frac{8}{20} = 0.4$$

$$\text{miss ratio} = \frac{\text{Total no. of miss}}{\text{Total no. of pages}} = \frac{12}{20} = 0.6$$

Optimal page replacement algorithm

most frequently used (MFU), most Recent used (MRU)

2, 3, 0, 1, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2, 1, 2, 0, 1, 7, 0, 1

NE	0	1	2	0	3	0	4	2	3	0	3	2	1	2	0	1	7	0	1
1	7	7	7	2	2	2	2	2	2	2	2	2	2	2	2	2	7	7	7
2	0	0	0	0	0	4	4	4	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	3	3	3	3	3	3	3	3	3	1	1	1	1	1	1	1
m	m	m	m	H	m	H	m	H	H	m	H	H	m	H	H	m	H	H	H

here 7 is lessing after a long time so, 7 will be placed by 2. whereas 0 and 1 will be using in next upcoming

Optimal page replacement gives best solution among FIFO and LRU

$$\text{Hit} = 11, \text{ Hit Ratio} = \frac{\text{Hit}}{\text{Total pages}} = \frac{11}{20} = 0.55$$

$$\text{miss} = 9, \text{ miss Ratio} = \frac{\text{miss}}{\text{Total pages}} = \frac{9}{20} = 0.45$$

We can see more no. of hits is optimal page replacement

Frame allocation:-

There are two important tasks in virtual memory management: a page replacement strategy and frame allocation strategy.

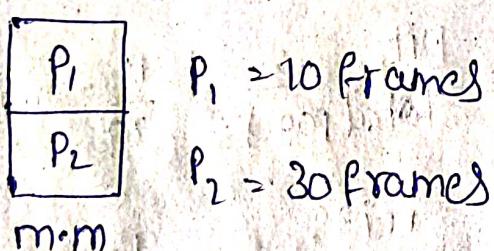
There will be no wastage of main memory and CPU.

→ Frame allocation can be divided into 5 types:

- ① equal allocation
- ② proportional allocation
- ③ priority allocation
- ④ global allocation
- ⑤ local allocation.

① equal allocation:-

The main memory size = 40



→ Here we have taken only 2 processes, so memory is divided into 2 parts i.e

P_1 can take 20 frames

P_2 can take 20 frames

So we can allocate to 10 frames in P_1 , where as 10 will be wasted and in P_2 , it can allocate only

o frames remaining 10 frames cannot allocate in memory.

If we do equal allocation then we have a disadvantage so we move on to next allocation.

Propositional allocation

The propositional allocation we have to take the frames based upon the process size. ~~the process it will~~ size of the process. Ex) $P_1 = 10, P_2 = 30$

$$\text{Propositional allocation } (P_1) = \frac{P_1}{P_1 + P_2} \times \text{(memory allocation size)}$$

$$P_1 = \frac{10}{(10+30)} \times 40 \Rightarrow \frac{10}{40} \times 40 = 10$$

$$\text{Propositional allocation } (P_2) = \frac{P_2}{(P_1 + P_2)} \times \text{(memory allocation size)}$$

$$P_2 = \frac{30}{(10+30)} \times 40 = \frac{30}{40} \times 40 = 30$$

Priority allocation

The frame allocation is depending on the priority which process will be executed first and which process will be executed next.

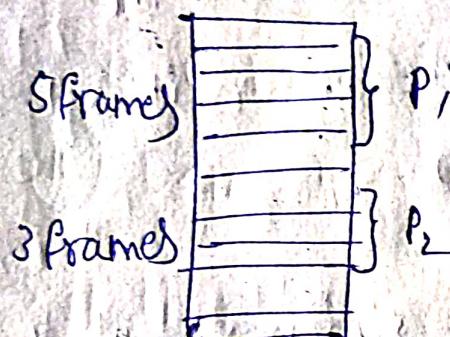
Global replacement allocation

This is mainly useful for performing so we know what is replacement. i-e

Whenever we want to use a page, mm checks, if the page is available in main memory, then it will

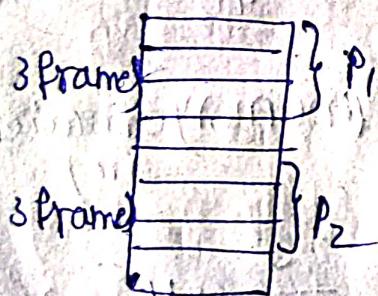
execute, and if the page is not occurred in mm, then we say page fault has occurred, then it has to get that particular page from secondary memory to main memory.

for global allocation:



In this, we can allocate a page from secondary memory to main memory in any of the process.

for local allocation:



In this we has to allocate a page which is necessary for P_1 , in that particular space itself.

Trashing:

If a process spends more time on paging rather than executing is called as an trashing.

Cause of Trashing:

- If a process does not have "enough" pages, the page fault rate is very high. This leads to: o Low CPU utilization.
- Trashing is a process is busy swapping pages in and out.
- OS thinks that it needs to increase the degree of multiprogramming.

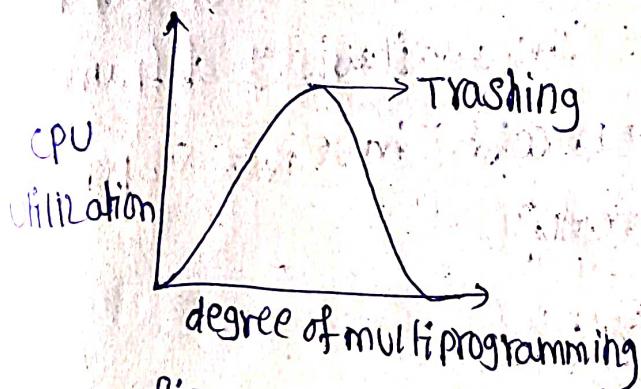


Fig: Trashing

- If Trashing occurs, CPU performance goes down.
- Techniques used to handle the Trashing:
 - ① Working set model
 - ② Page fault frequency.

Working set model:

The accuracy of the working set is depending on the source of parameters.

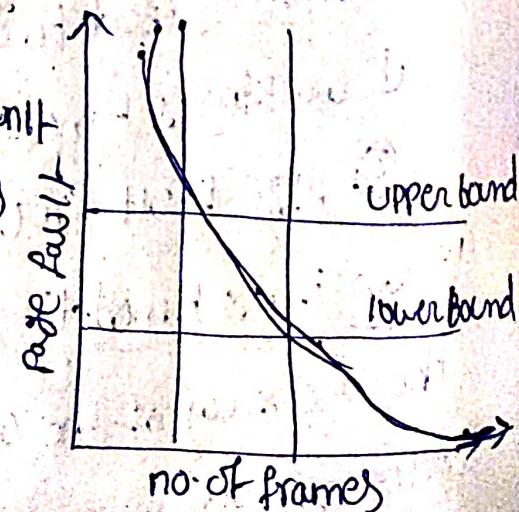
- if 'D' is the total "demand for frames", now, &
 if 'm' is the "no. of frames available in the
 memory";
 These are 2 possibilities.
- (i) $D > m$ i.e. total demand for frames exceeded
 the no. of frames, then "thrashing will occur" as
 some processes would not get enough frames
 - (ii) $D \leq m$, then there would be "no thrashing".

- if there are enough extra frames, then some
 more processes can be loaded in the memory.
- the other hand if the summation of working
 set sizes exceeds the availability of frames,
 then some of the processes have to be suspended
 (swapped out of memory)

② page fault frequency

more direct approach to handle thrashing is the
 one that uses page fault frequency concept.

- if the page fault rate is:
 too high it indicates that upper unit
 the process has too few frames
 allocated to it.
- if the page fault is below
 the lower unit, frames can
 be removed from the process.



memory mapped files

→ according to accessing a file on disk will require the standard system calls:

- ① open()
- ② read()
- ③ write()

→ every time the file is accessed requires a system call and disk access.

• memory mapped file I/O allows file I/O to be treated as routine memory access by mapping a disk block to page in memory.

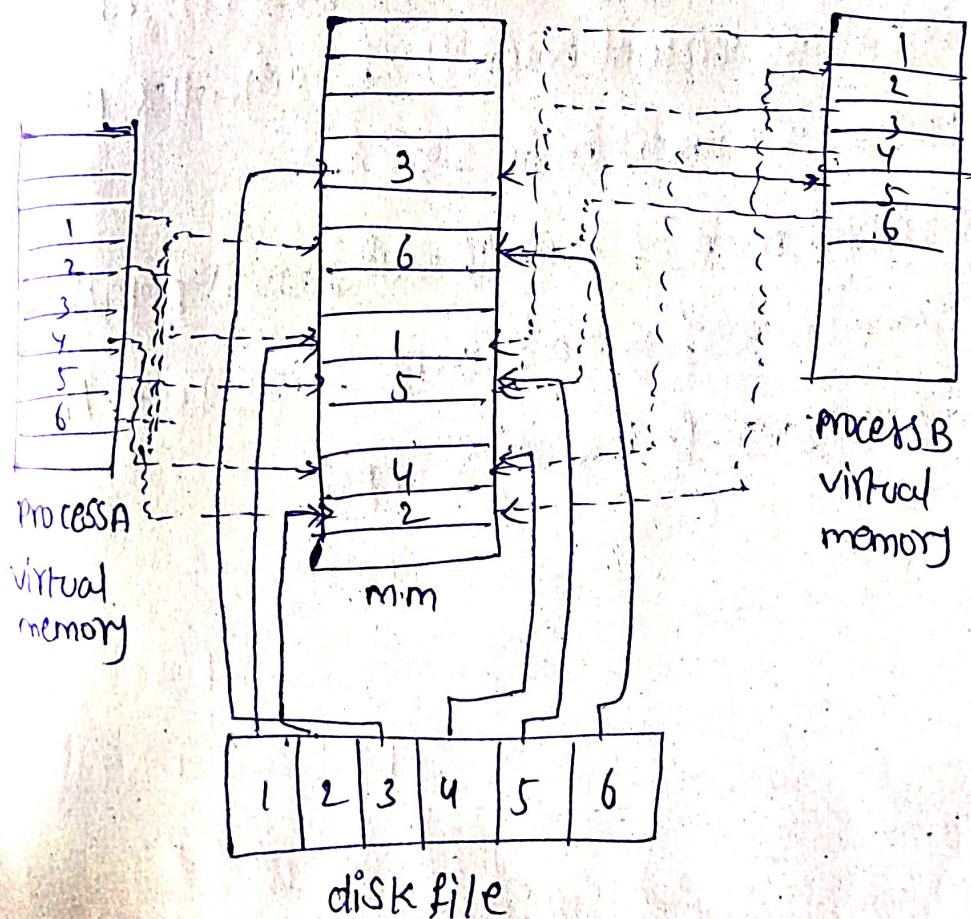


Fig:- memory mapped files

→ a file is initially read using demand paging, resulting in a page fault. i.e

→ we are having page 7 & it is not matched with the memory at that time page fault is occur

- Page sized portion of the file or pages is read from the file system into a physical page (or) memory page.
- The OS periodically checks if the page in the memory mapping the file have been modified, then update the physical file and (or) memory file on disk.
- When the file is closed, all the memory mapped data are written back to disk and removed from the virtual memory of the process.
- also allows several processes to map the same file into the virtual memory of each, allowing the pages in memory to be shared.
- The virtual map of each shared process points to the same page of main memory — The page that holds a copy of the disk block.