



- Q 4 what is the scheduling criteria explain any two primitive scheduling algorithm with example.

In multiprogramming environment there may be the situation where one or more processes are simultaneously in ready state & if only one CPU is available, then a choice has to be made as to which process should execute next, the part of OS that makes this choice is called scheduler & the algorithms it uses are called scheduling algorithms.

A major issues related to scheduling when to make scheduling decisions there are varieties of situations where scheduling is required.

- 1> when a new process is created, then decision regarding whether parent or child process should be executed.
- 2> when a process exits then some other process must be selected from set of ready states. If no process is in ready state a system supplied idle process is execute.
- 3> when a process blocks on Input/Output or a semaphore or for some other reason another process has to be selected to run.
- 4> when an I/O interrupt occurs a scheduling decision may be made.

Scheduling algorithms can be divided into two categories, depending on how they deal with clock interrupt.

- 1> A non-preemptive
- 2> A preemptive

* non preemptive :-

A non preemptive scheduling algorithm pi-



process to run & then just lets it run until it blocks or until voluntarily releases its CPU in other no scheduling decisions are made during clock interrupts after clock interrupt processing has been completed the process that was running before the interrupt is always resumed.

2) A preemptive:-

A preemptive scheduling algorithm picks a process to run & lets it run for a maximum of some fixed time. If it is still running at the end of time interval it is suspended & the scheduler picks another process to run.

Different scheduling algorithms are needed in different environment & different application areas. Three distinguishable environments are -

- 1) Batch system
- 2) Interactive system
- 3) Real time system.

1) Batch system:-

No user are impatiently waiting for quick response. This approach reduces process switches thus improving performance.

2) Interactive system:-

Here preemption is essential to keep one process from hogging the CPU & denying service to others. One process might shut out all others indefinitely. Hence preemption is needed to prevent this behaviour.

3) Real time system:-



Here preemption is sometimes not needed since the processes know that they may not run for long periods of time & usually do their work quickly.

* Goal of scheduling algorithms

1) All systems:-

- Fairness: giving each process a fair share of CPU
- Policy enforcement: seeing that stated policy is carried out
- Balance: keeping all parts of system busy

2) Batch system:-

- Throughput: maximum jobs per hour
- Turn around time: minimum time between submission & termination
- CPU utilization - keep CPU busy all the time

3) Interactive system:-

- Response time: respond to request quickly
- Proportionality: meet user expectations.

4) Real time systems:-

- meeting deadlines: avoid losing data
- predictability: avoid quality degradation in multimedia systems.

* Scheduling criteria:-

Different scheduling algorithms have different properties & are applicable in different situations selecting which algorithm to use in which situation.



i) CPU

is to keep CPU as busy as possible. It may range from 0% to 100%.

ii) throughput:

is the number of processes completed per time unit. It may range from 10 processes/second upto 1 process.

iii) turnaround time:

is the time interval between submission of a process & completion of a process. It is the sum of the periods spent waiting to get into memory, waiting in ready queue executing or CPU & doing I/O.

iv) waiting time:

is the amount of time a process spends waiting in ready queue. It is the sum of the periods spent waiting in ready queue.

v) Response time:

It is the time from submission of a request until first response is produced. It is the amount of time taken to start responding & not the time taken to output the response.

We want to maximize CPU utilization & throughput & minimize turnaround waiting & response time. For interactive systems minimizing the variance in response time is important then minimizing average response time.

Scheduling algorithm:

CPU scheduling deals with problems of deciding which of the processes in ready queue is to be allocated the CPU.

There are many different CPU scheduling algorithms

- 1) First-come, first-served scheduling
- 2) Shortest job first scheduling
- 3) Priority scheduling
- 4) Round Robin scheduling
- 5) Fair share scheduling.

1) First-come First served scheduling-

- It is a non preemptive algorithm
- In FCFS scheme, the process that request the CPU first is allocated the CPU
- The FIFO queue is used to the FCFS implementation.
- There is a single queue of ready processes & process are assigned to the CPU in the order they request it.
- The code of FCFS scheduling is simple to write & understand.
- The average waiting time under the FCFS policy is often quite long

2) shortest job first scheduling:-

- The SJF scheduling algorithm associated with each process the length of the processes next CPU burst.
- When the CPU is available, it is assigned to the process that has the smallest next CPU burst.
- If the next CPU bursts of two processes are the same FCFS scheduling is used to break the tie.
- Thus instead of SJF most appropriate term is shortest-next-CPU-burst algorithm because scheduling depends on the length of the next CPU burst of a process rather than its total length.

This approach takes into account all tasks in such a way as to ensure that each task is allocated some fraction of the time.

- It is non preemptive as well as preemptive algorithm.
- optimal when all the jobs are available simultaneously.
- shorter jobs do not have to wait for lengthy jobs hence increasing average turn around time & reducing average waiting time
- 2. Priority scheduling:
 - priority scheduling algorithm is associated with SJF.
 - A priority is associated with each process & the CPU is allocated to the process with the highest priority equal priority process are scheduled in FCFS order.
 - An SJF algorithm is simply a priority algorithm where the priority (P_j) is the inverse of the next CPU burst. The larger the CPU burst the lower the priority & vice versa
- 3. Round Robin algorithm/scheduling:
 - The RR scheduling algo is designed especially for time sharing system.
 - It is similar to FCFS scheduling, but preemption is added to switch between processes.
 - A small unit of time called a time quantum or a time slice is defined. A time quantum is generally from 10 to 100 millisecond.
 - The ready job queue is treated as a circular queue.
 - To implement RR scheduling the we keep the ready queue as a FIFO queue of process.
 - Each process is assigned a time interval called quantum, in which it can run. If process is still running at the end of the quantum, the CPU is preempted & given to another process.
 - If process has blocked or finished before quantum has elapsed CPU switches the next process in queue.
- 4. Fair share scheduling:

1) Fair time system:



- This approach takes into account who owns a process before scheduling it.
- Each user is allocated some fraction of CPU scheduler picks a process in such a way as to enforce it.
- Thus if 2 users have each between promised 50% of CPU, will each get that, irrespective of number of processes owned by each one of them.

* Non preemptive scheduling algorithms.

Process	Arrival Time	burst time
P1	0	9
P2	1	4
P3	2	3
P4	3	5

0	1	5	10	15	20	25
P1	P2	P4	P1	P3		

$$\Delta_{ij} = \text{start time}_j - \text{Arrival time}_i + \text{new start time}_j - \text{old finish time}_i$$

$$\Delta_1 = 0 - 0 + 10 - 9 = 0 + 9 = 9$$

$$\Delta_2 = 1 - 1 = 0$$

$$\Delta_3 = 15 - 2 = 13$$

$$\Delta_4 = 5 - 3 = 2$$

$$\text{avg} = \frac{9+0+13+2}{4} = \frac{24}{4} = 6$$

6.5

$$\Delta_i = (\text{Finish time}_i - \text{Arrival time}_i) + (\text{new finish time}_i - \text{old finish time}_i) + \text{idle time}$$

$$\Delta_1 = 1 - 0 + 10 - 1 + 9 = 1 + 16 + 9 = 26$$

$$\Delta_2 = 5 - 1 = 4$$

$$\Delta_3 = 16 - 2 = 14$$

$$\Delta_4 = 10 - 3 = 7$$

$$\text{avg} = \frac{26+4+14+7}{4} = \frac{51}{4} = 12.75$$

preemptive priority scheduling algorithm.

process burst time Arrival time priority

P1	8	0	4
P2	6	1	6 highest
P3	7	3	3
P4	9	3	1

P1	P2	P3	P4	P5
0	1	7	14	21

$$AWT = (ST - AT) + NST + OFT$$

$$P1 = 0 - 0 + 7 - 1 = 6$$

$$P2 = 1 - 1 = 0$$

$$P3 = 14 - 3 = 11$$

$$P4 = 21 - 3 = 18$$

$$Avg = \frac{18+11+0+6}{4} = \frac{35}{4} = 8.75$$

$$TT = (FT - AT) + (NST - OFT) + bt$$

$$P1 = (0-0) + 7 - 1 + 8 = 0 + 6 + 8 = 14$$

$$P2 = 7 - 1 = 6$$

$$P3 = 21 - 3 = 18$$

$$P4 = 30 - 3 = 27$$

$$Avg = \frac{27+18+6+14}{4} = \frac{65}{4} = 16.25$$

Q. What is deadlock? Explain deadlock concept.
 Many request units need multiple resources if the resource are not available at a time a wait state will occur at that stage because the resources are required by other waiting processes they have to wait.



Q 2 what is deadlock explain deadlock avoidance

concept - in a multiple environment several process may request finite no of resources A process request resources if the resource are not available at that time, the process enters a wait state within processes may never again change state because the resources they have requested are held by other waiting processes this situations is called deadlock deadlock occurs when some process makes a request that cannot be granted immediately.

deadlock may involve some or different resource types. there are various features that characterise deadlock i.e a deadlock function can arise if following 4 conditions hold simultaneously in a system.

1) mutual Exclusion:-

At least one resources must be held in a non sharable mode i.e only one process at a time can use the resource if another process requests that resource if another process request that must be delayed until the resource has been released.

2) circular wait:-

To impose total ordering of all resources type & to require that each process request resources in an increasing order enumeration whenever a process request an instance of resource type with higher order it has to release resource to lower order if acquired.

* Resource allocation:-

A directed edges from resource type $j \rightarrow p$ called as assignment edges. signifies that an instances of resour-

ce type R_j has been allocated to process p_i . A directed edge $p_i \rightarrow R_j$ called request edges. In addition we introduce new type edges called claim edge. $l_i \rightarrow R_j$ which indicates that process p_i may request resource R_j at some time in future. This edge resembles a request edge in direction but is represented by dashed line.

Suppose that p_i request R_j the request can be granted only if connecting the request edge $f_i \rightarrow R_j$ to an assignment edge $R_j \rightarrow p_i$ does not result in formation of cycle in resource allocation graph if no cycle in resource allocation graph.

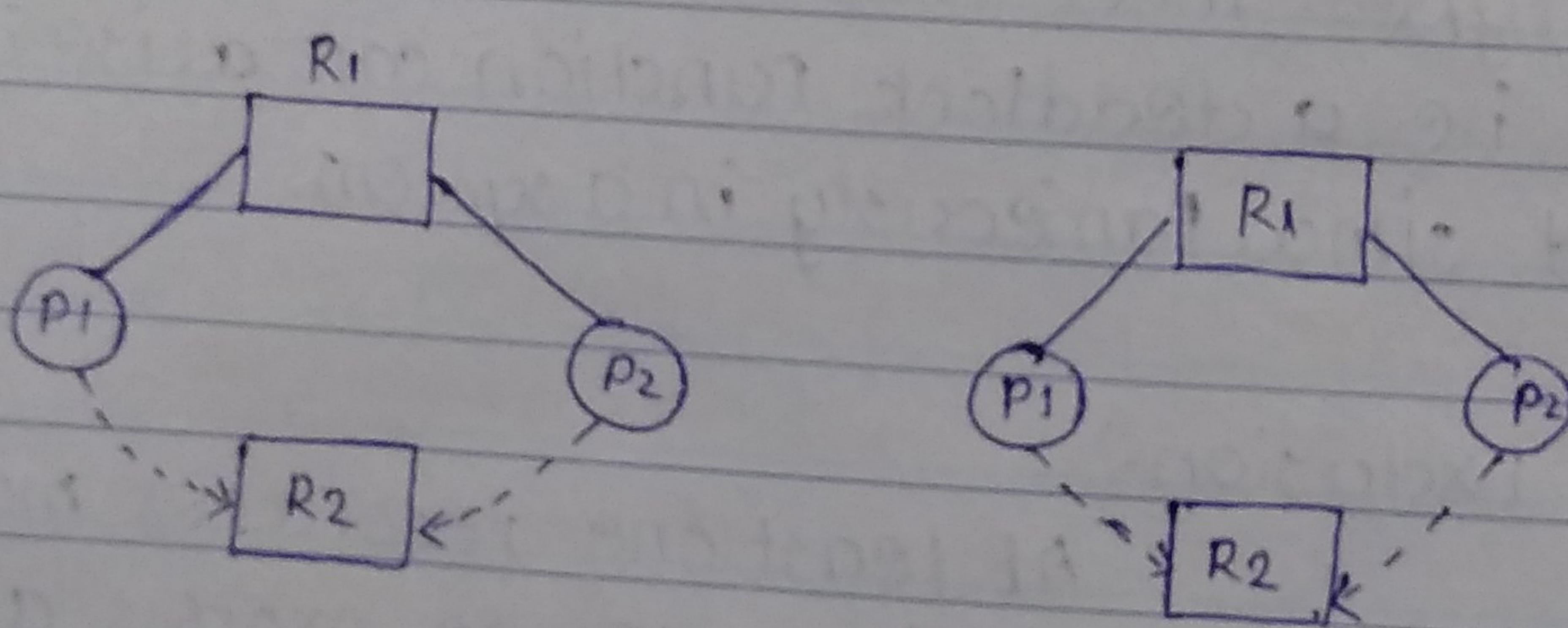


Fig:- Resource allocation graph
deadlock avoidance

1) Hold & wait:

A process must be holding at least one resource waiting to acquire additional resource that are currently being hold by other processes

2) No preemption:

Resource can not be preempted i.e. resource can be released only voluntarily by the process holding it after that process has completed its task.

3) Circular wait:



A set of $\{p_0, p_1, \dots, p_n\}$ of waiting processes must exist such that p_i is waiting for a resource that is held by p_j , p_j is waiting for resources that is held by p_i .

We can deal with deadlock problems in one of the following 3 ways:

- 1) We can use protocols to prevent or avoid deadlock ensuring that the system will never enter a deadlock state.
- 2) We can allow the system to enter a deadlock state & detect & it recovery.
- 3) We can ignore the problem altogether & pretend that deadlock never occur in the system.

* Banker's Algorithms:-

The resource allocation graph algorithm is not applicable to a resource allocation system with the multiple instances of each resource type. Hence banker's algorithm was introduced. The name was chosen as this algorithm could be used in a banking system to ensure that the bank never allocates its available cash such that it can no longer satisfy the needs of all its customers.

When a new process enters the system, it must declare the maximum number of instances of each resources type that it may need. This number may not exceed the total number of resource in the system. When a user request a set of resources the system must determine whether the allocation of these resources will leave the system in safe state.

Let, n = no of process in system

m = no of resource type

* We need following data structure:-

- 1) Available:-



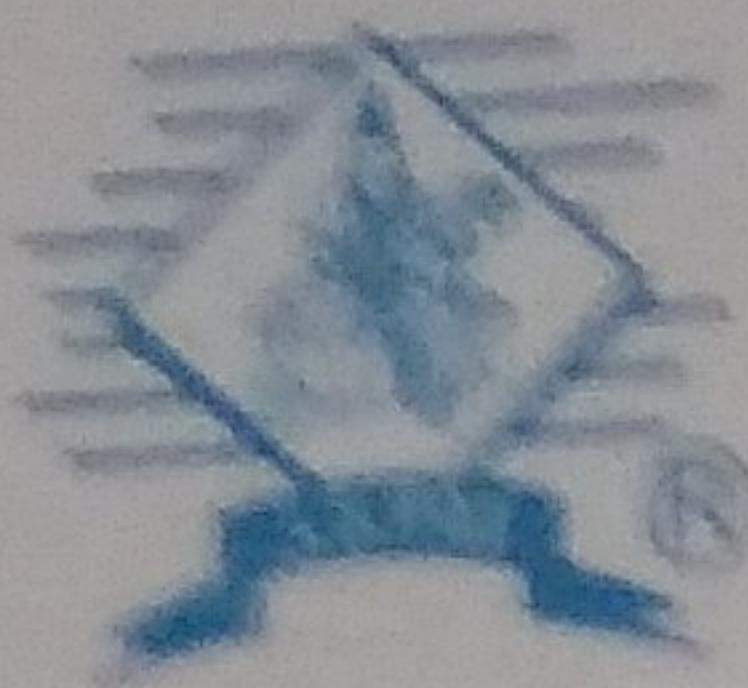
Indicate no of available resource of each type if
available[j] = k there are k instances of resource type Rj avail.

- 2) Max: define maximum demand of each process if max[i][j] = k then process p_i may request at most k instances of resource type R_j.
- 3) Allocation: defines the number of resources of each types currently allocated k instances of Resource type R_j:
- 4) Need: indicates remaining resource need of each process if need[i][j] = k then process p_i may need k more instances of resource type.

* Safe state:

A safe state is there if system can allocate resources to each process in some order & still avoid a deadlock. If the resource that process p_i needs is not immediately available then p_i can wait until all p_i have finished when they have finished p_i can obtain all of its needed resource complete its designed task, then it allocate resources & terminate when p_i terminated p_{i+1} can obtain its needed resources & so on, if no such sequence exists then system state is said to be unsafe.

A safestate is not a deadlock state. deadlock state is unsafe state but not all unsafe state are deadlock. initially a system is in safe state whenever process request a resource that is currently available that system must decide whether the resource can be allocated immediately or whether the process must wait. the request is generated only if the allocation leaves the system in a safe state.



Q 1 Explain segmentation in memory management

concept:

We can improve CPU utilization & speed of computers in response to its users by allocating various process to share CPU. But for this we must be able to share memory as well. Selection of memory management method for specific system depends on many factors especially on hardware design of system.

Memory is central of the operation of modern computer system. It consists of large array of words or bytes each with its own address. CPU fetches instructions from memory according to the value of program counter. The instruction is decoded & may cause operands to be fetched from memory after the instruction has been executed on the operands result may be stored back in memory. The memory units see only a stream of memory address does not know how they are generated.

The normal procedure is to select one of the processes in the ready queue & to load that process in memory as the process is executed it access instruction & data from memory. Eventually the process terminates & memory space is declared available.

* contiguous & non-contiguous memory allocation.

The main memory has to accommodate both OS & various user processes as may be placed either in low memory high memory generally interrupt vector tables are placed in low memory.

Various user process resides in memory at the same time the issues is how to allocate available memory to the process in ready queue in contiguous



memory allocation each process is contained in single contiguous section of memory.

The simplest method of memory allocation is to divide memory into several fixed sized partition where each partition may contain exactly one process. A generalization to fixed partition scheme must be generally used in batch environment initially all memory is in batch environment available for user processes & considered as one large enough for this process if we find such we allocate only has much memory has needed keeping the rest available satisfy to future request.

At any given time we have a list of available block size & ready queue. as can order the ready queue according to some scheduling algorithm generally a set of holes of various size is scattered throughout the memory at any given time. If the hole is too large it is split into 2 or more parts is allocated to arriving process other the return to be set of holes. When a process terminates, it releases its block of memory which then placed block in the set of holes. If the new holes is adjacent to other holes, these adjacent holes are merged to form one large hole. The system can use dynamic storage allocation methods - first fit, best fit or worst fit for allocating memory. But these algorithm suffer from external fragmentation as process are loaded & removed from memory free memory space is broken into little pieces.

Internal fragmentation occurs when memory is allocated in unit of blocks & if process gets a block slightly larger than request memory. This difference is called internal fragmentation memory that is internal to a partition but is not being used. The solution is external fragmentation is compaction. Compaction is possible only if relocation is dyno.



Segmentation is memory management scheme that supports this user view of memory. A logical address space is a collection of segments. Each segment has a name & a length. The addresses specify both the segment name & the offset within the segment.

The user therefore specifies each address by two quantities: a segment name & offset. For simplicity of implementation, segments are numbered & are referred to by a segment number rather than by a segment name. Thus a logical address consists of a two tuple:

(segment number, offset)

Normally, the user program is compiled & the compiler automatically constructs segments reflecting the input program.

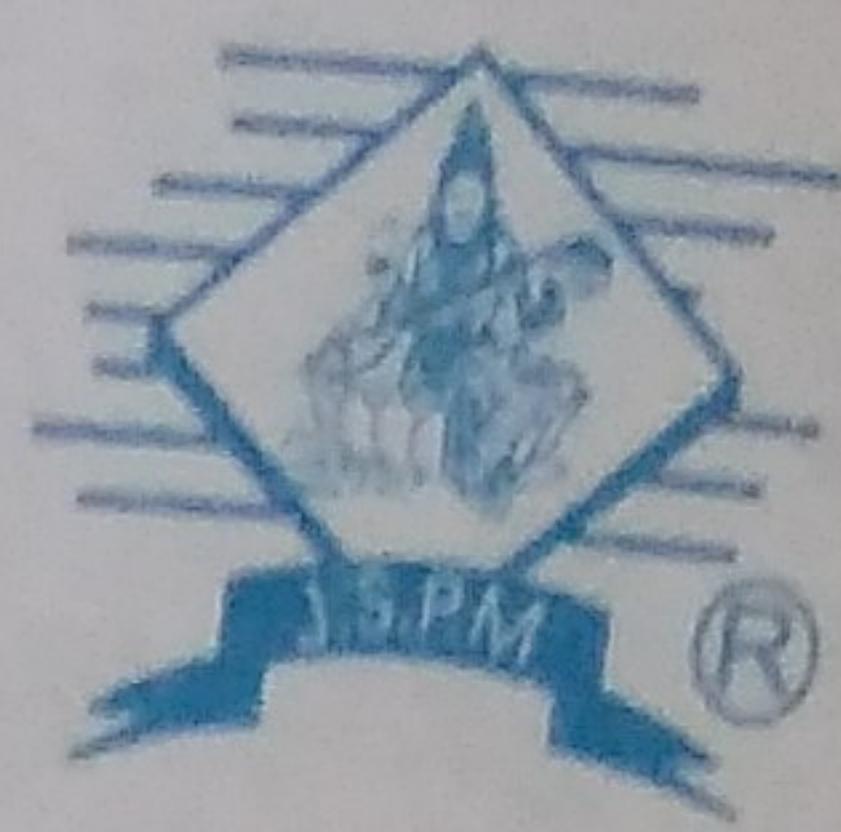
A C compiler might create separate segments for the following:

- 1) The code
- 2) Global variables
- 3) The heap from which memory is allocated
- 4) The stack used by each thread
- 5) The standard C library

* MAPPING:

User can have access to object in the program by a two dimensional address. In the actual physical memory is still of course a one dimensional sequence of bytes. Thus we must define an implementation to map two dimensional physical address.

The mapping is effected by the segment tables. Each entry in the segment table has a segment base & segment limits. The segment base contains starting physical address where the segment resides in memory whereas the segment limits specify the



it's done as the execution time.

* Segmentation.

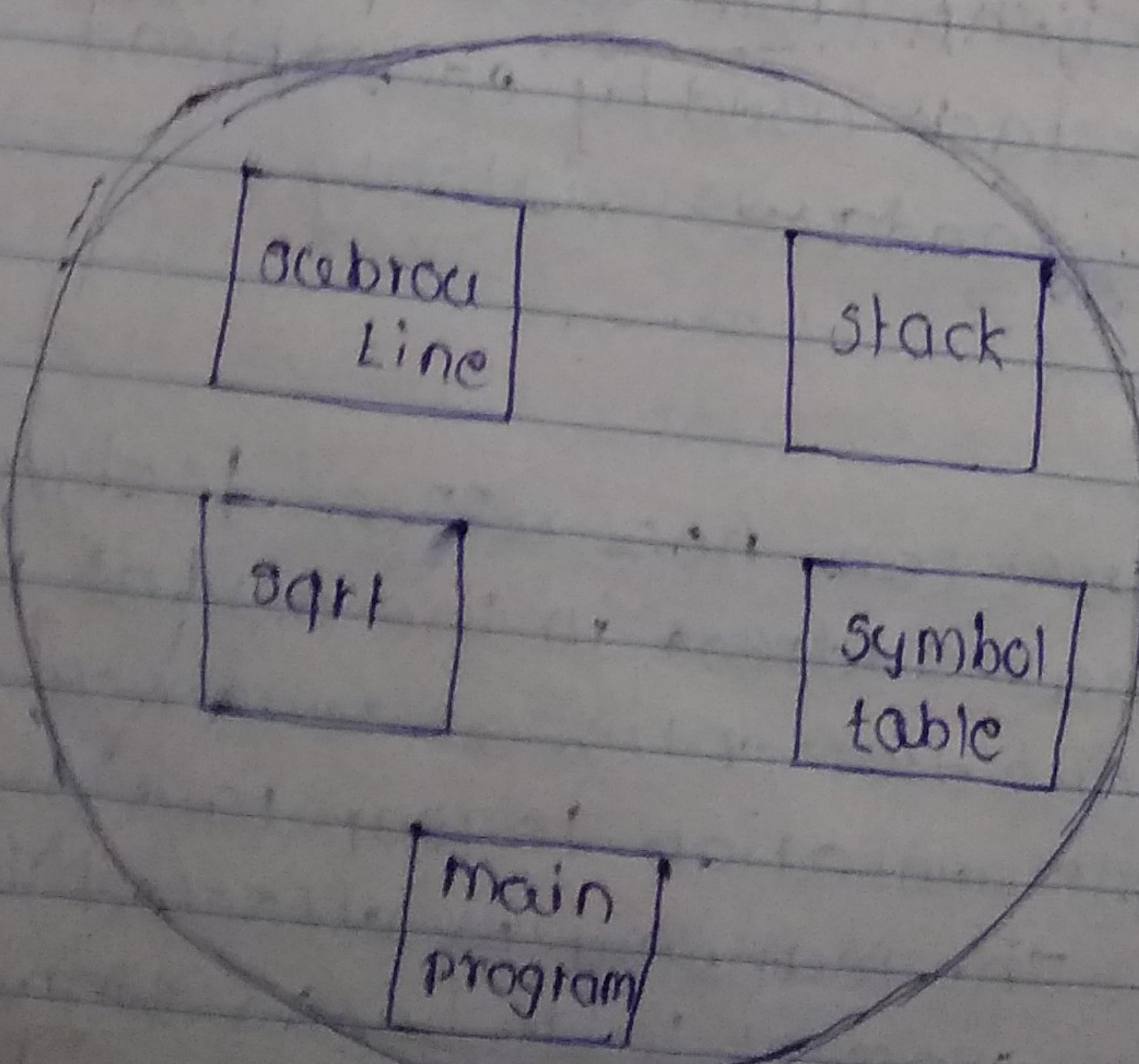
An important aspect of paging is the separation of the user view of memory & actual physical memory. The user view is mapped into physical memory. This mapping allows differentiation between logical & memory & physical memory.

* Basic Methods.

uses view memory as the allocation of the collection of variable-sized segment with no necessary ordering among segments.

* User view of programs.

User view of program you think of it as main program with a set of methods procedures or function it may also include various data structures object, array, stack, variable etc



length of the segments.

- A logical address consist of the two parts, a segment number s & an offset into a segment.
- The offset d is logical address must between 0 & segment limit.
- If it is not we trap to the operating system when an offset is legal it is added to the segment base & produce an address in physical memory of the desired bytes. the segment table is thus essentially an array of base limit register pairs.

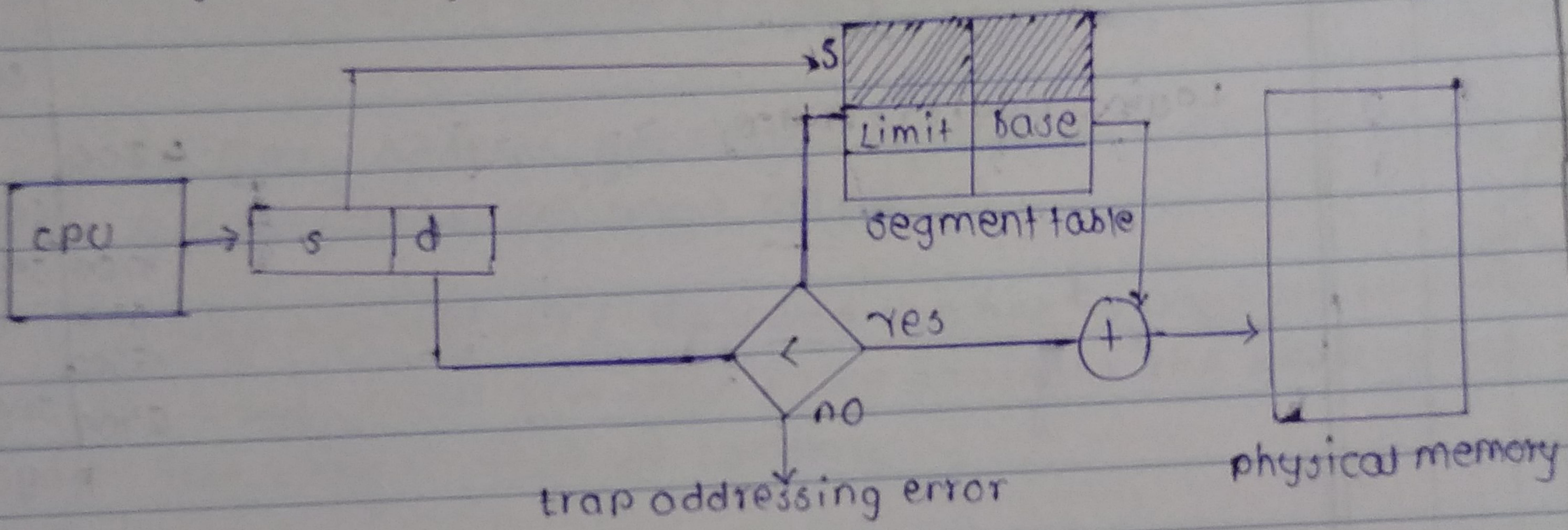


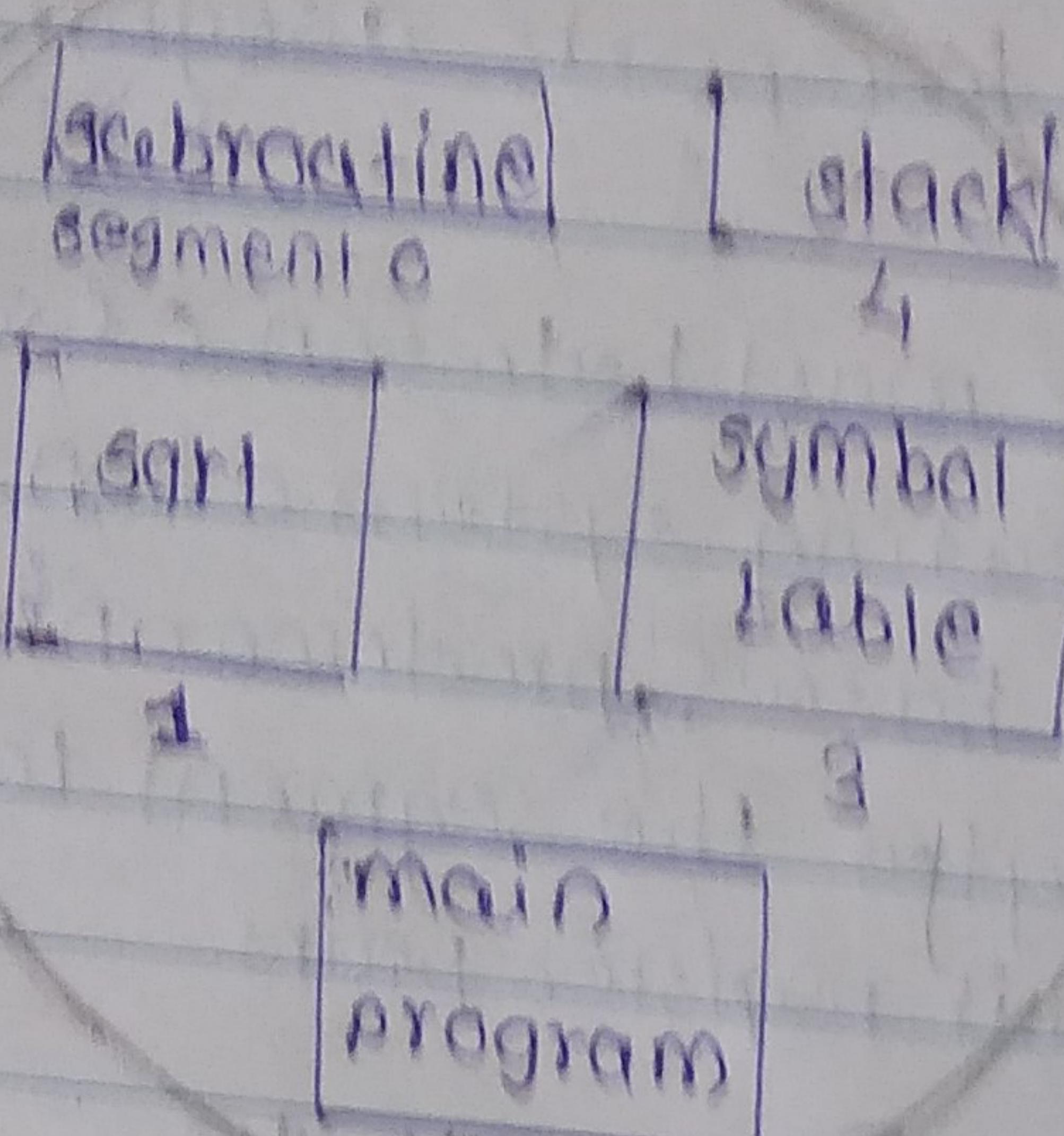
fig : segmentation hardware

* consider the situation shown the following figure

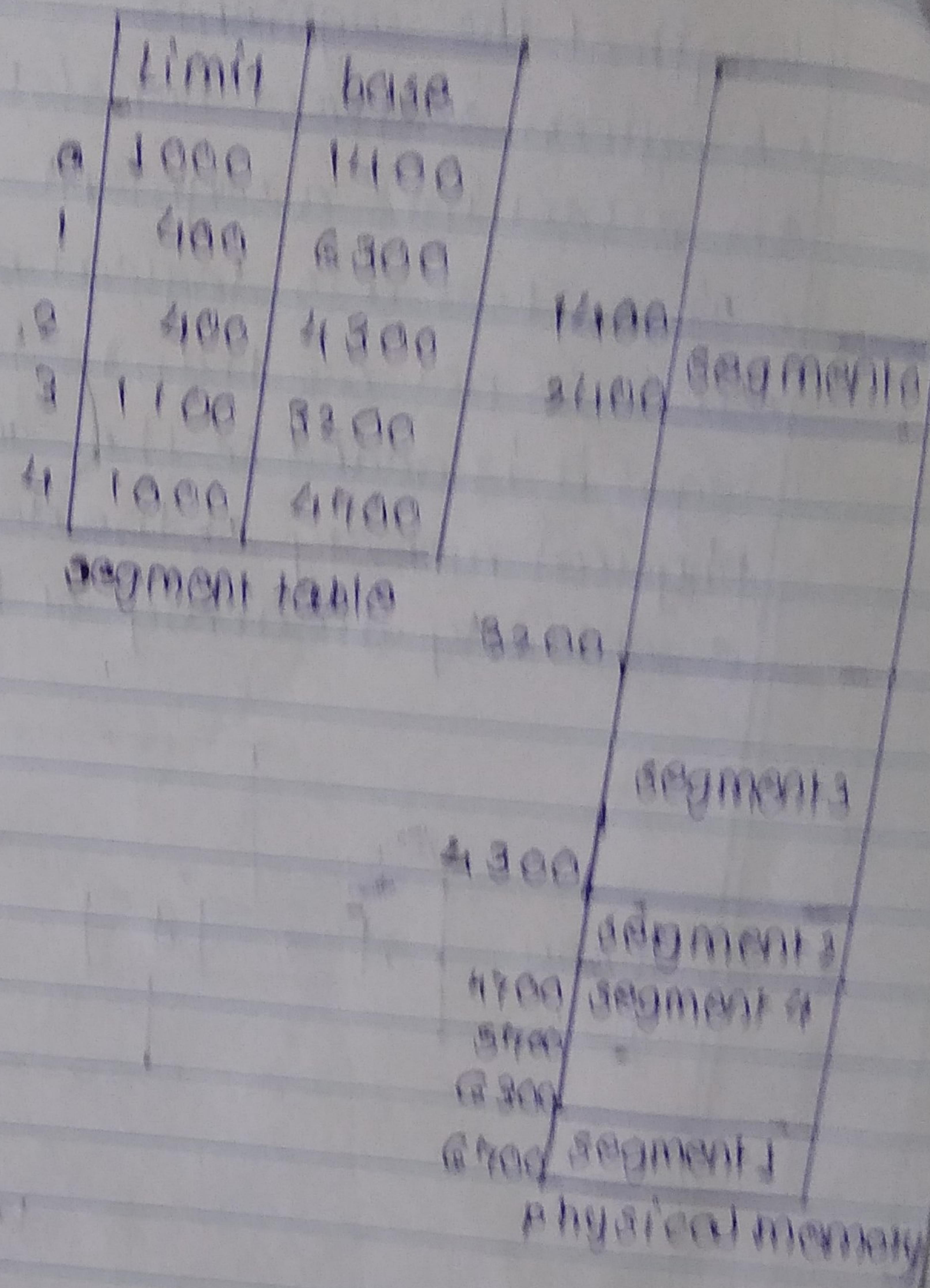
- we have segment numbered from 0 through
- The segment are stored in physical memory
- The segment table has a separate entry for each segments in physical memory (base) & the length of that segment (limit)

example :-

segment 2 is 400 bytes long & begins at location 4300 thus a reference to byte 30 of segment 2 is mapped onto location 4300 + 53 = 4353 a reference to segment 3 byte 852 is mapped of 8200 + 852 = 8282 a reference to byte 1222 of segment 0 would result in a trap to the operating system as this segment is only 1000 bytes long



Logical address space



* **Paging:** permits a physical address space of a process to be non-contiguous.

* Basic Method:

Physical memory is broken into fixed size blocks called frames. Logical memory is also broken into blocks of some size called pages. When a process is to be executed, its pages are mapped into any available memory frames from the backing store. The backing store is divided into fixed size blocks that are of the same size as the memory frames.

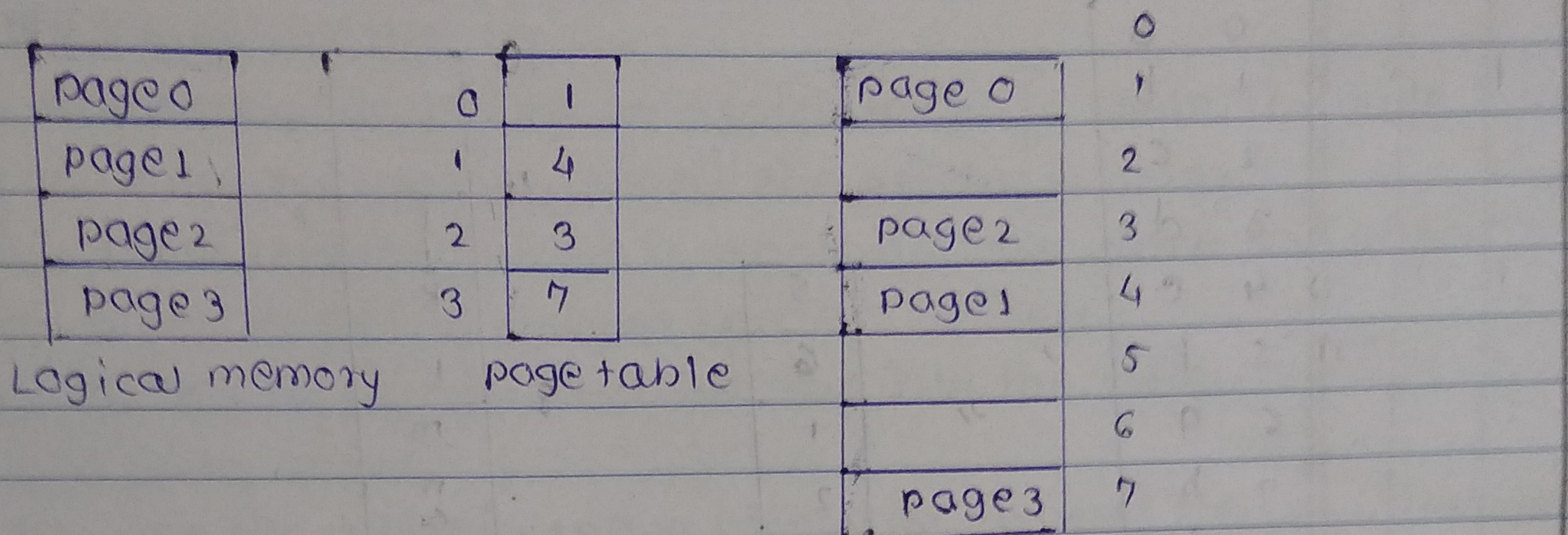


Fig:- paging model of logical & physical memory.

Every address generated by the CPU is divided into two parts 1) page number(p) 2) page offset(d)

1) page-number(p) :- the page number is used as an index into page table. The page table contains the base address of each page in physical memory

2) page offset(d) :- This base address is combined with page offset to define the physical memory address that is sent to the memory unit.

Consider the example given below.

Logical address 0 is page 0, offset 0

Indexing into page table we find that page 0 is in frame 5 thus logical address 0 maps to physical address

$$20 = (5 \times 4) + 0$$

- Logical address 3 (page 0 offset 3) maps to physical address

$$(5 \times 4) + 3 = 23$$

- Logical address 4 (page 1 offset 0) maps to physical address

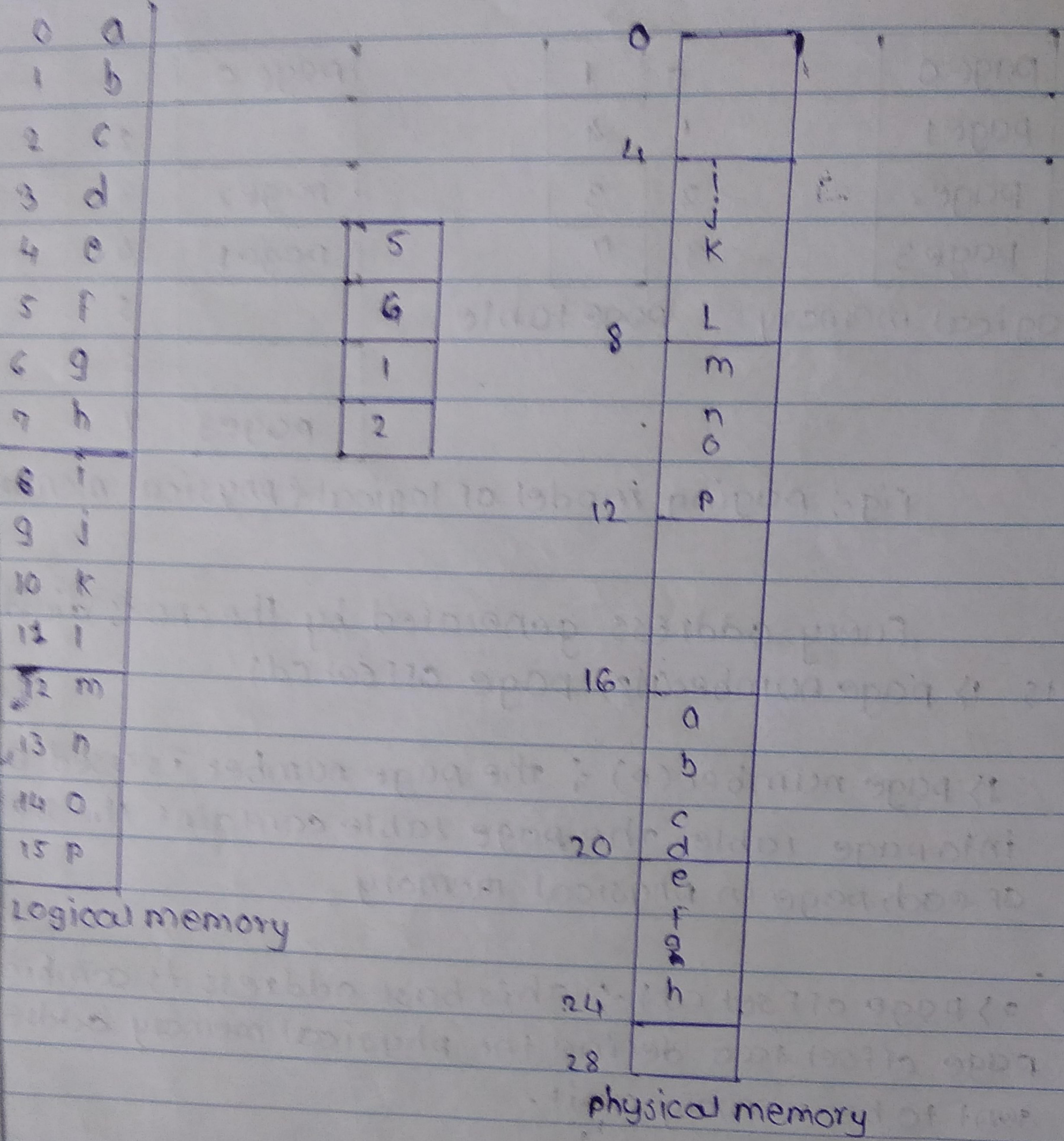
$$(6 \times 4) + 0 = 24$$

- Logical address 13 (page 3 offset 1) maps to physical address

$$(2 \times 4) + 1 = 9$$



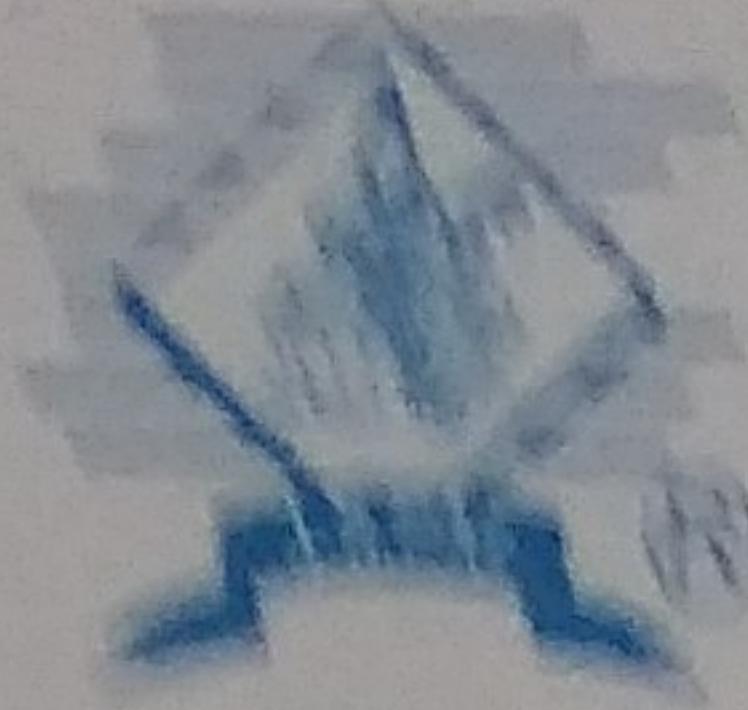
2) Invented page look
- usually
has no



* Structure of the Page Table

1) Hierarchical paging:

- most modern computer systems support a large logical address space (2^{32} to 2^{64}). In such an environment, the page table itself excessively large.
- so to avoid contiguous allocation of page table in main memory we divide the page table into smaller pieces.
- one way is to use a two level paging algorithm in which the page table itself is also paged.



Q. Invented page tables.

- usually each process has an associated page table. the page table has one entry for each page that the process is using
- this table representation is a natural one since process references pages through the pages virtual address.
- the operating system must be translate the reference into physical memory address
- starting the address space identifier ensures that logical page for a particular process is mapped to the corresponding physical page frame

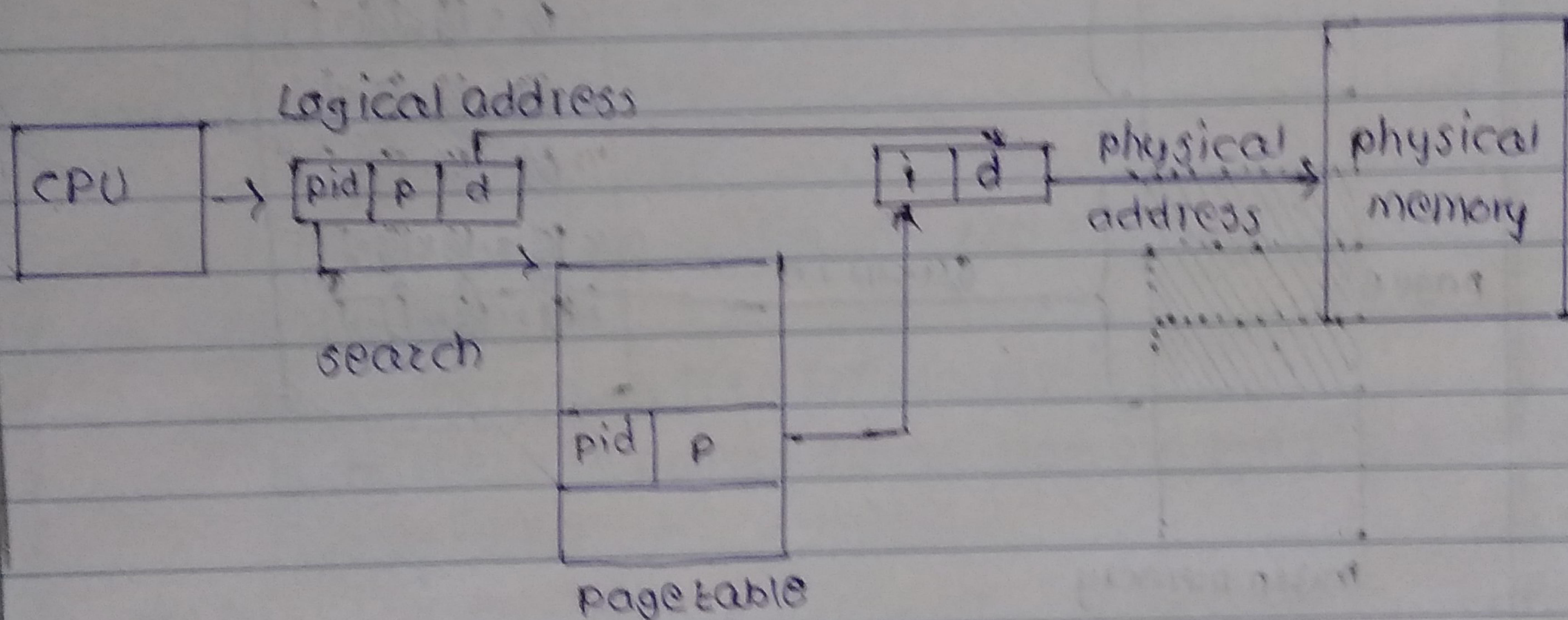


fig - Invented page tables.

Q. 4 Explain demand paging with page fault with details.

An executable might be loaded from disk into physical memory at program execution time.

- But user selects option available in lists this result in the loading the executable code for all options regardless of option is ultimately selected by the user or not.
- An alternative strategy is to initially load page only as they are needed this technique is known as demand paging & is



- commonly used in virtual memory system.
- with demand paged virtual memory pages are only loaded when they are demand during program execution pages that are never accessed are thus never loaded in physical memory.
 - A demand paging system is similar to a paged system with swapping where process reside in secondary memory.

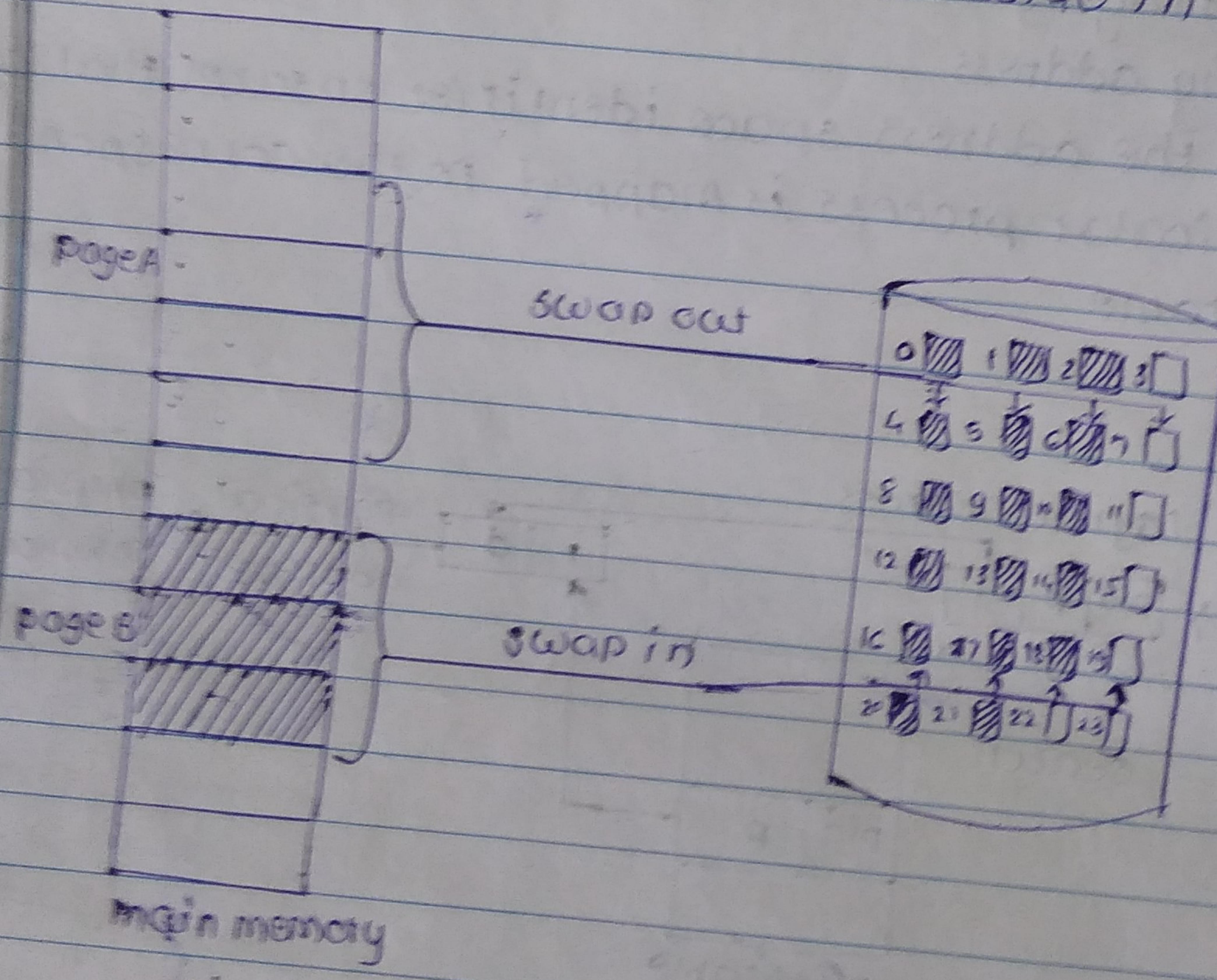
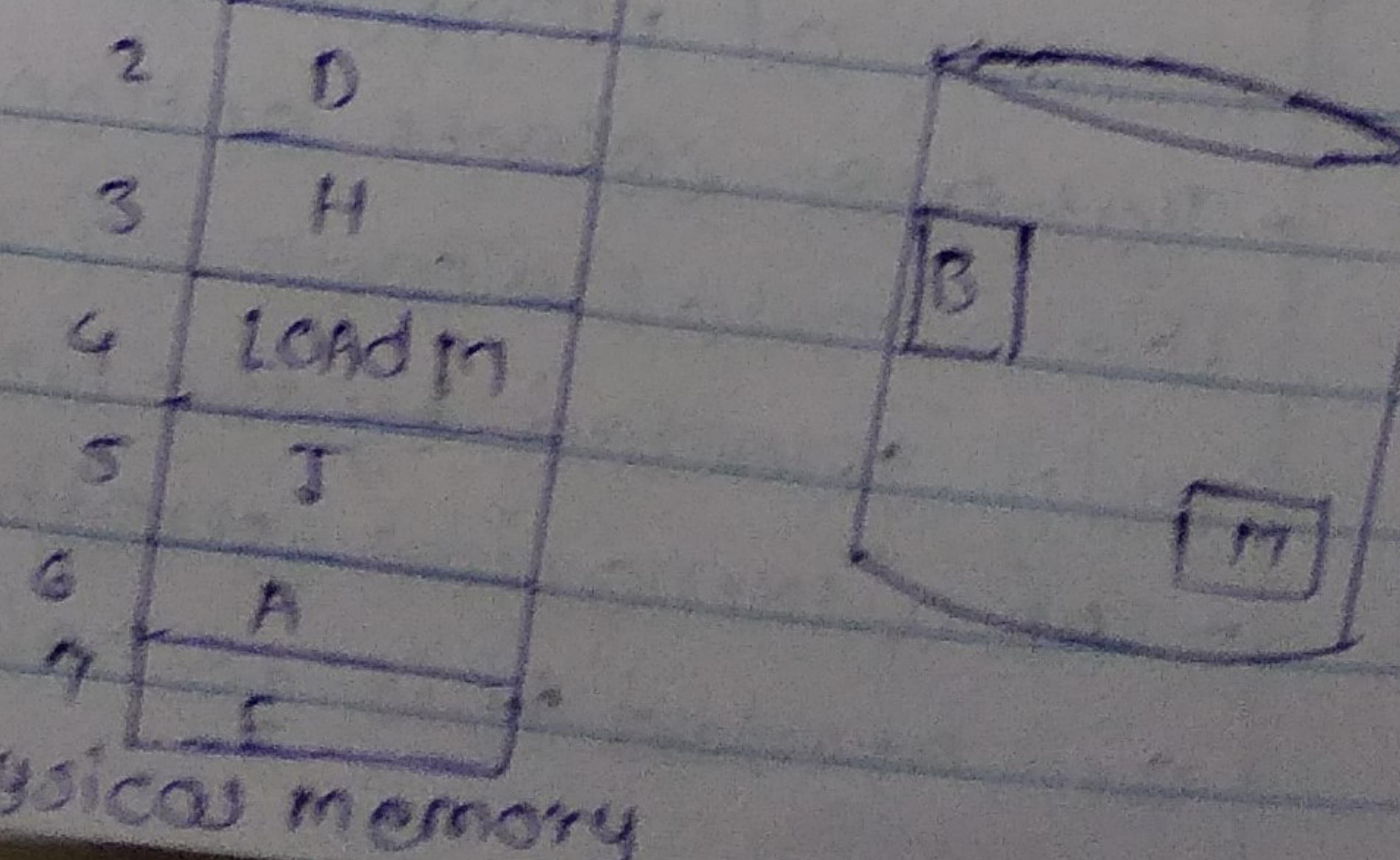


Fig : transfer of a paged memory to contiguous disk space

* Page Replacement :

H		I		J		K		L		M		N		O		P		Q		R		S		T		U		V		W		X		Y		Z					
P → 100dM		3 V		4 V		5 V		6 V		7 V		8 V		9 V		10 V		11 V		12 V		13 V		14 V		15 V		16 V		17 V		18 V		19 V		20 V					
Logical memory for user 1		frame																																							
0 A		1 B		2 C		3 D		4 E		5 F		6 G		7 H		8 I		9 J		10 K		11 L		12 M		13 N		14 O		15 P		16 Q		17 R		18 S		19 T		20 U	
logical memory for page table for user 1		valid-invalid bit																																							

Fig : Page Replacement





If we increase our degree of multiprogramming we run out of allocating memory itself as follows. while a user process is executing a page fault occurs. The OS determines where the desired page is residing on the disk but then finds that there are no free frames. On the free format list, all memory is in use.

- The operating system has several options at this point is could terminate the user process or uses the demand page but it is not the best choice
- The OS could instead swap out process freeing all its frames & reducing the level of multiprogramming. but page replacement is most common solutions.

★ Basic Page Replacement:

The page fault service routine to include page replacement.

- 1) Find the location of the desired page on the disk
- 2) Find a free frame.
 - a) If there is a free frame, use it.
 - b) If there is no free frame, use a page replacement algorithm to select victim frame
- 3) Read the desired page into the newly freed frame changing the page & frame tables
- 4) Restart the user process.

- Page replacement is basic to demand page it computes the separation between logical & physical memory
- If we have multiple process in memory we must decide how many frame to allocate to each process.
- For implementing demand paging we must develop a frame allocation algorithm & a page replacement algorithm

2) FIFO Page Replacement:

- The simplest page replacement algorithm is first-in-first-out (FIFO) algo



JSPM

Page No. :

-use replace the
22 clock

rithm. A FIFO replacement algorithm associates with each page the time when that page was brought into memory.

- when a page must be replaced the oldest page is chosen.
- we can create a FIFO queue to hold all page in memory, we replace the page at the head of the queue.
- when a page is brought into memory we insert it at the queue

3) Optimal page Replacement:

- The optimal page replacement algorithm has the lowest page fault rate of all algorithm & will never suffer from belay's anomaly
- such algorithm does exist & has been called OPT or MIN it is simply this.
 - Replace the page that will not be used for the longest period of time
 - use of this page replacement algorithm guarantees the lowest possible page fault rate a fixed number of frame

4) LRU Page Replacement:

- IF the optimal algorithm is not feasible an approximation of the optimal algorithm is possible
- This approach is the least-recently used (LRU) algorithm.
- The LRU policy is often used as a page replaced algorithm & is considered to be guid.
- The LRU algorithm is somewhat difficult to implement the problem is determine an order for the frames defined by the time of last use their implementation are feasible.

5) Counters:

In the simplest case, we associate with each page table entry a time of use field & add to the CPU logical counter. The clock is incremented for every memory references.

- In this way, we always have the time of the last reference each page



- we replace the page with the smallest LRU.

2) Stack:

Another approach to implementing LRU replacement is to keep a stack of page numbers.

- Removing a page & putting it on the top of the stack them requires changing six pointers at worst.
- LRU replacement does not suffer from belady's anomaly.

3) LRU - Approximation page Replacement:

Many system provides a reference bit to help to page replacement & set by the hw & those bits are associated with each page. Initially, all bits are cleared by o.s. As a user process executes the bit associated with each page.

★ page Fault:

processes resides in secondary memory when we want to execute a process will swap it into memory. Rather than swapping entire process into memory, we use a lazy swapper never swap a page into memory unless that page will be needed. thus it avoid reading into memory pages that will not be used any way decreasing the swap time & the amount of physical memory needed. To distinguish between page in memory, & those on disk valid & invalid bit schema can be used when this bit is set to valid, the associated page is both legal in memory.

This trap is result of the operating system failures to bring the desired page into memory. the procedure for handling this page fault is follows.

- 1) we check an internal table for this process to determine whether the reference was a valid or invalid memory access
- 2) If the reference was invalid we terminate process if it was



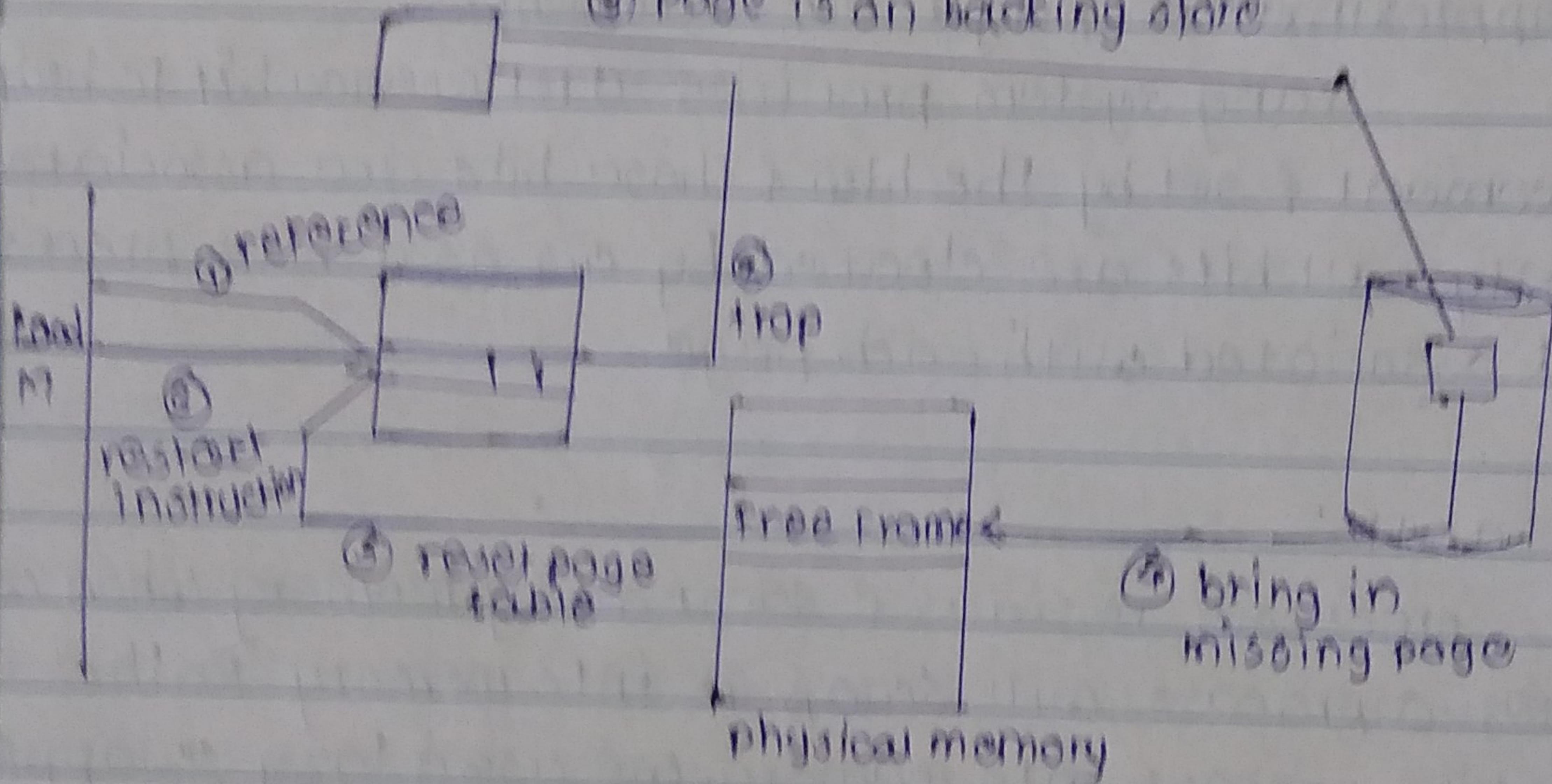
valid but are have not yet brought in that page we now page if we find a free frame

we schedule a disk operation to read the desired page into the newly allocated frame.

when the disk read is complete we modify the internal table kept with the process & the page table to indicate that the page is now in memory.

we restart the instruction that was interrupted by the trap. the process can now access page as though it has always in memory

④ page is an backing store



7 steps in handling a page fault

the hardware to support demand paging is the same as the hardware for paging & swapping

④ page table: this table has the ability to mark an entry invalid through valid-invalid bit or special value of protection bit

⑤ secondary memory: This memory holds those page that are not present in main memory. the secondary memory is usually a high speed disk. It is known as the swap device & the section of disk used for this purpose is known as swap space.



a) i) write short note

ii) deadlock Recovery

There are 2 options for breaking a deadlock

- 1) Abort one or more processes to break the circular wait.
- 2) To preempt some resources from one or more of the deadlocked processes.

i) Process Termination

one of the following 2 methods is used in which the system reclaims all resources allocated to the terminated processes.

a) Abort all deadlocked process:

It clearly breaks the deadlock but a great overhead. The processes may have computes for a long term time & result of these partial computations must be declared & probably recomputed later. This may also cause database or file to be in inconsistent state.

b) Abort one process at a time until deadlock cycle is eliminated.

It incurs less overheads since after each process is discarded, deadlock detection algorithm is invoked to determine whether any process are still deadlocked.

Here we have to decide which processes must be terminated. The answer is termination of process which will incur minimum cost. The factors to be considered while taking this decision are:

1) what the priority of process is

2) How long the process has computed & how much longer the process will compute before completing its designated task.



- (a) How many & what type of resources the process has used
- (b) How many processes will be needed to be terminating
- (c) whether the process is interactive or batch.

(d) Resource Preemption:

The idea is the preemption some resources processes & given these resources to other processes until deadlock cycle is broken.

* Issues need to be address if preemption is required to dead with deadlock.

(e) Selecting a victim:

which process & resource are to be preempted we must determine order of preemption to minimize cost. cost factor include parameters such as number of resources a deadlocked process is holding & the amount of time a deadlocked process has consumed during its execution.

(f) Rollback:

We must rollback the process to some safe state & restart it from that state we can either rollback a process totally or only as far as necessary to break the deadlock.

(g) starvation:

It may happen that the same process is always selected as a victim. As a result this process never completes its designated task. To solve this we must ensure that a process can be picked up as a victim only finite number of times. The most common solution is to include number of times the most rollback in the cost factor.

Resource
assign
see



1.2) Resource Allocation Graph:

A directed edges from resource type $J \rightarrow P$ called as assignment edges signifies that an instance of resource type R_j has been allocated to process P_i . A directed edge $P_i \rightarrow R_j$ is called request edge. In addition we have one type of edges called claim edges, $P_i \rightarrow R_j$ which indicates that process P_i may request resource R_j at some time in future. This edge resembles a request edge in direction but is represented by dashed line.

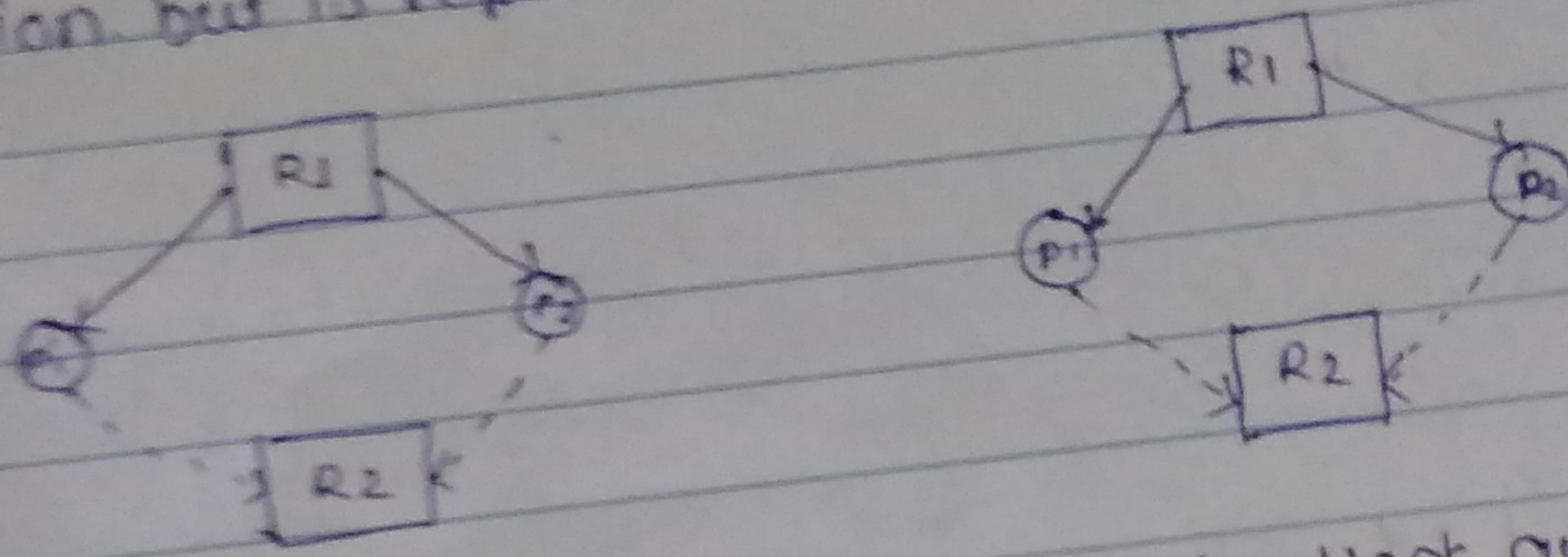


Fig - resource allocation graph deadlock avoidance

- 1) Hold & Wait: A proposed must be holding at least one resource waiting to acquire additional resources that are currently being held by other processes

- 2) No preemption: Resource can not be preempted i.e. resource can be released only voluntarily by the process holding it after that process has completed its task

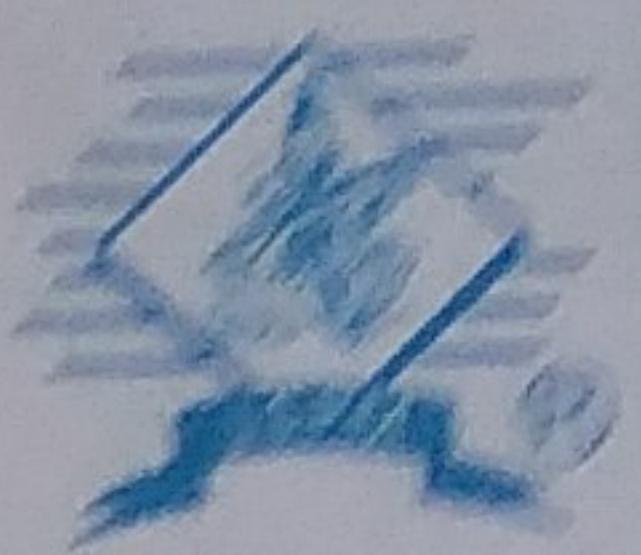
- 3) Circular wait: A set $\{P_0, P_1, \dots, P_n\}$ of waiting process must be exist that P_0 is waiting for a resource that is held by P_1 , P_1 is waiting for a resource that is held by P_2 .

We can deal with deadlock problem in one of the following 3 ways -

1) We can use protocols to prevent or avoid deadlock ensuring that the system will never enter a deadlock state

2) We can allow the system to enter a deadlock state detect it & recover

3) We can ignore the problem altogether & pretend that deadlocks never



JSPM

Page No.:

occur in the system.

dead storage occurs due to fragmentation. In addition, fragmentation may also be caused by other factors such as overwriting can result in fragmentation. When a file is deleted, its space is freed up and can be used by another file. This results in fragmentation.

can be resolved by using fragmentation management.

fragmentation management is a process of reorganizing files and free space in a disk to increase its efficiency.

fragmentation management is a process of reorganizing files and free space in a disk to increase its efficiency.

fragmentation management is a process of reorganizing files and free space in a disk to increase its efficiency.

fragmentation management is a process of reorganizing files and free space in a disk to increase its efficiency.

fragmentation management is a process of reorganizing files and free space in a disk to increase its efficiency.