

Ch#: Virtual memory

What is Virtual memory?

⇒ Virtual memory is a storage technique that provides user an illusion of having a very big main memory.

⇒ This is done by considering a part of secondary memory as the main memory used to execute process that are not in main memory.

⇒ Virtual memory is a separation of logical memory from physical memory.

⇒ In this method, we keep only a part of the process in the memory and other part on the secondary disk.

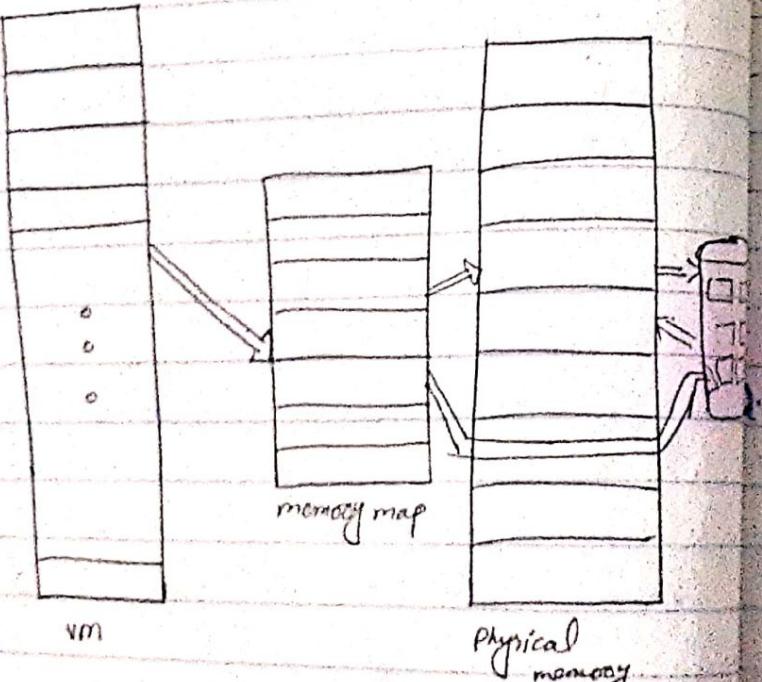
Why virtual memory needed?

- Because main memory has limited size we can do swapping for process replacement, but swapping has a drawback, it decreases the degree of multi-programming.

To overcome this problem, we use the concept of virtual memory.

brake & allow instead of locking
big power in the main memory, the
load the different parts of more than
process in the main memory
o by doing so, the degree of
multi-programming will be increased,
and CPU utilization will also be increased

Diagram:



⇒ It abstracts physical memory into our
extremely large uniform array of
pages

Pros:

- o Degree of multi-programming increases.
- o Bigger sizes programs can easily be executed.
- o Reduces the external fragmentation.
- o Allows only part of the program execution.
- o Allows processes to share files easily and to implement shared memory.
- o Less number of I/O would be needed to load or swap each user program into memory.

Cons:

- o System becomes more slower because swapping takes time.
- o The user will have to concern hard disk space for its use.

How VM implemented?

- o VM in OS is commonly implemented by "demand paging".
- o Can also be implemented in a Segmentation system.

Demand Paging

Definition:-

- o Demand paging is a technique where pages are loaded from secondary memory to main memory only when they are demanded during program execution. Pages that are never demanded.
- o It is a combination of paging & swapping.
- o Also known as Lazy swapping because swapping is done only when CPU requires.
- o Complete program should be present in the backing storage in form of pages.

⇒ Virtual memory is commonly implemented in demand paging.

⇒ It decreases the swap time and the amount of physical memory needed.

⇒ The demand paging system depends on the page table implementation because the page table helps map logical memory to physical memory.

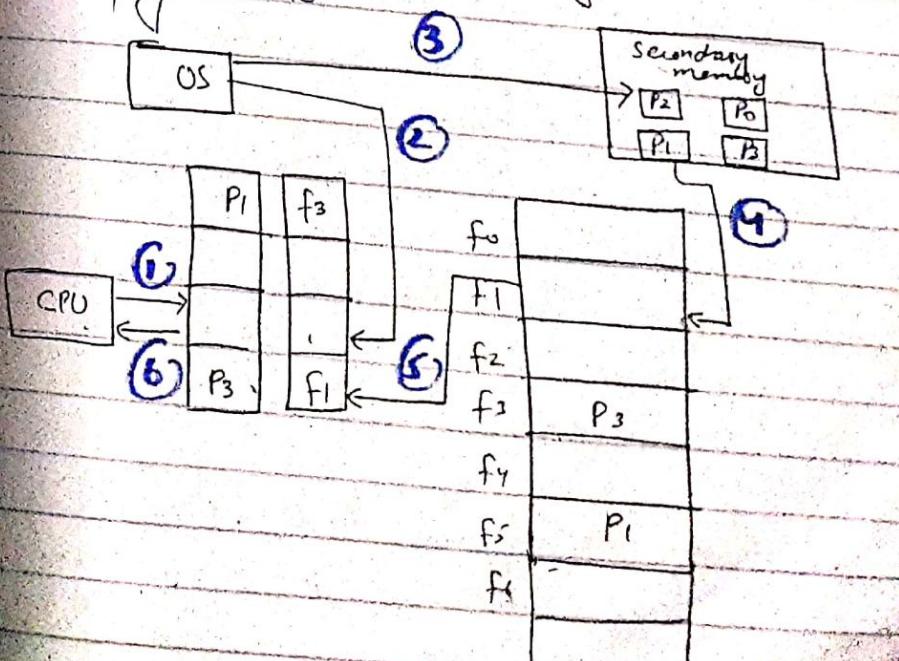
⇒ Bitwise operators are implemented in the page table to indicate that pages are valid or invalid.

⇒ All valid pages exist in primary memory.

⇒ All invalid pages are exist in secondary memory.

Example:-

- o Suppose we have to execute a process P having four pages P_0, P_1, P_2 and P_3 . Currently we have Pages P_1 and P_3 in page table.



Effective memory access time:

⇒ When the page fault rate is 'P' while executing any process, then the effective memory access time is calculated as follows:

$$EMA = (P) * (S) + (1-P) * m$$

where:

- o P is the page fault rate.
- o S is the page fault service time.
- o m is the main memory access time.

⇒ EAT \propto page fault rate.

Pros

- o memory is utilized efficiently.
- o Avoids external fragmentation.
- o Less I/O needed for demand paging.
- o This process is not constrained by the size of physical memory.
- o It becomes easier to share the page.

⇒ Hence, the OS follows these steps when a page fault occurs and the required page is brought into memory.

- o with this technique, pages of the process that are never called are never loaded.
- o no compaction is required in demand paging

Cons:

- o There is an increase in overheads due to interrupts and page tables.
- o memory access time is demand paging is longer.

Q. What is page fault?

- o page fault dominates more like error. It mainly occurs when any program tries to access the data or the code that is in the address space of the program, but that data is not currently located in the RAM of the system.

- o So basically when the page referenced by the CPU is not found in the main memory then this situation is termed as page fault.

whenever any page fault occurs, then the required page has to be fetched from the secondary memory into the main memory.

(Page Replacement)

- o Page replacement is a process of swapping out an existing page from the frame of a main memory and replacing it with the required page.
- o Page replacement is required when:-
- All the frames of main memory are already occupied.

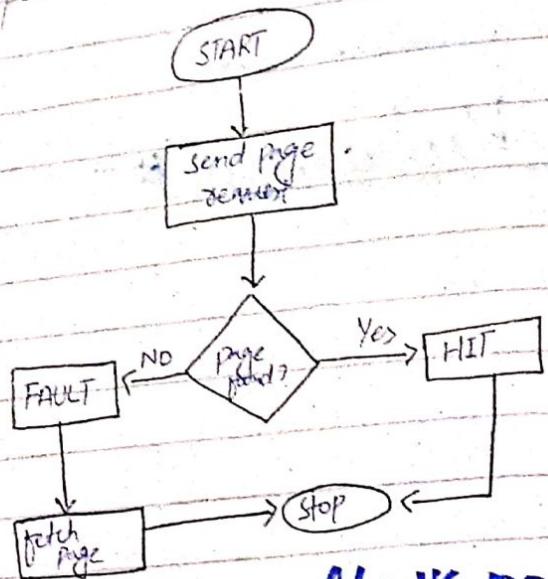
Page replacement algorithm

- o The algorithms used to select the page to be replaced are called page replacement algorithms.
- o The efficiency lies in the least time that is required for a page to be paged in.

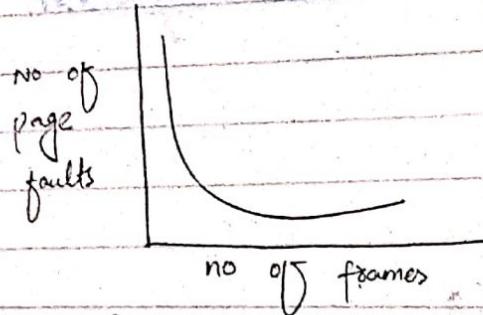
Types:

- o Local page replacement
- o Global page replacement strategy

Why we need?
o The main goal of page replacement alg.
is to provide lowest page fault rate.



No of page faults vs no of frames:



Algorithms:

- o FIFO
- o LRU (Counting based algorithm)
- o Second chance
- o Optimal replacement
- o MRU

i) FIFO:

- o FIFO stands for first in first out page replacement algorithm.
 - o This is the simplest page replacement algorithm.
 - o In this OS keeps track of all pages in the memory in a queue, oldest page is in the front of the queue.
 - o When a page needs to be replaced page in the front of the queue is selected for removal.
 - o Not effective way.
- Pros:
- o Simple and easy to use.
 - o FIFO does not cause overhead.

Cons:

- o Poor performance
- o Does not use the frequency of the least used time and just replaces the oldest page.
- o Supposes from Belady's anomaly.
- o Inefficient use of memory

Example:

- o Reference String: 1, 3, 0, 3, 5, 6
- o no of page slots: 3

P	1	3	0	3	5	6
f ₁	1	1	1	1	5	5
f ₂	-	3	3	3	3	6
f ₃	-	-	0	0	0	0
	X	X	X	V	X	X

no of page faults = 5

no of hit = 1

page fault probability = 5/6

Optimal Replacement:

- o In this algorithm, pages are replaced which are not used for the longest duration of time in the future.

- o This algorithm replaces the page which will not be referred for so long in future.

Pros:

- o Excellent efficiency
- o less complexity
- o Simple data structure can be used to implement.
- o Less no of page faults.
- o Can be used to measure the performance of other algorithms.

Cons:

- o More time consuming
- o Difficult for error handling
- o Need future awareness of the program, which is not possible every time.

Example:

- o Reference String: 7, 0, 1, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2
- o No of page slots: 4

P	7	0	1	8	0	3	0	4	2	3	0	2
f ₁	7	7	7	7	7	3	3	3	3	3	3	3
f ₂	-	0	0	0	0	0	0	0	0	0	0	0
f ₃	-	-	1	1	1	1	1	4	4	4	4	4
f ₄	-	-	-	2	2	2	2	2	2	2	2	2
	X	X	X	X	V	X	V	V	V	V	V	V

no of page faults = 6

no of hit = 7

total pages = 13

Pros:

- o Does not suffer from Belady's anomaly.
- o it is open for full analysis.
- o more efficient

Cons:

- o EXPensive and more Complexity.
- o Need for additional data structure.

Example:

Reference String: 7 0 1 2 0 3 0 4 2 3 0 8
No. of page slots: 4

P	7	0	1	2	0	3	0	4	2	3	0	3	2
f ₁	7	7	7	7	7	3	3	3	3	3	3	3	3
f ₂	-	0	0	0	0	0	0	0	0	0	0	0	0
f ₃	-	-	1	1	1	1	1	1	4	4	4	4	4
f ₄	-	-	-	2	2	2	2	2	2	2	2	2	2
	X	X	X	X	V	X	V	V	V	V	V	V	V

no of page fault = 6

no of hit = 7

(3) LRU

o page will be replaced which is least recently used.

o if replace the page that has not been referenced for the longest time.

o LRU works on the idea that pages have been most heavily used in the past few instutions are most likely to be used

heavily in the next few instutions too.

o we are looking in the past which has not been used longest time.

(4) LFU:-

- In LFU, the page with the minimum count frequency is selected for replacement with the page that needs to enter into the system.
- A counter is assigned to every page that is loaded into the memory.
- LFU algorithm is sometimes also combined with LRU replacement algorithm & then implemented.
- If the frequency of pages remains constant, the page that comes first is loaded first.

Example:

no of page slot = 3

Reference String: 7 0 1 2 0 3 0 4 2 3 0 3 2

P.	7	0	1	2	0	3	0	4	2	3	0	3
f ₁	7	7	7	2	2	2	2	4	4	3	3	1
f ₂	-	0	0	0	0	0	0	0	0	0	0	0
f ₃	-	-	1	1	1	3	3	3	2	2	2	2
	X	X	X	X	V	X	V	X	X	V	X	X

Frequency: $0=2$ $1=10$ $2=10$ $3=10$ $4=1$

(5) MRU

- Also known as fetch-and-displace algorithm.
- Deal with the case such as sequential scanning access pattern.

Example:

RS : 7 0 1 2 0 3 0 4 2 3 0 3 2 1 2
frame: 3

P.	7	0	1	2	0	3	0	4	2	3	0	3	2	1
f ₁	7	7	7	7	7	7	7	7	7	7	7	7	7	7
f ₂	-	0	0	0	0	0	0	3	0	4	4	4	4	4
f ₃	-	-	1	2	2	2	2	2	2	2	3	0	3	2
	X	X	X	X	V	X	X	V	X	X	V	X	X	X

no of page fault: 13
no of hit = 2

(b) Second Chance Algorithm

- o Modification of FIFO replacement algorithm.
- o Drawback in FIFO is that even if page is in use, it may be replaced due to its arrival time.
- o A reference bit or use bit is used to give information about regarding whether page has been used.
- o Initially, reference bit for all pages are 0.
- o As the page is referenced, bit is set to 1 to indicate that the page is in use.
- o A page whose reference bit is 1 will be replaced, and will be given second chance.

Example:

Frame = 3

Reference string: 2 3 2 1 5 2 4 5 3

frame	2	3	2	1	5	2	4	5	3
f ₁	2°	2°	2°	2°	2°	2°	2°	2°	3°
f ₂	3°	3°	3°	5°	5°	5°	5°	5°	
f ₃	-	-	1°	1°	1°	4°	4°	4°	
	X	X	✓	X	X	✓	X	✓	X

PF = 6 Hit = 3 PRR = 6/9

Belady's Anomaly

Also called FIFO anomaly.

- o Belady's anomaly is the phenomenon of increasing the number of page faults on increasing the number of frames in main memory. It is not always happen but for some reference string.
- o Following algorithms suffer from Belady's anomaly:

o FIFO

o Random page replacement

o Second chance

- o An algorithm suffers from Belady's anomaly if it does not follow stack property.

Why stack based algorithm do not suffer from Belady's anomaly?

- o Because these type of algorithms assign a priority to a page that is independent of the number of page frames.

o A stack-based algorithm is one for which it can be shown that the set of pages in memory for N frames is always a subset of the set of pages that would be in memory with N+1 frames.

- o if the number of frames increases then in pages will still be the most recently referenced pages.

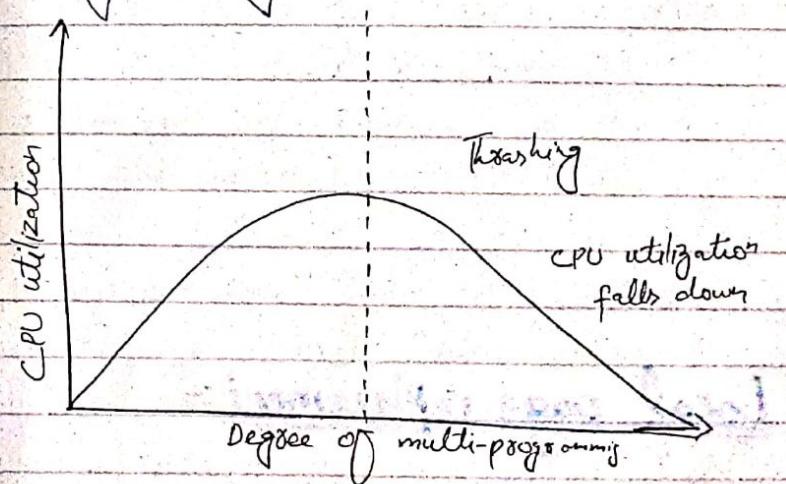
- o while in FIFO, if a page b comes in physical memory before a page - a then priority of replacement of b is greater than a, but this is not independent of the number of page frames.

\Rightarrow Because stack based algorithms make their page replacement decisions based on future references rather than simply relying on past references like the FIFO.

Thrashing

\Rightarrow A situation when system is spending most of the time in servicing page faults (swapping), than executing the actual process.

\Rightarrow In degree of multiprogramming, thrashing occurs because page faults are more in degree of multi-programming, then CPU spends most of his time in swapping the pages between main memory and the secondary memory.



Causes:

- o when more processes in memory as compared to the available no. of frames in memory. o Lack of frames.
- o high degree of multiprogramming.
- o Page replacement policy.

Effect of Thrashing:-

- ⇒ At the time, when thrashing starts then the operating system tries to apply either the
- o Global page replacement algorithm
 - o Local page replacement algorithm

Global page replacement algorithm

- o Global page replacement can bring any page, if tries to bring more pages when looking is found.
- o But what actually will happen is that no process gets enough frames, and as a result, the thrashing will increase more and more.
- o Therefore, the global page replacement algorithm is not suitable when thrashing happens.

Local page replacement:-

- o Unlike the global page algorithm, local page replacement will select pages which only belong to that process.
- o So there is a chance to reduce the thrashing.

- o But it is proven that there are many disadvantages if we use local page replacement.
- o Therefore, local page replacement is just an alternative to global page replacement in a thrashing scenario.

Techniques to Handle Thrashing

We already have local replacement that is better than global replacement to avoid thrashing.

- o But it also has disadvantages and not suggestible. Some more techniques are:
 - o Working set model
 - o Page fault frequency

Working set model:-

- o This model is based on locality model.
- o A locality is a set of pages that are actively used together.
- o The locality model states that as a process executes, it moves from one locality

Def: working set model is a conceptual model proposed by Peter Denning to prevent thrashing.

To analyse:

- o A program is generally composed of several different localities which may overlap.
- o Working set model is a dynamic page replacement algorithm that allocates frames to a particular process assuming that the nearest future of pages will be used is a closed approximation of the recent past.

page in memory.

- o The model uses a parameter Δ to define the working-set window.

Working set: The set of pages in the most recent page reference.

(working set) Δ : No of pages in the most recent page reference.

Example:

(i) Page references: 2 6 1 5 7 7 7 5 1 / t_1

$$\Delta = 10$$

$$WS(t_1) = \{1, 2, 5, 6, 7\}$$

(ii) Page references = 3 4 4 3 4 3 4 4 4 / t_2

$$\Delta = 10$$

$$WS(t_2) = \{3, 4\}$$

o if a page is in active use, it will be in working set.

o if it is no longer used, it will be dropped from the working set.

Example:

0 2 6 1 5 7 7 7 5 1, 6 2 3 4 1 2 3
4 4 3 4 3 4 4 4 1 3 2 3

$$\Delta = 10$$

$$WS(t_1) = \{1, 2, 3, 4, 5, 6, 7\}$$

\Rightarrow Through working set we can suspension and resume of the process can be done and prevent from thrashing.

Example:-

Assume working set of window set is.

4, initially pages in working set are d, a

{Assumeearliest 8-5 is ... e, e, d, a, ...}

what is the no of hits and misses

n of frames used by working set

of next 10 ref pg page is

{e, d, a} c, c, d, b, c, e, c, e, a, d

$t_1 - \{e, d, a, c\} - \text{miss} - 4$

$t_2 - \{d, a, c\} - \text{Hit} - 3$

$t_3 - \{a, c, d\} - \text{Hit} - 3$

- t₄ - [a, d, b] → miss - 3
- t₅ - [d, b, c] - H - 3
- t₆ - [d, b, c, e] - m - 4
- t₇ - [b, c, e] - H - 3
- t₈ - [c, e] - H - 2
- t₉ - [a, c, e] - m - 3
- t₁₀ - [c, e, a] - m - 3

(1) No of hits = 5

(2) No of frames on average = $\frac{32}{10} = 3.2$

$$D = ws(t_1) + ws(t_2) + \dots$$

Working set size :

Disadvantages:

- o Monitoring the working set is a very clumsy work.

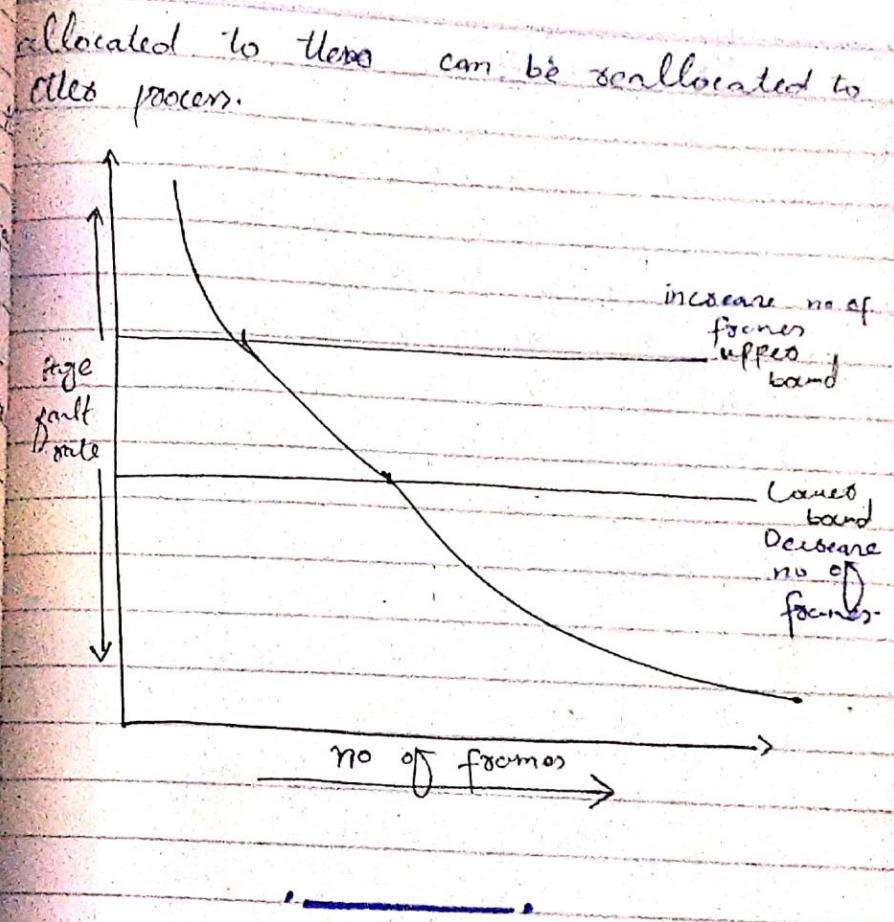
Note:

- o if enough empty frames in physical memory then initial address process.
- o $D >$ empty frames then suspend process

- o In working set model processes are divided into two groups:
 - o Active
 - o Inactive
- o When a process is active its entire working set must always be in memory never execute a thread whose working set is not resident.
- o When a process becomes inactive, its working set can migrate to disk. Threads from inactive processes need scheduled for execution.
- o balance set:
 - o The collection of active processes is called the balance set.
 - o The system must have a mechanism for gradually moving processes into and out of the balance set.
 - o As working sets change, the balance set must be adjusted.

2) Page fault frequency:

- o it is direct approach for handling thrashing.
- o In this method we monitor the page fault rate of each process.
- o High page fault rate:
 - indicates that process has very few frames allocated to its.
- o Low page fault rate:
 - indicates that process has large no. of frames allocated.
- o In order to prevent thrashing, upper limits have been set on the desired page fault rate.
 - o if page fault rate < lower limit.
 - frames can be removed from the process.
 - o if page fault rate > upper limit.
 - more frames allocate to this process.
 - o if there is no free frames then of the process to be suspended and fix.



- o What is belady's anomaly? (repeat)
- o FIFO and LRU both use previous information in page replacement policy? How is the one different from another than? (repeat)
- o Differentiate b/w LFU and MFU algorithm
- o What is demand paging? (repeat)
- o What is virtual memory? (repeat)
- o Under what circumstances do page fault occurs?
 - o How page fault frequency can be used as method of hashing?
 - o What is the cause of hashing?
 - o Why page table is needed?
 - o What is page fault? (repeat)
 - o Define hit ratio. (repeat)
 - o What is reference bit or dirty bit?
 - o What is hashing?
 - o What causes page faults?
 - o Discuss different threading issues in detail.

Q#1) Apply LRU and optimal algorithm.
 7, 0, 1, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2, 1, 2, 0, 1, 7, 0, 1
 frame size is 3.

Q#2) Apply LRU & optimal algorithm
 4, 3, 1, 2, 1, 1, 1, 3, 4, 2, 1, 5, 6, 3, 1, 2, 3, 7, 6, 3, 2, 1, 2, 3, 6
 frame size is a) 5 b) 6

Q#3) Apply LRU, FIFO & Optimal.
 6, 7, 0, 2, 1, 3, 3, 4, 2, 1, 5, 6, 2, 1, 2, 3, 7, 6, 3, 2, 1, 2, 3, 6
 frame size a) 3 b) 4 c) 6

Q#4) Apply LRU, LFU & FIFO.
 4, 3, 2, 1, 5, 1, 2, 3, 2, 3, 2, 3, 4, 5, 5, 5
 frame size 3.

Q#5) Apply FIFO & optimal.
 1, 3, 4, 2, 4, 5, 6, 4, 2, 3, 1, 7, 8
 frame size 3.

Q#6) Apply LFU and optimal
 1, 2, 3, 4, 1, 2, 3, 5, 1, 3, 2, 1, 4, 5, 6, 3, 2
 frame size is 3.

Q#7) Apply LRU, FIFO & optimal
 8, 2, 5, 1, 9, 8, 4, 6, 1, 7, 2, 1, 6, 1, 8, 7, 8, 2, 1, 7
 frame size a) 3 b) 4 c) 5

Q#8) Apply Optimal Page replacement
1, 2, 1, 1, 3, 4, 2, 1, 5, 6, 3, 1, 2, 3, 7, 6, 3, 2, 1, 2
frame size 4 and 7.

Q#9) Apply LRU and FIFO
3, 1, 9, 2, 4, 8, 6, 9, 2, 1, 3, 7, 5, 8, 1, 9, 3, 4
frame size 3.

Q#: What is reference bit or dirty bit?

o It is used in page replacement algorithm.

o If no frames are free, two pages transfer are required (one out and one in).

o This situation effectively doubles the page fault service time and increases the effective access time.

o We can reduce this overhead by using a modify bit or dirty bit.

o It is associated with each page of frame.

o It indicates whether a page has been modified or not.

o It is set by the hardware and used to show the page is modified after being loaded into the cache.

What is reference string?

o An algorithm can be evaluated by running it on a particular string of memory references and computing the number of page faults. The string of memory references is called a reference string.

o It can be generated artificially.

What is equal allocation?

o The easiest way to split m frames among n processes is to give everyone an equal share, m/n frames.

o For instance, if there are 93 frames and five processes, each process will get 18 frames. The three leftover frames can be used as a free frame buffer pool.

o This is called equal allocation.

What is proportional allocation?

o In proportional allocation, we allocate available memory to each process according to its size.

o If total no. of frames are m , we allocate a_i frames to process p_i , where a_i is approximately:

$$a_i = s_i / S * m$$

What is pure demand paging?

- o In the extreme case, we can start executing a process with no pages in memory.
- o When the OS sets the instruction pointer to the first instruction of the process, which is on a non-memory resident page, a page fault occurs immediately.
- o Pages are loaded with on-demanding of the process. This is called PDP.

What is Copy-on-Writer

- o Copy-on-write is a benefit of virtual memory.
- o COW allows both parent and child processes to initially share the same page in memory.
- o If either process modifies a shared page, only then is the page copied.
- o COW allows more efficient process copy as only modified pages are copied.
- o Free pages are allocated from a pool of zeroed-out pages.