# Answer to the question number:1

**a) What do you mean by an Operating System? What are the three main purposes of an operating system? Write a short note on, "An operating system is a resource allocator".**

An Operating System (OS) is a software that acts as an intermediary between computer hardware and user applications. It manages computer hardware resources and provides common services for computer programs. Here are the three main purposes of an operating system:

1. **Resource Management:** The OS manages computer hardware resources such as CPU (Central Processing Unit), memory (RAM), storage devices, and input/output (I/O) devices. It allocates these resources efficiently among various running processes, ensuring they have fair and timely access to the resources they need.

2. **Abstraction:** The OS provides a layer of abstraction to shield user applications from the complexities of the underlying hardware. It presents a simplified interface, allowing programmers to interact with hardware using standardized commands and system calls without needing to understand the intricate details of hardware architecture.

3. **Control and Coordination:** The OS controls the execution of computer programs, ensuring they run smoothly and securely. It coordinates various system activities, including process scheduling, memory management, file system operations, and device I/O, to maintain system stability, reliability, and security.

Now, focusing on the aspect of an operating system as a "resource allocator":

An operating system acts as a resource allocator by efficiently distributing computer resources among competing processes or tasks. It ensures that each process receives a fair share of resources such as CPU time, memory space, and I/O bandwidth according to its priority and requirements.

Resource allocation involves making decisions about which processes should run, how much CPU time each process should receive, and how memory should be allocated among different processes. The OS employs various scheduling algorithms, memory management techniques, and I/O control mechanisms to optimize resource utilization and maximize system throughput.

By effectively managing and allocating resources, the operating system enhances system performance, improves responsiveness, and ensures the efficient utilization of hardware resources, thereby enabling multiple concurrent processes to execute smoothly without interference or resource contention.

**b) List five services provided by an operating system and explain how each creates convenience for users.**

The five services are:

**a. Program execution.** The operating system loads the contents (or sections) of a file into memory and begins its execution. A user-level program could not be trusted to properly allocate CPU time.

**b. I/O operations.** Disks, tapes, serial lines, and other devices must be communicated with at a very low level. The user need only specify the device and the operation to perform on it, while the system converts that request into device- or controller-specific commands. User-level programs cannot be trusted to access only devices they should have access to and to access them only when they are otherwise unused.

**c. File-system manipulation.** There are many details in file creation, deletion, allocation, and naming that users should not have to perform. Blocks of disk space are used by files and must be tracked. Deleting a file requires removing the name file information and freeing the allocated blocks. Protections must also be checked to assure proper file access. User programs could neither ensure adherence to protection methods nor be trusted to allocate only free blocks and deallocate blocks on file deletion.

**d. Communications.** Message passing between systems requires messages to be turned into packets of information, sent to the network controller, transmitted across a communications medium, and reassembled by the destination system. Packet ordering and data correction must take place. Again, user programs might not coordinate access to the network device, or they might receive packets destined for other processes.

**e. Error detection.** Error detection occurs at both the hardware and software levels. At the hardware level, all data transfers must be inspected to ensure that data have not been corrupted in transit. All data on media must be checked to be sure they have not changed since they were written to the media. At the software level, media must be checked for data consistency; for instance, whether the number of allocated and unallocated blocks of storage match the total number on the device. There, errors are frequently process-independent (for instance, the corruption of data on a disk), so there must be a global program (the operating system) that handles all types of errors. Also, by having errors processed by the operating system, processes need not contain code to catch and correct all the errors possible on a system.

**c) Give two reasons why caches are useful. What problems do they solve? What problems do they cause? If a cache can be made as large as the device for which it is caching (for instance, a cache as large as a disk), why not make it that large and eliminate the device?**

**Why Are Caches Useful?**

The main benefit of a cache comes from the principle of **program locality**, which is the tendency for programs to access the same data or code repeatedly within a short time span. This principle has two main types:

1. **Temporal Locality**: Data or instructions that were recently accessed are likely to be accessed again soon.
2. **Spatial Locality**: Data or instructions that are near each other in memory are likely to be accessed together.

Most programs spend the majority of their execution time accessing a small subset of code and data. For example, a common rule of thumb is that a program spends **90% of its time executing just 10% of its code**. This means if this frequently used data and code can be stored in a smaller, faster memory (the cache), the program can run much faster.

By storing this subset in the cache, the CPU can access it more quickly than if it had to reach out to the main memory (RAM) each time. Caches are therefore used to bridge the **performance gap** between the CPU and the main memory. Without caches, the CPU would spend a lot of time waiting for data, leading to inefficient processing.

## What Problems Do Caches Solve?

1. **Speed Discrepancy**: There is a significant speed gap between the CPU and main memory. CPUs are much faster, and caches help bridge this gap by providing a small, high-speed memory that the CPU can access much more quickly than main memory.
2. **Efficient Access to Frequently Used Data**: Caches store the most frequently accessed data, reducing the need to fetch it repeatedly from slower memory, thereby increasing the overall speed of data access.
3. **Improved Performance**: By storing commonly used instructions and data close to the CPU, caches significantly reduce the average time it takes to access memory, leading to improved program performance.

## What Problems Do Caches Cause?

1. **Coherency Issues**: In a system with multiple processing units or cores, each might have its own cache. When one core updates a value in its cache, it can create inconsistency if other cores are working with outdated values in their caches. This is called a **cache coherency problem**.
2. **Increased Complexity**: Cache management adds complexity to the system. The CPU needs mechanisms to decide what data to store in the cache, what data to evict when the cache is full, and how to synchronize data across multiple caches in multi-core systems.
3. **Multitasking Overhead**: When the CPU switches between processes (context switching), the cached data might no longer be relevant for the next process, causing more cache misses. This means the CPU has to fetch new data from slower memory, reducing the benefits of caching during heavy multitasking.
4. **Storage Limitations**: Caches are relatively small because of their cost and size limitations. If data is not in the cache (cache miss), the CPU has to go to main memory, which is slower.

## Why Not Make Cache as Large as the Device Being Cached?

There are two main reasons why caches are not made as large as the device (such as main memory or a disk) they are caching for:

1. **Cost**: Cache memory, typically made from SRAM (Static RAM), is much more expensive than DRAM (Dynamic RAM), which is used in main memory. It is also significantly more expensive than storage devices like hard disks or SSDs. Making a large cache would be financially impractical due to the high cost of fast, small-sized cache memory.
2. **Speed Advantage**: The very fast access speed of a cache is partly due to its small size. If the cache were to be as large as main memory, it would lose some of its speed advantage, making it less effective as a high-speed memory layer. Large memory caches would also introduce more complexity in managing and searching the cache, which could slow down data retrieval.

# Answer to the question number:2

**a) Direct memory access is used for high-speed I/O devices in order to avoid increasing the CPU's execution load.**

**i) How does the CPU interface with the device to coordinate the transfer?**
**ii) The CPU is allowed to execute other programs while the DMA controller is transferring data.**

a. How does the CPU interface with the device to coordinate the transfer?
b. How does the CPU know when the memory operations are complete?
c. The CPU is allowed to execute other programs while the DMA controller is transferring data. Does this process [the DMA controller's data transfers] interfere with the execution of the user programs? If so, describe what forms of interference are caused.
d. Since multicore processors are so common now, what would be the effect of using one of the cores as a replacement for a separate DMA controller?
e. It's not quite right to say that DMA should be used with "high-speed I/O devices". How else should the class of devices be described, to help explain why and when DMA is useful?

a. Commands and status requests are sent from the CPU to the device via some special registers (I/O ports, or memory-mapped I/O registers, depending on the system design). The device notices when these registers are written-to, and will take the appropriate action (perhaps after a delay, if it is occupied with some other operation).
b. The device could write to a special register, or generate an interrupt. In either case, the processor notices what has happened (perhaps after a delay, if it is occupied with some other operation).
c. The CPU and the device will compete for cycles on the memory bus. The memory controller will try to fairly allocate bus cycles between the CPU and the device. So, any CPU program that would be capable of using all the memory bus cycles could run more slowly while the DMA is active.
d. The purpose of a DMA controller is to move data between memory and a device, and not much more. If you have only two cores, it seems wasteful to take a general-purpose thing (one core) and use it as a special-purpose thing (a DMA controller). If you have 12 cores, maybe the impact (one of 12, vs. 1 or 2) won't be so bad. You could also consider the effects on cache, since the DMA controller by-passes cache, while the core probably does not (that depends on details of the processor design).
e. You also need to consider the volume of data, not just the rate at which it is moved. Moving one byte by DMA has a high overhead compared to programmed I/O, for example. Given the time-cost of generating and handling interrupts, if only a small amount of data is moved

at high speed, the overhead time dominates the transfer time. It would be better to say "block-oriented I/O devices" (or something like that) to emphasize "lots of data for a long time".

a) In computer systems, communication between the CPU and devices is handled through specialized registers, which can be either I/O ports or memory-mapped I/O registers, depending on the system architecture. The CPU sends commands and status requests to these registers, and the device responds accordingly. When the CPU writes to these registers, the device detects the operation and carries out the appropriate task. However, if the device is busy with another task, there might be a slight delay before it acts on the command.

b) Once the device completes the requested action, it can notify the CPU in two ways: by writing to a special register or by generating an interrupt. The processor, when available, detects these notifications and takes appropriate action, though there could be a delay if it is engaged with another process.

c) The CPU and devices share access to the memory bus, which they both use to transfer data. The memory controller's role is to manage and fairly allocate bus cycles between the CPU and devices, ensuring neither monopolizes the bus. However, when direct memory access (DMA) is active, a CPU program that is capable of utilizing all available memory bus cycles might experience a slowdown as the bus resources are shared with the DMA controller.

d) The main function of a DMA controller is to move data between memory and devices efficiently. The advantage of DMA is that it offloads the task of data transfer from the CPU, allowing the CPU to focus on other tasks. However, if a system has limited processing power, such as with only two cores, using one core solely for data transfer through programmed I/O might seem inefficient compared to leveraging a dedicated DMA controller. In systems with many cores (like 12 cores), dedicating one core to handle data transfers may have less of an overall performance impact, but factors such as cache usage must also be considered, since DMA typically bypasses the cache, unlike the CPU, which may rely on it.

e) In terms of data transfer volume, DMA is particularly efficient for handling large blocks of data. For smaller transfers, the overhead associated with initiating DMA operations, including generating and handling interrupts, might outweigh the speed benefits. In such cases, programmed I/O, which directly involves the CPU in data transfer, could be more effective for transferring small amounts of data quickly. This is why DMA is generally more suited for "block-oriented I/O devices," where large volumes of data are transferred over a longer duration. The efficiency of DMA shines when handling large data transfers, reducing the need for the CPU to be heavily involved in the process.

**b) Describe the differences among short-term, medium-term, and long-term scheduling.**

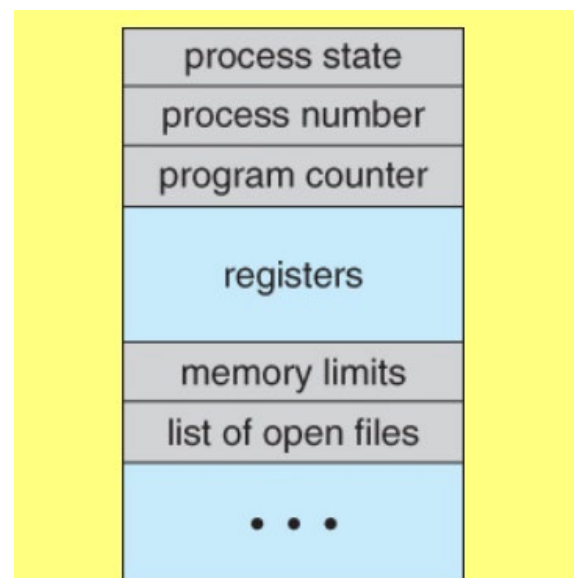Here's a tabular comparison of short-term, medium-term, and long-term scheduling:

| Aspect | Short-term Scheduling | Medium-term Scheduling | Long-term Scheduling |
|---|---|---|---|
| **Purpose** | Select which process to execute next on CPU. | Manage processes in memory, swapping if necessary. | Admit new processes into the system. |
| **Frequency** | Occurs frequently (milliseconds/microseconds). | Occurs less frequently than short-term scheduling. | Occurs infrequently (when new processes arrive). |
| **Decision Criteria** | Based on process priority, CPU burst, etc. | Based on memory management, swapping processes. | Based on system load, resource availability. |
| **Scope** | Operating at CPU level. | Operating at memory management level. | Operating at system level. |

These distinctions highlight how each type of scheduling focuses on different aspects of managing processes and resources within the operating system, optimizing system performance and resource utilization.

**c) What do you know about process and process control block (PCB)? What are the states of process? Explain with necessary diagrams.**

A process is an instance of a program in execution. It represents the execution context of a program, including the program code, data, execution state, and system resources allocated to it. In modern operating systems, processes are fundamental units of resource allocation and management, allowing multiple tasks to run concurrently on a single processor.

A Process Control Block (PCB) is a data structure used by the operating system to manage information about each process in the system. It contains various pieces of information required for process management, including process state, program counter, CPU registers, memory management information, and scheduling parameters.



Figure 3.3 - Process control block ( PCB )

For each process there is a Process Control Block, PCB, which stores the following ( types of ) process-specific information, as illustrated in Figure 3.1. ( Specific details may vary from system to system. )

Process State - Running, waiting, etc., as discussed above.

Process ID, and parent process ID.

CPU registers and Program Counter - These need to be saved and restored when swapping processes in and out of the CPU.

CPU-Scheduling information - Such as priority information and pointers to scheduling queues.

Memory-Management information - E.g. page tables or segment tables.

Accounting information - user and kernel CPU time consumed, account numbers, limits, etc.

I/O Status information - Devices allocated, open file tables, etc.

# Answer to the question number:3

**a) What are two differences between user-level threads and kernel-level threads? Under what circumstances is one type better than the other?.**

Here's a comparison of user-level threads and kernel-level threads in tabular form:

| Aspect | User-Level Threads | Kernel-Level Threads |
|---|---|---|
| **Creation and Management** | Managed entirely by user-level libraries. | Managed by the operating system kernel. |
| **Responsiveness** | Less responsive due to reliance on user-level thread library for scheduling and management. | More responsive as scheduling and management are handled directly by the kernel. |
| **Resource Usage** | Requires fewer system resources since all thread management is handled at the user level. | Requires more system resources as kernel-level threads involve additional overhead and kernel intervention. |

| | | |
|---|---|---|
| **Synchronization** | Synchronization mechanisms (e.g., mutexes, semaphores) are implemented at the user level, leading to potentially lower overhead but also to less efficient coordination between threads. | Synchronization mechanisms can take advantage of kernel support, leading to more efficient coordination between threads but potentially higher overhead due to kernel involvement. |
| **Fault Isolation** | A crash in one user-level thread does not necessarily affect other threads in the same process. | A crash in one kernel-level thread can potentially affect other threads and processes due to shared kernel resources. |
| **Portability** | More portable across different operating systems since they rely on user-level libraries that can be implemented consistently across platforms. | Less portable as kernel-level thread implementations may vary significantly between operating systems. |
| **Context Switching Overhead** | Generally lower context switching overhead compared to kernel-level threads since it occurs within user space. | Generally higher context switching overhead due to the involvement of the kernel in thread management. |
| **Concurrency** | Limited by the number of available CPU cores since scheduling is handled by user-level libraries. | Can take advantage of multi-core systems more effectively since the kernel can schedule threads across multiple CPU cores. |

**Under what circumstances is one type better than the other:**

1. **User-Level Threads:**

   o Better suited for applications requiring lightweight threading and fine-grained control over thread management, such as event-driven programming models or applications with a large number of threads.

   o Ideal for environments where portability across different operating systems is a priority or where resource efficiency is critical.

2. **Kernel-Level Threads:**

   o More suitable for applications requiring high responsiveness, scalability, and efficient utilization of multi-core systems, such as server applications, real-time systems, or compute-intensive tasks.

- o Preferred in environments where advanced synchronization mechanisms or fault isolation between threads and processes are necessary.
- o Useful when applications need to take advantage of kernel-level features or when interaction with hardware resources is required.

Overall, the choice between user-level threads and kernel-level threads depends on the specific requirements of the application, including factors such as responsiveness, scalability, resource efficiency, portability, and the need for advanced synchronization mechanisms or fault isolation.

## b) What are the benefits and the disadvantages of synchronous and asynchronous communication? Consider both the system level and the programmer level.

A.  **Synchronous and asynchronous communication**—A benefit of synchronous communication is that it allows a rendezvous between the sender and receiver. A disadvantage of a blocking send is that a rendezvous may not be required and the message could be delivered asynchronously. As a result, message-passing systems often provide both forms of synchronization.

B.  **Automatic and explicit buffering**—Automatic buffering provides a queue with indefinite length, thus ensuring the sender will never have to block while waiting to copy a message. There are no specifications on how automatic buffering will be provided; one scheme may reserve sufficiently large memory where much of the memory is wasted. Explicit buffering specifies how large the buffer is. In this situation, the sender may be blocked while waiting for available space in the queue. However, it is less likely that memory will be wasted with explicit buffering.

C.  **Send by copy and send by reference**—Send by copy does not allow the receiver to alter the state of the parameter; send by reference does allow it. A benefit of send by reference is that it allows the programmer to write a distributed version of a centralized application. Java's RMI provides both; however, passing a parameter by reference requires declaring the parameter as a remote object as well.

D.  **Fixed-sized and variable-sized messages**—The implications of this are mostly related to buffering issues; with fixed-size messages, a buffer with a specific size can hold a known number of messages. The number of variable-sized messages that can be held by such a buffer is unknown. Consider how Windows 2000 handles this situation: with fixed-sized messages (anything $< 256$ bytes), the messages are copied from the address space of the sender to the address space of the receiving process. Larger messages (i.e. variable-sized messages) use shared memory to pass the message.

## c) Explain the difference between preemptive and non-preemptive scheduling.

Here's a comparison of preemptive and non-preemptive scheduling in tabular form:

| Aspect | Preemptive Scheduling | Non-Preemptive Scheduling |
| --- | --- | --- |

| | | |
|---|---|---|
| **Interrupt Mechanism** | Operating system can interrupt a running process and allocate CPU to another process without cooperation from the running process. | Processes continue executing until they voluntarily release the CPU or complete execution. |
| **Control Over CPU** | Operating system has control over CPU allocation and can forcefully allocate CPU time to different processes. | Processes have control over CPU allocation and voluntarily relinquish CPU when they have completed their task or need to wait for an event. |
| **Responsiveness** | Generally offers better responsiveness as higher-priority tasks can be scheduled quickly, interrupting lower-priority tasks if necessary. | May have lower responsiveness as tasks can only be scheduled when the currently running task voluntarily yields the CPU. |
| **Fairness** | Ensures fairness by preventing long-running processes from monopolizing CPU time, thereby improving overall system performance. | May lead to unfairness if a long-running process prevents lower-priority processes from accessing CPU time, potentially impacting system performance. |
| **Resource Utilization** | Optimizes resource utilization by dynamically allocating CPU time to processes based on priority, improving system throughput. | May lead to suboptimal resource utilization if long-running processes monopolize CPU time, causing other processes to wait longer for execution. |
| **Complexity** | Generally involves more complex implementation due to the need for mechanisms to handle process preemption and context switching. | Generally simpler to implement as there is no need for mechanisms to forcibly interrupt running processes. |
| **Overhead** | May incur higher overhead due to frequent context switches and the need to manage process preemption. | May have lower overhead compared to preemptive scheduling as there are fewer forced context switches. |

In summary, preemptive scheduling offers better responsiveness, fairness, and resource utilization by allowing the operating system to interrupt processes as needed, while non-preemptive scheduling offers simplicity and potentially lower overhead at the cost of potentially longer response times and less efficient resource allocation. The choice between preemptive and

non-preemptive scheduling depends on the specific requirements and performance goals of the system.

**d) A multithreaded web server wishes to keep track of the number of requests it services (known as hits). Consider the two following strategies to prevent a race condition on the variable hits.**

**Using basic mutex lock**
**int hits;**
**mutex lock hit lock;**
**hit lock.acquire();**
**hits++;**
**hit lock.release();**

**Using an atomic integer**
**atomic t hits;**
**atomic inc (&hits);**

**Explain which of these two strategies is more efficient.**

In a multithreaded web server, multiple threads may need to update a shared variable, `hits`, which keeps track of the total number of requests the server has handled. Since multiple threads might try to modify `hits` at the same time, this situation could lead to a **race condition**, where the variable's value could become inconsistent due to simultaneous updates.

To prevent this race condition, two strategies can be used:

1. **Using a Mutex Lock**
2. **Using an Atomic Integer**

## Strategy 1: Using a Mutex Lock

In the first strategy, the server uses a mutex (mutual exclusion) lock to control access to `hits`. Here's how it works:

1. **Acquiring the Lock**: When a thread wants to update `hits`, it first attempts to acquire the lock (`hit_lock.acquire();`). If another thread already holds the lock, this thread has to wait until the lock is released.
2. **Updating the Variable**: Once the lock is acquired, the thread increments the `hits` variable (`hits++`).
3. **Releasing the Lock**: After updating, the thread releases the lock (`hit_lock.release();`), allowing other threads to acquire it.

**Drawbacks of Mutex Lock:**

- **Context Switching Overhead**: If a thread has to wait for the lock, the operating system may put it in a suspended state. When the lock becomes available, the operating system has to wake up the thread, which involves context switching. Context switching is costly because it requires saving and restoring the thread's state, which consumes CPU time.
- **Scalability Issues**: As the number of threads increases, the contention for the lock also increases. This can create a bottleneck, as multiple threads wait for the lock to access `hits`, thereby reducing efficiency.
- **Kernel Mode Operations**: Mutexes may involve kernel-level operations, which can slow down the process because the system has to switch between user mode (where the application runs) and kernel mode (where the lock management happens). This adds extra latency.

## Strategy 2: Using an Atomic Integer

In the second strategy, the server uses an atomic integer for `hits`. The `atomic_t` data type provides a way to perform operations on `hits` in a way that guarantees atomicity, meaning that the increment operation (`atomic_inc(&hits);`) will complete without interruption.

Here's how it works:

1. **Atomic Increment**: When a thread wants to update `hits`, it uses the `atomic_inc()` function, which increments `hits` in a single, indivisible operation. This means no other thread can access or modify `hits` while this operation is being performed.
2. **No Lock Needed**: Since `atomic_inc()` is atomic, there's no need for locks, so threads don't need to wait or be suspended.

**Advantages of Atomic Integer:**

- **Reduced Overhead**: Atomic operations are generally implemented at the hardware level, which means they don't require the thread to be suspended or involve context switching. This reduces the overhead significantly.
- **Scalability**: Since there's no lock, multiple threads can update `hits` without waiting. This makes the system more scalable and efficient, especially as the number of threads increases.
- **User Mode Operation**: Atomic operations work entirely in user mode, without needing to interact with the kernel. This eliminates the latency associated with user-kernel transitions, making atomic operations much faster.

## Conclusion: Which Strategy Is More Efficient?

Using an atomic integer (`atomic_t hits`) is more efficient in this case because:

1. **No Context Switching**: Since atomic operations don't require locks, they avoid the context-switching overhead that mutex locks introduce.
2. **Direct Access**: Atomic operations can be performed directly at the hardware level, allowing for faster execution without the need to involve the operating system.
3. **Reduced Complexity**: By eliminating the need for locks, atomic operations simplify the code and make it more efficient.

## b) Regarding the Producer-Consumer problem, explain the following term:

**i) Producer process ii) Consumer Process iii) Unbounded-buffer iv) Bounded-buffer**

Sure, let's break down each of these terms related to the Producer-Consumer problem:

i) **Producer process**: In the context of the Producer-Consumer problem, a producer process is a concurrent process or thread responsible for generating data or items and placing them into a shared buffer or queue. These items are produced at some rate and are intended for consumption by other processes.

ii) **Consumer process**: This is another concurrent process or thread in the Producer-Consumer problem. The consumer process is responsible for removing items from the shared buffer or queue and processing them. The consumer consumes the items produced by the producer process.

iii) **Unbounded-buffer**: An unbounded buffer, also known as an infinite buffer, is a shared data structure used in the Producer-Consumer problem where there is no restriction on the number of items it can hold. In other words, the buffer can grow dynamically as items are produced without any predefined limit.

iv) **Bounded-buffer**: A bounded buffer, on the other hand, has a fixed capacity, meaning it can only hold a limited number of items at a time. If the buffer becomes full, the producer process may need to wait until there is space available in the buffer before it can produce more items. Similarly, if the buffer becomes empty, the consumer process may need to wait until new items are produced and added to the buffer before it can consume more items.

**c) The following processes are being scheduled using a Priority Scheduling-Preemptive algorithm. Each process is assigned a numerical priority, with a higher number indicating a higher relative priority.**

**i) Show the scheduling order of the processes using a Gantt chart. ii) Fill the following table for exit time, turn-around time, and waiting time:**

same as six batch solve!!!!!!!!!

# Answer to the question number:5

**a) Suppose that a system is in an unsafe state. Show that it is possible for the processes to complete their execution without entering a deadlocked state.**

An unsafe state in a system does not necessarily mean that a deadlock will occur; it only indicates that the system cannot guarantee deadlock avoidance. This means there is a potential risk, but it is still possible for the processes to complete successfully without encountering deadlock.

To illustrate this, consider a system with a total of 12 resources distributed among three processes P0,P1,P_0, P_1,P0,P1, and P2P_2P2. The table below provides the maximum resources required by each process, the resources currently allocated to them, and their remaining need to complete execution:

| Process | Max | Current | Need |
|---------|-----|---------|------|
| P0P_0P0 | 10 | 5 | 5 |
| P1P_1P1 | 4 | 2 | 2 |
| P2P_2P2 | 9 | 3 | 6 |

The system currently has 2 resources available.

**Explanation:**

The system is considered to be in an unsafe state because there is no guaranteed sequence in which processes can execute without leading to deadlock. However, this does not mean that deadlock will inevitably occur. Let's explore the possible scenarios:

1. **Execution of P1P_1P1**: If process P1P_1P1 completes its execution, it will release the 2 resources it currently holds, adding them to the available resources. This would increase the total number of available resources to 4.

2. **Next Steps**:

   o With 4 available resources, process P0P_0P0 could then proceed to complete since its remaining need is 5. After completion, P0P_0P0 would release all its 5 allocated resources, resulting in a total of 9 available resources.

   o Following this, process P2P_2P2 could safely execute as its need of 6 resources can now be met with the available 9 resources. Upon completion, it would release all its resources.

b)

Consider the following snapshot of a system:

| Process | Allocation | | | | Max | | | | Available | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | A | B | C | D | A | B | C | D | A | B | C | D |
| P0 | 0 | 0 | 1 | 2 | 0 | 0 | 1 | 2 | 1 | 5 | 2 | 0 |
| P1 | 1 | 0 | 0 | 0 | 1 | 7 | 5 | 0 | | | | |
| P2 | 1 | 3 | 5 | 4 | 2 | 3 | 5 | 6 | | | | |
| P3 | 0 | 6 | 3 | 2 | 0 | 6 | 5 | 2 | | | | |
| P4 | 0 | 0 | 1 | 4 | 0 | 6 | 5 | 6 | | | | |

Answer the following questions using the banker's algorithm:

i) What is the content of the matrix *Need*?

ii) Define safe and unsafe state. Is the system in a safe state? If so, show a safe order in which the processes can execute.

iii) If a request from process *P1* arrives for (0,4,2,0), can the request be granted immediately?

Here's a more detailed explanation of the answer to each part of the question regarding the banker's algorithm:

---

## a. What is the content of the matrix Need?

The Need matrix represents the remaining resources each process requires to complete its execution. It is calculated as:

$$\text{Need} = \text{Max} - \text{Allocation}$$

Given the processes $P_0$ to $P_4$, their Need matrix is as follows:

- $P_0$: (0, 0, 0, 0)
- $P_1$: (0, 7, 5, 0)
- $P_2$: (1, 0, 0, 2)
- $P_3$: (0, 0, 2, 0)
- $P_4$: (0, 6, 4, 2)

This matrix indicates the resources still needed by each process to reach its maximum required allocation.

## b. Is the system in a safe state?

To determine if the system is in a safe state, we need to check if there is an order of execution where all processes can complete without leading to deadlock. The available resources at the beginning are:

$$Available = (1,5,2,0)$$

From this state:

- Processes $P_0$ or $P_3$ can run because their Need requirements (both (0, 0, 0, 0) and (0, 0, 2, 0), respectively) can be satisfied by the current available resources.
- Once process $P_3$ runs, it releases its allocated resources, making them available for other processes. This allows subsequent processes to run until all processes have completed.

Thus, there exists an execution sequence that allows all processes to complete, proving that the system is in a safe state.

## c. If a request from process $P_1$ arrives for (0, 4, 2, 0), can the request be granted immediately?

To determine if the request can be granted immediately, we need to check if the available resources after the request would still leave the system in a safe state. The current Available resources are:

$$Available = (1,5,2,0)$$

If process $P_1$ requests (0, 4, 2, 0), the new available resources would be:

$$Available\ after\ request = (1,1,0,0)$$

The request can be granted if the system remains in a safe state with these available resources. We can see that:

- Processes $P_0$, $P_2$, and $P_3$ can run with these available resources.
- Once any of these processes complete, they release their allocated resources, allowing $P_1$ and $P_4$ to run subsequently.

Therefore, the request from $P_1$ for (0, 4, 2, 0) can be granted immediately. One valid sequence of process completion is:
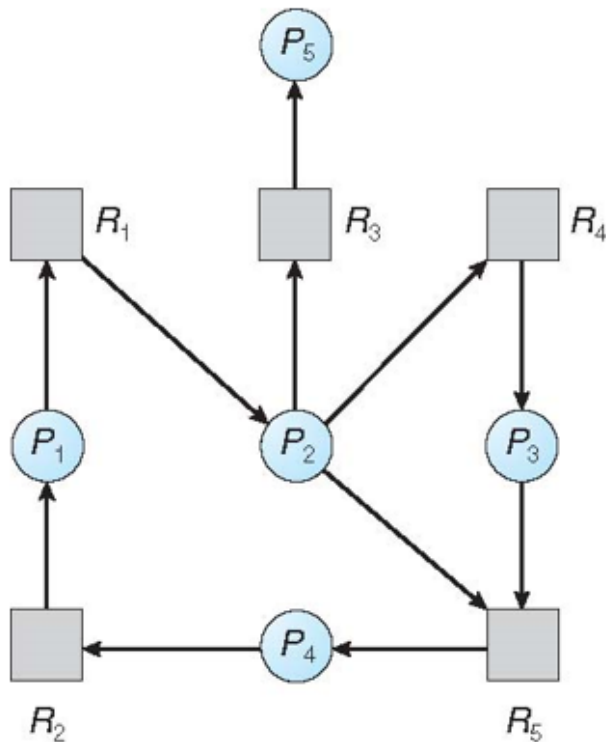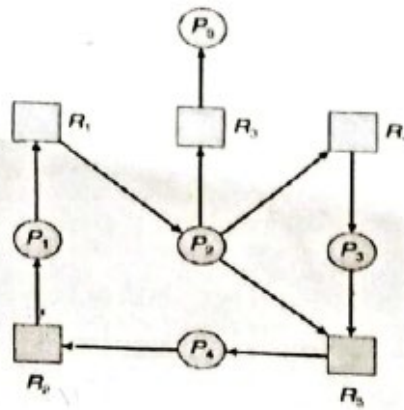
$$P_0, P_2, P_3, P_1, P_4$$

This confirms that after granting the request, the system can still avoid deadlock and complete all processes in a safe order.
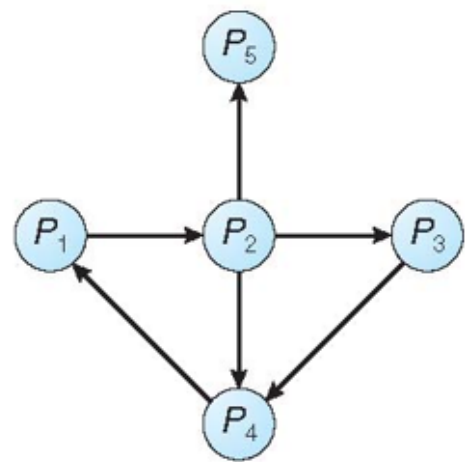
c)

List two deadlock detection mechanisms. Construct the wait-for graph that corresponds to the following resource allocation graph and determine whether or not there is deadlock.





(a)

(b)

Resource-Allocation Graph     Corresponding wait-for graph

Two common deadlock detection mechanisms are:

Resource Allocation Graph: This method involves creating a graph where nodes represent processes and resources. Edges represent requests and allocations. A cycle in this graph indicates a deadlock.

Wait-For Graph: This is a simplified version of the resource allocation graph where only the wait-for relationships between processes are represented. A cycle in this graph also indicates a deadlock.

# Answer to the question number:6

**a) In a system with a TLB, assume the following:**

**i) the memory access time is 150 nsec.**

**iii) the TLB hit rate is 80%.**

**ii) the TLB access time is 25 nsec.**

**Compute the effective access time for memory in this scenario.**

## Data:

TLB hit ratio = p = 0.8

TLB access time = 15 nanoseconds

Memory access time = m = 150 milliseconds

## Formula:

EMAT = p × (t + m) + (1 − p) × (t + m + m)

## Calculation:

EMAT = 0.8 × (15 + 150) + (1 − 0.8) × (15 + 150 + 150)

EMAT = 195 nanoseconds.

## Important points:

**b) Given memory partition of 100 KB, 500 KB, 200 KB and 600 KB (in order), how would each [ of the first fit, best fit and worst fit algorithms place processes of 212 KB,417 KB,112 KB and 426 KB (in order)? Rank the algorithms in terms of how efficiently they use memory.**

Given memory partitions of 100K, 500K, 200K, 300K, and 600K (in ord how would each of the First-fit, Best-fit, and Worst-fit algorithms place processes of 212K, 417K, 112K, and 426K (in order)? Which algorithm makes the most efficient use of memory?

### Solution:

#### First-Fit:
212K is put in 500K partition.
417K is put in 600K partition.
112K is put in 288K partition (new partition 288K = 500K - 212K).
426K must wait.

#### Best-Fit:
212K is put in 300K partition.
417K is put in 500K partition.
112K is put in 200K partition.
426K is put in 600K partition.

#### Worst-Fit:
212K is put in 600K partition.
417K is put in 500K partition.
112K is put in 388K partition.
426K must wait.

**In this example, Best-Fit turns out to be the best.**

**c) What is paging? Consider a logical address space of 64 pages of 1,024 words each, mapped onto a physical memory of 32 frames.**

**i) How many bits are there in the logical address?**

**ii) How many bits are there in the physical address?**

Paging is a memory management scheme used by computer operating systems to manage memory allocation more efficiently. In paging, the logical address space of a process is divided into fixed-size blocks called pages, and the physical memory is divided into blocks of the same size called frames. This allows the operating system to allocate memory in smaller, more manageable units.

Given the scenario:

Number of pages in logical address space = 64 pages
Size of each page = 1,024 words
Number of frames in physical memory = 32 frames

We can calculate the number of bits required for the logical and physical addresses as follows:

i) **Number of bits in the logical address**:

Since there are 64 pages, and each page contains 1,024 words, the logical address space can address $64 \times 1024$ words. To address this range, we need log2 of the total number of words, which gives us the number of bits required for the page offset.

$$\text{Page size} = 1024 \text{ words} = 2^{10} \text{ words}$$

$$\text{Number of pages} = 64 \text{ pages} = 2^6 \text{ pages}$$

So, we need 6 bits to address each page, which gives us:

$$\text{Bits for page offset} = \log_2(\text{Page size}) = 10$$

Since we have 64 pages, we need $\log_2(64) = 6$ bits for the page number.

Therefore, the total number of bits in the logical address is the sum of the bits for the page offset and the bits for the page number:

$$\text{Total bits in logical address} = \text{Bits for page number} + \text{Bits for page offset} = 6 + 10 = 16 \text{ bits}$$

ii) **Number of bits in the physical address**:

Since there are 32 frames, and each frame contains 1,024 words, the physical memory can address

$32 \times 1024$ words. Similar to the logical address, we calculate the number of bits required for the frame offset.

$$\text{Frame size} = 1024 \text{ words} = 2^{10} \text{ words}$$

$$\text{Number of frames} = 32 \text{ frames} = 2^5 \text{ frames}$$

So, we need 5 bits to address each frame, which gives us:

$$\text{Bits for frame offset} = \log_2(\text{Frame size}) = 10$$

Since we have 32 frames, we need $\log_2(32) = 5$ bits for the frame number.

Therefore, the total number of bits in the physical address is the sum of the bits for the frame offset and the bits for the frame number:

$$\text{Total bits in physical address} = \text{Bits for frame number} + \text{Bits for frame offset} = 5 + 10$$
$$= 15 \text{ bits}$$

**d) Explain the difference between internal and external fragmentation.**

Sure, here's a comparison between internal and external fragmentation in tabular form:

| Aspect | Internal Fragmentation | External Fragmentation |
|---|---|---|
| Definition | Unutilized space within allocated memory blocks. | Unutilized space between allocated memory blocks. |
| Cause | Occurs when allocated memory is larger than needed. | Occurs when total free memory is not contiguous. |
| Occurrence | Happens within a single memory block or partition. | Happens across multiple memory blocks or partitions. |
| Impact | Wastes memory resources due to unused allocated space. | Reduces overall available memory for allocation. |
| Mitigation | Adjusting allocation sizes or using memory compaction. | Using memory allocation algorithms to reduce gaps. |
| Example | A process allocated 100 KB but only needs 90 KB. | Free memory blocks scattered across memory space. |

In summary, internal fragmentation deals with wasted space within allocated memory blocks, while external fragmentation involves wasted space scattered across multiple memory blocks or partitions.

# Answer to the question number:7

**a) Define virtual memory. How does it empower computing?**

Virtual memory is a memory management technique used by modern computer operating systems to provide an abstraction of main memory (RAM). It extends the available memory beyond the physical RAM by using disk storage as an extension of physical memory.

In virtual memory, the operating system divides the physical memory into fixed-size blocks called pages and divides the logical memory (address space seen by the processes) into blocks of the same size called pages as well. Pages of memory are loaded into physical memory as needed, and when not in use, they can be swapped out to disk storage.

Virtual memory empowers computing in several ways:

1. **Enables Larger Programs**: Virtual memory allows running programs larger than the available physical memory by using disk space as an extension of RAM. This enables users to run more complex and memory-intensive applications without encountering "out of memory" errors.

2. **Provides Memory Protection**: Virtual memory provides memory protection by isolating processes from each other. Each process has its own virtual address space, preventing one process from accessing or corrupting the memory of another process.

3. **Facilitates Multitasking**: Virtual memory enables multitasking by allowing multiple processes to run simultaneously, even if the total memory requirements exceed the physical memory capacity. The operating system can swap pages in and out of memory as needed, allowing the CPU to switch between processes efficiently.

4. **Enhances Performance**: Virtual memory improves overall system performance by optimizing memory usage. Frequently accessed pages can be kept in physical memory, while less frequently accessed pages can be swapped out to disk. This minimizes the time spent on disk I/O operations and maximizes the availability of frequently accessed data in RAM.

5. **Simplifies Memory Management**: Virtual memory simplifies memory management for both the operating system and application developers. It provides a uniform and consistent view of memory to applications, regardless of the underlying hardware configuration. Application developers can write programs without needing to consider the specifics of physical memory allocation.

Overall, virtual memory plays a crucial role in modern computing by enabling efficient memory management, facilitating multitasking, and allowing for the execution of larger and more complex programs.

**b) What does TLB stands for? Briefly explain its usage in operating system.**

TLB stands for Translation Lookaside Buffer.

The Translation Lookaside Buffer (TLB) is a hardware cache used in modern computer processors to improve the efficiency of virtual memory operations. It works in conjunction with the memory management unit (MMU) of the CPU.

Here's a brief explanation of its usage in operating systems:

1. **Translation of Virtual Addresses**: When a program accesses memory, it uses virtual addresses. These addresses need to be translated to physical addresses so that the CPU can access the corresponding data in physical memory. The TLB caches the most recently used virtual-to-physical address translations, speeding up this translation process.

2. **Reducing Memory Access Time**: Without the TLB, the CPU would need to consult the page table stored in main memory for each memory access, resulting in significant overhead. By caching frequently used translations in the TLB, the CPU can avoid accessing the page table for every memory access, reducing memory access time and improving overall system performance.

3. **TLB Hit and TLB Miss**: When the CPU needs to translate a virtual address, it first checks the TLB. If the translation is found in the TLB (TLB hit), the corresponding physical address is retrieved directly from the TLB cache. If the translation is not found in the TLB (TLB miss), the CPU needs to consult the page table in main memory to retrieve the translation. In this case, the translation is also loaded into the TLB cache for future use.

4. **TLB Management**: The TLB has a limited capacity, so it needs to be managed efficiently. When the TLB becomes full, entries may need to be replaced using various replacement algorithms, such as least recently used (LRU). Additionally, TLB entries may need to be invalidated or updated when the page table is modified by the operating system.

Overall, the TLB plays a critical role in improving memory access performance in modern computer systems by caching virtual-to-physical address translations and reducing the overhead of accessing the page table for every memory access.

c)

What is fragmentation? Consider the following segment table:

| Segment | Base | Length |
|---------|------|--------|
| 0 | 219 | 600 |
| 1 | 2300 | 14 |
| 2 | 90 | 100 |
| 3 | 1327 | 58096 |
| 4 | 1952 | |

What are the physical addresses for the following logical addresses?

| i) 0, 430 | ii) 1, 10 |
|-----------|-----------|
| iii) 2, 500 | iv) 3, 400 |

Answer in the six batch solve sheet!!!!!!!!!!!!!!!!

# Answer to the question number:8

a) Consider the following page reference string:
2, 3, 2, 1, 5, 2, 4, 5, 3, 2, 5, 2
How many page faults would occur for the following replacement algorithms, assuming three frames? Remember that all frames are initially empty, so your first unique pages will cost one fault each.
i) Least Recently Used (LRU)
ii) Optimal Page Replacement
iii) FIFO

not solved yet!!!!!!!!!!!!!!!

b) Suppose that a disk drive has 200 cylinders, numbered 0 to 199. The drive is currently serving a request at cylinder 50. The queue of pending requests, in FIFO order, is:

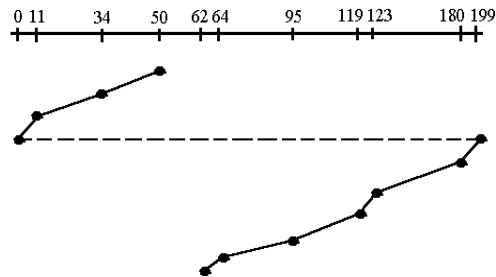95, 180, 34, 119, 11, 123, 62, 64

Starting from the current head position, what is the total distance (in cylinders) that the

**disk arm moves to satisfy all the pending requests? You should also provide a graphical view of head movement for each of the following disk-scheduling algorithms.**
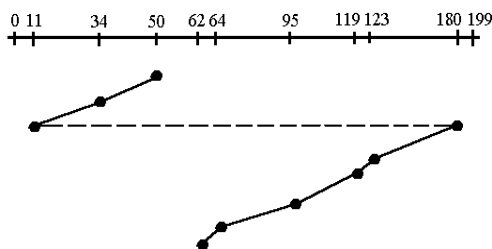
**i) C-SCAN**

**ii) LOOK**



*4. Circular Scan (C-SCAN)*
*Circular scanning works just like the elevator to some extent. It begins its scan toward the nearest end and works it way all the way to the end of the system. Once it hits the bottom or top it jumps to the other end and moves in the same direction. Keep in mind that the huge jump doesn't count as a head movement. The total head movement for this algorithm is only 187 track, but still this isn't the mose sufficient.*



*5. C-LOOK*
*This is just an enhanced version of C-SCAN. In this the scanning doesn't go past the last request in the direction that it is moving. It too jumps to the other end but not all the way to the end. Just to the furthest request. C-SCAN had a total movement of 187 but this scan (C-LOOK) reduced it down to 157 tracks.*