

**a) Operating System (OS):** An operating system (OS) is a software program that manages computer hardware and software resources and provides common services for computer programs. The OS acts as an intermediary between the computer hardware and the applications running on it. It controls and coordinates the use of the hardware among the various application programs for different users. The primary objectives of an operating system are to provide a user interface, manage computer resources, and ensure system security and stability.

**b) Multiprocessing System:** A multiprocessing system is a computer system that utilizes multiple processors (also known as CPUs or cores) to execute multiple tasks or processes simultaneously. In a multiprocessing system, tasks are divided among the available processors, allowing for parallel processing and potentially improving overall system performance. This approach contrasts with single-processor systems, where only one processor is available to execute tasks sequentially. Multiprocessing systems can range from simple setups with two or more processors sharing a common memory to complex distributed systems with multiple interconnected computers working together.

**c) Process:** In the context of an operating system, a process can be defined as a program in execution. It represents a single unit of work or a task that is being executed on a computer system. A process typically consists of the program code, data, and resources (such as memory and CPU time) allocated to it during execution. Processes are managed by the operating system's kernel, which schedules them for execution, allocates resources, and provides mechanisms for communication and synchronization between processes.

**d) Degree of Multiprogramming Scheduler:** The scheduler responsible for controlling the degree of multiprogramming in a system is often referred to as the "admission scheduler" or the "long-term scheduler." This scheduler determines how many processes can be simultaneously active in the system at any given time. It controls the admission of new processes into the system, based on factors such as available memory, CPU utilization, and system load. By regulating the number of processes allowed to run concurrently, the degree of multiprogramming scheduler helps maintain system stability, prevent resource exhaustion, and optimize overall system performance.

**e) Burst Time:** Burst time refers to the amount of time taken by a process to execute on the CPU without any interruptions. It is the duration between when a process starts executing on the CPU and when it completes its execution or is preempted. Burst time can vary for different processes and is a critical factor in CPU scheduling algorithms for determining the order in which processes are executed.

**f) Semaphores:** Semaphores are synchronization primitives used in concurrent programming to control access to shared resources. They are integer variables that can be accessed and modified only through atomic operations. Semaphores can have two main operations: "wait" (P) and "signal" (V).

- The "wait" operation decrements the semaphore value. If the resulting value is negative, the process executing the "wait" operation is blocked until the semaphore value becomes non-negative.
- The "signal" operation increments the semaphore value. If there are processes blocked on the semaphore, one of them is unblocked.

Semaphores are used to prevent race conditions and ensure that critical sections of code are executed mutually exclusively, thus avoiding issues like data corruption or inconsistency in multi-threaded or multi-process environments. They are a fundamental synchronization mechanism in operating systems and concurrent programming.

**g) Rollback:** Rollback refers to the process of reverting a system or a transaction to a previous state. In the context of databases, rollback typically refers to undoing changes made by a transaction that has encountered an error or has been aborted. When a transaction fails or is terminated, the system can rollback the changes made by that transaction to ensure data consistency and integrity. Rollback involves reversing the effects of the incomplete transaction, restoring the data to its state before the transaction began.

**h) Address Binding:** Address binding is the process of associating symbolic addresses used in a program with physical addresses in memory or other resources at runtime. It involves translating symbolic addresses, which are used by programmers to refer to variables, functions, or other memory locations in their code, into actual memory addresses that the hardware can understand and access. There are different types of address binding, including compile-time binding, load-time binding, and runtime binding, each of which determines when the association between symbolic and physical addresses occurs during the program's lifecycle. Address binding is crucial for enabling programs to access and manipulate data and resources effectively during execution.

**i) Operations on Files:** Files in an operating system can undergo various operations, including:

1. **Creation:** The process of creating a new file.
2. **Opening:** Accessing an existing file for reading, writing, or both.
3. **Reading:** Retrieving data from a file.
4. **Writing:** Storing data into a file.
5. **Closing:** Ending access to a file after reading or writing.
6. **Deletion:** Removing a file from the file system.
7. **Seeking:** Moving the file pointer to a specific position within the file.

8. **Renaming:** Changing the name of a file.
9. **Copying:** Creating a duplicate of a file.
10. **Appending:** Adding data to the end of an existing file.

These operations allow users and programs to manipulate files stored on disk or other storage devices effectively.

**j) Seek Time in Disk Scheduling:** Seek time in disk scheduling refers to the time taken by the disk arm to move from its current position to the desired position on the disk to read or write data. Disk drives consist of multiple platters, each with tracks divided into sectors where data is stored. When the operating system requests data from the disk, the disk arm must seek to the appropriate track before the data can be read or written. Seek time is a crucial factor in disk performance because it directly affects the overall time required to access data. Disk scheduling algorithms aim to minimize seek time by optimizing the order in which disk requests are serviced, reducing unnecessary movement of the disk arm, and improving overall disk efficiency.

## a) List and explain advantages of Multiprocessor system.

### a) Advantages of Multiprocessor Systems:

1. **Increased Processing Power:** Multiprocessor systems can execute multiple tasks concurrently, leading to a significant increase in processing power compared to single-processor systems. This allows for faster execution of tasks and improved system performance.
2. **Improved Throughput:** By distributing tasks among multiple processors, a multiprocessor system can handle a higher workload and achieve better throughput. This is particularly beneficial in environments with high-demand applications or heavy multitasking scenarios.
3. **Enhanced Reliability:** Multiprocessor systems often incorporate redundancy and fault-tolerance mechanisms to ensure system reliability. If one processor fails, the system can continue to operate using the remaining processors, reducing the risk of system downtime.
4. **Resource Sharing:** Multiprocessor systems enable efficient resource sharing among processes. Multiple processes can access shared resources simultaneously, leading to better resource utilization and reduced contention for resources.
5. **Scalability:** Multiprocessor systems are highly scalable, allowing additional processors to be added to the system as workload demands increase. This scalability makes multiprocessor systems suitable for accommodating growing computational requirements over time.
6. **Parallel Processing:** Multiprocessor systems support parallel processing, where tasks are divided into smaller sub-tasks and executed concurrently on different processors. This parallelism can lead to significant performance gains for parallelizable workloads.
7. **Flexibility:** Multiprocessor systems offer flexibility in resource allocation and task scheduling. Processes can be dynamically assigned to different processors based on system load and resource availability, optimizing system performance and responsiveness.

8.

## **b) Process Control Block (PCB):**

The Process Control Block (PCB) is a data structure used by operating systems to manage information about each running process. It contains various pieces of information required for process management, scheduling, and resource allocation. Here's a detailed explanation of the components typically found in a PCB:

1. **Process ID (PID):** A unique identifier assigned to each process in the system.
2. **Process State:** Indicates the current state of the process (e.g., running, ready, blocked).
3. **Program Counter (PC):** Stores the address of the next instruction to be executed by the process.
4. **CPU Registers:** Includes registers such as accumulator, index registers, and general-purpose registers, which store the process's current execution context.
5. **Memory Management Information:** Contains information about the process's memory allocation, including base and limit registers for memory protection.
6. **Priority:** Indicates the priority level of the process for scheduling purposes.
7. **Pointers to Other Data Structures:** Points to other data structures related to the process, such as the process's parent, child, or sibling processes.
8. **I/O Status Information:** Stores information about I/O devices requested by the process, including device identifiers, status, and buffer pointers.
9. **Accounting Information:** Tracks resource usage statistics such as CPU time consumed, memory usage, and I/O operations performed by the process.
10. **Scheduling Information:** Contains data used by the scheduler to determine the process's scheduling priority and execution order.

Diagram:

## **c) Explain different method for recovery from a deadlock.**

### **c) Methods for Recovery from Deadlock:**

#### **1. Deadlock Detection and Recovery:**

- o Deadlock detection algorithms periodically check the system's state to identify deadlock situations. Once a deadlock is detected, the operating system can initiate recovery mechanisms.
- o Recovery methods include killing processes involved in the deadlock, preempting resources from processes, or rolling back transactions to release resources.

#### **2. Resource Preemption:**

- o In situations where deadlock avoidance is not possible, the operating system can preempt resources from one or more processes to break the deadlock.
- o Preemption involves forcibly removing resources from processes and allocating them to other processes to resolve the deadlock.
- o Preemption strategies should consider factors such as process priority, resource utilization, and fairness.

#### **3. Process Termination:**

- o Another approach to deadlock recovery is to terminate one or more processes involved in the deadlock.
- o The operating system can select processes for termination based on criteria such as process priority, resource usage, and deadlock duration.
- o Termination should be performed carefully to minimize disruption and ensure system stability.

#### **4. Resource Rollback:**

- o In systems where transactions are involved, resource rollback can be used to release resources held by processes in a deadlock.
- o Rollback involves undoing the effects of transactions to return resources to their original states.
- o Rollback mechanisms should ensure consistency and integrity of data while resolving the deadlock.

#### **5. Combined Approaches:**

- o Some systems employ a combination of deadlock detection, prevention, and recovery techniques to effectively manage deadlocks.
- o By utilizing multiple approaches, systems can achieve better deadlock handling and minimize the impact of deadlocks on system performance and availability.

## **d) What is Fragmentation? Explain types of fragmentation in details.**

### **d) Fragmentation:**

Fragmentation refers to the phenomenon where available memory or disk space becomes divided into small, non-contiguous blocks, making it challenging to allocate large contiguous blocks of memory or disk space efficiently. Fragmentation can occur in both main memory (RAM) and secondary storage (disk).

Types of Fragmentation:

#### **1. Internal Fragmentation:**

- o Internal fragmentation occurs in memory allocation when a portion of allocated memory remains unused due to the allocation of memory in fixed-size blocks.
- o It typically occurs in systems that allocate memory in fixed-size units, leading to wastage of memory space.
- o Internal fragmentation is inherent in memory allocation techniques such as fixed partitioning and dynamic partitioning.

#### **2. External Fragmentation:**

- o External fragmentation occurs when free memory blocks are scattered throughout the memory space, making it difficult to allocate contiguous blocks of memory to processes.
- o It occurs over time as processes are allocated and deallocated memory, leaving small gaps between allocated blocks that are too small to be used effectively.
- o External fragmentation is a significant concern in dynamic memory allocation systems such as those employing techniques like best fit, worst fit, or next fit.

#### **3. Disk Fragmentation:**

- o Disk fragmentation refers to the scattering of files and data across non-contiguous sectors on a disk.
- o As files are created, modified, and deleted, their data may become fragmented, leading to inefficient disk access and reduced performance.
- o Disk defragmentation utilities are used to rearrange fragmented data on disks, consolidating fragmented files and improving disk performance.

Overall, fragmentation can degrade system performance and efficiency, making it essential for operating systems to employ strategies to minimize fragmentation and optimize resource utilization.



## **a) List and explain system calls related to Process and Job control.**

### **a) System Calls Related to Process and Job Control:**

1. **fork():**
  - o The fork() system call creates a new process by duplicating the calling process. The new process, called the child process, is an exact copy of the calling process, called the parent process.
  - o After forking, the parent and child processes have separate memory spaces but share the same code, data, and open file descriptors. They continue execution from the point of the fork() call, with different process IDs (PIDs).
2. **exec():**
  - o The exec() family of system calls replaces the current process's memory image with a new program. It loads and executes a new program file into the current process's address space, replacing the current program's code and data.
  - o The exec() system calls come in different variants, such as execve(), execl(), execv(), etc., allowing different ways to specify the new program and its arguments.
3. **wait() and waitpid():**
  - o The wait() system call suspends the execution of the calling process until one of its child processes terminates. It allows the parent process to wait for the completion of child processes.
  - o The waitpid() system call is similar to wait() but provides more options and flexibility for waiting on specific child processes based on their process IDs or exit statuses.
4. **exit():**
  - o The exit() system call terminates the calling process and releases all its resources. It allows a process to gracefully exit and return an exit status to its parent process, indicating the reason for termination.
5. **kill():**
  - o The kill() system call sends a signal to a specified process or a group of processes. It is commonly used to terminate or manipulate the execution of processes, such as terminating a process forcefully or requesting it to perform specific actions.

## **b) Explain multilevel Feedback queue Algorithm.**

### **Multilevel Feedback Queue Algorithm:**

The multilevel feedback queue (MLFQ) algorithm is a CPU scheduling algorithm that manages processes by assigning them to multiple queues with different priorities. Each queue has its own scheduling algorithm, and processes are promoted or demoted between queues based on their behavior and resource requirements. Here's an explanation of how the MLFQ algorithm works:

1. **Multiple Queues:**
  - o MLFQ maintains multiple queues, each with a different priority level. Typically, there are several queues, with higher-priority queues having shorter time quantum or higher scheduling priority.



**2. Scheduling Algorithm:**

- o Each queue uses its own scheduling algorithm to determine which process to execute next. Common scheduling algorithms include round-robin, shortest job first, or priority scheduling.

**3. Promotion and Demotion:**

- o Processes are initially placed in the highest-priority queue when they enter the system. As they execute, they may be promoted or demoted between queues based on their behavior.
- o Processes that use their entire time quantum without blocking or voluntarily yielding CPU time are demoted to lower-priority queues to give other processes a chance to execute.
- o Processes that block frequently or have I/O-bound behavior are promoted to higher-priority queues to ensure responsiveness.

**4. Aging:**

- o Aging is a mechanism used in MLFQ to prevent starvation of low-priority processes. As processes wait in lower-priority queues for an extended period, their priority gradually increases (ages), allowing them to eventually reach higher-priority queues and receive CPU time.

**5. Priority Boosting:**

- o Periodically, the system may perform a priority boost, where all processes are moved to the highest-priority queue. This helps prevent aging processes from being indefinitely starved in lower-priority queues.

By dynamically adjusting process priorities based on their behavior and resource requirements, the MLFQ algorithm aims to achieve good throughput, fairness, and responsiveness in CPU scheduling, making it suitable for systems with diverse workload characteristics.

## c) Describe in detail the 'Dinning Philosopher Problem' Synchronization problem.

### c) Dining Philosophers Problem:

The Dining Philosophers Problem is a classic synchronization problem in computer science that illustrates challenges in resource allocation and deadlock avoidance. It involves a group of philosophers sitting around a dining table, with each philosopher alternating between thinking and eating. There are only enough chopsticks (resources) for each philosopher to eat, and they must pick up two adjacent chopsticks to eat. The problem arises when multiple philosophers attempt to pick up chopsticks simultaneously, leading to potential deadlock.

Here's a detailed description of the problem and its solution:

#### 1. Problem Statement:

- o There are N philosophers sitting around a circular dining table, each with a plate of spaghetti in front of them.
- o Between each pair of adjacent philosophers, there is a single chopstick.
- o A philosopher requires two chopsticks to eat. They can only pick up adjacent chopsticks and must release them when finished eating.
- o The philosophers alternate between thinking and eating, but they may need to wait if the necessary chopsticks are not available.

#### 2. Challenge:

- o The challenge in the Dining Philosophers Problem is to design a synchronization solution that prevents deadlocks, where each philosopher is indefinitely blocked waiting for a chopstick held by another philosopher.

#### 3. Solution:

- o Several solutions exist to prevent deadlocks in the Dining Philosophers Problem, including resource hierarchy, resource ordering, and arbitration through a central authority (like a waiter).
- o One common solution is the "Chandy/Misra solution," which involves assigning a unique identifier (number) to each chopstick and having philosophers request chopsticks in a predefined order.
- o The algorithm ensures that philosophers never hold a chopstick if they cannot obtain its neighboring chopstick, thereby preventing deadlock.

#### 4. Pseudocode:

scssCopy code

```
do {  
    think(); // Philosopher thinks  
    pickup(left_chopstick);  
    pickup(right_chopstick); // Philosopher picks up two adjacent chopsticks  
    eat(); // Philosopher eats  
    putdown(left_chopstick);  
    putdown(right_chopstick); // Philosopher puts down both chopsticks  
} while (true);
```

#### 5. Deadlock Avoidance:

- o Deadlock is avoided by ensuring that philosophers always pick up chopsticks in a predefined order, such as picking up the lower-numbered chopstick first.

- o Additionally, philosophers can be programmed to release chopsticks if they cannot acquire both, allowing other philosophers to eat.

## **d) Write a note on interrupts**

### **d) Interrupts:**

Interrupts are signals sent by hardware or software to the CPU to notify it of an event that requires immediate attention. Interrupts temporarily suspend the execution of the current program and transfer control to an interrupt handler (also known as an interrupt service routine or ISR), which is responsible for handling the event. Here are key points about interrupts:

#### **1. Types of Interrupts:**

- o **Hardware Interrupts:** Generated by external hardware devices, such as timers, keyboards, or disk drives, to signal events such as data arrival or completion of a task.
- o **Software Interrupts:** Generated by software instructions, such as system calls or software exceptions, to request services from the operating system or handle exceptional conditions.

#### **2. Interrupt Handling Process:**

- o When an interrupt occurs, the CPU saves the current state of the program, including the program counter and processor registers.
- o It then jumps to the address of the interrupt handler associated with the interrupting event.
- o The interrupt handler executes to service the interrupt, performing tasks such as handling the event, updating system state, and resuming the interrupted program.
- o After handling the interrupt, the CPU restores the saved state and resumes execution of the interrupted program.

#### **3. Interrupt Priorities:**

- o Interrupts can have different priorities, allowing the system to handle critical events promptly. Higher-priority interrupts are serviced before lower-priority ones.
- o Interrupt priority levels can be managed by hardware or software mechanisms, such as interrupt controllers or interrupt handling routines.

#### **4. Benefits of Interrupts:**

- o **Improves System Responsiveness:** Interrupts allow the CPU to respond quickly to external events, ensuring timely processing of critical tasks.
- o **Supports Asynchronous I/O:** Interrupt-driven I/O enables efficient handling of input/output operations without requiring the CPU to continuously poll devices.
- o **Facilitates Multi-tasking:** Interrupts enable multitasking by allowing the CPU to switch between multiple tasks in response to various events, enhancing system efficiency.

Overall, interrupts play a crucial role in modern computer systems by enabling efficient handling of asynchronous events and supporting real-time processing requirements.

# What is meant by Free Space Management? Define Bit vector and Grouping

## **Free Space Management:**

Free space management refers to the process of tracking and managing available space on storage devices, such as disks or memory, to efficiently allocate storage resources to new data or files. It involves maintaining information about which storage blocks or sectors are currently unused and available for allocation. Effective free space management is essential for optimizing storage utilization and performance in computer systems.

## **Bit Vector:**

A bit vector, also known as a bit array or bitset, is a data structure that represents a fixed-size sequence of binary values, typically stored compactly in memory. In a bit vector, each bit corresponds to a single value, either 0 or 1, indicating the presence or absence of a particular attribute or condition.

In the context of free space management, a bit vector is often used to represent the allocation status of storage blocks or sectors within a storage device. Each bit in the bit vector corresponds to a storage block, with a value of 1 indicating that the block is allocated (in use) and a value of 0 indicating that the block is free (available for allocation). By storing this information compactly as a sequence of bits, bit vectors enable efficient querying and manipulation of storage allocation status.

## **Grouping:**

Grouping, in the context of free space management, refers to a technique used to organize and manage storage blocks or sectors into groups or clusters. Instead of managing individual blocks separately, grouping combines multiple contiguous blocks into larger units, simplifying free space management and reducing overhead.

In a grouped allocation scheme, the storage device is divided into fixed-size groups or clusters, each containing a predefined number of contiguous blocks. The allocation status of entire clusters is tracked using a bit vector or similar data structure, rather than tracking individual blocks. This reduces the number of entries in the free space management data structure, leading to improved storage utilization and efficiency.

Grouping helps minimize fragmentation and overhead associated with managing individual blocks, especially in storage systems with large capacities. However, it may also lead to internal fragmentation if allocated space within a group is not fully utilized. Various grouping strategies, such as block grouping or cylinder grouping, can be employed based on the characteristics of the storage device and the requirements of the system.

**i) Logical Address:** A logical address, also known as a virtual address, is a reference to a location in memory that is generated by a program or process. It represents the memory location used by the program during its execution, independent of the underlying hardware configuration. Logical addresses are typically generated by the CPU and passed to the memory management unit (MMU) for translation into physical addresses. In virtual memory systems, logical addresses are mapped to physical addresses through address translation mechanisms, allowing processes to access a larger address space than physically available memory.

**ii) Physical Address:** A physical address is the actual location in physical memory (RAM) where data is stored. It represents a unique location within the memory hardware, typically identified by a memory cell's specific row and column address. Unlike logical addresses, which are generated by the CPU and represent an abstract memory space, physical addresses directly correspond to specific memory locations in the hardware. The memory management unit (MMU) translates logical addresses into physical addresses, enabling the CPU to access the corresponding data stored in physical memory. Physical addresses are essential for memory management and data retrieval operations in computer systems.

### 3) Explain the Resource Allocation Graph in detail

The Resource Allocation Graph (RAG) is a graphical representation used in operating systems to illustrate the allocation of resources and the relationships between processes and resources in a system. It is particularly useful for analyzing and managing resource allocation in concurrent systems, such as those with multiple processes competing for shared resources. The Resource Allocation Graph consists of nodes and edges representing processes and resources, respectively, along with various types of arcs indicating different relationships. Here's a detailed explanation of the components and features of the Resource Allocation Graph:

1. **Nodes:**

- o Nodes in the Resource Allocation Graph represent two types of entities: processes and resources.
- o Process nodes depict individual processes or tasks in the system that require access to resources to execute.
- o Resource nodes represent shared resources such as printers, disk drives, or communication channels that processes may request and use.

2. **Edges:**

- o Edges in the Resource Allocation Graph represent the allocation or request of resources by processes.
- o An edge from a process node to a resource node indicates that the process currently holds or has been allocated the resource.
- o Conversely, an edge from a resource node to a process node indicates that the process is requesting the resource but has not yet been allocated it.

3. **Types of Arcs:**

- o The Resource Allocation Graph includes three types of arcs to represent different types of relationships:
  - **Allocation Arcs:** These arcs indicate that a process has been allocated a resource. They connect a process node to a resource node.
  - **Request Arcs:** These arcs indicate that a process is currently requesting a resource but has not yet been allocated it. They connect a process node to a resource node.
  - **Assignment Arcs:** These arcs represent the assignment of resources to processes in a way that avoids deadlocks. They connect a resource node to a process node that currently holds the resource.

4. **Deadlocks:**

- o One of the primary uses of the Resource Allocation Graph is to identify and prevent deadlocks, which occur when processes are unable to proceed because they are waiting for resources held by other processes, creating a circular waiting pattern.
- o Deadlocks can be detected in the Resource Allocation Graph by identifying cycles involving request and assignment arcs.
- o To prevent deadlocks, resource allocation strategies can be designed to ensure that processes never reach a state where they are waiting indefinitely for resources.

#### 5. **Cycle Detection and Resolution:**

- o A cycle in the Resource Allocation Graph indicates a potential deadlock situation, as it implies a circular waiting pattern among processes.
- o When a cycle is detected, deadlock resolution strategies can be applied to break the cycle and allow processes to proceed.
- o Common deadlock resolution techniques include resource preemption, where resources are forcibly removed from processes to break the deadlock, or process rollback, where processes are rolled back to a previous state to release resources.

Overall, the Resource Allocation Graph provides a visual representation of resource allocation and dependencies in a system, facilitating analysis, deadlock detection, and resource management in operating systems and concurrent systems.

## **d) What are the difference between Preemptive and Non-Preemptive Scheduling.**

### **Preemptive Scheduling:**

#### **1. Definition:**

- o In preemptive scheduling, the operating system has the authority to interrupt a currently executing process and move it out of the CPU, even if the process has not yet completed its execution.
- o The CPU can be allocated to another process with higher priority or a more urgent task, even if the currently executing process has not voluntarily relinquished the CPU.

#### **2. Priority Management:**

- o Preemptive scheduling allows the operating system to manage process priorities dynamically. Higher-priority processes can preempt lower-priority ones, ensuring that critical tasks are executed promptly.
- o The CPU scheduler constantly evaluates the priority of active processes and may preempt lower-priority processes to give CPU time to higher-priority ones.

#### **3. Response Time:**

- o Preemptive scheduling typically leads to shorter response times for high-priority tasks since they can preempt lower-priority tasks, even if those tasks are in the middle of execution.
- o This ensures that time-sensitive or interactive tasks receive prompt attention and do not experience delays due to lower-priority tasks.

### **Non-Preemptive Scheduling:**

#### **1. Definition:**

- o In non-preemptive scheduling, a process continues to execute on the CPU until it voluntarily relinquishes control by either completing its execution or entering a waiting state (e.g., waiting for I/O).
- o The operating system cannot interrupt the currently executing process and must wait for it to finish before allocating the CPU to another process.

#### **2. Priority Management:**

- o Non-preemptive scheduling does not involve dynamic priority management. Once a process starts executing, it retains control of the CPU until it completes or blocks.
- o Processes are typically scheduled based on their arrival time or assigned static priorities that do not change during execution.

#### **3. Response Time:**

- o Non-preemptive scheduling may lead to longer response times for high-priority tasks since they must wait for lower-priority tasks to complete or yield the CPU voluntarily.
- o Time-sensitive or interactive tasks may experience delays if they are queued behind long-running processes that monopolize the CPU.

### **Key Differences:**

#### **1. Interrupt Handling:**

- o Preemptive scheduling allows the operating system to interrupt and switch between processes at any time, whereas non-preemptive scheduling relies on processes voluntarily giving up the CPU.

#### **2. Priority Management:**

- o Preemptive scheduling dynamically manages process priorities and allows higher-priority tasks to preempt lower-priority ones, while non-preemptive scheduling typically uses static priorities or arrival time for scheduling.

#### **3. Response Time:**

- o Preemptive scheduling generally offers shorter response times for high-priority tasks since they can preempt lower-priority tasks, whereas non-preemptive scheduling may result in longer response times for high-priority tasks due to waiting for lower-priority tasks to complete.



## a) Solution for Critical Section Problem:

The critical section problem arises in concurrent computing when multiple processes or threads share a common resource or section of code, and there is a possibility of conflicting accesses leading to race conditions and data inconsistency. Various synchronization mechanisms have been developed to address the critical section problem and ensure mutual exclusion, progress, and bounded waiting. Here are some common solutions:

### 1. **Locks and Mutexes:**

- o Using locks or mutexes to enforce mutual exclusion is a widely used solution for the critical section problem.
- o Processes or threads acquire a lock before entering the critical section, ensuring that only one process can execute the critical section at a time.
- o Other processes attempting to enter the critical section while the lock is held are blocked or put in a wait state until the lock is released.

### 2. **Semaphores:**

- o Semaphores are synchronization primitives that can be used to control access to critical sections.
- o They provide two atomic operations: "wait" (P) and "signal" (V).
- o By properly initializing and managing semaphores, processes can coordinate access to shared resources and ensure mutual exclusion.

### 3. **Atomic Operations:**

- o Some platforms provide atomic operations such as compare-and-swap (CAS) or test-and-set, which can be used to implement lock-free synchronization algorithms.
- o Atomic operations ensure that certain operations on shared variables are performed atomically, preventing race conditions without the need for explicit locks.

### 4. **Mutex Algorithms:**

- o There are various mutex algorithms, such as Peterson's algorithm, Dekker's algorithm, and the Bakery algorithm, designed to achieve mutual exclusion in a distributed or shared memory environment.
- o These algorithms typically rely on atomic operations or other synchronization primitives to coordinate access to critical sections.

### 5. **Transactional Memory:**

- o Transactional memory is a higher-level synchronization mechanism that allows multiple operations to be grouped into atomic transactions.
- o Processes can execute transactions concurrently, and conflicts are automatically detected and resolved to ensure consistency.

These solutions aim to provide safe and efficient access to shared resources in concurrent systems while avoiding race conditions and data corruption.

## **b) Medium-term Scheduler:**

The medium-term scheduler, also known as the swap scheduler or memory scheduler, is a component of the operating system's memory management subsystem. Unlike the short-term and long-term schedulers, which primarily focus on CPU scheduling and process management, respectively, the medium-term scheduler is responsible for managing memory resources and alleviating memory pressure in the system.

Key characteristics and functions of the medium-term scheduler include:

### **1. Memory Management:**

- o The medium-term scheduler manages memory by swapping out inactive or idle processes or parts of processes from main memory (RAM) to secondary storage (usually disk).
- o Swapping frees up memory space for other processes that are ready to execute, thereby improving overall system performance and responsiveness.

### **2. Memory Pressure Management:**

- o The medium-term scheduler monitors memory usage and intervenes when memory pressure becomes high.
- o When available memory becomes scarce due to a large number of active processes or high memory allocation demands, the medium-term scheduler selects processes or pages for swapping to secondary storage to alleviate memory pressure.

### **3. Swapping Algorithms:**

- o The medium-term scheduler employs various swapping algorithms to select processes or memory pages for swapping based on criteria such as least recently used (LRU), page fault frequency, or process priority.
- o These algorithms aim to minimize the impact on system performance while efficiently managing memory resources.

### **4. Page Fault Handling:**

- o The medium-term scheduler is involved in handling page faults generated when a process attempts to access memory that is not currently in main memory.
- o It may initiate swapping operations to bring required memory pages into main memory or decide to terminate or suspend processes to free up memory resources.

Overall, the medium-term scheduler plays a crucial role in maintaining memory balance and managing memory resources effectively in operating systems, ensuring optimal performance and efficient utilization of memory resources.

# Explain Indexed Allocation briefly.

Indexed allocation is a file allocation method used in operating systems to manage and allocate disk space for files efficiently. In indexed allocation, each file has its own index block, which serves as an index or table containing pointers to the actual data blocks where the file's contents are stored on the disk. Here's a brief overview of indexed allocation:

## 1. **Index Block:**

- o Each file in indexed allocation is associated with an index block, also known as an index node or i-node.
- o The index block contains an array or list of pointers, with each pointer corresponding to a disk block where a portion of the file's data is stored.
- o The size of the index block determines the maximum number of pointers it can hold, thus limiting the maximum file size supported by indexed allocation.

## 2. **Direct Indexing:**

- o In some implementations of indexed allocation, the index block contains pointers to data blocks directly, known as direct indexing.
- o Each pointer in the index block points to a specific data block on the disk, allowing for efficient direct access to file data.

## 3. **Indirect Indexing:**

- o To support larger files, indexed allocation may employ indirect indexing, where some pointers in the index block point to additional index blocks rather than directly to data blocks.
- o These additional index blocks, called single indirect blocks, double indirect blocks, or triple indirect blocks, contain pointers to data blocks, enabling the storage of more extensive files with a larger number of data blocks.

## 4. **Advantages:**

- o Indexed allocation offers fast access to file data since the index block provides direct pointers to data blocks.
- o It supports random access to file data, allowing efficient retrieval of data from any part of the file without the need for sequential scanning.
- o Indexed allocation supports dynamic file growth, as additional index blocks can be allocated when needed to accommodate larger files.

## 5. **Disadvantages:**

- o Indexed allocation requires additional space overhead for maintaining index blocks, which can increase the storage overhead for small files.
- o It may suffer from fragmentation issues, especially if file sizes fluctuate frequently, leading to scattered allocation of data blocks.
- o The maximum file size supported by indexed allocation is limited by the size of the index block and the number of pointers it can hold.