

Operating System

Theory Content

INDEX

Index Structure for Operating System Content

1. Introduction to Operating Systems

- Definition and Purpose
- History and Evolution
- Types of Operating Systems
 - Batch Operating Systems
 - Time-Sharing Operating Systems
 - Distributed Operating Systems
 - Real-Time Operating Systems
 - Embedded Operating Systems

2. Operating System Architecture

- Kernel
 - Monolithic Kernels
 - Microkernels
 - Hybrid Kernels
- System Calls
- User Interface
 - Command-Line Interface (CLI)
 - Graphical User Interface (GUI)
- Device Drivers

3. Process Management

- Process Concept
 - Definition and Characteristics
 - Process States
 - Process Control Block (PCB)
- Process Scheduling
 - CPU Scheduling Algorithms
 - First-Come, First-Served (FCFS)
 - Shortest Job Next (SJN)
 - Priority Scheduling
 - Round Robin (RR)
 - Multilevel Queue Scheduling
- Operations on Processes
 - Creation and Termination
 - Inter-process Communication (IPC)
 - Shared Memory
 - Message Passing

4. Threads and Concurrency

- Thread Concept
 - Definition and Types
 - Multithreading Models
- Benefits of Threads
- Thread Libraries
- Concurrency Issues

- Race Conditions
- Deadlock
- Starvation
- Synchronization Mechanisms
 - Mutexes
 - Semaphores
 - Monitors
 - Locks

5. Memory Management

- Memory Hierarchy
- Memory Allocation
 - Contiguous Allocation
 - Non-contiguous Allocation
 - Paging
 - Segmentation
- Virtual Memory
 - Concept and Benefits
 - Page Replacement Algorithms
 - FIFO
 - LRU
 - Optimal
 - Thrashing

6. Storage Management

- File System
 - File system Structure
 - File operations
 - File system implementation
 - Disk management and optimize
- File Allocation Methods
 - Contiguous Allocation
 - Linked Allocation
 - Indexed Allocation
- Disk Scheduling
 - FCFS
 - SSTF
 - SCAN
 - C-SCAN
 - LOOK
- RAID Levels

7. Input /Output (I/O) Management

- I/O Device and controller
- Device driver
- Kernel I/O Subsystem
- I/O technique(polling ,interrupts, DMA)
- I/O Scheduling

8. LINUX

- Introduction

- Linux History
- Key Feature of Linux
- Linux Architecture
- Detail Explanation of Linux Component
- Benefits of Using Linux
- Use Case of Linux

9. **Shell Programming**

- Introduction
- What is a Shell
- Key Concept in Shell Programming
- Writing a Shell Scripts
 - Variables
 - Control Structure
 - Function
 - Command Shell Commands
 - Input and Output Redirection
- Common Command in Shell Programming
 - File and Directory Operation
 - Text processing
 - System Information and Management
 - Networking
 - Process Control
 - Redirection and Pipelines
 - Shell Script Example

Chapter 1

Introduction to Operating Systems

An **operating system** is software that manages a computer's hardware. It also provides a basis for application programs and acts as an intermediary between the computer user and the computer hardware. An amazing aspect of operating systems is how they vary in accomplishing these tasks in a wide variety of computing environments. Operating systems are everywhere, from cars and home appliances that include "Internet of Things" devices, to smart phones, personal computers, enterprise computers, and cloud computing environments. In order to explore the role of an operating system in a modern computing environment, it is important first to understand the organization and architecture of computer hardware. This includes the CPU, memory, and I/O devices, as well as storage. A fundamental responsibility of an operating system is to allocate these resources to programs. Because an operating system is large and complex, it must be created piece by piece. Each of these pieces should be a well-delineated portion of the system, with carefully defined inputs, outputs, and functions. In this chapter, we provide a general overview of the major components of a contemporary computer system as well as the functions provided by the operating system. Additionally, we cover several topics to help set the stage for the remainder of the text: data structures used in operating systems, computing environments, and open-source and free operating systems.

History and Evolution

The development of operating systems has evolved significantly since the early days of computing:

- **Early Systems:** Initially, computers had no operating systems. Each program needed full control of the hardware, which was highly inefficient.
- **Batch Systems (1950s-1960s):** Programs were batched together and executed one after another without interaction from the user.
- **Time-Sharing Systems (1960s-1970s):** Allowed multiple users to interact with the computer simultaneously, giving the illusion of exclusive access.
- **Personal Computers (1980s-1990s):** The rise of PCs brought operating systems like MS-DOS and early versions of Windows and macOS, focusing on user-friendly interfaces.
- **Modern Systems (2000s-Present):** Emphasis on networking, security, and multi-user capabilities, with advanced operating systems like Windows 10, macOS, Linux, and mobile OSes like Android and iOS.

Types of Operating Systems

Batch Operating Systems

- **Definition:** Processes jobs in batches without user interaction during the processing.
- **Characteristics:** Efficient for large computations, minimal interaction, jobs are processed in the order they are submitted.
- **Example:** Early mainframe systems.

Time-Sharing Operating Systems

- **Definition:** Allows multiple users to share system resources simultaneously.
- **Characteristics:** Each user has a small time slice of the CPU, providing the impression of concurrent execution.
- **Example:** UNIX.

Distributed Operating Systems

- **Definition:** Manages a group of distinct computers and makes them appear to be a single computer.
- **Characteristics:** Shares computational tasks across multiple machines, enhancing performance and reliability.
- **Example:** Google's internal infrastructure.

Real-Time Operating Systems (RTOS)

- **Definition:** Designed to process data as it comes in, typically within a strict time constraint.
- **Characteristics:** High reliability and predictability, used in environments where timing is critical.
- **Example:** Systems used in medical devices, industrial control systems, and automotive applications.

Embedded Operating Systems

Definition: Tailored for use in embedded systems, which are part of a larger device.

- **Characteristics:** Highly specialized and optimized for the specific hardware, limited user interface, real-time capabilities.

- **Example:** Embedded Linux, VxWorks used in routers, smart appliances, and automotive systems.

Chapter 2

Operating System Architecture

Kernel

The kernel is the core component of an operating system, responsible for managing system resources, communication between hardware and software, and ensuring efficient operation. It operates in a privileged mode, providing essential services like process management, memory management, and device control.

Monolithic Kernels

- **Definition:** A single large process running entirely in a single address space.
- **Characteristics:** All OS services (file system, device drivers, memory management, etc.) run in kernel space, offering high performance but potentially less stability and security.
- **Example:** Linux, traditional UNIX systems.

Microkernels

- **Definition:** A minimalistic approach where only essential services like basic inter-process communication (IPC) and minimal process management run in kernel mode.
- **Characteristics:** Other services (e.g., device drivers, file systems) run in user space, improving modularity, stability, and security but possibly reducing performance due to increased context switching.
- **Example:** Minix, QNX.

Hybrid Kernels

- **Definition:** Combines aspects of monolithic and microkernel designs.
- **Characteristics:** Core services run in kernel space (as in monolithic kernels), while other services run in user space (as in microkernels), aiming to balance performance, flexibility, and stability.
- **Example:** Windows NT, macOS.

System Calls

System calls are the programming interface between a program and the operating system. They allow user-space applications to request services from the kernel, such as file manipulation, process control, and communication.

User Interface

The user interface is the layer through which users interact with the operating system. It can be categorized into two main types:

Command-Line Interface (CLI)

- **Definition:** A text-based interface where users input commands via a keyboard.
- **Characteristics:** Efficient for experienced users, scriptable, and powerful for administrative tasks. Requires knowledge of command syntax.
- **Example:** Bash in Linux, Command Prompt in Windows.

Graphical User Interface (GUI)

- **Definition:** A visual-based interface using graphical elements like windows, icons, and buttons.
- **Characteristics:** User-friendly, intuitive, and accessible to a wider audience. Provides a more engaging experience but requires more system resources.
- **Example:** Windows desktop environment, macOS, GNOME on Linux.

Device Drivers

Device drivers are specialized programs that allow the operating system to communicate with hardware devices. They act as translators between the OS and the hardware, ensuring that peripheral devices like printers, graphics cards, and network adapters function correctly. Device drivers can be part of the kernel or run in user space, depending on the operating system architecture.

Chapter 3

Operating System Process Management

Process Concept

Definition and Characteristics

A process is an instance of a program in execution. It is a fundamental concept in operating systems, encapsulating the program's code, its activity, and the resources it uses. Key characteristics of a process include:

- **Executable Program:** The program code itself.
- **Program Counter (PC):** A register that holds the address of the next instruction to be executed.
- **Stack:** Contains temporary data such as function parameters, return addresses, and local variables.
- **Data Section:** Includes global variables.
- **Heap:** Dynamically allocated memory during the process runtime.
- **Process State:** Reflects the current status of the process.

Process States

A process undergoes several states during its execution cycle, including:

- **New:** The process is being created.
- **Ready:** The process is waiting to be assigned to a processor.
- **Running:** Instructions are being executed.
- **Waiting:** The process is waiting for some event (such as I/O completion) to occur.
- **Terminated:** The process has finished execution.

These states are part of the process lifecycle, enabling the operating system to manage processes efficiently.

Process Control Block (PCB)

The Process Control Block (PCB) is a data structure used by the operating system to store all the information about a process. It is essential for process management and contains:

- **Process State:** Current state of the process (new, ready, running, waiting, terminated).
- **Program Counter:** Address of the next instruction to execute.
- **CPU Registers:** Contents of all process-centric registers.
- **Memory Management Information:** Details of the process's memory allocation, such as page tables or segment tables.
- **Accounting Information:** Information like CPU usage, process priority, and elapsed execution time.
- **I/O Status Information:** List of I/O devices allocated to the process and the status of I/O requests.
- **Process ID (PID):** Unique identifier for the process.

The PCB is crucial for context switching, allowing the operating system to save and restore the state of a process when switching between processes.

Process Scheduling

CPU Scheduling Algorithms

Process scheduling is a key function of the operating system, responsible for determining which process runs at any given time. Effective scheduling maximizes CPU utilization, throughput, and responsiveness while minimizing waiting time and response time. Various algorithms are used to achieve these goals:

First-Come, First-Served (FCFS)

- **Description:** The simplest scheduling algorithm, where processes are scheduled in the order they arrive in the ready queue.
- **Characteristics:**
 - Non-preemptive: Once a process starts, it runs to completion.
 - Simple to implement and understand.
 - Can lead to long waiting times, especially if a short process follows a long one (convoy effect).

- **Example:** If processes arrive in the order P1, P2, P3 with burst times of 24, 3, and 3 milliseconds, respectively, P1 will run first, followed by P2 and P3.

Shortest Job Next (SJN)

- **Description:** Also known as Shortest Job First (SJF), this algorithm selects the process with the smallest burst time to execute next.
- **Characteristics:**
 - Can be preemptive or non- preemptive.
 - Optimal in minimizing average waiting time.
 - Requires knowledge of future burst times, which is often not possible, leading to approximation.
- **Example:** Given processes with burst times 6, 8, 7, and 3 milliseconds, the order would be 3, 6, 7, 8 milliseconds.

Priority Scheduling

- **Description:** Processes are assigned a priority, and the CPU is allocated to the process with the highest priority (lowest numerical value).
- **Characteristics:**
 - Can be preemptive or non- preemptive.
 - Lower priority processes may suffer from starvation.
 - Can be mitigated with aging, gradually increasing the priority of waiting processes.
- **Example:** If processes P1, P2, P3 have priorities 3, 1, 2 respectively, P2 will run first, followed by P3 and then P1.

Round Robin (RR)

- **Description:** Each process is assigned a fixed time slice or quantum and is executed in a circular order.
- **Characteristics:**
 - Preemptive: Each process runs for a maximum of one time quantum.
 - Provides a good balance between responsiveness and throughput.
 - Performance depends on the length of the time quantum.
- **Example:** With a time quantum of 4 milliseconds and processes P1, P2, P3 with burst times 6, 8, and 7 milliseconds, the execution order would cycle through P1, P2, P3, then back to P1, and so on until all processes are complete.

Multilevel Queue Scheduling

- **Description:** Processes are divided into different queues based on their priority or type, and each queue has its own scheduling algorithm.
- **Characteristics:**
 - Supports different types of processes (e.g., interactive, batch) with different needs.
 - Processes are permanently assigned to a queue.
 - Scheduling can be hierarchical, with higher-priority queues preempting lower-priority ones.
- **Example:** A system with interactive processes in a high-priority queue using Round Robin and batch processes in a low-priority queue using FCFS.

Operations on Processes

Creation and Termination

Creation

- **Process Creation:** New processes are created by existing processes through system calls like `fork` in UNIX/Linux or `CreateProcess` in Windows. The creating process is known as the parent process, and the new process is the child process.
- **Reasons for Creation:**
 - **User Request:** User launches a new application.
 - **System Initialization:** System starts background services.
 - **Batch Job:** Execution of a batch job.
 - **Parent Request:** Parent process creates a child for parallel task execution.
- **Steps in Process Creation:**
 - Assigning a unique process ID (PID).
 - Allocating memory for the process.
 - Setting up process control block (PCB).
 - Loading the program into the allocated memory.
 - Passing initialization parameters to the process.

Termination

- **Process Termination:** Processes are terminated by themselves, by their parent processes, or by the operating system.
- **Reasons for Termination:**
 - **Normal Completion:** The process completes its execution successfully.
 - **Errors:** The process encounters a fatal error.
 - **External Request:** A user or parent process requests termination.
 - **Resource Unavailability:** Insufficient resources or deadlock.
- **Steps in Process Termination:**
 - Releasing the process's allocated resources.
 - Updating the process state to terminated.
 - Removing the process control block from the process table.
 - Optionally, notifying the parent process of termination.

Inter-Process Communication (IPC)

IPC mechanisms allow processes to communicate and synchronize their actions. This is crucial in multi-process systems to ensure coordinated execution and data sharing.

Shared Memory

- **Definition:** A segment of memory that multiple processes can access. It is the fastest IPC mechanism due to direct memory access.
- **Characteristics:**
 - Processes must ensure mutual exclusion to avoid race conditions.
 - Operating system provides mechanisms to create and manage shared memory regions.
- **Example:** In UNIX, shared memory can be created and managed using `shmget`, `shmat`, and `shmdt` system calls.
- **Use Case:** High-speed data exchange between processes, like in producer-consumer scenarios.

Message Passing

- **Definition:** Processes communicate by sending and receiving messages. This can be synchronous or asynchronous.
- **Characteristics:**

- Provides a higher level of abstraction compared to shared memory.
- Synchronization can be built into the communication mechanism.
- More suitable for distributed systems.
- **Example:** UNIX provides message passing through `msgsnd` and `msgrcv` system calls. In Windows, named pipes and mail slots can be used.
- **Use Case:** Communication between processes that may not share the same memory space or are running on different machines.

Chapter-4

Threads and Concurrency

Thread Concept

Definition and Types

Definition

A thread is the smallest unit of processing that can be performed by an operating system. It is a component of a process, which can run independently and concurrently with other threads of the same process. Threads share the process's resources, such as memory and file descriptors, but each thread has its own execution context, including a unique program counter, stack, and set of registers.

Types of Threads

1. User Threads

- **Managed by User-Level Libraries:** User threads are created and managed by user-level libraries rather than the operating system kernel.
- **Invisible to the Kernel:** The kernel is unaware of user threads, and as a result, it schedules the entire process rather than individual threads.
- **Examples:** Pthreads (POSIX threads), Java threads.
-

2. Kernel Threads

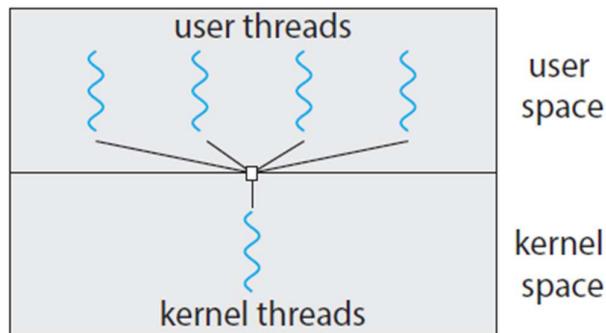
- **Managed by the Kernel:** Kernel threads are created and managed by the operating system kernel.
- **Visible to the Kernel:** The kernel is aware of each thread and can schedule them individually.
- **Examples:** Windows threads, Linux kernel threads.

Multithreading Models

1. Many-to-One Model

- **Description:** Multiple user threads are mapped to a single kernel thread.
- **Advantages:** Simple and efficient.

- **Disadvantages:** If one thread performs a blocking operation, the entire process blocks.
- **Use Case:** Not commonly used due to its limitations in modern systems.



- Fig 1. Many-to-One Model

2. One-to-One Model

- **Description:** Each user thread is mapped to a separate kernel thread.
- **Advantages:** Provides better concurrency, as the kernel can schedule multiple threads on multiple processors.
- **Disadvantages:** Can be resource-intensive, as each thread requires a kernel thread.
- **Use Case:** Used in Windows and Linux systems.

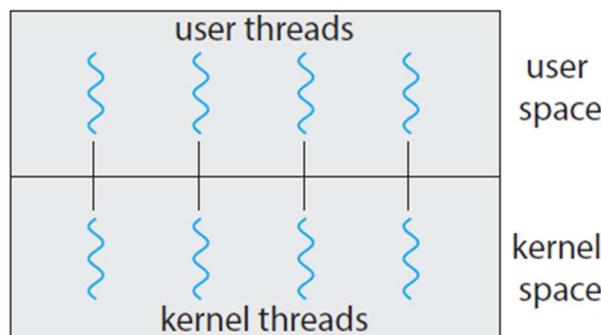


Fig 2. One to One Model

3. Many-to-Many Model

- **Description:** Multiple user threads are mapped to multiple kernel threads.
- **Advantages:** Combines the benefits of the above two models, allowing efficient management and better concurrency.
- **Disadvantages:** More complex to implement.
- **Use Case:** Used in systems that require high scalability and concurrency.

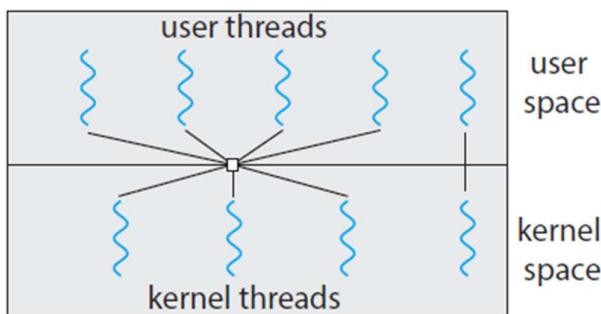


Fig 3 .Many-to-Many Model

Benefits of Threads

1. Responsiveness

- Threads can make applications more responsive by performing background tasks without freezing the user interface.
- Example: A web browser can use one thread to render a web page while another handles user input.

2. Resource Sharing

- Threads within the same process share memory and resources, allowing for efficient communication and data sharing.
- Example: Multiple threads in a server application can handle different client requests while sharing a common cache.

3. Economy

- Creating and managing threads is less resource-intensive compared to processes.
- Example: A server handling multiple connections can spawn threads more efficiently than processes, reducing overhead.

4. Scalability

- Threads can take advantage of multiprocessor architectures by running concurrently on multiple CPUs.
- Example: A parallel computing application can distribute tasks across multiple threads to speed up processing.

Thread Libraries

1. POSIX Threads (Pthreads)

Pthreads refers to the POSIX standard (IEEE 1003.1c) defining an API for thread creation and synchronization. This is a **specification** for thread behavior, not an **implementation**. Operating-system designers may implement the specification

- A standard API for creating and managing threads, widely used in UNIX-like systems.
- Example: `pthread_create`, `pthread_join` functions.

```
#include <pthread.h>
#include <stdio.h>
#include <stdlib.h>
int sum; /* this data is shared by the thread(s) */
void *runner(void *param); /* threads call this function */
int main(int argc, char *argv[])
{
pthread t tid; /* the thread identifier */
pthread attr t attr; /* set of thread attributes */
/* set the default attributes of the thread */
pthread_attr_init(&attr);
/* create the thread */
pthread_create(&tid, &attr, runner, argv[1]);
/* wait for the thread to exit */
pthread_join(tid,NULL);
printf("sum = %d\n",sum);
}
/* The thread will execute in this function */
void *runner(void *param)
{
int i, upper = atoi(param);
sum = 0;
for (i = 1; i <= upper; i++)
sum += i;
pthread_exit(0);
}
```

- Figure : Multithreaded C program using the Pthreads API.

2. Windows Threads

- The Windows API for thread management, providing functions to create, synchronize, and terminate threads.
- Example: CreateThread, WaitForSingleObject functions.

```
#include <windows.h>
#include <stdio.h>
DWORD Sum; /* data is shared by the thread(s) */

/* The thread will execute in this function */
DWORD WINAPI Summation(LPVOID Param)
{
    DWORD Upper = *(DWORD*)Param;
    for (DWORD i = 1; i <= Upper; i++)
        Sum += i;
    return 0;
}

int main(int argc, char *argv[])
{
    DWORD ThreadId;
    HANDLE ThreadHandle;
    int Param;

    Param = atoi(argv[1]);
    /* create the thread */
    ThreadHandle = CreateThread(
        NULL, /* default security attributes */
        0, /* default stack size */
        Summation, /* thread function */
        &Param, /* parameter to thread function */
        0, /* default creation flags */
        &ThreadId); /* returns the thread identifie

    /* now wait for the thread to finish */
    WaitForSingleObject(ThreadHandle, INFINITE);

    /* close the thread handle */
    CloseHandle(ThreadHandle);

    printf("sum = %d\n", Sum);
}
```

- Above Figure: Multithreaded C program using the Windows API
- .

3. Java Threads

Threads are the fundamental model of program execution in a Java program, and the Java language and its API provide a rich set of features for the creation and management of threads. All Java programs comprise at least a single thread of control—even a simple Java program consisting of only a main () method runs as a single

thread in the JVM. Java threads are available on any system that provides a JVM including Windows, Linux, and mac OS. The Java Thread API is available for Android applications as well.

There are two techniques for explicitly creating threads in a Java program. One approach is to create a new class that is derived from the Thread class and to override its run () method. An alternative—and more commonly used technique is to define a class that implements the Runnable interface. This interface defines a single abstract method with the signature public void run(). The code in the run() method of a class that implements Runnable is what executes in a separate thread. An example is shown below:

```

import java.util.concurrent.*;
class Summation implements Callable<Integer>
{
    private int upper;
    public Summation(int upper) {
        this.upper = upper;
    }
    /* The thread will execute in this method */
    public Integer call() {
        int sum = 0;
        for (int i = 1; i <= upper; i++)
            sum += i;
        return new Integer(sum);
    }
}
public class Driver
{
    public static void main(String[] args) {
        int upper = Integer.parseInt(args[0]);
        ExecutorService pool = Executors.newSingleThreadExecutor();
        Future<Integer> result = pool.submit(new Summation(upper));
        try {
            System.out.println("sum = " + result.get());
        } catch (InterruptedException | ExecutionException ie) { }
    }
}

```

Figure . Illustration of Java Executor framework API.

Concurrency Issues

Race Conditions

- **Definition:** Occurs when the outcome of processes or threads depends on the non-deterministic ordering or timing of their execution.

- **Example:** Multiple threads updating a shared variable simultaneously without proper synchronization.

Deadlock

Definition

Deadlock is a situation in operating systems where a set of processes are blocked because each process is holding a resource and waiting for another resource held by another process in the set. In simpler terms, a deadlock occurs when two or more processes cannot proceed because each is waiting for the other to release a resource.

Conditions for Deadlock

For a deadlock to occur, the following four conditions must hold simultaneously in the system:

1. **Mutual Exclusion:** At least one resource must be held in a non-shareable mode; only one process can use the resource at a time.
2. **Hold and Wait:** A process holding at least one resource is waiting to acquire additional resources held by other processes.
3. **No Preemption:** Resources cannot be forcibly taken away from a process holding them; they must be released voluntarily by the holding process.
4. **Circular Wait:** A set of processes are waiting for each other in a circular chain. Each process is waiting for a resource held by the next process in the chain.

Example Scenario:

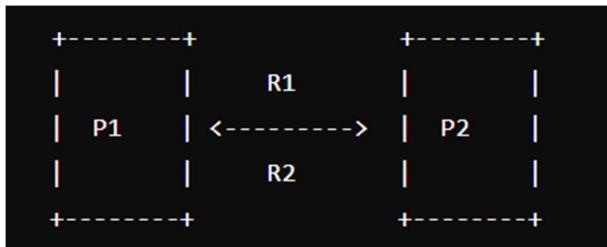
Let's consider a simple scenario with two processes (P1 and P2) and two resources (R1 and R2):

1. Process P1 holds resource R1 and waits for resource R2.
2. Process P2 holds resource R2 and waits for resource R1.

This scenario creates a circular wait, leading to a deadlock.

Diagram:

Below is a visual representation of this deadlock scenario:



In this diagram:

- P1 holds R1 and waits for R2.
- P2 holds R2 and waits for R1.

Both processes are unable to proceed because they are each waiting for a resource held by the other, creating a circular dependency and resulting in a deadlock.

6Deadlock Prevention

Deadlock prevention ensures that at least one of the necessary conditions for deadlock cannot hold. This can be achieved through several strategies:

- **Mutual Exclusion:** Make resources shareable wherever possible.
- **Hold and Wait:** Require processes to request all resources at once, ensuring that they hold no resources while waiting for others.
- **No Preemption:** Allow the operating system to preempt resources from processes (e.g., suspending a process and releasing its resources).
- **Circular Wait:** Impose an ordering of resource types and require processes to request resources in an increasing order of enumeration.

Deadlock Avoidance

Deadlock avoidance involves dynamically examining the resource-allocation state to ensure that a circular wait condition can never exist. The most well-known deadlock avoidance algorithm is the **Banker's Algorithm**, which works as follows:

- Each process must declare the maximum number of resources it may need.
- The system must determine whether the state of allocation will remain safe after granting the requested resources.

- A state is safe if the system can allocate resources to each process in some order and still avoid a deadlock.

Deadlock Detection

In cases where deadlock prevention and avoidance are not feasible, deadlock detection algorithms are used to detect the presence of deadlocks. The system must periodically check for deadlocks and take action if one is detected. Key steps include:

- **Resource Allocation Graphs:** A directed graph used to represent the state of resources and processes. A cycle in the graph indicates a possible deadlock.
- **Deadlock Detection Algorithm:** Examines the state of resource allocation and identifies cycles in the graph to detect deadlocks.

Deadlock Recovery

Once a deadlock is detected, the system must recover from it using one of the following methods:

- **Process Termination:** Terminate one or more processes involved in the deadlock until the cycle is broken.
 - **Abort all deadlocked processes:** Can be costly as it leads to loss of work.
 - **Abort one process at a time:** Requires rerunning the deadlock detection algorithm each time a process is terminated.
- **Resource Preemption:** Temporarily take resources away from processes and reallocate them to break the deadlock.
 - **Selecting a victim:** Choose the process that can be aborted or preempted with the least cost.
 - **Rollback:** Roll back the process to a safe state and restart it with the possibility of acquiring all necessary resources.

Starvation

Definition

Starvation is a situation in operating systems where a process or thread is perpetually denied necessary resources to proceed, despite the resource being available. It occurs when a scheduling algorithm unfairly allocates resources,

prioritizing certain processes or threads over others, leading to some processes or threads being consistently deprived of resources.

Causes of Starvation

1. **Priority Inversion:** Lower-priority processes or threads may starve if higher-priority processes or threads monopolize the resources they need.
2. **Resource Contention:** If a resource is heavily contested by multiple processes or threads, those with lower priority or less frequent access may starve due to constant competition.
3. **Poor Scheduling Policies:** Inefficient scheduling algorithms may unfairly allocate resources, causing some processes or threads to be consistently overlooked in favor of others.

Effects of Starvation

1. **Reduced Performance:** Starvation can lead to degraded system performance as critical processes or threads may not receive the resources they need to execute efficiently.
2. **Deadlocks:** In extreme cases, starvation can contribute to deadlocks if a process or thread is unable to acquire the resources it needs to release resources held by others, leading to a circular dependency.
3. **Fairness Issues:** Starvation can result in fairness issues, where certain processes or threads are unfairly denied resources, impacting system reliability and user experience.

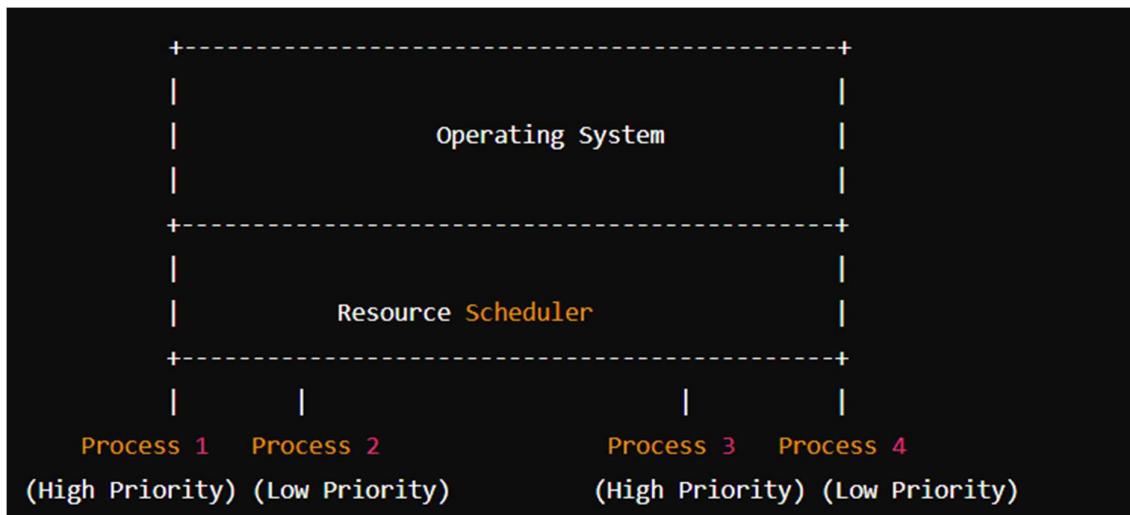
Prevention and Mitigation

1. **Fair Scheduling Policies:** Use scheduling algorithms that prioritize fairness and prevent resource starvation. For example, round-robin scheduling ensures that all processes or threads receive equal access to resources over time.
2. **Priority Inversion Avoidance:** Implement techniques to prevent priority inversion, such as priority inheritance or priority ceiling protocols, which ensure that lower-priority processes or threads can temporarily inherit higher priority when accessing critical resources.
3. **Resource Reservation:** Reserve resources for critical processes or threads to ensure they are not starved by lower-priority tasks. This can be achieved through techniques such as resource partitioning or quality of service (QoS) guarantees.
4. **Dynamic Resource Allocation:** Dynamically adjust resource allocation based on system load and process or thread priorities to ensure that critical tasks are not starved during peak demand periods.

Example

Consider a multi-user operating system where multiple users are running CPU-bound tasks. If the scheduling algorithm favors short tasks over long ones, long-running tasks may be continuously delayed, leading to starvation. As a result, these long-running tasks may experience significantly longer wait times and reduced performance compared to short tasks, despite both types of tasks being equally important.

Below diagram illustrating the concept of starvation in an operating system:



In this diagram:

- **Operating System:** Manages the allocation of resources among processes.
- **Resource Scheduler:** Component responsible for scheduling processes and allocating resources.
- **Process 1, Process 2, Process 3, Process 4:** Represent processes in the system.
 - **High Priority:** Processes with high priority.
 - **Low Priority:** Processes with low priority.

In a scenario where the resource scheduler prioritizes high-priority processes over low-priority ones, low-priority processes may experience starvation. Even though resources are available, they are consistently allocated to high-priority processes, leaving low-priority processes waiting indefinitely for access to resources.

Synchronization Mechanisms

1. Mutexes

- **Definition:** Mutual exclusion locks that ensure only one thread accesses a critical section at a time.
- **Characteristics:** Simple to implement, but can lead to contention and priority inversion.

2. Semaphores

Definition

Semaphores are a synchronization mechanism used in operating systems to control access to shared resources and coordinate the execution of multiple processes or threads. They act as counters that can be incremented or decremented atomically and are typically used to enforce mutual exclusion and manage access to critical sections of code.

Types of Semaphores

1. Binary Semaphores (Mutexes)

- Binary semaphores have only two possible values: 0 and 1.
- They are commonly used to implement mutual exclusion, where the semaphore is locked (set to 1) when a resource is being accessed and unlocked (set to 0) when the resource is released.

2. Counting Semaphores

- Counting semaphores can have an integer value greater than or equal to zero.
- They are used to control access to a finite number of resources, allowing multiple processes or threads to access the resources concurrently up to a certain limit.

Operations on Semaphores

1. Wait (P) Operation

- Decrements the value of the semaphore by 1.
- If the semaphore's value becomes negative, the process or thread is blocked until the semaphore's value becomes non-negative.

2. Signal (V) Operation

- Increments the value of the semaphore by 1.

- If there are processes or threads blocked on the semaphore, one of them is woken up.

Example Use Cases

1. Mutual Exclusion

- Binary semaphores are used to enforce mutual exclusion, ensuring that only one process or thread can access a critical section of code at a time.
- Example: Protecting a shared variable from concurrent access by multiple threads.

2. Resource Allocation

- Counting semaphores are used to control access to a finite number of resources, such as a pool of shared memory buffers or a fixed-size thread pool.
- Example: Limiting the number of concurrent connections to a network server.

Implementation Considerations

1. Atomicity

- Semaphore operations must be implemented atomically to prevent race conditions and ensure correct synchronization.

2. Deadlock Avoidance

- Care must be taken to avoid deadlocks when using semaphores by ensuring that processes or threads do not acquire multiple semaphores in different orders.

3. Initialization and Cleanup

- Semaphores must be properly initialized before use and released (if necessary) when they are no longer needed to prevent resource leaks.

Advantages of Semaphores

1. Flexibility

- Semaphores can be used to solve a wide range of synchronization problems, from simple mutual exclusion to complex resource allocation scenarios.

2. Efficiency

- Semaphore operations are typically fast and efficient, making them suitable for use in performance-critical applications.

3. Portability

- Semaphores are a standard synchronization primitive supported by most operating systems, making code that uses semaphores highly portable across different platforms.

Limitations of Semaphores

1. Complexity

- Semaphores can be more complex to use correctly compared to other synchronization mechanisms like mutexes or condition variables.

2. Potential for Deadlocks

- Improper use of semaphores can lead to deadlocks, where processes or threads are blocked indefinitely waiting for resources that are held by other processes or threads.

•

3. Monitors

Definition

Monitors are a higher-level synchronization mechanism used in operating systems to manage access to shared resources and coordinate the execution of concurrent processes or threads. Unlike low-level synchronization primitives like semaphores, monitors encapsulate both data and synchronization operations within a single construct, making them easier to use and less error-prone.

Components of Monitors

1. Data Variables

- Monitors contain shared data variables that represent the state of the resource being managed.
- These variables can be accessed and modified by multiple processes or threads.

2. Procedures or Methods

- Monitors contain procedures or methods that provide operations to manipulate the shared data.
- These procedures are executed within the context of the monitor and have exclusive access to the shared data.

3. Condition Variables

- Monitors contain condition variables that allow processes or threads to wait for specific conditions to become true before proceeding.
- Condition variables are associated with predicates that represent the conditions being waited for.

Operations on Monitors

1. Enter Monitor

- A process or thread must acquire exclusive access to the monitor before it can execute any of its procedures or access its data.
- If the monitor is currently in use by another process or thread, the entering process or thread is blocked until it can acquire access.

2. Execute Procedure

- Once inside the monitor, a process or thread can execute any of its procedures, which have exclusive access to the shared data.
- Other processes or threads are blocked from entering the monitor or executing procedures until the current procedure completes and releases control.

3. Wait (Condition Variable)

- A process or thread can wait on a condition variable associated with a specific predicate.
- Waiting causes the process or thread to release control of the monitor and be blocked until another process or thread signals or broadcasts on the condition variable.

4. Signal (Condition Variable)

- A process or thread can signal or notify one waiting process or thread that the condition it was waiting for has become true.
- Signaling allows the waiting process or thread to re-enter the monitor and continue execution.

5. Broadcast (Condition Variable)

- A process or thread can broadcast or notify all waiting processes or threads that the condition they were waiting for has become true.
- Broadcasting allows multiple waiting processes or threads to re-enter the monitor and continue execution.

Example Use Cases

1. Producer-Consumer Problem

- Monitors can be used to synchronize access to a shared buffer between multiple producer and consumer processes or threads.
- Procedures within the monitor provide operations to add items to the buffer (producer) and remove items from the buffer (consumer), while condition variables are used to signal when the buffer is empty or full.

2. Readers-Writers Problem

- Monitors can be used to synchronize access to a shared resource between multiple reader and writer processes or threads.
- Procedures within the monitor provide operations to read from the resource (reader) and write to the resource (writer), while condition variables are used to control access based on the number of readers or writers currently active.

Advantages of Monitors

1. Simplicity

- Monitors encapsulate both data and synchronization operations within a single construct, making them easier to use and understand compared to low-level synchronization primitives like semaphores.

2. Safety

- Monitors ensure mutual exclusion and prevent race conditions by providing exclusive access to shared data and controlling the execution of procedures within the monitor.

3. Expressiveness

- Monitors provide high-level synchronization primitives like condition variables, which allow processes or threads to wait for specific conditions to become true before proceeding.

Limitations of Monitors

1. Blocking

- Monitors can lead to potential blocking if a process or thread holding the monitor is preempted while executing a procedure, causing other processes or threads to be blocked from entering the monitor.

2. Overhead

- Monitors may introduce additional overhead compared to low-level synchronization primitives due to the need for mutual exclusion and condition variable management.

4. Locks

Definition

Locks are a fundamental synchronization mechanism used in operating systems to control access to shared resources and ensure mutual exclusion between multiple processes or threads. A lock allows only one process or thread at a time to acquire exclusive access to a resource, preventing concurrent access and potential data corruption.

Types of Locks

1. Mutex (Mutual Exclusion)

- Mutex locks are the simplest form of locks and provide mutual exclusion by allowing only one thread to acquire the lock at a time.
- A thread attempting to acquire a locked mutex will be blocked until the lock is released by the thread holding it.

2. Spinlock

- Spinlocks are lightweight locks that busy-wait until they can acquire the lock.
- Rather than blocking, a thread attempting to acquire a locked spinlock will repeatedly check the lock's status until it becomes available.

3. Read-Write Lock (Reader-Writer Lock)

- Read-write locks allow multiple readers to acquire the lock simultaneously but only one writer to acquire it exclusively.
- Readers can access the resource concurrently without blocking each other, while writers must wait for exclusive access.

Operations on Locks

1. Acquire (Lock) Operation

- A process or thread attempts to acquire the lock to gain exclusive access to the resource.
- If the lock is available, the process or thread acquires it and proceeds to access the resource.
- If the lock is already held by another process or thread, the acquiring process or thread is blocked or spins until the lock becomes available.

2. Release (Unlock) Operation

- A process or thread releases the lock to allow other processes or threads to acquire it.
- Once the lock is released, one of the waiting processes or threads (if any) is allowed to acquire it.

Advantages of Locks

1. Simplicity

- Locks provide a simple and intuitive mechanism for controlling access to shared resources, making them easy to use and understand.

2. Efficiency

- Locks are generally lightweight and efficient, with low overhead compared to other synchronization mechanisms like semaphores or monitors.

3. Flexibility

- Locks can be used in a wide range of synchronization scenarios, from protecting critical sections of code to managing access to shared data structures.

Limitations of Locks

1. Deadlocks

- Improper use of locks can lead to deadlocks, where multiple processes or threads are blocked indefinitely waiting for resources that are held by others.

2. Priority Inversion

- Priority inversion can occur when a low-priority process holds a lock required by a high-priority process, causing the high-priority process to wait unnecessarily.

3. Performance Degradation

- Spinlocks can cause performance degradation on multiprocessor systems due to busy-waiting, especially if locks are held for extended periods.

Example Use Cases

1. Critical Section Protection

- Mutex locks are commonly used to protect critical sections of code that access shared resources, ensuring that only one thread can execute the critical section at a time.

2. Resource Allocation

- Spinlocks can be used to protect access to memory pools or other finite resources, allowing efficient allocation and deallocation by multiple threads.

Chapter 5

Memory Management

Memory Hierarchy

1. Introduction to Memory Hierarchy

- Memory hierarchy organizes computer memory into different levels based on speed, cost, and capacity.
- It ensures faster access to frequently used data and instructions while optimizing resources.

2. Levels of Memory Hierarchy

- **Registers:** Located inside the CPU, registers are the fastest form of memory, used to store data and instructions currently being processed.
- **Cache Memory:** Situated between registers and main memory, cache memory stores frequently accessed data and instructions, speeding up CPU operations.
- **Main Memory (RAM):** Holds data and instructions actively used by programs. It's slower than cache but faster than secondary storage.
- **Secondary Storage:** Devices like hard drives and SSDs provide non-volatile storage for data and programs not currently in use.

3. Memory Management Tasks

- **Address Mapping:** The OS maps logical addresses to physical addresses in memory.
- **Allocation and Deallocation:** It allocates memory space to processes and deallocates it when no longer needed.
- **Protection:** Ensures processes do not interfere with each other's memory space, preventing unauthorized access.
- **Sharing and Fragmentation:** Facilitates memory sharing among processes and manages memory fragmentation to optimize usage.

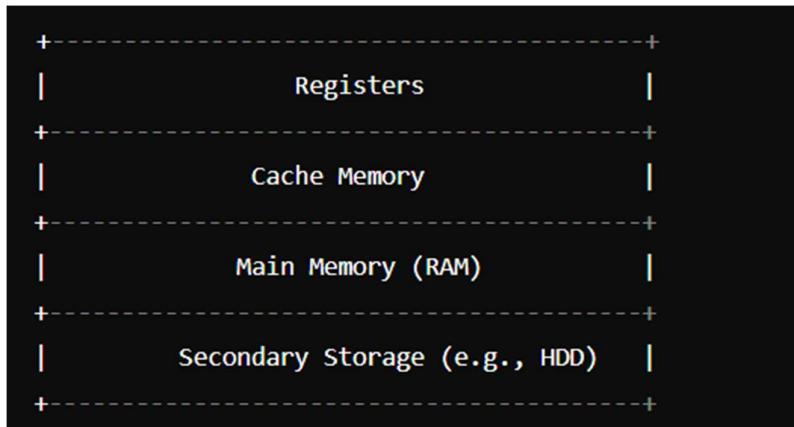
4. Memory Management Techniques

- **Virtual Memory:** Uses part of secondary storage as an extension of main memory, swapping data between RAM and disk as needed.
- **Paging and Segmentation:** Divides physical memory into manageable chunks for efficient allocation.
- **Memory Protection:** Hardware and software mechanisms prevent processes from accessing unauthorized memory locations.

5. Performance Considerations

- Efficient memory management is crucial for system performance.
- Factors like access time, throughput, and latency influence memory hierarchy effectiveness.
- Techniques such as caching, prefetching, and compression enhance performance.

Diagram:



This diagram represents the typical memory hierarchy in an operating system. Registers are the fastest memory, followed by cache, main memory, and secondary storage. The operating system manages memory allocation, address mapping, protection, and optimization across these levels to ensure efficient system performance.

Memory Allocation Strategies

Contiguous Allocation:-

1. Introduction

- Memory allocation strategies determine how processes are assigned memory space during their execution.
- Efficient memory allocation is crucial for optimizing system performance and resource utilization.

2. Fixed Partitioning

- **Definition:** Divides physical memory into fixed-size partitions, each capable of holding one process.
- **Key Concepts:**
 - Partitions may be of equal or unequal sizes.
 - Processes are assigned to partitions based on their size, with smaller processes allocated to smaller partitions.
- **Advantages:**
 - Simple and easy to implement.
 - Minimizes external fragmentation.
- **Challenges:**
 - May lead to internal fragmentation if partition sizes do not match process sizes.
 - Inflexible and inefficient for varying process sizes.

3. Dynamic Partitioning (Variable Partitioning)

- **Definition:** Allocates memory dynamically to processes based on their size.
- **Key Concepts:**
 - Memory is divided into variable-sized partitions, with each partition assigned to a process as needed.
 - When a process terminates, its partition is reclaimed and added to the free memory pool.
- **Advantages:**
 - Flexible and efficient for varying process sizes.
 - Reduces internal fragmentation compared to fixed partitioning.
- **Challenges:**
 - May suffer from external fragmentation if memory becomes fragmented over time.
 - Requires efficient memory management algorithms to allocate and deallocate partitions.

4. Memory Allocation Algorithms

- **First Fit:** Allocates the first available partition that is large enough to accommodate the process.
- **Best Fit:** Allocates the smallest available partition that is large enough to accommodate the process.
- **Worst Fit:** Allocates the largest available partition, leaving the smallest possible hole.
- **Next Fit:** Similar to first fit, but starts searching for a suitable partition from the location of the last allocation.

Example and Diagram

Let's consider an example where a system has 1000 MB of memory, and it uses dynamic partitioning to allocate memory to processes. Here's a step-by-step illustration:

1. **Initial State:** All 1000 MB of memory is free.
2. **Process A (200 MB):** Allocates the first 200 MB of memory.
3. **Process B (300 MB):** Allocates the next 300 MB of memory.
4. **Process C (100 MB):** Allocates the next 100 MB of memory.
5. **Process B terminates:** 300 MB is freed between Process A and Process C.

Here's a simple diagram to illustrate contiguous allocation and fragmentation:

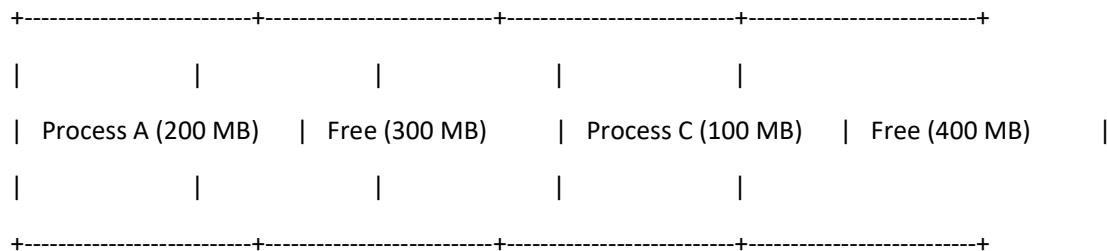
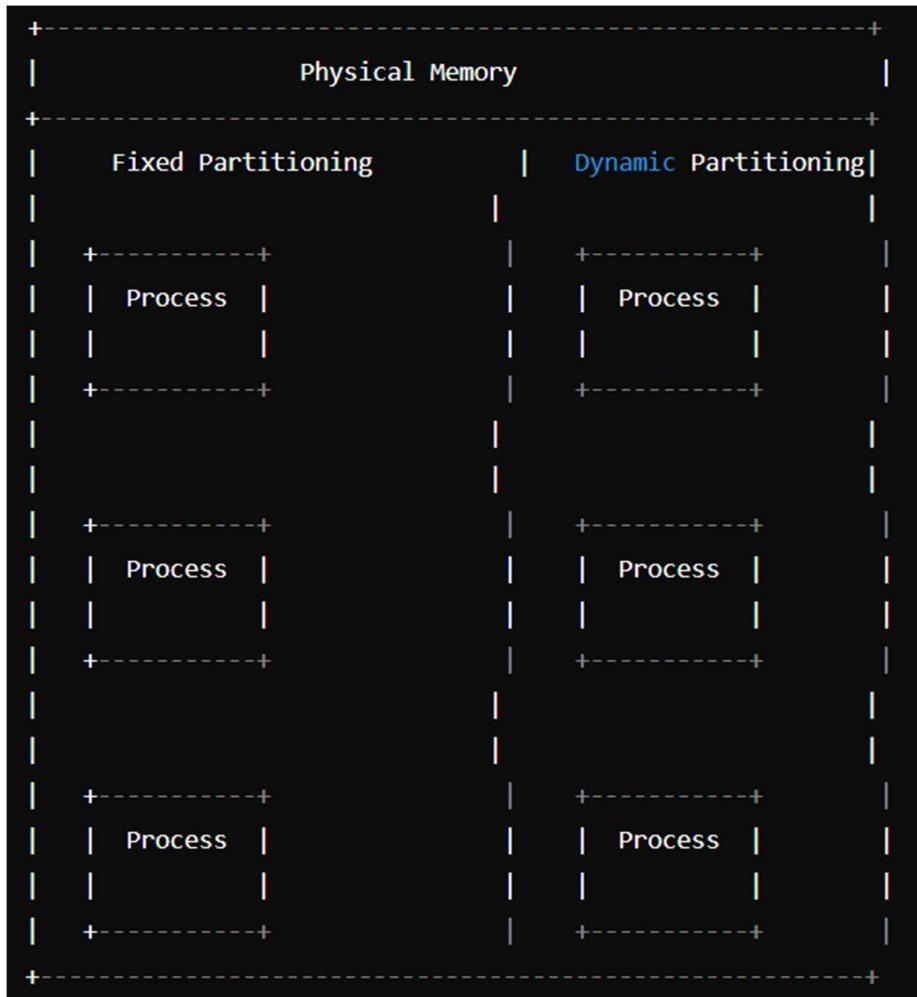


Diagram:



In the diagram, physical memory is represented, showing both fixed partitioning and dynamic partitioning approaches. In fixed partitioning, memory is divided into fixed-sized partitions, each assigned to a process. In dynamic partitioning, memory is divided into variable-sized partitions, which are allocated to processes as needed. Memory allocation strategies play a crucial role in optimizing memory usage and system performance.

Non Contiguous Allocation:-

Paging and Segmentation

1. Introduction

- Paging and segmentation are memory management techniques used by operating systems to manage memory allocation efficiently.

2. Paging

- **Definition:** Paging divides the physical memory into fixed-size blocks called "pages" and the logical memory into blocks of the same size called "frames".
- **Key Concepts:**
 - Each process's address space is divided into fixed-size pages.
 - Pages are stored in physical memory frames, with a page table mapping logical pages to physical frames.
 - Pages can be loaded into any available frame in physical memory, allowing for efficient memory allocation.
- **Advantages:**
 - Simplifies memory allocation and management.
 - Eliminates external fragmentation.
- **Challenges:**
 - Internal fragmentation may occur if a page is not fully utilized.
 - Overhead of maintaining page tables.

3. Segmentation

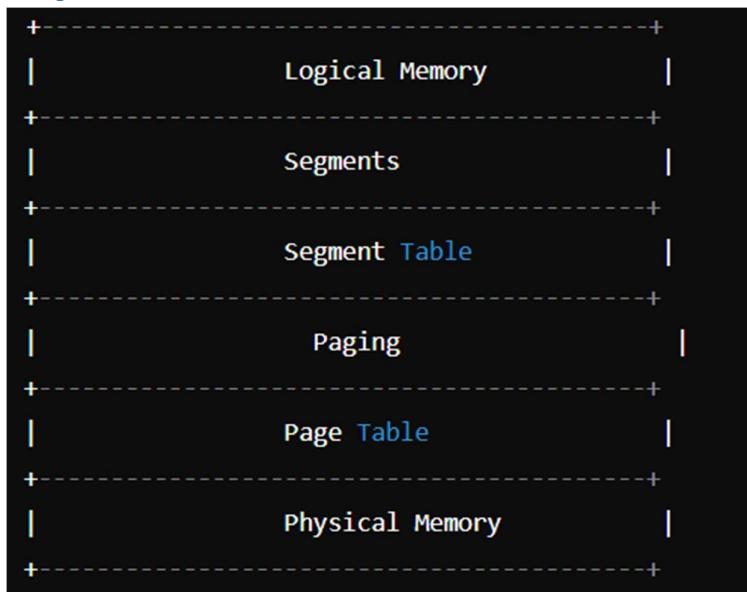
- **Definition:** Segmentation divides the logical memory into variable-sized segments that represent logical units of a program.
- **Key Concepts:**
 - Each segment corresponds to a particular type of data (e.g., code, stack, heap).
 - Segments are of variable sizes, depending on the program's requirements.
 - Each segment is mapped to a contiguous region of physical memory.
- **Advantages:**
 - Provides a more flexible memory allocation scheme.
 - Allows better protection and sharing of memory segments.

- **Challenges:**
 - External fragmentation may occur when variable-sized segments are allocated and deallocated.
 - Overhead of managing segment tables and ensuring protection.

4. Combining Paging and Segmentation

- Many modern operating systems use a combination of paging and segmentation to leverage the advantages of both techniques.
- This approach, known as "paged segmentation" or "segmented paging," combines the flexibility of segmentation with the simplicity of paging.

Diagram:



In the diagram, the logical memory of a process is first divided into segments, representing different parts of the program (e.g., code, data, stack). Each segment is then mapped to a contiguous region of physical memory using a segment table. Within each segment, paging is used to further divide the memory into fixed-size pages, which are mapped to physical frames using a page table. This combination of segmentation and paging allows for efficient memory management while providing flexibility in memory allocation.

Virtual Memory

Concept and Benefits

Concept

Virtual memory is a memory management technique that provides an "idealized abstraction of the storage resources that are actually available on a given machine" which creates the illusion to users of a very large (main) memory. The primary mechanism behind virtual memory is the mapping of virtual addresses to physical addresses using a combination of hardware and software, allowing an operating system to use more memory than what is physically available.

Benefits

1. Increased Memory Utilization

- Virtual memory allows systems to run larger applications than the physical memory would allow by using disk space as an extension of RAM.

2. Isolation and Protection

- Each process has its own virtual address space, which prevents one process from interfering with another's memory, thus increasing security and stability.

3. Simplified Memory Management

- Programmers can work with a large, contiguous memory space without worrying about the underlying physical memory constraints.

4. Efficient Multitasking

- Virtual memory enables multiple processes to share the physical memory efficiently, improving the system's overall performance and responsiveness.

Page Replacement Algorithms

Page replacement algorithms are used in virtual memory systems to decide which memory pages to swap out when a new page needs to be loaded into memory and there is no free space available.

FIFO (First-In, First-Out)

- **Concept:** The oldest page in memory (the one that was loaded first) is the first to be replaced.
- **Advantages:** Simple to implement.

- **Disadvantages:** Does not consider the page's usage frequency, which can lead to poor performance (Belady's anomaly).

[LRU \(Least Recently Used\)](#)

- **Concept:** The page that has not been used for the longest time is replaced.
- **Advantages:** Generally provides better performance than FIFO by considering usage history.
- **Disadvantages:** Can be more complex and costly to implement, as it requires keeping track of the order of page accesses.

[Optimal](#)

- **Concept:** Replaces the page that will not be used for the longest period in the future.
- **Advantages:** Provides the best possible performance.
- **Disadvantages:** Impossible to implement in practice because it requires future knowledge of the page requests.

[Thrashing](#)

[Definition](#)

Thrashing occurs when a system spends more time swapping pages in and out of memory than executing actual tasks. This happens when there is insufficient physical memory to support the working set of active processes, leading to a high rate of page faults and constant paging.

[Causes](#)

1. **Insufficient Physical Memory:** When the combined working sets of all processes exceed the physical memory.
2. **Overcommitment:** Running too many processes concurrently.
3. **Inadequate Page Replacement Algorithm:** Poor choices in page replacement can exacerbate thrashing.

[Effects](#)

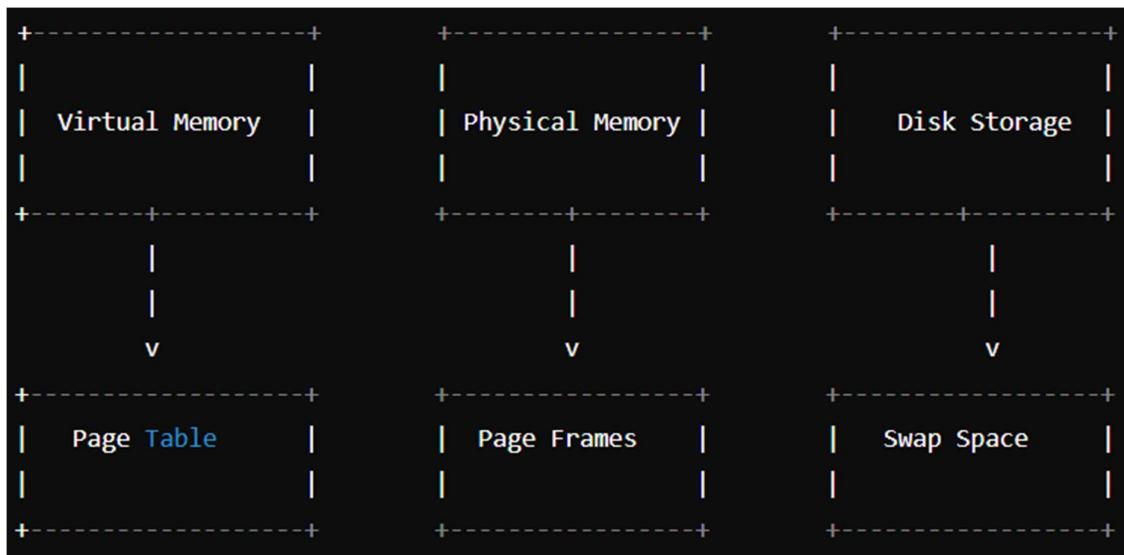
1. **Performance Degradation:** The system becomes slow and unresponsive as it spends most of its time swapping pages.
2. **Increased CPU Usage:** High CPU time spent on handling page faults.

Mitigation Strategies

1. **Increasing Physical Memory:** Adding more RAM can help reduce thrashing.
2. **Reducing the Multiprogramming Level:** Limiting the number of concurrent processes.
3. **Efficient Page Replacement Algorithms:** Using more intelligent algorithms like LRU to minimize unnecessary paging.

Diagram for Virtual Memory Concept

Here's a simple diagram to illustrate the virtual memory concept:



In this diagram:

- **Virtual Memory:** Logical view of memory space for a process.
- **Page Table:** Maps virtual addresses to physical addresses.
- **Physical Memory:** Actual RAM.
- **Disk Storage:** Used to extend RAM via swap space.

Chapter-6

Storage Management

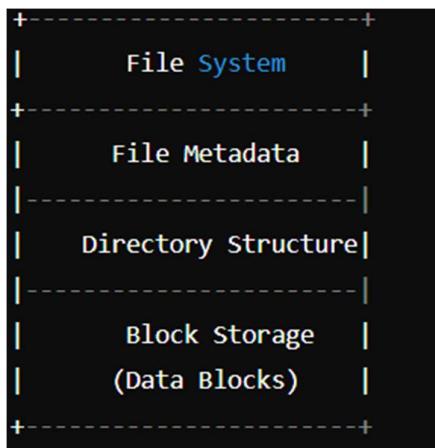
File Systems:

In computing, a file system is a method and data structure used by an operating system to manage and organize files on a storage medium (typically a disk drive). It provides a way to store, retrieve, and manipulate data stored on the disk. File systems manage the allocation of space on the disk and provide a logical structure for organizing files and directories.

File System Structure:

1. **File:** A named collection of related information stored on disk.
2. **Directory (Folder):** A special type of file that contains a list of file names and their corresponding metadata (such as size, permissions, etc.). Directories provide a hierarchical structure for organizing files.
3. **Metadata:** Information about a file, including its name, size, type, permissions, creation date, and last modification date.
4. **Blocks:** The disk space allocated to store file data. Files are divided into fixed-size blocks, and each block is assigned a unique address.
5. **File Allocation Table (FAT) or Inode Table:** A data structure used by the file system to keep track of the allocation of disk space to files. It maintains a mapping between file names and their corresponding blocks on the disk.

Basic File System Diagram:



In the diagram above:

- **File Metadata:** Stores information about each file, such as its name, size, permissions, etc.
- **Directory Structure:** Organizes files into directories (folders) and maintains the hierarchy of the file system.
- **Block Storage (Data Blocks):** Stores the actual data of the files on the disk. Each file is divided into blocks, and these blocks are allocated on the disk to store the file data.

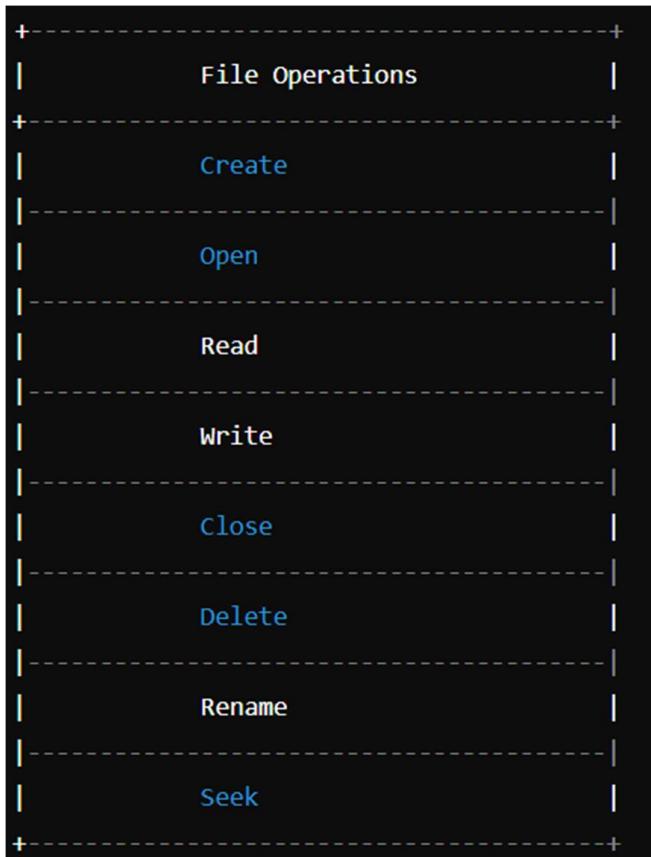
File Operations:

File operations refer to the actions or activities that can be performed on files within a file system. These operations include creating, opening, reading, writing, closing, deleting, and modifying files. File systems provide a set of system calls or APIs (Application Programming Interfaces) that applications and users can use to perform these operations.

Common File Operations:

1. **Create:** Create a new file in the file system.
2. **Open:** Open an existing file for reading, writing, or both.
3. **Read:** Read data from a file into memory.
4. **Write:** Write data from memory to a file.
5. **Close:** Close an opened file, releasing any associated resources.
6. **Delete:** Delete a file from the file system.
7. **Rename:** Change the name of a file.
8. **Seek:** Move the file pointer to a specific position within the file.

Basic File Operations Diagram:



In the diagram above:

- Each box represents a file operation that can be performed on a file within the file system.
- Arrows indicate the flow or sequence of these operations.
- Applications and users interact with the file system by invoking these operations through system calls or APIs.

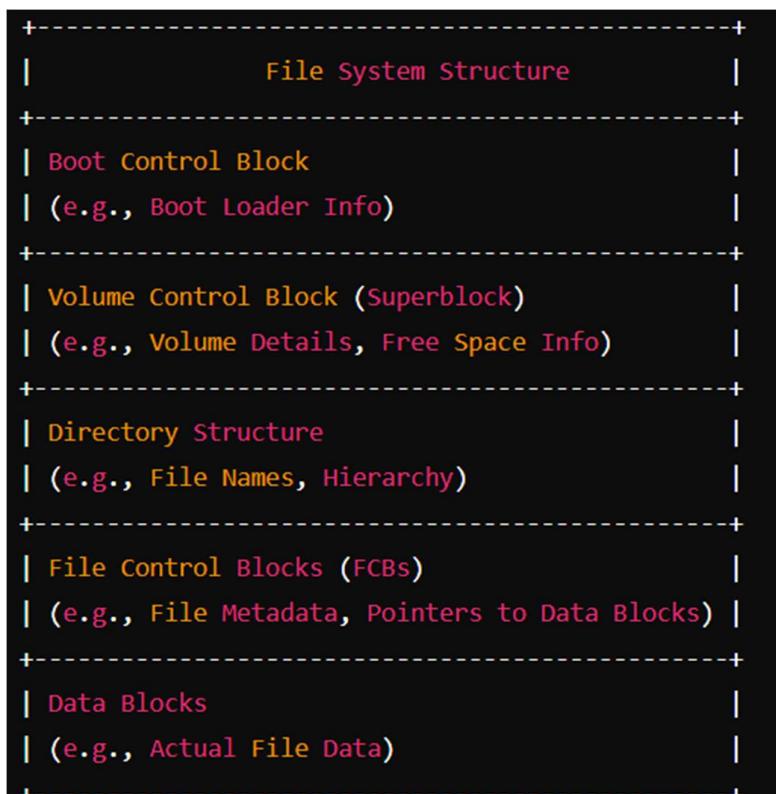
File System Implementation:

File system implementation involves the methods and data structures used by an operating system to manage files and directories on storage devices. This implementation ensures efficient data storage, retrieval, and management, providing a logical structure for organizing and accessing files.

Key Components of File System Implementation:

1. **Boot Control Block:** Contains information needed to boot the operating system from the volume. It includes details like the location of the boot loader.
2. **Volume Control Block (Superblock):** Contains volume details such as the number of blocks, block size, free block count, and the location of the free block management.
3. **Directory Structure:** Stores information about files, including their names, sizes, and locations on the disk. It organizes files hierarchically.
4. **File Control Block (FCB):** Contains details about a file, including its name, permissions, size, and pointers to the data blocks on the disk.
5. **Data Blocks:** The actual storage locations where file data is kept. Files are divided into blocks, and each block has a unique address.

Basic File System Implementation Diagram:



Detailed Explanation:

1. **Boot Control Block:**
 - Typically located at the beginning of the disk.
 - Contains information necessary for the system to boot from the disk.
2. **Volume Control Block (Superblock):**

- Located after the boot block.
- Contains critical information about the file system, such as the total number of blocks, block size, free block count, and the location of key file system structures.

3. Directory Structure:

- A hierarchical structure that keeps track of all files and directories in the file system.
- Each directory entry points to a file control block (FCB) for detailed information about the file.

4. File Control Blocks (FCBs):

- Stores metadata about files, such as file names, access permissions, sizes, and pointers to the blocks where the file data is stored.
- Each file has a corresponding FCB.

5. Data Blocks:

- The actual locations on the disk where file data is stored.
- Files are divided into fixed-size blocks, and each block can be addressed individually.

Example Workflow:

- When a file is created, the directory structure is updated with an entry for the new file, and an FCB is created.
- When a file is opened, the FCB is accessed to get metadata and pointers to data blocks.
- Reading or writing to a file involves accessing the data blocks pointed to by the FCB.
- When a file is deleted, its directory entry and FCB are removed, and its data blocks are marked as free in the volume control block.

Disk Management and Optimization:

Disk management involves organizing and maintaining data on storage devices such as hard drives and SSDs. The goal is to maximize efficiency, performance, and reliability of data storage and retrieval. Optimization techniques ensure that data is accessed quickly and efficiently, reducing wear and tear on the disk and improving overall system performance.

Key Components of Disk Management:

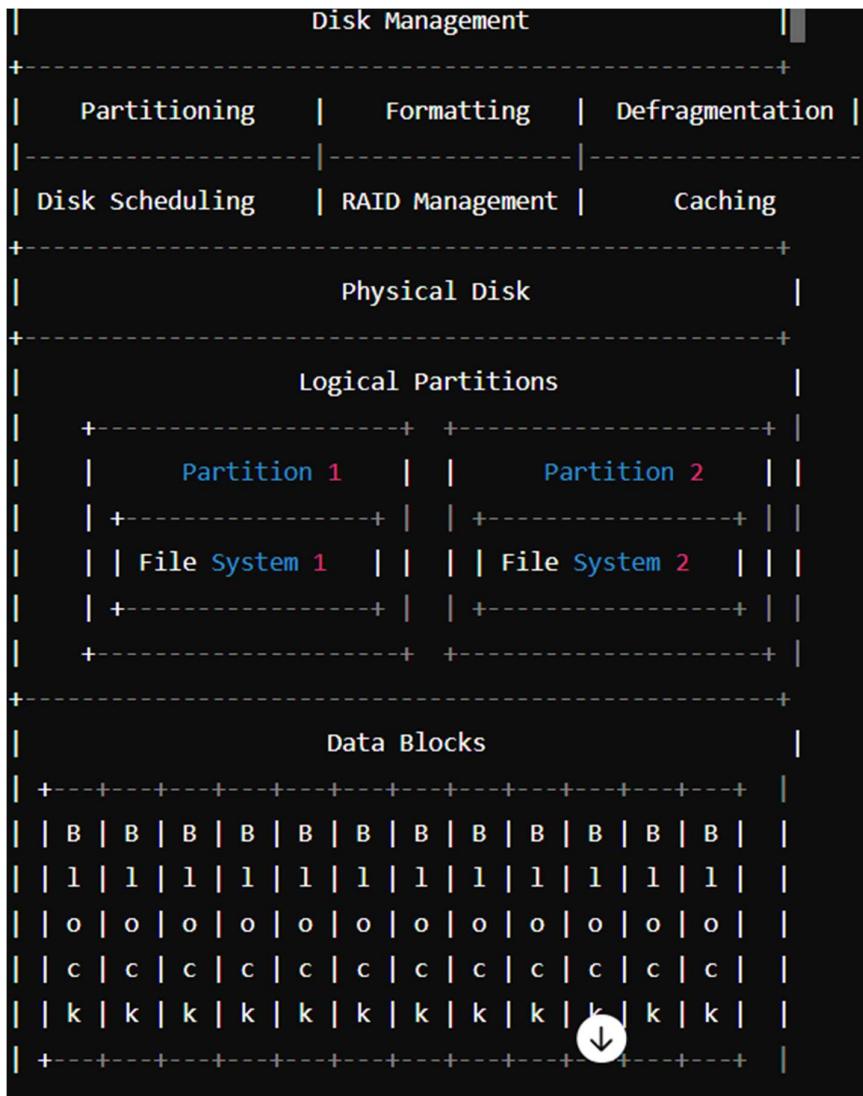
1. **Partitioning:**
 - Dividing a disk into separate regions called partitions.
 - Each partition can be managed as a distinct logical disk, often with its own file system.
2. **Formatting:**
 - Preparing a partition to hold files by creating a file system structure within it.
 - Includes creating a boot sector, superblock, inodes, and data blocks.
3. **Disk Scheduling:**
 - Managing the order in which disk I/O requests are processed.
 - Common algorithms include First-Come, First-Served (FCFS), Shortest Seek Time First (SSTF), and Elevator (SCAN) algorithms.
4. **Defragmentation:**
 - Reorganizing the layout of files on disk to ensure that file data blocks are contiguous.
 - Reduces the time needed to read or write files by minimizing disk head movement.
5. **RAID (Redundant Array of Independent Disks):**
 - Combining multiple physical disks into a single logical unit for redundancy and performance improvement.
 - Different levels (e.g., RAID 0, RAID 1, RAID 5) offer varying balances of performance, redundancy, and storage capacity.
6. **Caching:**
 - Storing frequently accessed data in faster storage (e.g., RAM) to speed up read/write operations.
 - Helps reduce latency and improve overall system performance.

Basic Disk Management Diagram:

Detailed Explanation:

1. **Partitioning:**
 - **Purpose:** Allows multiple operating systems or file systems on a single disk.
 - **Example:** A disk can have partitions for Windows, Linux, and a separate data storage partition.
2. **Formatting:**
 - **Purpose:** Initializes the file system on a partition, preparing it for data storage.

- **Example:** Formatting a partition with NTFS for Windows or ext4 for Linux



3. Disk Scheduling:

- **Purpose:** Optimizes the order of disk I/O operations to improve performance.
- **Example:** SSTF algorithm selects the disk I/O request with the shortest seek time next.

4. Defragmentation:

- **Purpose:** Reorganizes fragmented files to ensure data blocks are contiguous.
- **Example:** Windows Disk Defragmenter tool.

5. RAID:

- **Purpose:** Increases performance and/or provides fault tolerance.
- **Example:** RAID 1 mirrors data across two disks for redundancy.

6. Caching:

- **Purpose:** Reduces access time by storing frequently accessed data in faster memory.
- **Example:** Disk cache stores recently read data in RAM.

Workflow Example:

- **Disk Setup:** The disk is partitioned into multiple logical units.
- **Formatting:** Each partition is formatted with a file system.
- **File Operations:** Disk scheduling algorithms optimize the read/write operations.
- **Performance Maintenance:** Defragmentation and caching improve data access speed.
- **Redundancy:** RAID configurations ensure data safety and performance.

File Allocation Methods

File allocation methods are techniques used by operating systems to allocate space for files on a disk. The goal is to efficiently manage disk space while ensuring fast access to files. There are three primary methods: contiguous allocation, linked allocation, and indexed allocation.

1. Contiguous Allocation

In contiguous allocation, each file occupies a set of contiguous blocks on the disk. The file system needs to know the starting block and the length of the file (number of blocks).

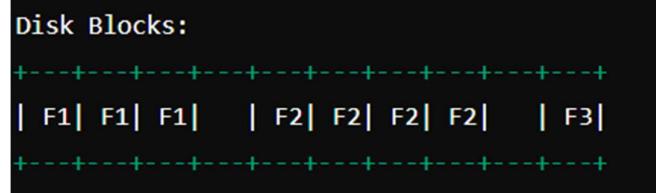
Advantages:

- Simple to implement.
- Fast sequential and direct access.

Disadvantages:

- Leads to external fragmentation.
- Difficult to find contiguous space for large files.

Diagram: Contiguous Allocation



2. Linked Allocation

In linked allocation, each file is a linked list of disk blocks. The directory contains a pointer to the first and last blocks of the file. Each block contains a pointer to the next block.

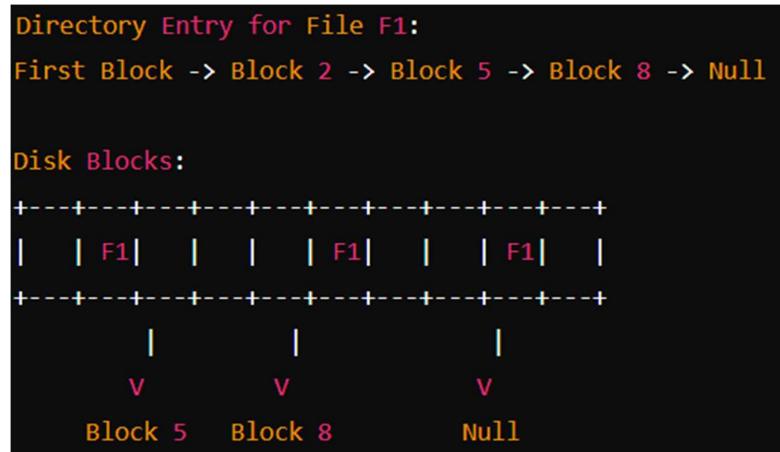
Advantages:

- No external fragmentation.
- Easy to grow files.

Disadvantages:

- Slow direct access.
- Extra space for pointers.
- Reliability issues if pointers are lost or corrupted.

Diagram: Linked Allocation



3. Indexed Allocation

In indexed allocation, each file has an index block which contains pointers to all the blocks occupied by the file. The directory contains the address of the index block.

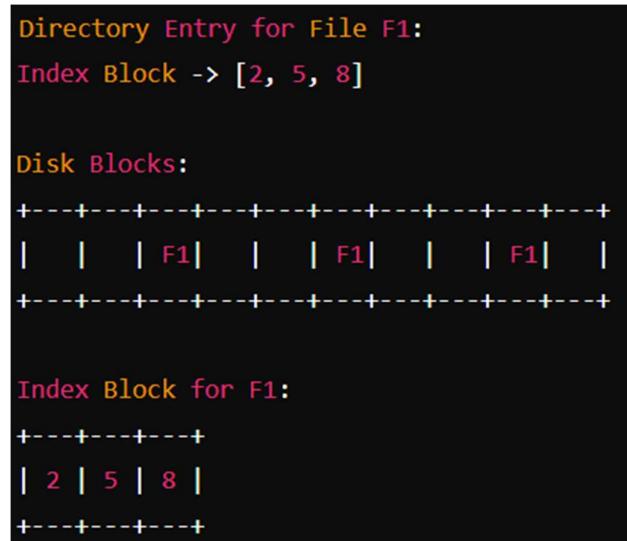
Advantages:

- Fast direct access.
- No external fragmentation.

Disadvantages:

- Overhead of index block.
- Limited file size if index block is too small.

Diagram: Indexed Allocation



Detailed Explanation:

1. Contiguous Allocation:

- **Implementation:** Allocate contiguous blocks for each file.
- **Usage:** Suitable for systems where fast access is crucial, and files do not change size frequently.
- **Example:** Storing a video file which needs fast sequential access.

2. Linked Allocation:

- **Implementation:** Allocate blocks anywhere on the disk, maintaining pointers from one block to the next.
- **Usage:** Suitable for systems with varying file sizes where space efficiency is crucial.

- **Example:** Storing log files which grow over time and do not require fast random access.

3. Indexed Allocation:

- **Implementation:** Use an index block to keep track of all blocks associated with a file.
- **Usage:** Suitable for systems requiring fast access to both sequential and random blocks of data.
- **Example:** Storing database files which require quick access to records.

Workflow Example:

Contiguous Allocation:

- **Step 1:** The OS finds a contiguous space large enough for the file.
- **Step 2:** The file is stored in consecutive blocks.
- **Step 3:** The directory entry records the starting block and length of the file.

Linked Allocation:

- **Step 1:** The OS allocates any free block for the file.
- **Step 2:** Each block contains a pointer to the next block.
- **Step 3:** The directory entry records the first and last blocks of the file.

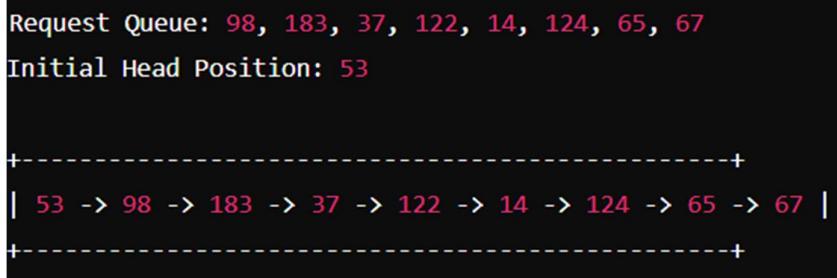
Indexed Allocation:

- **Step 1:** The OS allocates an index block.
- **Step 2:** The index block contains pointers to all other blocks used by the file.
- **Step 3:** The directory entry records the address of the index block

Disk Scheduling Algorithms

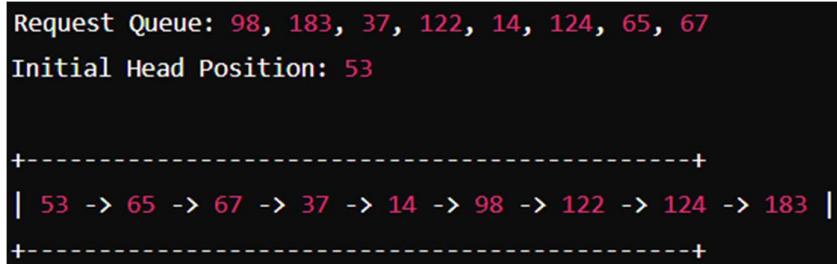
1. First-Come, First-Served (FCFS)

- **Description:** Processes I/O requests in the order they arrive.
- **Advantages:** Simple and fair.
- **Disadvantages:** Can lead to long wait times and poor performance if requests are not optimally ordered.

Diagram: FCFS

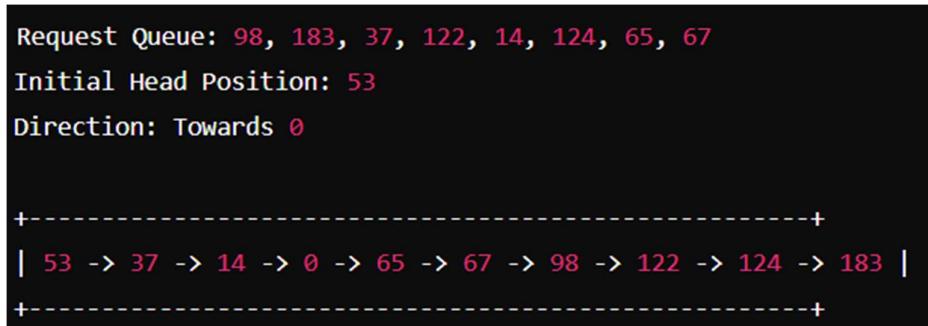
2 . Shortest Seek Time First (SSTF)

- Description:** Selects the I/O request that requires the least movement of the disk arm.
- Advantages:** Reduces seek time compared to FCFS.
- Disadvantages:** Can lead to starvation of some requests.

Diagram: SSTF

3 . SCAN (Elevator Algorithm)

- Description:** The disk arm moves in one direction, servicing requests until it reaches the end, then reverses direction.
- Advantages:** Provides a more uniform wait time compared to FCFS and SSTF.
- Disadvantages:** Can still lead to higher wait times for edge requests.

Diagram: SCAN

4 .Circular SCAN (C-SCAN)

- **Description:** Similar to SCAN, but after reaching the end, the disk arm returns to the beginning and starts again.
- **Advantages:** Provides more uniform wait time and treats all requests more equally.
- **Disadvantages:** Can be less efficient compared to SCAN in certain scenarios.

Diagram: C-SCAN

```
Request Queue: 98, 183, 37, 122, 14, 124, 65, 67
Initial Head Position: 53
Direction: Towards 0

+-----+
| 53 -> 37 -> 14 -> 0 -> 183 -> 124 -> 122 -> 98 -> 67 -> 65 |
+-----+
```

5 .LOOK

- **Description:** Similar to SCAN, but the disk arm only goes as far as the last request in each direction, then reverses.
- **Advantages:** Reduces unnecessary movements compared to SCAN.
- **Disadvantages:** More complex to implement.

Diagram: LOOK

```
Request Queue: 98, 183, 37, 122, 14, 124, 65, 67
Initial Head Position: 53
Direction: Towards 0

+-----+
| 53 -> 37 -> 14 -> 65 -> 67 -> 98 -> 122 -> 124 -> 183 |
+-----+
```

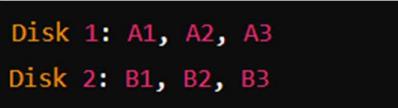
RAID Levels

RAID levels are configurations that combine multiple physical disk drives into a single logical unit for data redundancy and performance improvement.

1. RAID 0 (Striping)

- **Description:** Data is split across multiple disks.
- **Advantages:** High performance.
- **Disadvantages:** No redundancy, data loss if a single disk fails.

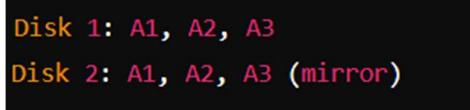
Diagram: RAID 0



2. RAID 1 (Mirroring)

- **Description:** Data is duplicated on two or more disks.
- **Advantages:** High redundancy, simple recovery.
- **Disadvantages:** High cost due to duplication.

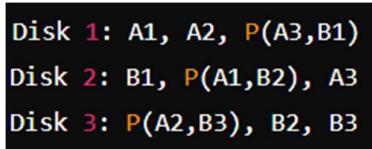
Diagram: RAID 1



3. RAID 5 (Striping with Parity)

- **Description:** Data and parity information are striped across three or more disks.
- **Advantages:** Good balance of performance and redundancy.
- **Disadvantages:** Complex to implement, slower writes due to parity calculation.

Diagram: RAID 5



4. RAID 6 (Striping with Double Parity)

- **Description:** Similar to RAID 5, but with additional parity for extra redundancy.
- **Advantages:** Can tolerate failure of two disks.
- **Disadvantages:** Reduced write performance, higher complexity.

Diagram: RAID 6

```
Disk 1: A1, A2, P(A3,B1), Q(A2,B3)
Disk 2: B1, P(A1,B2), A3, Q(A1,B2)
Disk 3: P(A2,B3), B2, B3, Q(A3,B1)
Disk 4: Q(A1,B2), Q(A3,B1), Q(A2,B3), B4
```

RAID 10 (1+0)

- **Description:** Combines RAID 1 and RAID 0 (mirrored sets in a striped set).
- **Advantages:** High performance and redundancy.
- **Disadvantages:** High cost, complex to implement.

Diagram: RAID 10

```
Stripe Set 1:
Disk 1: A1, A2, A3
Disk 2: A1, A2, A3 (mirror)

Stripe Set 2:
Disk 3: B1, B2, B3
Disk 4: B1, B2, B3 (mirror)
```

Detailed Explanation:

Disk Scheduling Algorithms:

- **FCFS:** Processes requests in the order they arrive, simple but can lead to inefficient disk arm movement.
- **SSTF:** Selects the nearest request, reducing seek time but can cause starvation.
- **SCAN:** Moves in one direction servicing requests, then reverses, providing a more uniform wait time.
- **C-SCAN:** Similar to SCAN but resets to the beginning after reaching the end, ensuring all requests are treated equally.
- **LOOK:** Like SCAN but only goes as far as the last request, reducing unnecessary movements.

RAID Levels:

- **RAID 0:** Stripes data across multiple disks for high performance but lacks redundancy.
- **RAID 1:** Mirrors data for high redundancy, providing simple and fast recovery.
- **RAID 5:** Stripes data with parity across three or more disks, balancing performance and redundancy.
- **RAID 6:** Adds an extra parity block to RAID 5, allowing for two disk failures.
- **RAID 10:** Combines mirroring and striping, offering high performance and redundancy but at a high cost.

Chapter-7

Input /Output Management:

I/O management is a critical function of an operating system that deals with the communication between the system and the outside world, typically through various I/O devices. It involves managing input from devices like keyboards and mice and output to devices like monitors and printers. Effective I/O management ensures smooth and efficient operation of these devices.

I/O Device and Controller

Key Components of I/O Management:

1. I/O Devices:

- Hardware components used for inputting data to or outputting data from the computer.
- Examples: Keyboards, mice, monitors, printers, disk drives.

2. Device Controllers:

- Hardware interfaces that connect I/O devices to the computer and manage the communication between the device and the system.
- Each controller manages a specific type of device (e.g., disk controller, display controller).

3. Device Drivers:

- Software modules that provide the operating system with an interface to communicate with the device controllers.
- Abstracts the specifics of the hardware, allowing the OS to use the device without needing to know its details.

4. I/O Ports:

- Communication endpoints used by the system to send and receive data to and from the I/O devices.
- Includes hardware ports (e.g., USB, serial, parallel) and logical ports (e.g., network ports).

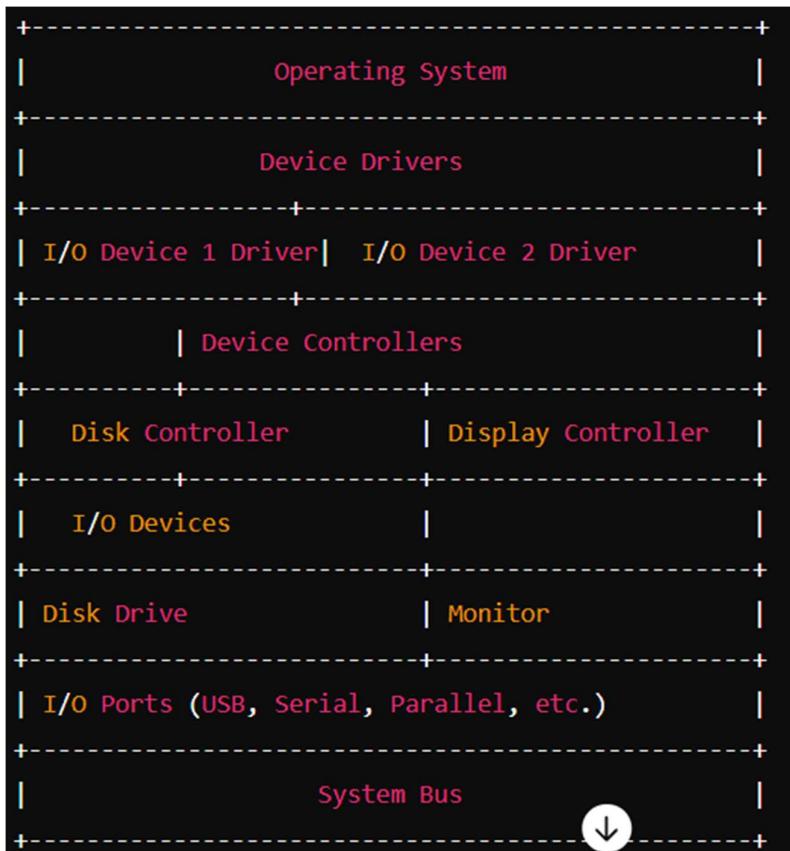
5. Interrupts and Polling:

- Mechanisms for managing I/O operations.
- Interrupts: The device signals the CPU to gain attention.
- Polling: The CPU periodically checks the status of the device.

6. DMA (Direct Memory Access):

- Allows devices to transfer data directly to/from memory without involving the CPU, enhancing performance.

Basic I/O Management Diagram:



Detailed Explanation:

1. I/O Devices:

- **Role:** Facilitate user interaction with the computer and provide data input/output.
- **Examples:** Keyboards (input), monitors (output), printers (output), disk drives (input/output).

2. Device Controllers:

- **Role:** Act as intermediaries between the I/O devices and the system.
- **Examples:** Disk controllers manage hard drives, display controllers manage monitors.

3. Device Drivers:

- **Role:** Provide a standard interface for the OS to interact with hardware.
- **Function:** Translate OS commands into device-specific operations.

4. I/O Ports:

- **Role:** Serve as communication gateways for data transfer.

- **Examples:** USB ports connect peripherals like mice and keyboards, network ports handle internet connections.

5. Interrupts and Polling:

- **Interrupts:** The device sends an interrupt signal to the CPU when it needs attention, allowing the CPU to perform other tasks until needed.
- **Polling:** The CPU regularly checks the status of a device to see if it needs attention.

6. DMA (Direct Memory Access):

- **Role:** Improves system efficiency by allowing data transfer directly between devices and memory.
- **Function:** Bypasses the CPU, freeing it to perform other tasks.

Workflow Example:

- **Input:** A user presses a key on the keyboard.
 - The keyboard controller detects the keypress and sends a signal to the CPU via the system bus.
 - The keyboard driver translates the signal into a format the OS can understand.
 - The OS processes the input and takes the appropriate action.
- **Output:** The system needs to display data on the monitor.
 - The OS sends data to the display driver.
 - The display driver translates the data into commands understood by the display controller.
 - The display controller updates the monitor with the new data.
- **Data Transfer:** A file is copied from a disk drive to memory.
 - The OS initiates the data transfer request.
 - The disk controller handles the transfer using DMA, moving data directly from the disk to memory without involving the CPU.

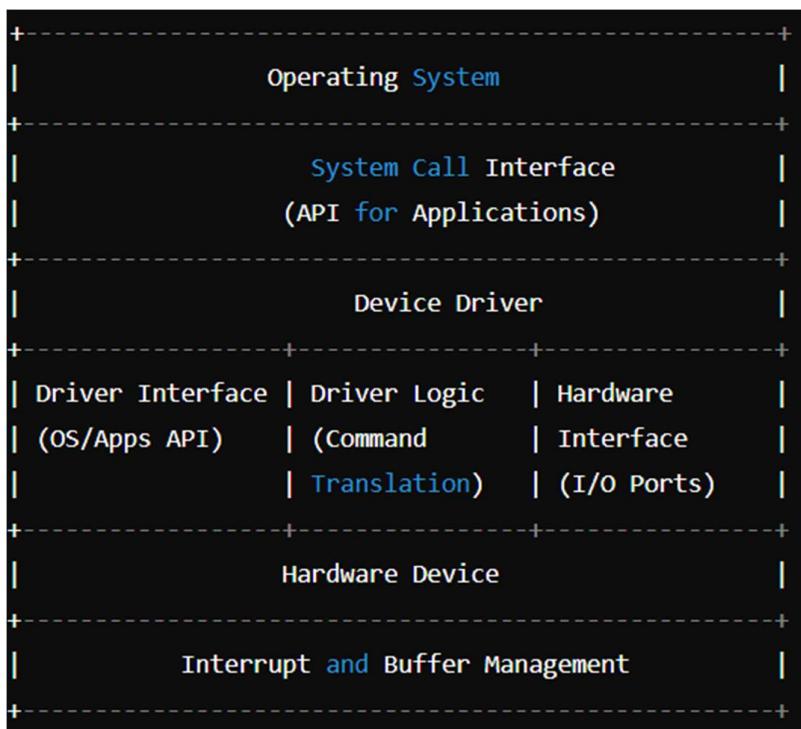
Device Drivers:

Device drivers are specialized software components that allow the operating system and applications to interact with hardware devices. They act as translators, converting high-level commands from the OS into low-level commands that the hardware can understand, and vice versa.

Key Components of Device Drivers:

1. **Driver Interface:** The API that allows the OS and applications to communicate with the driver.
2. **Driver Logic:** The core code that translates generic commands from the OS into specific commands for the hardware device.
3. **Hardware Interface:** The part of the driver that interacts directly with the device through I/O ports or memory-mapped I/O.
4. **Interrupt Handling:** Manages hardware interrupts, which are signals sent by devices to get the CPU's attention.
5. **Buffer Management:** Manages data buffers used for transferring data between the OS and the device.
6. **Configuration and Initialization:** Handles device configuration and initializes the device when it is first connected or powered on.

Basic Device Driver Diagram:



Detailed Explanation:

1. **Driver Interface:**
 - **Role:** Provides a standardized API for the OS and applications to interact with the hardware.
 - **Example:** Functions to read from or write to a device.
2. **Driver Logic:**
 - **Role:** Contains the logic to translate high-level commands from the OS into device-specific commands.

- **Example:** Converting a file write command into specific instructions for a printer.

3. Hardware Interface:

- **Role:** Communicates directly with the hardware device through I/O ports or memory-mapped I/O.
- **Example:** Sending control signals to a disk drive to initiate a read operation.

4. Interrupt Handling:

- **Role:** Manages interrupts sent by devices to signal events (e.g., completion of an I/O operation).
- **Example:** An interrupt from a keyboard indicating a key press.

5. Buffer Management:

- **Role:** Handles data buffering to ensure smooth data transfer between the OS and the device.
- **Example:** Using a buffer to store data temporarily while reading from a disk.

6. Configuration and Initialization:

- **Role:** Configures device settings and initializes the device when connected or powered on.
- **Example:** Setting up communication parameters for a network card.

Workflow Example:

1. Device Initialization:

- When a device is connected, the OS loads the appropriate driver.
- The driver initializes the device, setting necessary parameters and preparing it for use.

2. Data Transfer:

- An application requests data transfer (e.g., reading a file).
- The OS calls the corresponding function in the device driver.
- The driver translates the command and communicates with the hardware to perform the operation.
- Data is transferred through buffers managed by the driver.

3. Interrupt Handling:

- The device sends an interrupt signal to indicate completion of an operation.
- The driver handles the interrupt, processes the data, and notifies the OS.

Example of Device Driver Interaction:

- **Printing a Document:**
 - **Application:** Sends a print command.
 - **OS:** Calls the print function in the printer driver.
 - **Driver Logic:** Converts the print command into printer-specific instructions.
 - **Hardware Interface:** Sends instructions to the printer via I/O ports.
 - **Printer:** Executes the print command and sends an interrupt on completion.
 - **Driver:** Handles the interrupt and informs the OS that the print job is complete.

Kernel I/O Subsystem

The Input/Output (I/O) management subsystem in an operating system (OS) is a critical component responsible for managing communication between the computer's hardware and its software. The kernel, as the core part of the OS, plays a vital role in this management by providing an interface for I/O operations, ensuring efficient data transfer, and maintaining system stability and performance.

Components of the Kernel I/O Subsystem

1. Device Drivers

- **Function:** Device drivers act as intermediaries between the OS and hardware devices. Each device driver is tailored to a specific type of hardware, translating OS commands into device-specific operations.
- **Role in I/O:** They handle the specifics of data transfer, interrupt handling, and device control, enabling the kernel to perform I/O operations without needing to know the intricate details of the hardware.

2. Buffering

- **Function:** Buffers are temporary storage areas used to hold data while it is being transferred between the OS and a device.
- **Role in I/O:** Buffering helps to accommodate the speed differences between devices and the CPU, preventing data loss and improving efficiency by allowing the system to continue processing while waiting for I/O operations to complete.

3. Caching

- **Function:** Caching involves storing frequently accessed data in faster storage (like RAM) to reduce access times.
- **Role in I/O:** It improves performance by reducing the need to repeatedly access slower storage devices, thus speeding up read and write operations.

4. Spooling

- **Function:** Spooling is a mechanism for managing I/O operations by queueing data that needs to be processed.
- **Role in I/O:** Commonly used for devices like printers, spooling allows the system to queue multiple print jobs and process them sequentially, ensuring efficient resource usage and preventing conflicts.

5. I/O Scheduling

- **Function:** I/O scheduling determines the order in which I/O requests are processed.
- **Role in I/O:** Effective scheduling algorithms enhance system performance by optimizing the sequence of I/O operations, reducing wait times, and ensuring fair resource allocation among processes.

Kernel I/O Subsystem Operations

1. System Calls

- System calls are the interface through which user applications interact with the kernel. For I/O operations, common system calls include `read()`, `write()`, `open()`, and `close()`.
- These calls are abstracted to hide the complexity of hardware operations, providing a simple and consistent interface for application developers.

2. Interrupt Handling

- Interrupts are signals sent by hardware devices to alert the CPU of events like data arrival or completion of an I/O operation.
- The kernel's interrupt handler prioritizes and processes these signals, ensuring timely and efficient handling of I/O operations.

3. Direct Memory Access (DMA)

- DMA allows certain hardware subsystems to access main system memory independently of the CPU, enabling high-speed data transfers.
- The kernel manages DMA operations to ensure data integrity and coordinate between the CPU and peripheral devices.

Efficiency and Performance Considerations

1. Throughput and Latency

- Throughput refers to the amount of data processed in a given time, while latency is the time taken to complete an I/O operation.
- The kernel aims to maximize throughput and minimize latency by optimizing I/O scheduling, buffering, and caching strategies.

2. Concurrency and Parallelism

- Modern systems often handle multiple I/O operations concurrently. The kernel manages this through techniques like multi-threading and asynchronous I/O, ensuring efficient parallel processing and resource utilization.

3. Error Handling and Recovery

- The kernel implements robust error detection and handling mechanisms to manage I/O errors, such as bad sectors on a disk or communication failures with devices.
- Recovery strategies include retries, error correction codes, and failover mechanisms to maintain system reliability.

I/O Techniques :

Effective I/O management is essential for the smooth operation of a computer system. There are several techniques used to manage I/O operations: polling, interrupts, and Direct Memory Access (DMA). Each technique has its own advantages and is suitable for different scenarios.

1. Polling

Polling is a technique where the CPU repeatedly checks the status of an I/O device to determine if it is ready for data transfer. This method is straightforward but can be inefficient because it keeps the CPU busy.

2. Interrupts

Interrupts allow an I/O device to signal the CPU when it needs attention, freeing the CPU to perform other tasks until the device is ready. This method is more efficient than polling as it reduces CPU idle time.

3. Direct Memory Access (DMA)

DMA allows an I/O device to transfer data directly to/from memory without involving the CPU, significantly speeding up data transfer rates and freeing the CPU for other tasks.

Diagram: Direct Memory Access (DMA)

Diagram: Polling

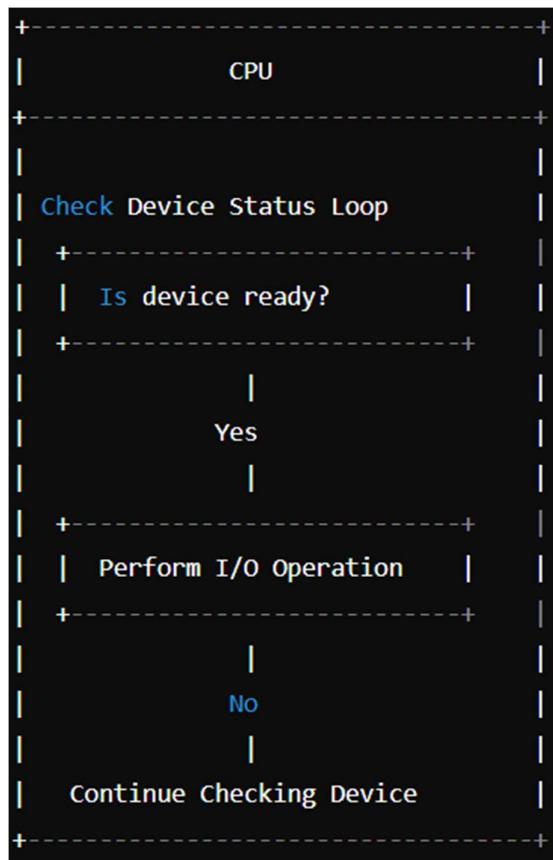
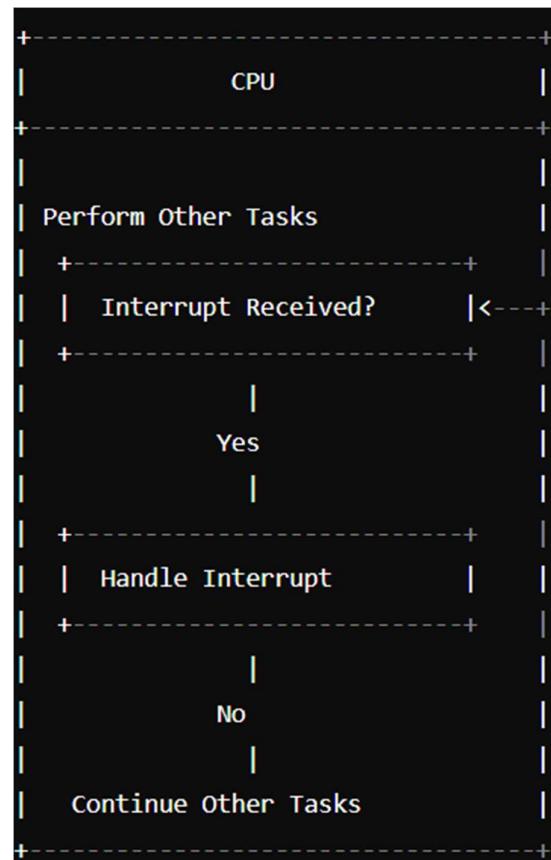
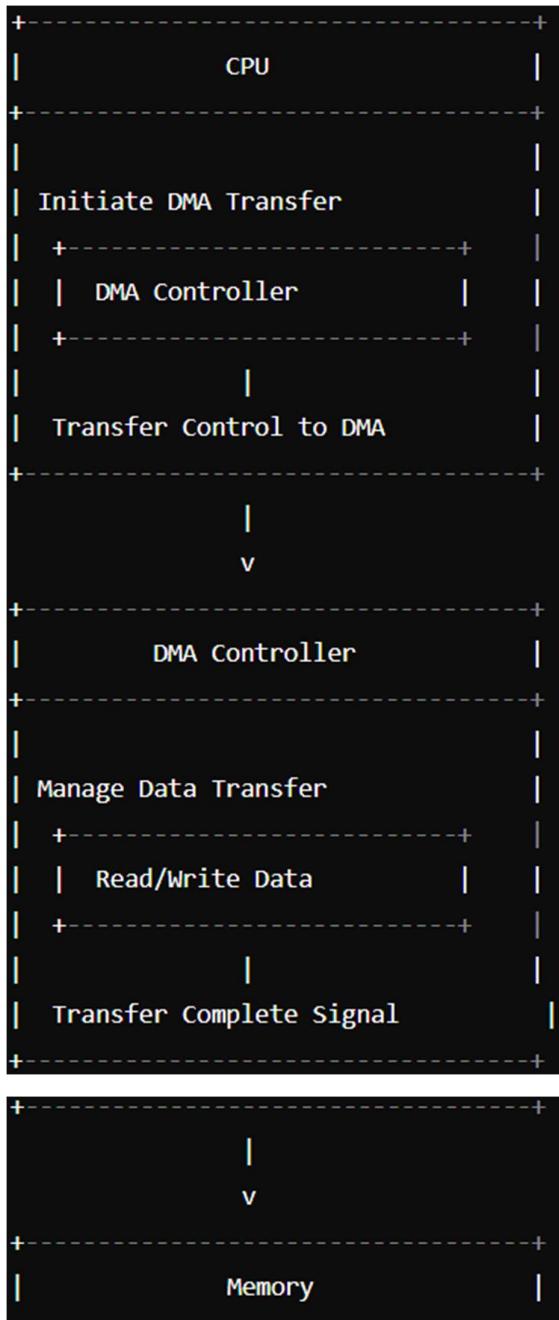


Diagram: Interrupts





I/O scheduling:

I/O scheduling refer **chapter 6 Disk Scheduling Algorithms Page No - 54, 55, 56**

Chapter-8

Linux:

Introduction to Linux

Linux is a powerful, open-source operating system based on the Unix architecture. Created by Linus Torvalds in 1991, it has grown to become one of the most widely used operating systems in the world, powering everything from personal computers and servers to smartphones and embedded devices. Linux is known for its stability, security, and flexibility.

Linux History

Linux looks and feels much like any other UNIX system; indeed, UNIX compatibility has been a major design goal of the Linux project. However, Linux is much younger than most UNIX systems. Its development began in 1991, when a Finnish university student, Linus Torvalds, began creating a small but self-contained kernel for the 80386 processor, the first true 32-bit processor in Intel's range of PC-compatible CPUs.

Early in its development, the Linux source code was made available free—both at no cost and with minimal distributional restrictions—on the Internet. As a result, Linux's history has been one of collaboration by many developers from all around the world, corresponding almost exclusively over the Internet. From an initial kernel that partially implemented a small subset of the UNIX system services, the Linux system has grown to include all of the functionality expected of a modern UNIX system. In its early days, Linux development revolved largely around the central operating-system kernel—the core, privileged executive that manages all system resources and interacts directly with the computer hardware. We need much more than this kernel, of course, to produce a full operating system.

We thus need to make a distinction between the Linux kernel and a complete Linux system. The **Linux kernel** is an original piece of software developed from scratch by the Linux community. The **Linux system**, as we know it today, includes a multitude of components, some written from scratch, others borrowed from other development projects, and still others created in collaboration with other teams.

The basic Linux system is a standard environment for applications and user programming, but it does not enforce any standard means of managing the available functionality as a whole. As Linux has matured,

a need has arisen for another layer of functionality on top of the Linux system. This need has been met by various Linux distributions. A **Linux distribution** includes all the standard components of the Linux system, plus a set of administrative tools to simplify the initial installation and subsequent upgrading of Linux and to manage installation and removal of other packages on the system. A modern distribution also typically includes tools for management of file systems, creation and management of user accounts, administration of networks, web browsers, word processors, and so on.

Key Features of Linux

1. **Open Source:** Linux is free to use, modify, and distribute.
2. **Multitasking:** Linux efficiently handles multiple tasks at the same time.
3. **Multiuser Capability:** Multiple users can access system resources simultaneously without interfering with each other.
4. **Security:** Linux provides robust security features including file permissions, encryption, and user authentication.
5. **Portability:** Linux can run on various hardware platforms.
6. **Customizability:** Users can modify and customize the system to meet their needs, thanks to its open-source nature.
7. **Community Support:** A vast and active community of developers and users who contribute to its development and provide support.

Linux Architecture

Linux follows a modular architecture that is structured in layers, with the kernel being the core component. The key components of Linux architecture include:

1. **Kernel:** The core part of Linux, responsible for managing system resources, hardware devices, memory, and processes.
2. **System Libraries:** Essential libraries that provide functions and interfaces for application development and system utilities.
3. **System Utilities:** Basic commands and utilities for file manipulation, system monitoring, and configuration.
4. **Shell:** A command-line interface that allows users to interact with the kernel and execute commands.
5. **Application Programs:** User-level programs that perform various tasks, ranging from simple text editing to complex software development.

Diagram: Linux Architecture



Detailed Explanation of Linux Components

1. Kernel:

- **Description:** The heart of the Linux operating system. It manages hardware resources and provides essential services to other parts of the system.
- **Responsibilities:** Process management, memory management, device management, system calls, and security.
- **Types:** Monolithic kernel (as in Linux), microkernel, hybrid kernel.

2. System Libraries:

- **Description:** Libraries of functions and routines that applications use to interact with the kernel and perform various tasks.
- **Examples:** GNU C Library (glibc), which provides standard C library functions.

3. System Utilities:

- **Description:** Programs and tools that provide system functionality. These can be simple commands or more complex software.
- **Examples:** ls, cp, mv, top, ps.

4. Shell:

- **Description:** A command interpreter that allows users to execute commands and scripts.
- **Types:** Various shells like Bash (Bourne Again Shell), Zsh (Z Shell), Ksh (Korn Shell).

- **Functionality:** Command execution, scripting, environment customization.
5. **Application Programs:**
- **Description:** Software that performs user-specific tasks and functions.
 - **Examples:** Web browsers (Firefox, Chrome), office suites (LibreOffice), development tools (GCC, Python).

Benefits of Using Linux

1. **Cost:** Linux is free, reducing software costs significantly.
2. **Flexibility:** Can be used on a variety of devices, from servers to desktops to embedded systems.
3. **Security:** Known for its strong security features, making it a preferred choice for servers and sensitive applications.
4. **Performance:** Efficient and can handle high loads, making it ideal for servers and high-performance computing.
5. **Community:** A large, active community provides extensive support and continuous development.

Use Cases of Linux

1. **Servers:** Most web servers and enterprise servers run on Linux due to its reliability and scalability.
2. **Development:** Preferred by developers for its powerful tools and flexibility.
3. **Embedded Systems:** Used in devices like routers, smart TVs, and IoT devices.
4. **Desktops:** Popular among enthusiasts and professionals who value customization and control.
5. **Cloud Computing:** Backbone of many cloud infrastructures like AWS, Google Cloud, and Microsoft Azure.

Chapter-9

Shell programming:

Introduction

Shell programming involves writing scripts for the command-line interpreter, known as the shell. These scripts automate tasks, manipulate files, execute commands, and control the flow of programs on Unix and Unix-like operating systems, including Linux.

What is a Shell?

A shell is a command-line interface that allows users to interact with the operating system by typing commands. It acts as an intermediary between the user and the kernel. There are several types of shells, with the most common being:

- **Bourne Shell (sh)**
- **Bourne Again Shell (bash)**
- **C Shell (csh)**
- **Korn Shell (ksh)**
- **Z Shell (zsh)**

Key Concepts in Shell Programming

1. **Shell Script:** A file containing a series of commands to be executed by the shell.
2. **Shebang (# !):** A special character sequence at the beginning of a script that specifies the interpreter to be used.
3. **Variables:** Used to store data that can be used and manipulated within the script.
4. **Control Structures:** Conditional statements and loops to control the flow of execution.
5. **Functions:** Blocks of code that can be reused multiple times within the script.

Writing a Shell Script

A basic shell script starts with a shebang followed by the path to the shell. Commands are written in sequence, and the script is made executable using the chmod command.

Example: Hello World Script

```
#!/bin/bash
echo "Hello, World!"
```

Variables

Variables in shell scripts are used to store values, which can be strings, numbers, or the output of commands.

Example: Using Variables

```
#!/bin/bash
name="John"
echo "Hello, $name!"
```

Control Structures

Control structures allow you to make decisions and repeat actions. The most common control structures are if statements and loops (for, while, until).

Example: if Statement

```
#!/bin/bash
if [ "$name" == "John" ]; then
    echo "Hello, John!"
else
    echo "Hello, stranger!"
fi
```

Example: for Loop

```
#!/bin/bash
for i in 1 2 3 4 5
do
    echo "Number: $i"
done
```

Functions

Functions allow you to group commands and reuse them. They can take arguments and return values.

Example: Function

```
#!/bin/bash
greet() {
    echo "Hello, $1!"
}

greet "John"
```

Common Shell Commands

1. **File Operations:** cp, mv, rm, touch, mkdir
2. **Text Processing:** cat, echo, grep, sed, awk
3. **System Information:** uname, df, top, ps, uptime
4. **Networking:** ping, curl, wget, ifconfig, netstat

Input and Output Redirection

Shell scripts can redirect input and output using symbols like >, >>, <, 2>, and |.

Example: Redirecting Output

```
#!/bin/bash
echo "This is a test" > output.txt
```

Example: Piping

```
#!/bin/bash
ls -1 | grep "txt"
```

Shell Script Example: Backup Script

Here's a simple shell script to back up a directory:

```

#!/bin/bash
SOURCE_DIR="/path/to/source"
BACKUP_DIR="/path/to/backup"

echo "Starting backup from $SOURCE_DIR to $BACKUP_DIR"

cp -r $SOURCE_DIR $BACKUP_DIR

if [ $? -eq 0 ]; then
    echo "Backup successful!"
else
    echo "Backup failed!"
fi

```



Best Practices in Shell Programming

- Use Comments:** Explain the purpose of the script and complex sections with comments.
- Error Handling:** Check the exit status of commands and handle errors gracefully.
- Modularity:** Use functions to break down the script into manageable and reusable sections.
- Portability:** Write scripts that are portable across different Unix-like systems.
- Security:** Validate inputs and avoid using hard-coded sensitive information.

Common Commands in Shell Programming

Shell programming leverages a variety of commands that allow users to interact with the file system, manipulate text, manage processes, and perform other system-related tasks. Below are some of the most commonly used commands in shell programming, categorized by their functionality.

File and Directory Operations

- ls:** List directory contents.

```
ls -l
```

- cd:** Change the current directory.

```
cd /path/to/directory
```

3. **pwd**: Print the current working directory.

```
pwd
```

4. **cp**: Copy files or directories.

```
cp source_file destination_file
```

5. **mv**: Move or rename files or directories.

```
mv old_name new_name
```

6. **rm**: Remove files or directories.

```
rm file_name
```

7. **mkdir**: Create a new directory.

```
mkdir new_directory
```

8. **rmdir**: Remove an empty directory.

```
rmdir empty_directory
```

Text Processing:

1. `echo`: Display a line of text.

```
sh
echo "Hello, World!"
```

2. `cat`: Concatenate and display file content.

```
sh
cat file_name
```

3. `grep`: Search for patterns in text.

```
sh
grep "pattern" file_name
```

4. `sed`: Stream editor for filtering and transforming text.

```
sh
↓
sed 's/old/new/g' file_name
```

5. `awk`: Pattern scanning and processing language.

```
sh
awk '{print $1}' file_name
```

6. `cut`: Remove sections from each line of files.

```
sh
cut -d":" -f1 /etc/passwd
```

7. `sort`: Sort lines of text files.

```
sh
sort file_name
```

8. `uniq`: Report or omit repeated lines.

```
sh  
uniq file_name
```

System Information and Management:

1. `uname`: Print system information.

```
sh  
uname -a
```

2. `df`: Report file system disk space usage.

```
sh  
df -h
```

3. `du`: Estimate file space usage.

```
sh  
du -sh directory_name
```

4. `top`: Display Linux tasks.

```
sh  
top
```

5. `ps`: Report a snapshot of current processes.

```
sh  
ps aux
```

6. `kill`: Terminate processes.

```
sh  
kill PID
```

7. `uptime`: Tell how long the system has been running.

```
sh  
uptime
```

8. `who`: Show who is logged on.

```
sh  
who
```



Networking:

1. `ping`: Send ICMP ECHO_REQUEST to network hosts.

```
sh  
ping example.com
```

2. `curl`: Transfer data from or to a server.

```
sh  
curl -o http://example.com/file
```

3. `wget`: Non-interactive network downloader.

```
sh  
wget http://example.com/file
```



4. `ifconfig`: Configure network interfaces.

```
sh  
ifconfig
```

5. `netstat`: Print network connections, routing tables, interface statistics, masquerading connections, and multicast memberships.

```
sh  
netstat -tuln
```

Process Control:

1. `&`: Run a command in the background.

```
sh  
command &
```

2. `jobs`: List active jobs.

```
sh  
jobs
```

3. `fg`: Bring a background job to the foreground.

```
sh  
fg %1
```

4. `bg`: Resume a suspended job in the background.

```
sh  
bg %1
```

5. `nohup`: Run a command immune to hangups.

```
sh  
nohup command &
```

Redirection and Pipelines:

1. Output Redirection:

- Overwrite: `>`.

```
sh  
echo "Hello, World!" > file.txt
```

- Append: `>>`.

```
sh  
echo "Hello again!" >> file.txt
```

2. Input Redirection:

```
sh  
wc -l < file.txt
```

3. Pipelines:

```
sh  
ls -l | grep "pattern"
```

Shell Script Example:

Here's an example of a shell script using some of the commands mentioned:

```
#!/bin/bash

# Print the current directory
echo "Current directory: $(pwd)"

# List files in the current directory
echo "Files in the directory:"
ls -l

# Create a new directory
mkdir my_new_directory
echo "Created my_new_directory"

# Copy a file
cp file1.txt my_new_directory/
echo "Copied file1.txt to my_new_directory"

# Search for a pattern in a file
grep "search_pattern" file1.txt
```

```
# Display system information
uname -a

# Show disk usage
df -h

# Print the current users
who
```