# Processes

OPERATING SYSTEMS

Sercan Külcü | Operating Systems | 16.04.2023

# Contents

# Chapter 3:
# Processes

## 1 Introduction

In this section, we will discuss the fundamental concepts of processes and their significance in an operating system. We will also provide an overview of the functions of processes in multi-tasking and concurrency.

A process is defined as a program in execution. It is a fundamental concept in an operating system that enables the system to execute multiple tasks concurrently. Processes are essential for the efficient use of system resources and play a crucial role in managing the system's overall performance.

In this section, we will explore the importance of processes in an operating system and provide a brief overview of the functions of processes in multi-tasking and concurrency. So, let's dive in and understand the critical role that processes play in an operating system!

### 1.1 Definition of a process

In an operating system, a process is defined as an instance of a program in execution. A process is a fundamental concept in operating systems, and it is responsible for executing instructions, allocating and releasing resources, and communicating with other processes. A process has its own memory space, execution context, and system resources that it uses to accomplish its tasks.

A process is created when a program is loaded into memory, and it is terminated when the program completes its execution or when it is terminated by the operating system. While a process is running, it may spawn child processes, communicate with other processes, and perform various operations on system resources.

Processes are essential components of an operating system, and they provide the foundation for the system's functionality. Without processes, an operating system would not be able to execute applications, manage system resources, or provide a platform for multi-tasking and concurrency.

In the following sections, we will explore the importance of processes in an operating system and provide an overview of their functions in multi-tasking and concurrency. We will also discuss the various attributes that define a process and the mechanisms that an operating system uses to manage processes.

## 1.2 Importance of processes in an operating system

Processes are fundamental components of any operating system. They are essential for the efficient and effective execution of tasks, which makes them critical to the performance and functionality of an operating system. In this chapter, we will discuss the importance of processes in an operating system.

### 1.2.1 Resource Allocation and Management

Processes play a vital role in resource allocation and management. A process can request resources such as memory, CPU time, and I/O devices. The operating system is responsible for allocating these resources to the requesting process. Each process is allocated a specific amount of resources, which helps to ensure that all processes receive a

fair share of resources. Proper resource allocation and management are essential to maintain the stability and reliability of the operating system.

### 1.2.2 Multitasking and Concurrency

The ability to run multiple processes simultaneously is known as multitasking. In a multitasking environment, processes share the CPU time, and the operating system manages the allocation of CPU time to each process. The operating system uses scheduling algorithms to ensure that all processes get a fair share of CPU time. This feature allows users to run multiple applications and perform multiple tasks simultaneously, enhancing the efficiency and productivity of the system.

Concurrency refers to the ability of a system to execute multiple tasks simultaneously. Processes enable concurrency by allowing multiple applications to run concurrently. This feature enables users to perform multiple tasks simultaneously, making the system more efficient and productive.

### 1.2.3 Security and Protection

Processes play a crucial role in maintaining the security and protection of the operating system. Each process runs in its address space, which prevents it from accessing the memory of other processes. This feature ensures that one process cannot interfere with the execution of another process. Additionally, the operating system assigns specific privileges and permissions to each process, which helps to ensure that processes only access the resources they are authorized to access.

### 1.2.4 Fault Isolation and Recovery

Processes provide fault isolation and recovery capabilities. Each process runs independently, which means that if one process fails, it does not affect the execution of other processes. The operating system can

terminate a faulty process without affecting the other running processes. This feature helps to maintain the stability and reliability of the system.

In conclusion, processes are critical components of any operating system. They play a crucial role in resource allocation and management, multitasking and concurrency, security and protection, and fault isolation and recovery. The efficient execution of tasks is only possible because of the underlying processes that run in the background. Therefore, it is important to understand the importance of processes in an operating system to maintain the stability, reliability, and efficiency of the system.

## 1.3 Overview of the functions of processes in multi-tasking and concurrency

In a modern operating system, it is common to have multiple applications running concurrently. This means that the operating system needs to manage the execution of multiple processes and ensure that they can coexist and operate without interfering with each other. This is where multi-tasking and concurrency come into play.

A process is an executing program with its own memory space, set of resources, and state. Each process has a unique identifier and can interact with other processes through inter-process communication mechanisms. In a multi-tasking environment, the operating system can manage the execution of multiple processes simultaneously.

The primary function of processes in multi-tasking environments is to allow multiple applications to execute concurrently. This is achieved by the operating system allocating time slices to each process, allowing them to execute for a specified period before suspending execution and allowing other processes to execute. This is known as time-sharing, and it enables the operating system to make the most efficient use of system resources.

Another function of processes in multi-tasking environments is to provide isolation and protection between applications. Each process has its own memory space, which means that it cannot access the memory of another process without explicit permission. This provides a level of security and protection between applications.

Concurrency is the ability of an operating system to manage multiple processes that execute simultaneously. This is achieved by the operating system dividing the system resources among the processes, allowing them to execute concurrently. The operating system provides mechanisms to ensure that the processes do not interfere with each other, and this is achieved through synchronization mechanisms such as semaphores and mutexes.

Processes can communicate with each other using inter-process communication mechanisms. These mechanisms allow processes to exchange data and synchronize their actions. This is an essential function of processes in multi-tasking and concurrent environments, as it enables different applications to work together and share resources.

In summary, processes are essential components of modern operating systems. They enable multi-tasking and concurrency, which allows multiple applications to execute concurrently, and they provide isolation and protection between applications. Inter-process communication mechanisms allow processes to communicate and synchronize their actions, enabling different applications to work together and share resources.

## 2   Process States and Transitions

A process is a fundamental concept in an operating system, and it is essential to understand the different states a process can be in and how it transitions between these states. The concept of process states is

crucial in multi-tasking and concurrency because it allows the operating system to manage and prioritize the execution of processes.

In this chapter, we will cover the different process states, including new, ready, running, blocked, and terminated. We will also discuss how a process transitions between these states and the importance of understanding process states in a multi-tasking and concurrent environment.

## 2.1  Process states:

Processes are the core of an operating system, and they have various states throughout their execution. In this chapter, we will discuss the different process states, which are new, ready, running, blocked, and terminated.

### 2.1.1  New State:

When a process is created, it is in the new state. At this point, the process is just an idea or a request, and the operating system has not yet allocated resources to it. Once the operating system assigns resources to the process, it moves to the next state.

### 2.1.2  Ready State:

When a process has been assigned resources, it moves to the ready state. In this state, the process is waiting to be allocated a processor by the operating system. The process is ready to execute, but the operating system has not yet assigned a processor to it.

### 2.1.3  Running State:

When the operating system assigns a processor to a process, it moves to the running state. In this state, the process is executing its instructions on the assigned processor.

### 2.1.4  Blocked State:

When a process is waiting for an event, such as I/O or a system resource, it moves to the blocked state. In this state, the process cannot execute until the event it is waiting for occurs. Once the event occurs, the process moves back to the ready state.

### 2.1.5  Terminated State:

When a process completes its execution, it moves to the terminated state. In this state, the process is no longer executing, and its resources have been deallocated by the operating system.

Understanding the different process states is crucial to the efficient operation of an operating system. The ability to manage and manipulate process states is a critical component of any operating system, allowing the operating system to allocate resources effectively and provide a seamless user experience.

## 2.2 Transitions between process states

Transitions between process states are crucial to the functioning of a modern operating system. A process can move between different states during its lifetime, depending on the type of operation it is currently performing. Understanding these transitions is key to understanding how a multi-tasking and concurrent system works.

When a process is first created, it enters the "new" state. In this state, the operating system has allocated resources such as memory and process control blocks to the process, but it is not yet ready to execute. From here, the process may transition to the "ready" state, where it is waiting for the CPU to be assigned to it.

Once the CPU is assigned to the process, it enters the "running" state. In this state, the process is actively executing its instructions. However, at any point, the process may be interrupted and transition back to the "ready" state. This can happen, for example, if the operating system needs to switch to another process that has become ready to execute.

A process can also transition to the "blocked" state, which occurs when the process is waiting for some external event, such as user input or data from a file. In this state, the process is not executing any instructions and is temporarily suspended.

Finally, a process may transition to the "terminated" state when it has completed its execution. In this state, the operating system releases any resources that were allocated to the process.

Understanding the transitions between process states is important for building a robust and efficient operating system. It enables the operating system to prioritize processes based on their state, ensuring that processes that are ready to execute are given access to the CPU. It also allows the operating system to manage resources effectively, by releasing resources when they are no longer needed.

## 2.3 Importance of process states in multi-tasking and concurrency

In a modern operating system, the ability to manage multiple processes simultaneously is a critical feature. This feature enables the system to

be more efficient and responsive, as well as providing the ability to execute multiple programs simultaneously.

The management of multiple processes is a complex task, requiring careful attention to detail and the use of sophisticated algorithms. One of the key components of process management is the concept of process states. The state of a process is a reflection of its current activity level and can be used to control its behavior.

There are five process states: new, ready, running, blocked, and terminated. Each of these states plays an important role in the management of processes and is used to control how processes are scheduled for execution.

When a process is first created, it is in the new state. In this state, the process has been created but has not yet been assigned any resources or executed by the system. Once the process has been assigned the necessary resources and is ready to be executed, it enters the ready state.

In the ready state, the process is waiting for its turn to be executed by the system. The operating system uses a scheduling algorithm to determine which process to execute next from the pool of ready processes.

Once a process has been selected for execution, it enters the running state. In this state, the process is actively executing its code and using system resources.

The blocked state is used to indicate that a process is waiting for some external event to occur before it can continue executing. For example, a process may be blocked while waiting for user input or while waiting for a file to be loaded from disk.

Finally, when a process has completed its execution or has been terminated by the system, it enters the terminated state. In this state, the process is no longer executing and its resources have been freed by the system.

The importance of process states in multi-tasking and concurrency cannot be overstated. These states provide a mechanism for the operating system to control how processes are executed and to ensure that system resources are used efficiently.

By carefully managing the state of each process, the system can ensure that all processes are executed fairly and that no single process monopolizes the system's resources. Additionally, the use of process states makes it possible to manage complex multi-tasking and concurrency scenarios, allowing the system to execute multiple processes simultaneously without conflicts or other issues.

In conclusion, the management of process states is a critical aspect of modern operating systems. By providing a mechanism for controlling the behavior of processes and managing system resources, process states make it possible to execute multiple processes simultaneously and ensure that the system is both efficient and responsive.

## 3  Process Control Block (PCB)

Welcome to the chapter on Process Control Block (PCB). In an operating system, managing processes is a critical task, and one of the key components used for this purpose is the Process Control Block. The Process Control Block is a data structure that contains essential information about a process, and it serves as a central point for the operating system to manage and control the process. In this chapter, we will cover the definition of a PCB, its contents, and its importance in process management. Understanding PCB is essential for anyone studying operating systems and its processes. So, let's dive into the world of PCBs and learn how they help in managing processes in an operating system.

In modern operating systems, process management is an essential component to ensure the efficient and effective execution of processes.

One of the key structures used in process management is the Process Control Block (PCB). In this chapter, we will discuss what a PCB is, its functions, and how it is used in process management.

## 3.1  Definition of a PCB

A Process Control Block (PCB) is a data structure used by the operating system to manage information about a process. It is also known as a task control block or a process descriptor. A PCB is created by the operating system when a process is created and is responsible for keeping track of important information about the process, such as its current state, register values, memory allocation, and other attributes that are necessary for process execution.

## 3.2  Contents of a PCB

The contents of a PCB may vary depending on the operating system, but it generally contains the following information:

Process state: The current state of the process, which can be running, ready, blocked, or terminated.

Process ID: A unique identifier assigned by the operating system to each process.

Program counter: A register that holds the address of the next instruction to be executed.

CPU registers: The values of the CPU registers that are associated with the process, such as the accumulator, stack pointer, and index register.

Memory management information: Information about the memory allocated to the process, including its base address, limit, and page table information.

Priority: A value that determines the relative importance of the process compared to other processes.

I/O status information: Information about any I/O devices that are associated with the process, including the device status and any pending I/O operations.

## 3.3 Importance of PCB in process management

The PCB is a critical data structure used by the operating system to manage processes effectively. The operating system uses the information stored in the PCB to make decisions about how to allocate resources to the process, such as CPU time, memory, and I/O devices. Without the PCB, it would be difficult for the operating system to manage multiple processes concurrently and efficiently. The PCB is also used by the operating system to switch between processes quickly, as it contains all the necessary information required to save the state of a process and restore it later.

In conclusion, the Process Control Block (PCB) is a vital component of process management in modern operating systems. It is responsible for storing and managing critical information about a process, including its current state, register values, memory allocation, and other attributes. The PCB allows the operating system to manage multiple processes efficiently and switch between them quickly, making it an essential part of any modern operating system.

## 4 Process Scheduling

In a multi-tasking and concurrent operating system, several processes compete for the same resources, such as CPU time and memory. The

scheduler is responsible for assigning these resources to the processes efficiently and fairly. This is where process scheduling comes into play. In this chapter, we will discuss the definition of process scheduling, popular scheduling algorithms such as First-Come-First-Serve (FCFS), Shortest Job First (SJF), Round Robin (RR), Priority Scheduling, and Multilevel Feedback Queue (MLFQ), and the importance of process scheduling in multi-tasking and concurrency.

Process scheduling is a fundamental task in the operating system, and its primary purpose is to allocate CPU time to multiple processes in an efficient and effective manner. The CPU is a valuable resource that must be allocated carefully, as it can significantly impact the performance of the entire system. The scheduler is responsible for selecting the next process to run based on a set of predefined criteria. These criteria can vary depending on the scheduling algorithm used.

We will explore the most popular scheduling algorithms in detail in this chapter, including their strengths and weaknesses. We will also discuss the contents of a process control block (PCB), which is an essential data structure used by the scheduler to store information about each process. Finally, we will discuss the importance of process scheduling in multi-tasking and concurrency.

## 4.1  Definition of process scheduling

Process scheduling is a critical aspect of operating systems that allows for efficient and effective utilization of system resources. It involves determining which process will be executed by the CPU and when. The objective of process scheduling is to optimize system performance by reducing the turnaround time, waiting time, and response time of processes. In this chapter, we will explore the concept of process scheduling and its importance in modern operating systems.

Process scheduling is the mechanism used by operating systems to determine which process will be executed by the CPU next. The process scheduler is responsible for managing the queue of processes waiting to be executed and allocating system resources to them. The scheduler must balance competing demands for resources while ensuring that each process receives a fair share of CPU time. Process scheduling can be preemptive or non-preemptive. In preemptive scheduling, the CPU can be taken away from a process at any time, while in non-preemptive scheduling, a process holds the CPU until it voluntarily relinquishes it.

Process scheduling is a key component of any operating system, and different scheduling algorithms have been developed to manage system resources efficiently. These algorithms are designed to prioritize certain processes over others based on various criteria, such as the length of time a process has been waiting, its priority level, or the amount of CPU time it has already consumed.

## 4.2 Scheduling algorithms:

### 4.2.1 First-Come-First-Serve (FCFS)

First-Come-First-Serve (FCFS) is one of the simplest CPU scheduling algorithms used in operating systems. In this algorithm, the process that arrives first is allocated the CPU first. The FCFS algorithm is also known as the First-In-First-Out (FIFO) scheduling algorithm because the process that comes first into the ready queue will be the first one to be executed.

When a process enters the ready queue, it is assigned the CPU based on the order in which it entered the queue. The process continues to use the CPU until it finishes its execution, blocks for I/O, or terminates.

The FCFS algorithm is simple to implement and understand, but it can cause long waiting times for the processes that arrive later. This is

because a process with a longer burst time can hold the CPU for a long time, causing other processes to wait in the ready queue.

Moreover, the FCFS algorithm can lead to a phenomenon known as convoy effect. The convoy effect occurs when a long process holds the CPU, causing other short processes to wait behind it. This can lead to poor resource utilization and long average waiting times.

Despite its drawbacks, the FCFS algorithm is still widely used in batch processing systems, where it is essential to execute processes in the order they were submitted. However, it is not suitable for interactive systems, where users expect a quick response time.

In summary, the FCFS algorithm is a simple and straightforward scheduling algorithm that can be used in batch processing systems. However, it can cause long waiting times and poor resource utilization in interactive systems.

### 4.2.2 Shortest Job First (SJF)

In process scheduling, the goal is to maximize the throughput, minimize response time, and minimize turnaround time. The Shortest Job First (SJF) scheduling algorithm is a non-preemptive algorithm that prioritizes the process with the shortest CPU burst time.

The SJF scheduling algorithm is a non-preemptive algorithm that selects the process with the smallest CPU burst time. The CPU burst time is the amount of time required by a process to complete its execution on the CPU. In this algorithm, the ready queue is maintained in order of the burst time of the processes.

The SJF scheduling algorithm can be either preemptive or non-preemptive. In preemptive SJF scheduling, if a new process with a shorter burst time enters the ready queue, the currently executing process is interrupted, and the new process is scheduled to execute. In

non-preemptive SJF scheduling, once a process starts executing, it continues until completion or until it enters the blocked state.

The SJF scheduling algorithm has the following advantages:

- It results in the shortest average waiting time for the processes in the ready queue.
- It is optimal in the sense that it minimizes the average waiting time for the processes, assuming that the CPU burst times are known in advance.
- It prioritizes short processes, which leads to a faster turnaround time.

The SJF scheduling algorithm has the following disadvantages:

- It requires knowledge of the CPU burst time of each process, which is not always available.
- It can lead to starvation for long processes if there are a large number of short processes in the system.

The SJF scheduling algorithm is a non-preemptive algorithm that selects the process with the shortest CPU burst time. It prioritizes short processes, resulting in a faster turnaround time and shorter average waiting time for processes in the ready queue. However, it requires knowledge of the CPU burst time of each process, which may not always be available, and it can lead to starvation for long processes if there are a large number of short processes in the system.

### 4.2.3 Priority Scheduling

Priority Scheduling is one of the most widely used scheduling algorithms in modern operating systems. As the name suggests, this algorithm assigns priorities to each process based on their importance

and then schedules them according to their priority. This approach ensures that the high-priority tasks get executed first, thus improving the overall performance and responsiveness of the system.

The basic idea behind Priority Scheduling is to assign a priority level to each process based on its importance. This priority level can be determined based on various factors, such as the amount of CPU time required by the process, the amount of memory it needs, or the urgency of the task. Once the priority levels are assigned, the scheduler can then schedule the processes based on their priority, giving the highest-priority process the CPU time first.

Priority levels can be assigned either statically or dynamically. In static priority scheduling, the priority level of a process is fixed at the time of creation and remains unchanged throughout the life of the process. In dynamic priority scheduling, on the other hand, the priority level of a process can be adjusted dynamically based on its behavior.

Priority Scheduling can be further divided into two different scheduling policies: Preemptive and Non-preemptive. In Preemptive Priority Scheduling, a higher-priority process can interrupt a lower-priority process while it is running. This allows the higher-priority process to start executing immediately, thus ensuring that the most important tasks get executed first. In Non-preemptive Priority Scheduling, however, a running process cannot be interrupted by a higher-priority process. In this case, the scheduler must wait for the currently running process to finish before scheduling the higher-priority process.

One of the key advantages of Priority Scheduling is that it allows the most important tasks to be executed first. This improves the overall performance and responsiveness of the system, especially in real-time systems where timely execution of critical tasks is essential. Priority Scheduling also allows for efficient utilization of system resources, as it ensures that the most important tasks get executed first, thus reducing wastage of CPU time.

One major disadvantage of Priority Scheduling is that it can lead to starvation. If a low-priority process is constantly preempted by higher-priority processes, it may never get a chance to execute, thus leading to starvation. Another disadvantage is that it can lead to low-priority processes getting neglected. If there are too many high-priority processes, the low-priority processes may never get a chance to execute, leading to poor performance and responsiveness.

Priority Scheduling is widely used in modern operating systems, such as Linux, Windows, and macOS. In these systems, priority levels are assigned based on various factors, such as the type of process, its behavior, and its importance. The scheduler then uses these priorities to schedule the processes, giving the highest-priority process the CPU time first.

### 4.2.4  Round Robin (RR)

Round Robin (RR) is a widely used CPU scheduling algorithm in operating systems. This algorithm is designed to schedule multiple processes concurrently in a fair and efficient manner. In this chapter, we will discuss the details of the Round Robin scheduling algorithm.

The Round Robin (RR) scheduling algorithm is a preemptive scheduling algorithm. In this algorithm, each process is assigned a fixed time interval, called a time quantum or time slice, to use the CPU. The time quantum is usually a small unit of time, typically between 10 to 100 milliseconds.

When a process is given the CPU, it is allowed to run for a time quantum. If the process completes its task before the end of the time quantum, it releases the CPU voluntarily. If the process does not complete its task before the end of the time quantum, it is preempted and moved to the back of the ready queue. The next process in the ready queue is then given the CPU to execute.

If a process arrives while the CPU is busy, it is placed at the end of the ready queue. This ensures that processes are executed in the order in which they arrived.

Round Robin (RR) scheduling algorithm offers the following advantages:

- Fairness: The Round Robin scheduling algorithm provides fairness to all processes by giving each process an equal opportunity to use the CPU.
- Time-sharing: The Round Robin scheduling algorithm is ideal for time-sharing systems where multiple users access the system simultaneously.
- Low response time: The Round Robin scheduling algorithm ensures that each process gets a turn to execute quickly, resulting in a low response time.

The Round Robin scheduling algorithm has the following disadvantages:

- Inefficiency: The Round Robin scheduling algorithm can be inefficient when the time quantum is too long or too short.
- Overhead: The Round Robin scheduling algorithm has a higher overhead compared to other scheduling algorithms, as it requires the scheduler to keep track of the time quantum for each process.

The Round Robin (RR) scheduling algorithm is a widely used scheduling algorithm in operating systems. It provides fairness to all processes and is ideal for time-sharing systems. However, it can be inefficient and has a higher overhead compared to other scheduling algorithms.

### 4.2.5 Multilevel Queue Scheduling (MLQS)

In the previous chapters, we discussed various CPU scheduling algorithms, such as FCFS, SJF, Priority Scheduling, and Round Robin. In this chapter, we will discuss another important algorithm called Multilevel Queue Scheduling (MLQS). The MLQS algorithm is widely used in modern operating systems, especially in systems that need to manage multiple types of processes with varying priority levels.

Multilevel Queue Scheduling is a scheduling algorithm that divides the ready queue into several separate queues, each with its own scheduling algorithm. Each queue has its own priority level, and the scheduling algorithm is applied to each queue based on the priority level. This approach allows the operating system to prioritize different types of processes based on their needs, without compromising the responsiveness of the system.

In Multilevel Queue Scheduling, the ready queue is divided into multiple queues, each with its own priority level. The operating system assigns each process to a queue based on its priority level. Each queue can have its own scheduling algorithm, such as FCFS, SJF, Priority Scheduling, or Round Robin.

In most implementations of MLQS, the highest priority queue is served first, followed by the next highest priority queue, and so on. Within each queue, the scheduling algorithm is applied to determine the order in which processes are executed. The processes in each queue are typically scheduled in a non-preemptive manner, meaning that a process must yield the CPU voluntarily before another process can be executed.

Multilevel Queue Scheduling is an important scheduling algorithm in modern operating systems. It allows the operating system to prioritize different types of processes based on their priority level and the needs of the system. This is important in systems that need to manage multiple types of processes with varying priority levels, such as real-time systems and systems that run both user-level and system-level processes.

The MLQS algorithm can also help improve the responsiveness of the system by ensuring that high-priority processes are given priority access to the CPU. By dividing the ready queue into multiple queues, the operating system can ensure that each type of process is given the appropriate amount of CPU time, without compromising the performance of the system.

Multilevel Queue Scheduling is an important scheduling algorithm in modern operating systems. It allows the operating system to prioritize different types of processes based on their priority level and the needs of the system. The MLQS algorithm can help improve the responsiveness of the system and ensure that high-priority processes are given priority access to the CPU.

### 4.2.6 Multilevel Feedback Queue Scheduling (MLFQS)

Multilevel Feedback Queue Scheduling (MLFQS) is a complex scheduling algorithm that dynamically adjusts the priority of processes based on their behavior over time. It is an extension of the Multilevel Queue Scheduling (MLQS) algorithm and is widely used in modern operating systems.

The MLFQS algorithm works by dividing the ready queue into multiple priority queues. Each queue is assigned a different priority level, with the highest priority queue being reserved for the most important processes, such as system processes or real-time processes. Each queue also has its own scheduling algorithm, with different algorithms being used for different priority levels.

When a process is created, it is placed in the highest priority queue. The process is then given a certain amount of time to execute, known as a time slice or quantum. If the process completes its execution before the time slice expires, it is removed from the queue. If the time slice expires before the process completes its execution, the process is preempted and moved to a lower priority queue.

The priority of a process in MLFQS is determined dynamically based on its behavior over time. Processes that use a lot of CPU time are given lower priority to prevent them from monopolizing the CPU. Conversely, processes that use less CPU time are given higher priority to ensure that they are executed quickly.

MLFQS also includes a mechanism for promoting and demoting processes between priority levels. When a process is blocked, it is moved to a lower priority queue. When it becomes unblocked, it is moved back to its original priority level. This ensures that long-running processes do not hold up the system and that important processes are executed as quickly as possible.

In summary, Multilevel Feedback Queue Scheduling (MLFQS) is a powerful scheduling algorithm that dynamically adjusts the priority of processes based on their behavior over time. It is designed to ensure that important processes are executed quickly while preventing long-running processes from monopolizing the CPU.

### 4.2.7 Lottery Scheduling

In operating systems, the lottery scheduling algorithm is a probabilistic scheduling algorithm used to allocate resources to processes. It is a unique scheduling algorithm because it provides equal opportunities for all processes to win a "lottery ticket" and acquire CPU time. The algorithm uses a lottery ticket metaphor to determine which process gets the CPU next, making the process selection entirely random.

The lottery scheduling algorithm assigns each process a set of lottery tickets. A lottery ticket represents the process's share of the CPU time. The more lottery tickets a process has, the higher the chances of it winning the lottery and acquiring the CPU.

To allocate CPU time, the lottery scheduling algorithm randomly selects a ticket from the ticket pool. The process that owns the ticket wins the lottery and gets to run on the CPU for a set time quantum. After the

time quantum expires, the process returns the CPU, and the lottery begins again.

The lottery scheduling algorithm uses a number of data structures to maintain the ticket pool and manage processes' ticket allocation. One such data structure is a list of processes, each with a number of lottery tickets associated with it. The algorithm also maintains a list of unused tickets and a counter that tracks the number of tickets in the pool.

The lottery scheduling algorithm has several advantages over other scheduling algorithms. One significant advantage is its ability to provide fairness to all processes. Since each process has an equal chance of winning the lottery, the scheduling algorithm ensures that no process is left behind or unfairly treated.

Another advantage of the lottery scheduling algorithm is its simplicity. The algorithm is easy to implement and does not require complex data structures or sophisticated algorithms. This simplicity translates to low overhead costs and makes it a good choice for systems with limited resources.

However, the lottery scheduling algorithm has some disadvantages. One significant disadvantage is that the algorithm is entirely random. There is no guarantee that a process with a large number of tickets will win the lottery or that a process with few tickets will not win the lottery several times in a row. This randomness can lead to inefficiencies in the system's performance and unpredictability in the scheduling results.

Additionally, the lottery scheduling algorithm may not be suitable for systems with strict scheduling requirements. The randomness of the algorithm may lead to scheduling delays and missed deadlines, which can be detrimental in real-time systems.

The lottery scheduling algorithm is a unique scheduling algorithm that uses a lottery ticket metaphor to allocate CPU time to processes. The algorithm provides fairness to all processes and is easy to implement, making it a good choice for systems with limited resources.

However, the algorithm's randomness can lead to inefficiencies and unpredictability in the scheduling results, making it unsuitable for systems with strict scheduling requirements. In such cases, other scheduling algorithms, such as Round Robin or Priority Scheduling, may be more appropriate.

## 4.2.8 Fair-Share Scheduling

In a multi-user system, it is important to ensure fairness and prevent any single user from monopolizing system resources. Fair-Share Scheduling is a scheduling algorithm that addresses this issue by allocating system resources fairly among all users. In this chapter, we will discuss the concept of Fair-Share Scheduling and how it works in an operating system.

Fair-Share Scheduling is a scheduling algorithm that dynamically allocates system resources based on the proportion of resources each user is entitled to. Each user is assigned a "share" of the system resources, and the scheduler ensures that each user gets their fair share. The concept of fair sharing can be applied to various system resources, including CPU time, memory, and I/O devices.

The Fair-Share Scheduling algorithm works by maintaining a record of each user's resource usage over time. This record is used to calculate each user's share of the resources based on a predetermined policy. The policy may be based on factors such as the number of active processes, the amount of CPU time used, or a combination of factors.

When a user requests a resource, such as CPU time or memory, the scheduler checks their entitlement to that resource based on the user's share. If the user's share has been used up, they may be placed in a waiting queue until their share becomes available again. This prevents any user from monopolizing system resources and ensures that all users are treated fairly.

The main advantage of Fair-Share Scheduling is that it ensures fairness in the allocation of system resources. This is particularly important in multi-user systems where resources are shared among multiple users. By dynamically adjusting each user's share based on their resource usage, the scheduler ensures that no user can monopolize the system resources.

Another advantage of Fair-Share Scheduling is that it allows administrators to set policies that reflect the organization's priorities. For example, a policy can be set to give higher priority to certain users or groups of users, based on their role within the organization.

One potential disadvantage of Fair-Share Scheduling is that it can be complex to implement and maintain. The scheduler needs to keep track of each user's resource usage over time, which requires additional system overhead. Additionally, the policies used to calculate each user's share can be complex and difficult to configure.

Another potential disadvantage of Fair-Share Scheduling is that it may not be suitable for all types of systems. For example, in a system where users are performing real-time tasks, such as video streaming or audio processing, Fair-Share Scheduling may not provide the necessary responsiveness.

Fair-Share Scheduling is a scheduling algorithm that ensures fairness in the allocation of system resources. It works by dynamically adjusting each user's share based on their resource usage, preventing any single user from monopolizing system resources. While Fair-Share Scheduling has its advantages, such as allowing administrators to set policies that reflect the organization's priorities, it also has potential disadvantages, such as increased complexity and possible lack of responsiveness in real-time systems.

### 4.2.9 Guaranteed Scheduling

In the context of Operating Systems, a Guaranteed Scheduling Algorithm is a type of scheduling algorithm that ensures that certain

processes are guaranteed a certain amount of CPU time, regardless of the presence of other processes in the system. This is particularly useful in situations where there are real-time processes that require a certain level of responsiveness from the system.

In a Guaranteed Scheduling Algorithm, the system sets aside a certain amount of CPU time for specific processes, known as guaranteed processes. These guaranteed processes are given a certain priority level, which determines the amount of CPU time they are allocated. Once a guaranteed process is scheduled to run, it is given the CPU until it either completes or reaches its maximum allocated time slice.

If there are no guaranteed processes ready to run, the system switches to a different scheduling algorithm to assign CPU time to non-guaranteed processes. This helps to ensure that non-guaranteed processes don't starve for CPU time.

One of the main advantages of a Guaranteed Scheduling Algorithm is that it ensures that certain processes receive the CPU time they need to function properly. This is particularly important in real-time systems, where a delay in processing a critical task could have serious consequences.

Another advantage of Guaranteed Scheduling Algorithm is that it allows for more precise control over system resources. By allocating specific amounts of CPU time to specific processes, system administrators can ensure that all processes receive an appropriate amount of resources, while also preventing any single process from monopolizing the CPU.

The main disadvantage of a Guaranteed Scheduling Algorithm is that it can be difficult to implement in a way that balances the needs of all processes in the system. For example, if too much CPU time is allocated to guaranteed processes, non-guaranteed processes may experience unacceptable levels of latency or may be starved for CPU time.

Additionally, a Guaranteed Scheduling Algorithm can be complex to implement, requiring careful tuning of the system parameters and a thorough understanding of the needs of each process in the system.

## 4.3 Importance of process scheduling in multi-tasking and concurrency

In multi-tasking and concurrent environments, it is essential to have a proper process scheduling mechanism in place. The scheduling algorithm plays a vital role in deciding which process gets to execute on the CPU and for how long. The process scheduling algorithm decides the efficiency of an operating system in handling multiple tasks simultaneously.

The primary goal of process scheduling is to improve system performance by reducing the CPU idle time and improving the response time of the system. The following are the key reasons why process scheduling is essential in a multi-tasking and concurrent environment:

Resource Utilization: In a multi-tasking environment, several processes compete for resources such as CPU, memory, and I/O devices. Process scheduling ensures that these resources are allocated efficiently and utilized to their maximum capacity.

Throughput: Process scheduling influences the throughput of the system. The throughput refers to the number of processes that the system can execute in a given period. A good process scheduling algorithm can significantly improve the throughput of the system.

Response Time: The response time of a system refers to the time taken by the system to respond to a user's input. A good process scheduling algorithm can ensure that the system responds to the user's input promptly.

Fairness: Process scheduling ensures that every process gets a fair share of the system resources. A fair process scheduling algorithm can ensure that no process is starved of system resources.

Prioritization: Process scheduling can prioritize processes based on their importance. The priority of a process determines its access to the system resources. A good process scheduling algorithm can ensure that critical processes get higher priority, ensuring that the system operates efficiently.

In conclusion, process scheduling is an essential aspect of any operating system, particularly in a multi-tasking and concurrent environment. The scheduling algorithm used by an operating system can significantly impact its performance, response time, throughput, fairness, and prioritization. Operating system designers need to consider these factors while designing the process scheduling algorithm to ensure that the operating system is efficient and responsive.

# 5 Interprocess Communication (IPC) and Synchronization

In this chapter, we will explore the various methods of IPC and synchronization, such as shared memory, message passing, semaphores, and monitors. These methods provide a means for different processes to exchange information and coordinate their activities. We will also discuss the importance of IPC and synchronization in multi-tasking and concurrency, and how they help in preventing race conditions, deadlocks, and other synchronization problems that may arise when multiple processes access shared resources simultaneously. Overall, this chapter will provide you with an understanding of how IPC and synchronization play a crucial role in the effective management of processes in operating systems.

## 5.1 Definition of IPC and synchronization

Interprocess Communication (IPC) and synchronization are two important concepts in operating systems that are necessary for effective multi-tasking and concurrency.

IPC refers to the mechanisms and techniques used by different processes to communicate with each other and share resources. In a multi-tasking environment, it is essential for different processes to communicate with each other to coordinate their activities, share data, and perform tasks collaboratively.

Synchronization, on the other hand, refers to the process of coordinating access to shared resources among different processes. In a multi-tasking environment, multiple processes may require access to the same resources simultaneously, and synchronization ensures that they do not interfere with each other.

IPC and synchronization are closely related concepts, as synchronization is necessary for proper communication and resource sharing between processes. There are various methods available for IPC and synchronization, each with its own advantages and disadvantages.

In the following chapters, we will explore different methods of IPC and synchronization in detail, along with their advantages, disadvantages, and real-world applications.

## 5.2 Methods of IPC and synchronization:

In a multi-tasking and concurrent environment, processes often need to communicate and synchronize with each other to achieve their intended goals. This is where Interprocess Communication (IPC) and synchronization techniques come into play. IPC and synchronization

allow processes to exchange information and coordinate their activities, ensuring that the system operates correctly and efficiently.

There are several methods of IPC and synchronization, each with its strengths and weaknesses. In this chapter, we will explore four of the most commonly used methods: shared memory, message passing, semaphores, and monitors.

### 5.2.1 Shared Memory

Shared memory is a technique that allows multiple processes to access the same region of memory. This region of memory is called a shared memory segment and is typically created by one process and then shared with other processes. Once a process has access to the shared memory segment, it can read from and write to it just like any other region of memory.

Shared memory is a fast and efficient way for processes to exchange large amounts of data because there is no need to copy the data between processes. However, it requires careful management to ensure that processes do not overwrite each other's data or access the shared memory segment at the same time.

### 5.2.2 Message Passing

Message passing is a technique that involves sending messages between processes. In this method, one process sends a message to another process, which receives and processes the message. Messages can be sent using either synchronous or asynchronous communication.

Synchronous communication means that the sender and receiver must synchronize their actions, such that the sender will not send another message until the receiver has processed the first message. Asynchronous communication, on the other hand, allows the sender to continue processing without waiting for the receiver to respond.

Message passing is a flexible and reliable method of IPC, but it can be slower and less efficient than shared memory, especially when large amounts of data need to be exchanged.

### 5.2.3 Semaphores

Semaphores are a synchronization technique that allows processes to coordinate their activities by controlling access to shared resources. A semaphore is essentially a counter that can be incremented and decremented by processes. When the counter reaches zero, the semaphore is considered to be locked, and any process attempting to access the shared resource must wait until the semaphore is unlocked.

Semaphores can be used to implement critical sections, where only one process at a time is allowed to access a shared resource. They can also be used to implement synchronization between processes, ensuring that one process completes its task before another process begins.

### 5.2.4 Monitors

Monitors are a synchronization technique that provides a higher-level abstraction than semaphores. A monitor is a module that encapsulates shared data and the procedures that operate on that data. Only one process can access a monitor at a time, and any other process that attempts to access the monitor is blocked until the first process completes its work.

Monitors are a powerful tool for IPC and synchronization because they simplify the development of concurrent programs. They provide a natural way to encapsulate shared data and procedures and ensure that processes do not interfere with each other.

IPC and synchronization are essential concepts in the field of operating systems. The methods we have discussed in this chapter provide ways for processes to communicate and coordinate their activities effectively, ensuring that the system operates correctly and efficiently. Shared

memory, message passing, semaphores, and monitors all have their strengths and weaknesses, and the choice of which method to use depends on the specific requirements of the system. As such, it is important for operating system developers to have a solid understanding of these concepts and their implementation.

## 5.3 Importance of IPC and synchronization in multi-tasking and concurrency

In a multi-tasking and concurrent operating system, several processes run simultaneously, accessing and manipulating shared resources, such as memory, files, or hardware devices. To ensure correct and safe operation, these processes must communicate and synchronize with each other through Interprocess Communication (IPC) and synchronization mechanisms.

IPC allows processes to exchange information and coordinate their activities. This is crucial for tasks such as data sharing, interlocking, or coordination, which can be accomplished using different methods, such as shared memory, message passing, semaphores, or monitors.

Synchronization mechanisms, on the other hand, ensure that processes access shared resources in an orderly and safe manner. For instance, mutual exclusion mechanisms, like semaphores or monitors, prevent two processes from accessing the same resource simultaneously, thus avoiding race conditions and data inconsistencies. Similarly, synchronization mechanisms like barriers, locks, or condition variables, ensure that processes wait for each other until a specific condition is met, allowing them to coordinate their activities.

In conclusion, IPC and synchronization are critical components of multi-tasking and concurrent operating systems, as they ensure that processes can communicate and coordinate with each other in a safe and efficient manner. By using these mechanisms, processes can access

and manipulate shared resources while avoiding data inconsistencies, race conditions, or deadlocks.

# 6 Case Study: Process Management in Linux Operating System

In this chapter, we will explore the process management in Linux, one of the most popular operating systems in the world. We will begin by giving an overview of the Linux process management, including its design principles and features.

Next, we will compare Linux process management with other operating systems, such as Windows and MacOS. This will help us to understand the strengths and weaknesses of Linux process management and how it differs from other systems.

Finally, we will discuss the impact of process management on the performance, reliability, and functionality of the Linux Operating System. We will analyze the various techniques and strategies employed by the Linux process management system to achieve these goals.

## 6.1 Overview of Linux process management

Linux is one of the most widely used operating systems in the world, powering everything from servers to mobile devices. One of the key reasons for its success is the robust process management capabilities it provides. In this chapter, we will take a closer look at the process management features of Linux.

Processes in Linux are managed using the process table, which is a data structure that holds information about all currently running processes. Each process is identified by a unique process ID (PID) and is associated

with other data, including its state, priority, parent process, and resource usage.

Linux provides a range of tools for managing processes, including the top command, which displays information about running processes, and the kill command, which is used to terminate a running process. The ps command is used to list processes and their attributes, while the nice and renice commands are used to adjust process priorities.

In Linux, processes are organized into a hierarchical structure, with each process having a parent process and the root process (init) as the ultimate parent. This structure helps to ensure that processes are properly managed and terminated when they are no longer needed.

Another key feature of Linux process management is the ability to create and manage threads. Threads are lightweight processes that share memory and other resources with their parent process. Linux provides a range of threading models, including POSIX threads, Native Posix threads library (NPTL), and LinuxThreads.

Overall, Linux's process management capabilities are a major strength of the operating system. They enable efficient multitasking and concurrency, ensuring that the system can handle multiple tasks simultaneously without becoming bogged down or crashing. Linux's process management features have also been influential in the development of other operating systems, making them an important area of study for anyone interested in operating system design and implementation.

## 6.2 Comparison with process management in other operating systems

Linux process management has several unique features and capabilities that set it apart from process management in other operating systems.

In this chapter, we will compare Linux process management with process management in other popular operating systems, including Windows and macOS.

6.2.1  Windows Process Management:

In Windows, the process management system is similar to that of Linux, where each process has its own virtual address space. However, there are some notable differences between the two. For example, in Windows, processes are assigned a priority value, which determines the order in which they are executed. This priority value can be changed by the operating system or the user, depending on the needs of the system. Additionally, Windows uses a system of "job objects" to group processes together and apply policies such as CPU time limits, memory limits, and more.

6.2.2  macOS Process Management:

Like Linux and Windows, macOS also uses a similar process management system. However, there are some notable differences between the three. For example, macOS uses a concept called "launchd" to manage processes. Launchd is a daemon that manages system services, user applications, and other processes. It provides a single point of control for starting, stopping, and monitoring processes on the system. Additionally, macOS uses a system of "sandboxing" to limit the resources available to individual processes, providing an additional layer of security.

6.2.3  Linux Process Management:

Linux process management is highly flexible and customizable. Each process has its own virtual address space and is managed by the kernel. The Linux kernel provides a variety of scheduling algorithms, including the Completely Fair Scheduler (CFS) and the Round Robin Scheduler.

Additionally, Linux supports a wide range of IPC mechanisms, including shared memory, message passing, and semaphores.

Linux also has a unique process management feature called "cgroups" (control groups). Cgroups allow processes to be organized into hierarchical groups, with each group having its own set of resource limits (CPU, memory, etc.). This feature is particularly useful for managing large-scale deployments such as web servers, where it is essential to limit resource usage.

In conclusion, Linux process management provides a high degree of flexibility and customization, allowing it to be adapted to a wide range of use cases. While other operating systems have similar process management systems, Linux stands out for its support for cgroups, its variety of scheduling algorithms, and its range of IPC mechanisms.

## 6.3 Impact on Linux Operating System's performance, reliability, and functionality

Linux is known for its excellent process management system, which allows for efficient multi-tasking and concurrency. The process management system in Linux is responsible for creating, managing, and terminating processes, as well as allocating resources to these processes.

The impact of Linux's process management system on its performance is significant. The kernel's scheduler is designed to be efficient and can quickly switch between processes, allowing for smooth multi-tasking. This means that users can run multiple programs simultaneously without any noticeable lag or slowdown.

Moreover, Linux's process management system is designed with reliability in mind. It includes several mechanisms for ensuring that processes run smoothly and without interruption. For example, Linux

uses signals to communicate with processes and notify them of events such as errors or resource availability.

In addition to performance and reliability, Linux's process management system also has an impact on the system's functionality. The system is designed to be flexible and can adapt to different requirements. For example, it allows for the creation of real-time processes that require immediate attention and prioritization.

Furthermore, Linux's process management system supports various synchronization and interprocess communication mechanisms, including shared memory, message passing, semaphores, and monitors. These mechanisms enable processes to communicate and synchronize their activities efficiently, which is essential for multi-tasking and concurrency.

Overall, Linux's process management system is a critical component of the operating system, and its impact on performance, reliability, and functionality cannot be overstated. It is a testament to the power and flexibility of open-source software development and community-driven innovation.

# 7  Conclusion

In conclusion, processes are the fundamental building blocks of operating systems. They allow users to run multiple tasks concurrently and efficiently use the resources of a computer system. Process management includes various activities such as process creation, scheduling, synchronization, and communication. Operating systems use a range of scheduling algorithms to manage processes effectively, depending on the needs of the system and the tasks being performed. Interprocess communication and synchronization play a crucial role in maintaining the integrity of processes and preventing errors or conflicts in multi-tasking and concurrent environments.

As we have seen, Linux operating system provides a robust and efficient process management system, allowing users to run multiple processes simultaneously and effectively utilize system resources. The Linux process management system is superior to other operating systems in terms of scalability, flexibility, and reliability.

Overall, understanding processes and their management is essential for anyone interested in operating systems and computer science. By understanding how processes work, how they communicate with each other, and how they interact with system resources, we can build more efficient and reliable operating systems that meet the demands of modern computing.