# Question & Answers

SCHEDULING

Sercan Külcü | Operating Systems | 10.04.2023

# Contents

## What is CPU scheduling?

CPU scheduling is the method by which the operating system decides which process or thread to execute from the ready queue. It assigns CPU time to the selected process based on various scheduling algorithms. The objective is to optimize system performance by ensuring efficient CPU utilization, improving responsiveness, and maintaining fairness across processes.

## What is the role of CPU scheduling in maximizing system throughput, responsiveness, and fairness?

CPU scheduling plays a crucial role in maximizing system throughput, responsiveness, and fairness. It ensures the CPU is efficiently utilized by prioritizing processes or threads based on their importance, which boosts throughput. By rapidly switching between tasks, the system remains responsive to user inputs. Additionally, CPU scheduling guarantees equitable allocation of resources, preventing any single process from monopolizing CPU time and promoting fairness.

## What are some common CPU scheduling algorithms?

First-Come, First-Served (FCFS): The process or thread that has been in the ready queue the longest is selected for CPU allocation.

Shortest Job First (SJF): The process with the shortest expected execution time is given CPU time next.

Priority Scheduling: Processes are assigned priorities, and the one with the highest priority is allocated CPU resources.

Round Robin: Each process is allocated a fixed time slice, or quantum, and is preempted once the time expires, allowing other processes to execute.

## How to choose CPU scheduling algorithm?

The choice of CPU scheduling algorithm directly influences system performance. For instance, Shortest Job First (SJF) can improve average turnaround and response times but may cause longer processes to experience increased waiting times. Priority Scheduling ensures high-priority processes are handled first, but it can lead to low-priority processes being starved of CPU time. Different algorithms offer trade-offs in terms of efficiency, fairness, and responsiveness.

## What is the relationship between CPU scheduling and process/thread management?

CPU scheduling and process/thread management are tightly interconnected in an operating system. The scheduler selects a process or thread from the ready queue to allocate CPU time, while the process/thread management system tracks the state of each process and handles state transitions, such as from ready to running or from waiting to ready. Both components work together to ensure efficient resource allocation and process execution.

## Briefly explain FCFS?

CFS (First Come, First Serve) is a non-preemptive CPU scheduling algorithm where processes are executed in the order they arrive in the ready queue. The first process in the queue is allocated CPU time,

followed by the second, and so on. Once a process starts, it runs to completion or until it performs I/O, blocking other processes. This can lead to longer waiting times for subsequent processes, especially if a long process is scheduled first.

## What is the RR scheduling algorithm?

The Round Robin (RR) scheduling algorithm allocates a fixed time quantum to each process, allowing it to run for a set duration before being preempted. Processes that arrive during a quantum are placed at the end of the queue. This ensures fair distribution of CPU time, preventing starvation since each process is given an equal opportunity. RR is essentially a cyclic version of FCFS, where no process has higher priority. It is also known as time-slicing scheduling.

## What is the difference between preemptive and non-preemptive scheduling?

Preemptive and non-preemptive scheduling differ in how the CPU is allocated to processes. In preemptive scheduling, the CPU is assigned for a limited time and can be interrupted if a higher-priority process arrives. This requires context switching between processes, introducing overhead. In contrast, non-preemptive scheduling allocates CPU time to a process until it terminates or enters a waiting state, without interruption. This results in lower overhead but lacks flexibility, as running processes are not preempted. Preemptive scheduling allows critical processes to gain immediate access to the CPU, while non-preemptive scheduling ensures that processes run to completion without interruption.

## What are starvation and aging in OS?

Starvation occurs in an operating system when a process is unable to access the resources it needs because other processes are continuously prioritized. This leads to the process being indefinitely delayed, reducing overall system efficiency and potentially causing deadlock. Aging is a technique used to prevent starvation by gradually increasing the priority of a process over time. As a process waits, its priority increases, ensuring that it will eventually be allocated resources and not be neglected. This mechanism helps maintain fairness and system stability.

## What is Context Switching?

Context switching is a fundamental process in multitasking operating systems. It occurs when the CPU switches from executing one process to another. During this switch, the operating system saves the current process's state—such as the program counter, registers, and other relevant data—into a Process Control Block (PCB). The PCB stores all necessary information about the process, including its state and memory allocation. After saving the state, the system loads the next process's information and resumes its execution. This enables efficient multitasking by allowing multiple processes to share the CPU.

## What are the goals of CPU scheduling?

The goals of CPU scheduling are to optimize system performance based on various factors. Maximizing CPU utilization ensures that the CPU is fully used, minimizing idle time. Fair allocation guarantees that no process monopolizes the CPU, and each process receives its fair share of time. Maximizing throughput aims to complete the most processes in a

given time period. Minimizing turnaround time reduces the total time a process takes to complete; while minimizing waiting time decreases the time a process spends in the ready queue. Lastly, minimizing response time is crucial for interactive systems, ensuring quick feedback to users. Different scheduling algorithms prioritize these goals based on system needs.

## How do operating systems handle priority-based scheduling and preemption?

Operating systems manage priority-based scheduling by assigning each process or thread a priority value. The scheduler selects the process with the highest priority for execution. If a higher-priority process becomes ready while another is running, it may preempt the current process. In this case, the scheduler interrupts the execution of the lower-priority process and allocates CPU time to the higher-priority one. This ensures that critical tasks are given precedence.

## How does the choice of time quantum impact CPU scheduling performance?

The time quantum, or time slice, determines how long a process runs on the CPU before being preempted. A longer time quantum reduces the scheduling overhead but may lower system responsiveness, as processes are interrupted less frequently. A shorter time quantum improves responsiveness by allowing processes to be preempted more often, but it increases the scheduling overhead due to more frequent context switches. The choice of time quantum affects the balance between system efficiency and responsiveness.

# What are some common techniques used in real-time CPU scheduling?

Real-time CPU scheduling relies on techniques like Earliest Deadline First (EDF), Rate Monotonic Scheduling (RMS), and Deadline Monotonic Scheduling (DMS). These methods prioritize processes based on their deadlines, ensuring critical tasks meet timing requirements. Unlike general-purpose scheduling algorithms, which focus on maximizing throughput and responsiveness, real-time scheduling guarantees that time-sensitive processes are executed within their specified time constraints.

# How do operating systems handle CPU scheduling in a multi-core/multi-processor environment?

In a multi-core or multi-processor environment, the operating system distributes tasks across multiple processors. Techniques like Symmetric Multiprocessing (SMP) and Non-Uniform Memory Access (NUMA) are commonly used for efficient scheduling. The scheduler must account for factors like cache affinity and processor affinity to optimize performance. Load balancing algorithms are also employed to ensure an even distribution of workload across all processors, preventing bottlenecks and enhancing system efficiency.

# What are some advanced techniques for improving the efficiency and fairness of CPU scheduling?

Multi-level feedback queue scheduling: This method dynamically adjusts process priorities based on CPU usage and behavior, improving both fairness and responsiveness.

Round-robin with dynamic time quantum: The time quantum is adjusted based on process behavior, optimizing CPU resource allocation.

Proportional-share scheduling: Processes are allocated CPU time in proportion to their assigned weight, ensuring fair resource distribution.

Gang scheduling: Groups of related processes are scheduled to run simultaneously on different processors, improving system resource utilization and performance.

Lottery scheduling: This technique assigns "tickets" to processes, where the number of tickets correlates with the process's priority. The scheduler randomly selects a process based on the number of tickets, allowing for probabilistic fairness and dynamic priority adjustments.

Fair Share Scheduling: This approach ensures that each user or group of processes gets a fair portion of the CPU time based on predefined allocations, which helps prevent resource monopolization by individual users or tasks.

Earliest Deadline First (EDF) with resource reservations: For real-time systems, EDF assigns priorities based on process deadlines. Coupled with resource reservation mechanisms, this technique guarantees that critical processes meet their timing constraints without starving other processes.

Hierarchical scheduling: This method organizes the scheduler into multiple levels, where each level manages a different group of processes. It allows for a more granular control over process management, enabling better optimization of system resources in complex environments.

# How do operating systems handle CPU scheduling in distributed computing and cloud computing environments?

In distributed and cloud computing environments, CPU scheduling is typically managed by a central scheduler or resource manager that oversees multiple nodes or instances. These systems often employ a mix of local and global scheduling strategies to ensure efficient workload distribution and optimal resource utilization. Virtualization technologies, such as containers and virtual machines, further aid in isolating and controlling resource allocation for individual applications or services, improving overall performance and scalability.

# What is thread-level speculation, and how can it be used to improve CPU scheduling performance?

Thread-level speculation is a technique that enhances CPU scheduling by allowing the processor to speculatively execute multiple threads in parallel, based on predicted outcomes. By overlapping thread execution, it makes better use of available resources, reducing idle times and improving overall performance and efficiency. This approach is especially beneficial in exploiting parallelism in workloads, leading to faster processing and resource utilization.

# How do operating systems handle dynamic workload changes?

Operating systems manage dynamic workload changes by adjusting CPU scheduling based on real-time system conditions and workload requirements. Key strategies include:

Load balancing: Distributing workloads across multiple processors or nodes to optimize resource usage and enhance performance.

Adaptive scheduling: Modifying scheduling parameters like time quantum and priority dynamically, based on the current workload and system state.

Predictive scheduling: Using historical data and predictive models to anticipate workload shifts, allowing for proactive scheduling adjustments.

Resource reservation: Allocating specific resources (e.g., CPU, memory) to high-priority tasks to ensure they meet deadlines or performance goals, particularly in real-time systems.

Elastic scaling: In cloud environments, systems can dynamically scale resources, such as adding more virtual machines or containers, to handle increased workload demands.

Priority-based adaptation: Continuously adjusting the priority of processes based on their importance or urgency, ensuring that critical tasks receive more CPU time while less critical tasks are deprioritized.

# What are some emerging trends and technologies in CPU scheduling?

Machine learning-based scheduling: Leveraging machine learning algorithms to adaptively adjust scheduling parameters, optimizing system performance based on real-time workload and environmental factors.

Energy-efficient scheduling: Implementing scheduling strategies that minimize power consumption, enhancing battery life in mobile devices and reducing overall energy usage in data centers.

Neuromorphic computing: Utilizing brain-inspired hardware and software to improve both performance and energy efficiency in CPU scheduling, enabling more efficient parallel processing and decision-making.

Heterogeneous scheduling: This approach focuses on efficiently managing workloads across diverse processing units, such as CPUs, GPUs, and specialized accelerators (e.g., FPGAs), to fully utilize the capabilities of heterogeneous systems.

Quantum computing: As quantum computing develops, new scheduling algorithms will be required to optimize quantum and classical resource usage, ensuring that quantum processors integrate seamlessly with traditional CPUs.

Virtualization-aware scheduling: Virtualization technologies, such as containers and hypervisors, are becoming more advanced, and CPU scheduling must take into account the dynamic nature of virtual machine workloads to optimize resource allocation across virtualized environments.

# What is the role of CPU scheduling in an operating system?

CPU scheduling is essential for efficient resource management in an operating system. It determines the order in which processes are assigned CPU time, ensuring fair distribution and optimizing system performance. The scheduler uses algorithms like priority scheduling, round-robin, or shortest job first to select the next process, balancing factors like responsiveness, throughput, and fairness. Through effective CPU scheduling, the operating system maximizes CPU utilization while minimizing waiting times and ensuring that all processes receive appropriate attention.

# What are the different types of CPU scheduling algorithms?

First-Come, First-Serve (FCFS):

Processes are scheduled in order they arrive. It is simple to implement, but it can lead to long waiting times for short processes if they arrive after longer ones. This issue is known as the "convoy effect."

Shortest Job First (SJF):

This algorithm prioritizes processes with the shortest burst time (execution time). It minimizes the average wait time, which is efficient for workloads where job lengths are predictable. However, it requires knowledge of the process runtime in advance, which is difficult to obtain. Additionally, it can lead to the starvation of longer processes if they are always preceded by shorter ones.

Priority Scheduling:

Each process is assigned a priority, and the CPU is allocated to the process with the highest priority. This can lead to starvation, where lower-priority processes might never get executed if higher-priority processes continuously arrive. It is also difficult to decide on appropriate priority levels, and the algorithm may become biased if priorities are not adjusted dynamically.

Round-Robin Scheduling (RR):

Round-robin scheduling allocates CPU time in equal time slices or quanta. It is fair since each process gets an equal chance to execute, which is suitable for time-sharing systems. However, frequent context switching can lead to high overhead, especially if the time quantum is too small. It might also not be efficient for processes with significantly different burst times.

Multilevel Feedback Queue Scheduling:

This algorithm uses multiple queues with varying priority levels, and processes are moved between these queues based on their behavior and CPU burst times. It adapts dynamically, providing a good balance between fairness and efficiency. It can prioritize processes that need immediate attention while allowing for flexibility in handling different types of tasks. However, it is complex to implement and may incur significant overhead due to queue management and priority adjustments.

Shortest Remaining Time First (SRTF):

A preemptive version of SJF, this algorithm always executes the process with the shortest remaining time. It can result in even lower wait times than non-preemptive SJF but requires accurate estimates of process execution time, which is difficult to obtain. Like SJF, it can also lead to starvation for longer processes.

Earliest Deadline First (EDF):

This real-time scheduling algorithm prioritizes processes based on their deadlines. The process with the earliest deadline gets the CPU. EDF is optimal for real-time systems where meeting deadlines is critical. However, it may not be feasible for systems where deadlines are difficult to predict, or dynamic workloads are involved.

Rate Monotonic Scheduling (RMS):

A fixed-priority algorithm used in real-time systems, RMS assigns priorities based on the frequency of the processes (i.e., processes that need to execute more frequently get higher priority). It's optimal for periodic tasks but assumes that all tasks have fixed, known periods, and may not be effective for tasks with varying execution times.

Lottery Scheduling:

This algorithm assigns a "lottery ticket" to each process, with a higher number of tickets corresponding to a higher priority. The scheduler randomly selects a ticket, and the process holding that ticket gets the

CPU. It offers a probabilistic method of allocating CPU time, balancing fairness and efficiency. However, it may introduce unpredictability and is less deterministic than other algorithms.

## How does CPU scheduling of multi-core processors differ compared to single-core processors?

In a multi-core processor environment, CPU scheduling must consider not only process priorities but also core affinity, which is the preference for assigning a process to the same core it previously ran on. This reduces cache misses and improves performance due to the locality of reference. Unlike single-core systems where only one process can run at a time, multi-core systems can execute multiple processes concurrently, so the scheduler must balance the load across the available cores. Strategies like load balancing, where tasks are dynamically redistributed among the cores to ensure even utilization, and processor affinity, which tries to keep processes on the same core to benefit from cache reuse, are key to optimizing performance in multi-core systems.

## How does the concept of "aging" help mitigate the issue of starvation in priority-based CPU scheduling?

"Aging" is a technique used in priority-based scheduling to prevent starvation, where processes with lower priority levels may never get a chance to execute. As time passes, the priority of processes waiting in the ready queue is gradually increased, ensuring that they eventually get CPU time. Aging can effectively prevent starvation by ensuring that no process is indefinitely delayed due to higher-priority tasks. However, one drawback of aging is that it introduces complexity in managing process priorities, and if the aging rate is too aggressive, it could lead to processes that should be prioritized being delayed unnecessarily.

Additionally, if the aging mechanism is not properly tuned, it may cause the system to oscillate between giving priority to new and old tasks, leading to suboptimal scheduling.

# How do Earliest Deadline First (EDF) and Rate Monotonic Scheduling (RMS) differ in their approach to handling deadlines?:

Earliest Deadline First (EDF) and Rate Monotonic Scheduling (RMS) are both widely used in real-time systems to handle task deadlines, but they differ in their approach. EDF is a dynamic priority scheduling algorithm, where tasks with the earliest deadline are given the highest priority. It guarantees optimal task scheduling, meaning that if a set of tasks can be scheduled, EDF will always find a feasible solution. However, EDF requires the ability to compute deadlines dynamically and is more computationally expensive. On the other hand, RMS is a static priority scheduling algorithm where tasks are assigned fixed priorities based on their periods (shorter periods receive higher priority). RMS is easier to implement but does not guarantee optimal scheduling for all task sets. For systems with tasks of varying periods, EDF is generally more suitable, as it dynamically adapts to changing deadlines and can handle tasks with different periods more efficiently. However, RMS may be preferred for systems where the task set is periodic and predictable, as it has a lower computational overhead.