# HOMEWORK 1

## 1. What are the three main purposes of an operating system?

The three main purposes of an operating system are:

1. It manages the computer hardware resources, such as the CPU, memory, and I/O devices, ensuring that programs have the necessary resources to run efficiently.
2. It provides a layer of abstraction and enforces a consistent interface so that application developers do not need to work with low-level hardware details. For instance, the OS offers system calls and libraries that streamline tasks like file I/O, memory allocation, and process scheduling.
3. An operating system is responsible for ensuring proper security and access control by isolating processes and protecting critical data structures from malicious or errant code.

Together, these core functions help maintain system stability, enable efficient program execution, and simplify the task of developing and running applications. Additionally, the OS manages concurrency through scheduling, which further enhances reliability and performance in multi-tasking environments. These purposes ensure robust computing experience.

## 2. What is the main difficulty that a programmer must overcome in writing an operating system for a real-time environment?

In real-time operating systems, the main difficulty lies in guaranteeing all critical tasks are completed within strict timing constraints. Unlike general-purpose operating systems, which focus on maximizing throughput or fairness, a real-time OS must ensure deterministic behavior under worst-case scenarios. This requires precise scheduling algorithms that consider deadlines, priorities, and resource contention. Programmers must handle interruptions, context switches, and shared data structures carefully, minimizing unpredictability or latency. Additionally, they must use priority inheritance protocols to avoid priority inversion. Limited resources in embedded environments further complicates development. Ensuring bounded response times demands rigorous testing, careful memory management, and efficient interrupt handling. The key challenge is creating a predictable, low-latency environment that meets deadlines despite real-world complexity and resource constraints. Failing to meet these constraints can result in catastrophic failures. Therefore, developers must design task scheduling and synchronization with meticulous attention to timing guarantees, ensuring the system remains stable under all conditions.

## 3. How does the distinction between kernel mode and user mode function as a rudimentary form of protection (security) system?

The distinction between kernel mode and user mode serves as a fundamental protection mechanism in modern operating systems.

In **Kernel Mode**, the OS has unrestricted access to all hardware and system resources, enabling it to manage critical tasks, such as process scheduling, memory management, and device control.

Conversely, **User Mode** restricts application code from directly accessing sensitive resources, such as kernel memory or privileged CPU instructions. This separation prevents user-level processes from accidentally or maliciously compromising the stability of the system.

When a program running in user mode requests services that require elevated privileges, it must execute a system call, which transitions the CPU to kernel mode. This controlled interface ensures that only validated requests can perform privileged operations. By enforcing strict mode transitions and limiting direct hardware access, the OS maintains integrity and protects processes from one another. This layered architecture underpins modern security, making kernel-user mode separation indispensable.

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

Caches are useful primarily because they reduce the effective access time to frequently used data and help mitigate latency between fast processors and slower storage systems.

1. Caching improves performance by storing copies of frequently accessed data in faster memory, decreasing wait times and enhancing overall throughput.
2. Caches reduce the load on underlying devices, prolonging their lifespan and conserving power. However, caches can introduce problems, including inconsistency (i.e., stale data when updates occur), and coherence overhead in multiprocessor environments.

If one were to make a cache as large as the device it caches, the distinction between cache and device would vanish, making management complex and expensive. Furthermore,

devices like disks serve as reliable permanent storage, while cache memory is typically volatile and costlier per unit capacity. Keeping a smaller, faster cache strikes a balance between performance gains and system affordability, ensuring rapid access without impractical replication of entire devices.

## 5. What is the purpose of interrupts? How does an interrupt differ from a trap? Can traps be generated intentionally by a user program? If so, for what purpose?

**Purpose:**
Interrupts allow hardware devices or events to signal the CPU that immediate attention is required, temporarily suspending the currently executing process to handle a critical situation. This mechanism ensures responsiveness and efficient handling of asynchronous events, such as keyboard input or network packets. By using interrupt service routines, the operating system can quickly manage device requests and resume normal execution once the event is handled.

**Trap:**
A trap, on the other hand, is a software-generated interrupt triggered by errors (like division by zero) or explicit system calls (like requesting OS services). Unlike hardware interrupts, traps originate within the CPU due to program execution rather than an external device. Traps can indeed be intentionally generated by user programs, typically through specific instructions that request operating system services. This allows applications to safely perform privileged operations via well-defined interfaces, thus maintaining system security and stability while granting access to critical resources as needed.