# Real Time Operating System

Unit - 3

Embedded System | Prepared By **Er. Prashant Pant**

# Operating System Basics

- For the embedded devices where response time for a task is not time critical, t**he Super loop based task execution model** for firmware execution can be used. Typical examples are electronic toys and gaming devices. Any response delay won't create any operational issues or potential hazards.
- However, certain applications demand critical response to tasks/events and any delay may become catastrophic. For example, Flight Control Systems, Air bag control and Anti Locking Brake (ABS) systems for vehicles, nuclear monitoring devices demand time critical task response.
- **It is addressed by adopting the following procedure:**
  - Assign priority to tasks & execute the high priority task when the task is ready to execute.
  - Dynamically change the priorities of tasks if required on a need basis.
  - Schedule the execution of tasks based on the priorities.
  - Switch the execution of task when a task is waiting for an external event or a system resource including I/O device operation.
- The above listed procedure is handled properly by the system software know as operating system.

# RTOS basics

- An operating system acts as a bridge between users and underlying system resources. It manages the system resources and makes them available to the user application/task on a need basis.
- **The primary functions of operating system are:**
  - To make the system convenient to use.
  - To organize & manage the system resources efficiently and correctly.
- It is designed with special characteristics to support real time applications.
- An RTOS must also be able to respond predictable outside events and process multiple events concurrently.
- **Types:** *A hard real-time operating system* has less jitter than *a soft real-time operating system.* The chief design goal is not high throughput, but rather a guarantee of a soft or hard performance category.
- **Difference:** An RTOS that can usually or generally meet a deadlineis a soft real-time OS, but if it can meet a deadline deterministically it is a hard real-time OS.

# RTOS Basics

- An RTOS has an advanced algorithm for scheduling.
- Scheduler flexibility enables a wider, computer-system orchestration of process priorities, but a real-time OS is more frequently dedicated to a narrow set of applications.
- **Key factors in a real-time OS**
  - minimal interrupt latency
  - minimal thread switching latency;
- A real-time OS is valued more for how quickly or how predictably it can respond than for the amount of work it can perform in a given period of time.

# Features of a RTOS

- **Reliability:**
  - Any RTOS must be reliable. This means that it operates for a reasonably long time without human interference.
  - Reliability also means the configuration to enable the system to choose the right or most profitable action for current operations.
  - For example, a system for the telecom or banking industries
- **Predictability:**
  - A system must execute actions within a known time frame and produce known results.
- **Performance:**
  - Real-time operating systems are designed to make work easier. Every system must solve a problem or reduce the workload.
- **Manageability:**
  - This means a system whose veracity or bulkiness is manageable. The software and hardware required to operate the RTOS must be of reasonable size. Technicians should also be easy to find and orient. The idea is to reduce the cost of implementation.
- **Scalability:**
  - The needs of any production or event environment change with time. This means that a system may require an upgrade or downgrade. Such provisions must be made during design and installation of any RTOS.
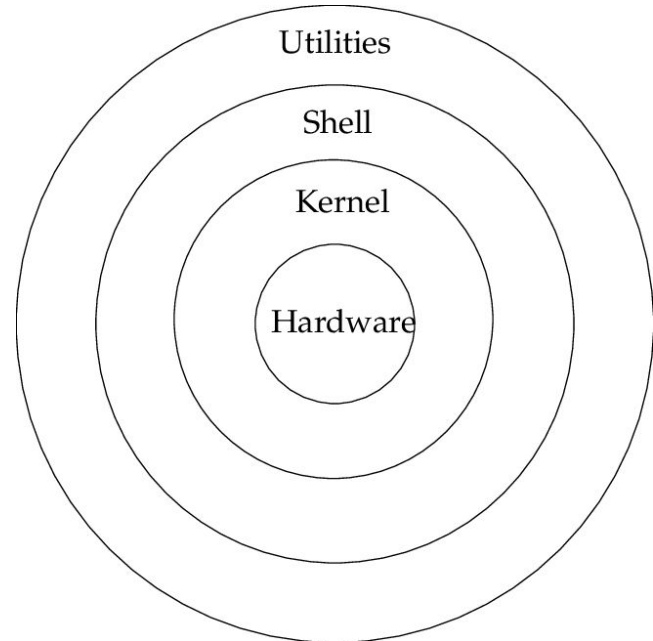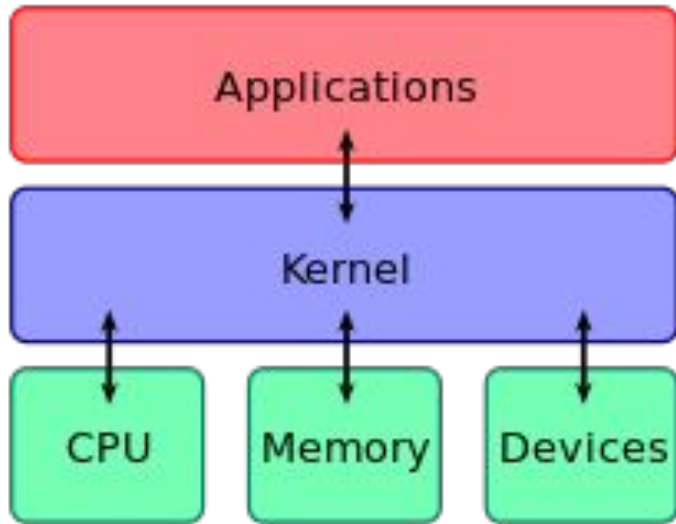
# The Kernel

- The kernel is a computer program that ís the core of a computer's operating system (OS), with complete control over everything in the system.
- It is responsible for managing the system resources and the communication among the hardware and other system services.
- It acts as the abstraction layer between system resources and user applications.
- It contains a set of system libraries and services.
- It is the first program loaded on startup.
- It handles the rest of startup as well as input/output requests from software translating them into data-processing instructions for the central processing unit.
- It handles memory and peripherals Like keyboards, monitors, printers, and speakers.
-

# The Kernel

- The kermel code is usually loaded into a protected area of memory to prevent it from being overwritten by programs or other parts of the operating system.
- The kernel performs its tasks, such as running processes and handling interrupts, in kernel space.
- In contrast, everything a user does is in user space.
- This separation prevents user data and kernel data from interfering with each other and causing instability and slowness.
- The kernel's interface is a low-level abstraction layer.
- When a process makes requests of the kernel, it is called a system call.
- Kernel designs differ in how they manage these system calls and resources.
- A monolithic kernel runs all the operating system instructions in the same address space, for speed.
- A micro-kernel runs most processes in user space, for modularity.

# Kernel Images

# Kernel Functions

- As kernel is the engine of any Operating System, so all the vital functions should be controlled and managed by kernel itself. There are various tasks and functions of a kernel but some of the important are:
- **Resource allocation:**
  - The kernel's primary function is to manage the computer's resources and allow other programs to run and use these resources.
  - These resources are-CPU, Memory and I/O devices.
- **Process Management:**
  - A process defines which memory portions the application can access.
  - The main task of a kernel is to allow the execution of applications and support them with features such as hardware abstraction.

# Kernel Functions

- **Primary Memory Management:**
  - The kernel has full access to the system's processes
  - It refers to the volatile memory (RAM) where processes are loaded and variables and shared data associated with each process are stored.
  - Memory Management unit (MMU) of the kernel is responsible for keeping track of which part of the memory area is currently used by which process.
  - Allocating and De-allocating memory space on a need basis (DMA)
- **I/O Device Management:**
  - To perform useful functions, processes need access to the peripherals connected to the computer, which are controlled by the kernel through device drivers.
  - A device driver is a computer program that enables the operating system to interact with a hardware device.
  - It provides the operating system with information of how to control and communicate with a certain piece of hardware.

# Kernel Functions

- **Inter-Process Communication:**
  - Communication between processes called Inter- Process Communication (IPC). There are various approaches of IPC say, semaphore, shared memory, message queue, pipe (or named FIFO), etc.
- **Scheduling:**
  - In a multitasking system, the kernel will give every program a slice of time and switch from process to process so quickly that it will appear to the user as if these processes were being executed simultaneously.
  - The kernel uses scheduling algorithms to determine which process is running next and how much time it will be given. The algorithm sets priority among the processes.

# Kernel Functions

- **System Calls and Interrupt Handling:**
  - A system call is a mechanism that is used by the application program to request a service from the operating system.
  - System calls include close, open, read, wait and write.
  - To access the services provided by the kernel we need to invoke the related kernel functions.
  - Most kernels provide a C Library or an API, which in turn invokes the related kernel functions.
- **Security or Protection Management:**
  - Kernel also provides protection from faults (error control) and from malicious behaviors (security).
  - One approach toward this can be language based protection system, in which the kernel will only allow code to execute which has been produced by a trusted language compiler.

# Kernel Functions

- System memory in Operating System can be divided into two, distinct regions: kernel space and user space.
- **Kernel Space:**
  - The program code corresponding to the kernel applications/ services are kept in a contiguous area (0S dependent) of primary memory and is protected from unauthorized access by their programs/applications.
  - The memory space at which kernel code is located is known as "Kernel Space'.
  - Kernel space is where the kernel (i.e., the core of the operating system) executes (i.e., runs) and provides its services.
  - Kernel space can be accessed by user processes only through the use of system calls.
  - System calls are requests in a Unix-like operating system by an active process for a service performed by the kernel, such as input/output (I/0) or process creation.
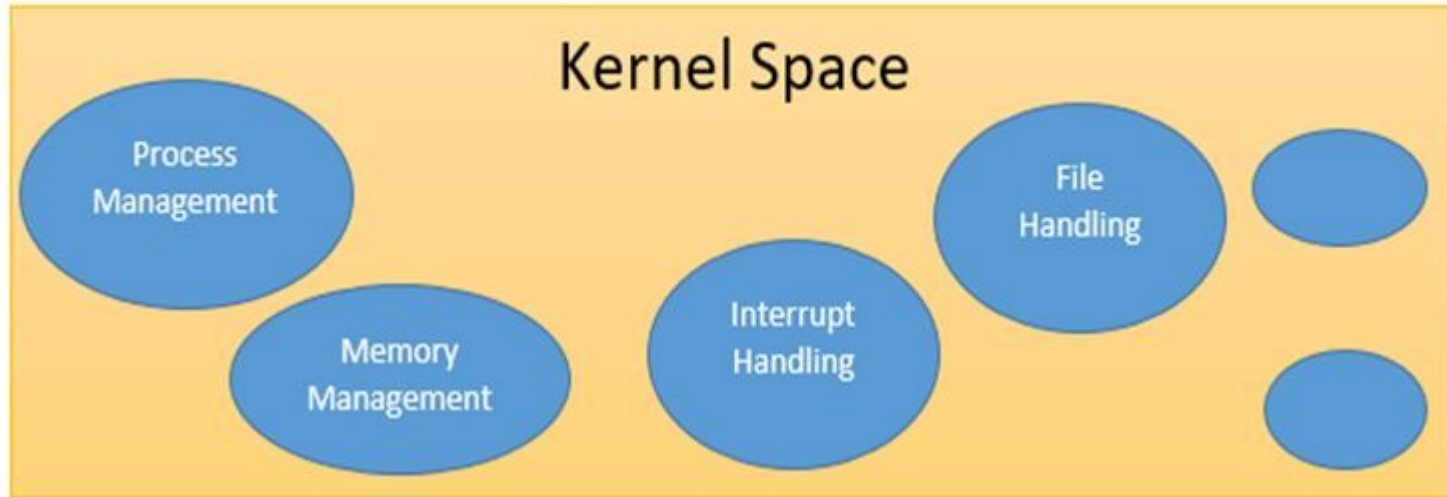
# Kernel Functions

- **User Space:**
  - All user applications are loaded to a specific area of primary memory "User space"
  - User space is that set of memory locations in which user processes (i.e. everything other than the kernel) run.
  - A process is an executing instance of a program.
  - One of the roles of the kernel is to manage individual user processes within this space and to prevent them from interfering with each other.
  - In an OS with virtual Memory support, the user applications are loaded into its corresponding

# Kernel Functions

- **Kernel Types: — The the design and architecture of the Kernel, it may be classified as:**
- **Monolithic Kernel:**
  - All kernel services run in the Kernel space.
  - All kernel modules run with, same Memory space under a single kernel thread. The tight internal integration of kernel may allows the effective utilization of the low-level features of the underlying system.
  - Earlier in this type of kernel architecture, all the basic system services like process and Mery, management, interrupt handling etc. were packaged into a single module in kernel space. This type of architecture led to some serious drawbacks like:
    - Size of kernel, which was huge.
    - Poor maintainability, which means bug fixing or addition of new features resulted ' recompilation of the whole kernel which could consume hours.
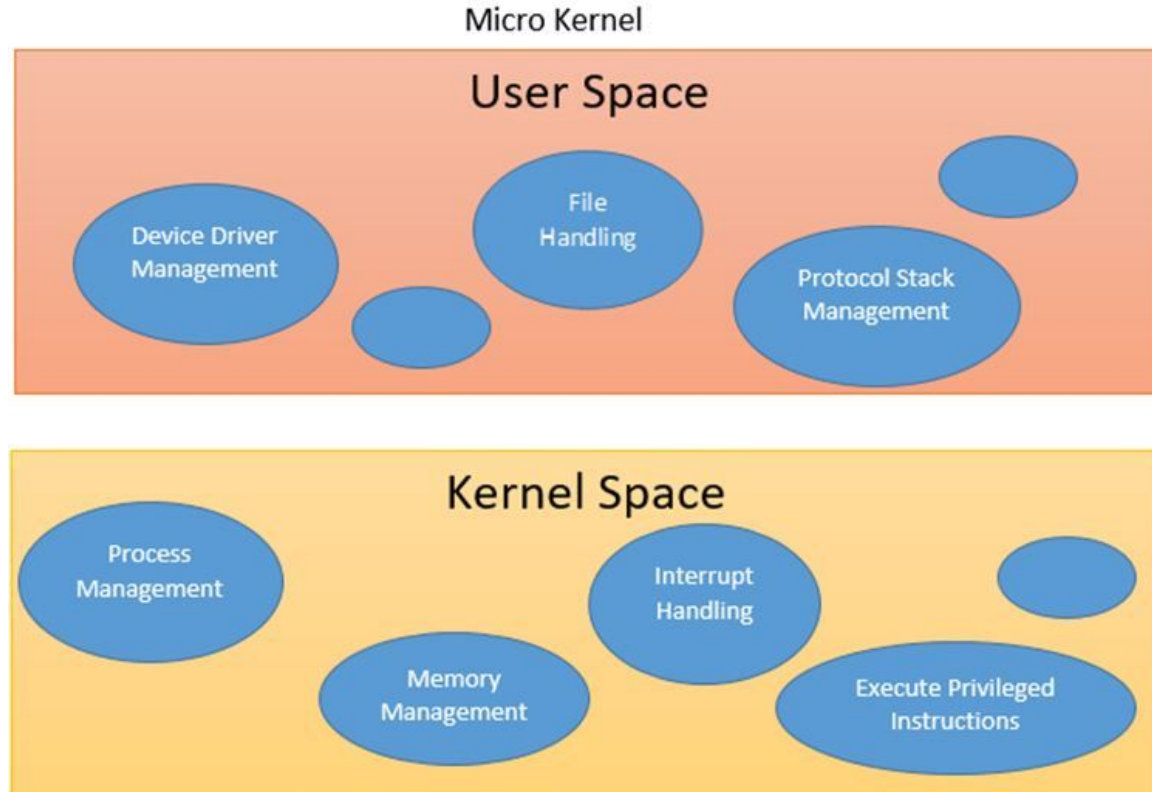
# Monolithic Kernel

# Kernel Functions

- Micro-kernels:
  - Microkernel incorporates only the essential set of OS services into the kemel.
  - The rest of the OS services are implemented in programs known as 'Servers' which runs in user space.
  - This provides a highly modular design and OS-neutral abstraction to the Kernel.
  - Memory management process management, timer systems and interrupt handlers are the essential services.
  - Examples: Mach, QNX, Minix 3 kernels.

# Micro Kernel



Micro Kernel

**User Space**

Device Driver Management

File Handling

Protocol Stack Management

**Kernel Space**

Process Management

Memory Management

Interrupt Handling

Execute Privileged Instructions

# Difference between GPOS and RTOS

- Depending on the type of kernel and kernel services, OS are classified into General Purpose Operating System (GPOS) and Real Time Operating System (RTOS).
- **General Operating System (GPOS):**
  - Operating Systems deployed in general computing Systems are GPOS.
  - The kernel of such OS is more generalized and contains all kinds of services required General Purpose for executing generic applications. It is often quite non-deterministic in behavior.
  - Their services can inject random delays into application software and may cause slow responsiveness of an application at unexpected times.
  - Usually deployed in computing systems where deterministic behavior is not an important criterion. For example: Windows Operating System.
- **Real Time Operating System (RTOS):**
  - RTOS supervises the application software tasks running on the hardware and organizes the access to system resources according to priorities & timing constraints.
  - 'Real Time' Implies deterministic timing behavior
  - It means the OS services consumes only known and expected amounts of time regardless the number of services.
  - It implements policies and rules concerning time critical allocation of a system's resources
  - It decides which applications should run in which order and how much time need to be allocated for each application.
  - Examples include Windows CE, QNX, VXWorks, RT Linux etc.

# The Real Time Kernel

A real-time kernel is a specialized operating system kernel designed to handle tasks with strict timing constraints, ensuring low latency and deterministic response times. Unlike standard kernels, which prioritize throughput and fair scheduling, real-time kernels focus on minimizing latency and providing consistent response times.

**Key Features and Mechanisms**

- **Priority Scheduling**: Real-time kernels use priority-based scheduling policies such as SCHED_FIFO (First In, First Out) and SCHED_RR (Round Robin). High-priority tasks are given preference for CPU execution, ensuring that critical tasks are completed within their deadlines.

# The Real Time Kernel

- **Low Latency**: Real-time kernels maintain low latency execution times, which is crucial for applications requiring immediate responses, such as industrial automation and robotics.
- **Deterministic Response**: The kernel ensures that tasks are executed within a predictable time frame, providing deterministic behavior essential for real-time applications.
- **Interrupt Handling**: Real-time kernels use dedicated CPUs to isolate processes from each other, reducing the impact of interrupts on task execution

# The Real Time Kernel

**Forecast**: By 2025, nearly 30% of the world's data will need real-time processing.

**Market Adoption**: Real-time computing is expected to accelerate in the coming years.

**Industrial Demand**: Increasing demand for industrial PCs, edge servers, PLCs, robots, and drones with real-time capabilities.

**Software Solutions**: Real-time computing via the Linux kernel is gaining popularity due to its robust support for hardware.

**PREEMPT_RT Patch**: Aims to enhance predictability and reduce latencies in the Linux kernel by modifying existing code.
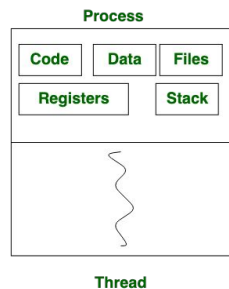
StackOver FLow

# Threads

**Threads** are a fundamental concept in modern operating systems. They represent the smallest unit of execution within a process. Essentially, a thread is a single sequence of instructions within a program that can be executed independently.

**Here's a breakdown:**

- **Process:** A process is an executing instance of a program. It encompasses the entire program code, data, and resources allocated to it by the operating system.
- **Thread:** A thread is a smaller unit within a process. Multiple threads can exist within a single process, each executing a different part of the program concurrently.
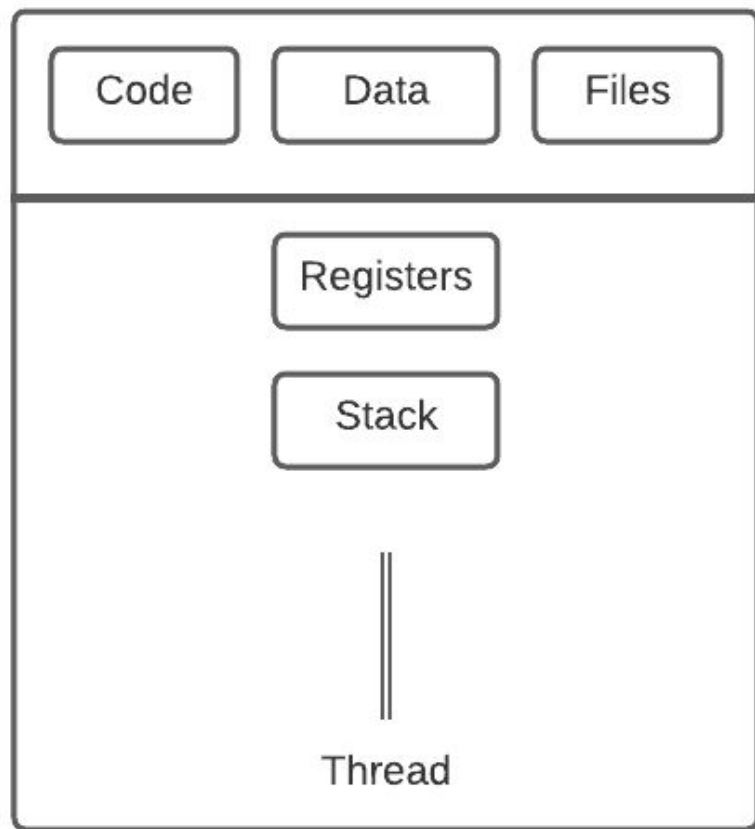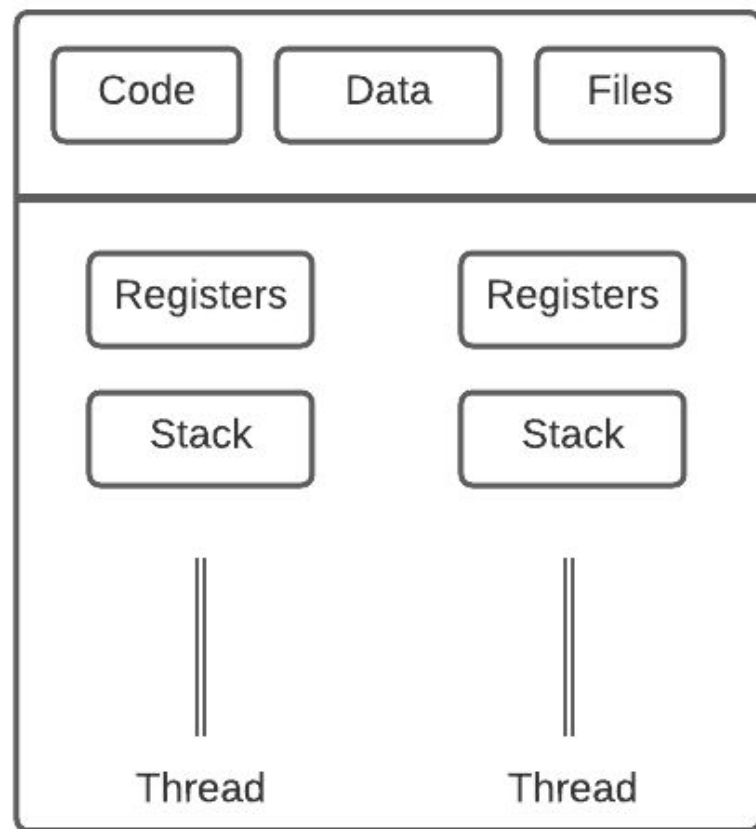
# Key Characteristics of threads

- **Lightweight:** Threads are "lightweight" compared to processes. Creating and managing threads typically requires fewer system resources.
- **Resource Sharing:** Threads within the same process share the same memory space, data, and open files. This allows for efficient communication and resource utilization.
- **Concurrency:** Threads can execute concurrently, either on a single-core processor through time-sharing or on a multi-core processor simultaneously.
- **Independent Execution:** Each thread has its own program counter, stack, and set of registers, allowing it to execute independently.

# Examples of Threads

- **Web Servers:** Use multiple threads to handle incoming requests from multiple clients concurrently.
- **Word Processors:** May use separate threads for spell checking, grammar checking, and background saving.
- **Multimedia Applications:** Can utilize threads for playing audio, decoding video, and rendering graphics simultaneously

| Code | Data | Files |
|------|------|-------|

| Registers |
|-----------|

| Stack |
|-------|

Thread

**Single-threaded Process**

| Code | Data | Files |
|------|------|-------|

| Registers | Registers |
|-----------|-----------|

| Stack | Stack |
|-------|-------|

Thread       Thread
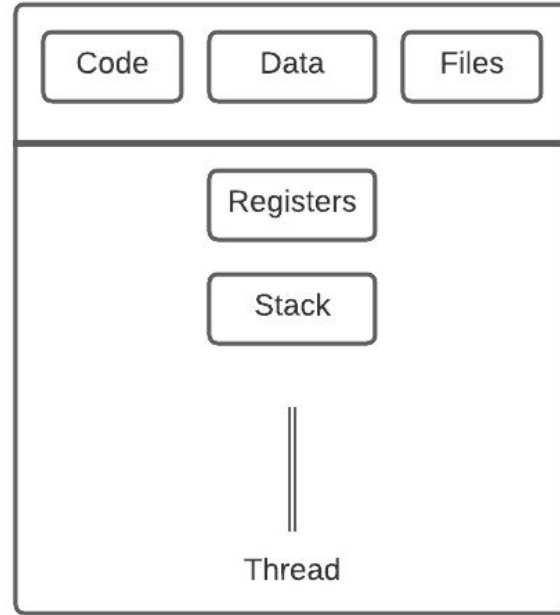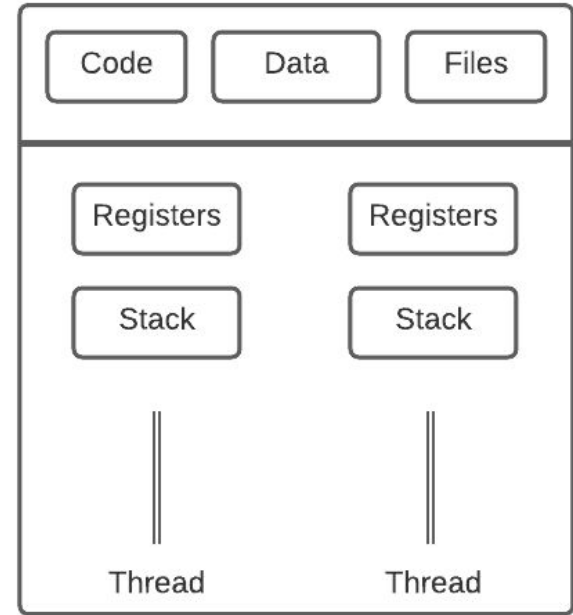
**Multi-threaded Process**

# Process Vs Thread



Single-threaded Process     Multi-threaded Process

# Multithreading

# Introduction

Multithreading is a powerful concept in computer science that allows a single program to execute multiple tasks concurrently.

**What is Multithreading?**

- **Threads:** Imagine a thread as a single flow of execution within a program. A process can have multiple threads, each independently executing a portion of the program's code.
- **Concurrency:** Multithreading enables concurrency, meaning multiple threads can appear to run simultaneously. This doesn't necessarily mean they are truly executing at the exact same time, especially on single-core processors. The operating system rapidly switches between threads, giving the illusion of simultaneous execution.

# Key Concepts

- **Shared Resources:** Threads within the same process share resources like memory, file handles, and open sockets. This allows for efficient communication and data sharing between threads.
- **Independent Execution:** Each thread has its own program counter, stack, and set of registers, allowing it to execute independently.
- **Concurrency vs. Parallelism:**
  - **Concurrency:** Multiple tasks appear to be executing simultaneously, but may actually be interleaved on a single processor.
  - **Parallelism:** Multiple tasks are truly executing simultaneously, typically on a multi-core processor.

# Introduction

**Benefits of Multithreading:**

- **Improved Performance:**
    - **Increased Throughput:** By executing multiple tasks concurrently, applications can process more work in a given time.
    - **Better Responsiveness:** Applications can remain responsive to user input while performing background tasks, such as downloading files or processing data.
- **Resource Utilization:** Threads can efficiently utilize multi-core processors, maximizing their processing power.
- **Simplified Programming:** Multithreading can simplify the design and implementation of concurrent tasks, making code more modular and easier to maintain.
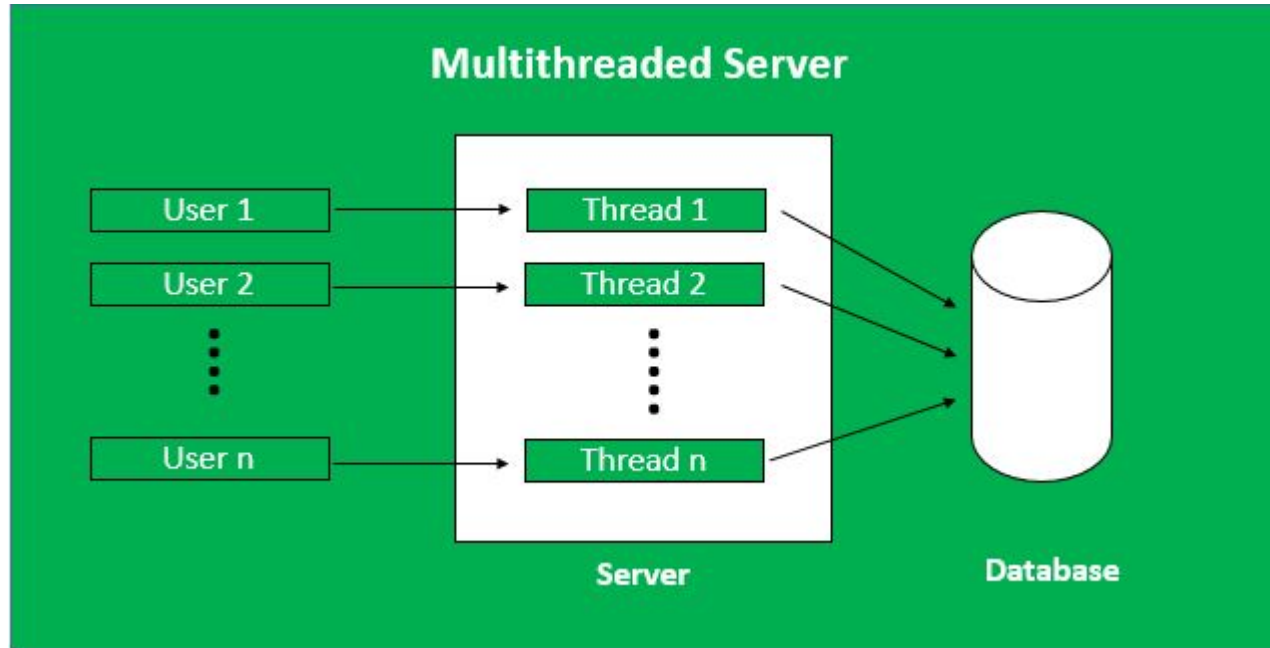
# Challenges of Multithreading

**Challenges of Multithreading:**

- **Synchronization:** Ensuring that multiple threads access and modify shared resources correctly can be challenging and requires careful synchronization mechanisms (e.g., mutexes, semaphores).
- **Deadlock:** A situation where two or more threads are blocked indefinitely, waiting for resources held by each other.
- **Race Conditions:** When the outcome of a computation depends on the relative order in which multiple threads access and modify shared data.

# Examples

- **Web Servers:** Handle multiple client requests concurrently.
- **Word Processors:** Perform background tasks like spell checking and grammar checking while the user continues to type.
- **Multimedia Applications:** Play audio, decode video, and render graphics simultaneously.
- **Operating Systems:** Utilize multiple threads for various system tasks, such as scheduling, memory management, and device handling.

# Example

# Multitasking

# Multitasking

- Multitasking in an RTOS (Real-Time Operating System) is a crucial aspect, allowing for the concurrent execution of multiple tasks or threads within a real-time environment.

**Key Characteristics:**

- **Real-Time Constraints:** RTOSs are specifically designed for systems with strict timing requirements. Multitasking in an RTOS must ensure that tasks are executed within their deadlines.
- **Deterministic Behavior:** RTOSs prioritize predictability and determinism. Multitasking mechanisms should minimize jitter and ensure timely task execution.
- **Resource Management:** RTOSs efficiently manage resources like CPU time, memory, and peripherals among multiple tasks.

# Multitasking Techniques

**Preemptive Scheduling:**

- The RTOS can interrupt a running task at any time and switch to another task with higher priority.
- This ensures that higher-priority tasks are executed promptly, meeting their deadlines.
- Requires careful consideration of interrupt handling and context switching overhead.

**Cooperative Scheduling:**

- Tasks voluntarily relinquish control of the CPU when they are finished with their current execution block or when they explicitly yield control.
- Simpler to implement but can lead to unpredictable behavior if one task monopolizes the CPU.

# Multitasking Techniques

**Priority-Based Scheduling:**

- Tasks are assigned priorities, and the RTOS executes the highest-priority task available.
- Can be preemptive or cooperative.
- Effective for systems with tasks of varying criticality.

**Round-Robin Scheduling:**

- Each task is allocated a fixed time slice (quantum) of CPU time.
- After the time slice expires, the RTOS switches to the next task in the ready queue.
- Suitable for systems with tasks of similar priority.
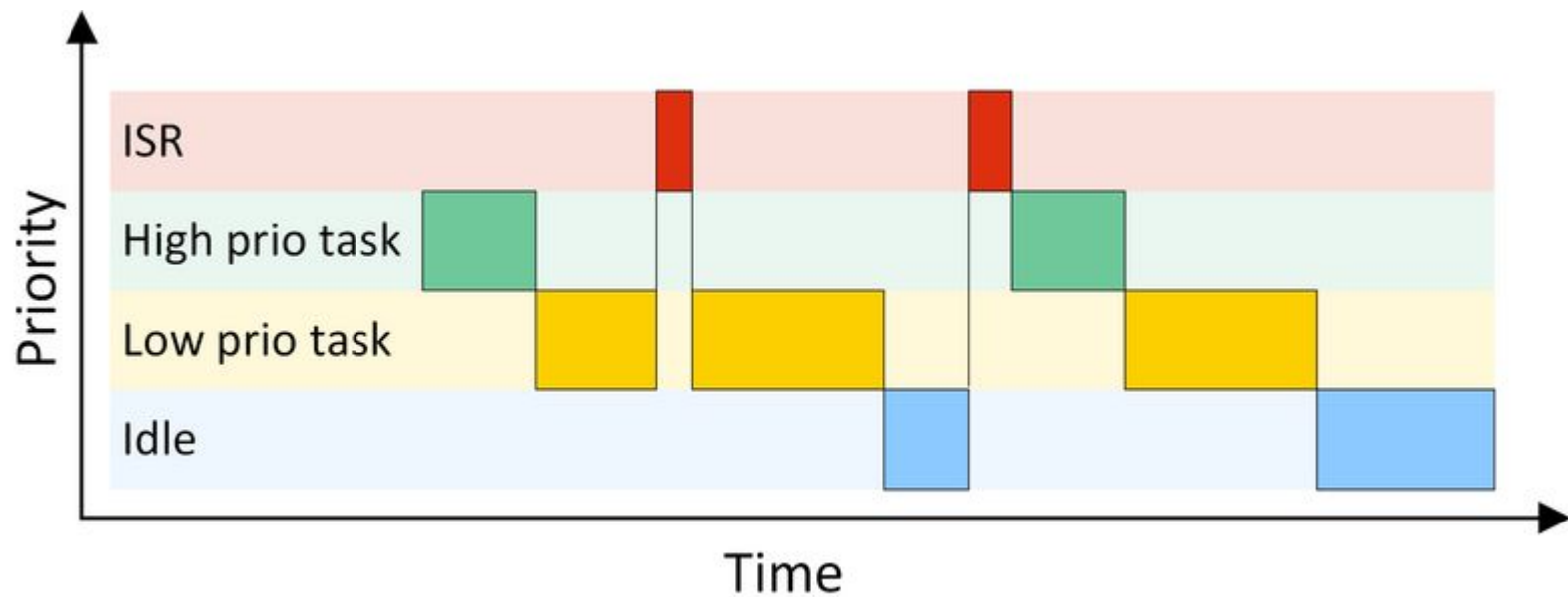
# Challenges in RTOS Multitasking

Challenges in RTOS Multitasking:

- **Real-Time Constraints:** Meeting deadlines is crucial in real-time systems. Multitasking mechanisms must ensure that tasks are executed within their time constraints.
- **Resource Contention:** Multiple tasks may compete for shared resources like memory and peripherals. Proper synchronization mechanisms (e.g., semaphores, mutexes) are essential to prevent race conditions and deadlocks.
- **Context Switching Overhead:** Frequent context switching between tasks can introduce overhead, potentially impacting system performance.

# Examples of RTOS

- **Industrial Automation:** Controlling robotic arms, manufacturing processes, and other real-time systems.
- **Automotive Systems:** Managing engine control, anti-lock brakes, and other safety-critical functions.
- **Aerospace and Defense:** Controlling aircraft flight systems, missile guidance systems, and other critical applications.

# Examples

# Interrupt Handling

# Interrupt Handling in RTOS

Interrupt processing in a Real-Time Operating System (RTOS) is crucial for handling asynchronous events efficiently. Here are some key points:

- **Interrupt Service Routine (ISR)**: When an interrupt occurs, the RTOS invokes an ISR to handle it. The ISR is a special function designed to respond quickly to the interrupt.
- **Minimal Work in ISR**: The ISR should perform minimal work, such as collecting data or resetting the interrupt source. Complex processing should be deferred to a task to avoid delays.
- **Task Scheduling**: After the ISR completes, the RTOS resumes the interrupted task or schedules another task based on priority.

# Interrupt Handling in RTOS

- **Interrupt Priorities**: Interrupts have priorities, and higher-priority interrupts can pre-empt lower-priority tasks. This ensures critical interrupts are handled promptly.
- **Deferred Interrupt Handling**: Some RTOS implementations use deferred interrupt handling, where the ISR records the interrupt reason and unblocks a task for further processing.
- **Interrupt Latency**: Minimizing interrupt latency is essential for real-time performance. This involves optimizing ISRs and using techniques like disabling interrupts for short critical sections.

# Interrupt Handling in RTOS (Terms)

- **Interrupt Vector Table (IVT)**: The IVT is a table that holds the addresses of all the ISRs. When an interrupt occurs, the CPU uses this table to jump to the appropriate ISR.
- **Types of Interrupts**:
  - **Hardware Interrupts**: Generated by hardware devices (e.g., keyboard, timer, network card) to signal the CPU.
  - **Software Interrupts**: Triggered by software instructions to request a system service or switch tasks.
- **Context Switching**: When an interrupt occurs, the current state (context) of the running task is saved, and the context of the ISR is loaded. Once the ISR completes, the saved context is restored to resume the interrupted task.
- **Nested Interrupts**: Higher-priority interrupts can pre-empt lower-priority ISRs. This requires careful handling to ensure the system remains stable and efficient.

# Clocking Communication in RTOS

Clocking communication in an RTOS involves synchronizing tasks based on timing events. This is crucial for ensuring that tasks execute in a timely and predictable manner. Here are some key points:

1. **Timers and Counters**: RTOS often uses hardware timers and counters to generate periodic interrupts. These interrupts can be used to trigger tasks at specific intervals.
2. **Scheduling Algorithms**: The RTOS scheduler uses algorithms to determine the order and timing of task execution based on clock events.
3. **Time-Triggered Scheduling**: Tasks are scheduled to run at precise times, ensuring deterministic behavior. This is common in safety-critical systems.

# Resource Sharing, Deadlock, Priority Inversion

# Resource Sharing

**Concept:** In a multi-threaded or multi-process environment, multiple entities (threads or processes) may need to access shared resources. These resources can include:

- **Memory:** Shared memory regions, global variables.
- **Peripherals:** Hardware devices like printers, network cards, sensors.
- **Files:** Data files accessed by multiple programs.
- **Databases:** Shared databases accessed by concurrent users.

**Challenges:**

**Data Corruption:** If multiple threads/processes access and modify shared data concurrently without proper synchronization, it can lead to data corruption and unpredictable behavior.

**Race Conditions:** A situation where the outcome of an operation depends on the relative order in which multiple threads access and modify shared data.

# Resource Sharing

**Solutions:**

- **Synchronization Mechanisms:** To ensure safe and correct access to shared resources, synchronization mechanisms are essential. These include:
  - **Mutexes:** Ensure only one thread can access a critical section at a time.
  - **Semaphores:** Control access to shared resources, allowing a limited number of threads to access them concurrently.
  - **Condition Variables:** Allow threads to wait for a specific condition to become true.

# Deadlock

**Definition:** A situation where two or more threads are blocked indefinitely, waiting for resources held by each other.

**Example:**

- **Two threads, A and B:**
  - Thread A holds resource 1 and is waiting for resource 2.
  - Thread B holds resource 2 and is waiting for resource 1.
- Neither thread can proceed, resulting in a deadlock.

# Conditions for Deadlock

**Conditions for Deadlock (Coffman Conditions):**

- **Mutual Exclusion:** Resources are held in an exclusive mode (only one thread can use a resource at a time).
- **Hold and Wait:** A process holding at least one resource is requesting additional resources that are currently held by other processes.
- **No Preemption:** Resources cannot be forcibly taken away from a process that holds them.
- **Circular Wait:** A circular chain of two or more processes exists, where each process in the chain is waiting for a resource held by the next process in the chain.

# Deadlock Prevention and Avoidance

**Deadlock Prevention and Avoidance:**

- **Resource Ordering:** Assign a unique number to each resource and require processes to request resources in ascending order.
- **Resource Preallocation:** Allocate all required resources to a process before it begins execution.
- **Deadlock Detection and Recovery:** Periodically check for deadlock situations and take corrective actions (e.g., preempting a process, rolling back transactions)

# Priority Inversion

- A high-priority task is blocked by a low-priority task that is holding a shared resource.
- A medium-priority task preempts the low-priority task, preventing it from releasing the resource and keeping the high-priority task blocked.

**Consequences:**

- Can lead to significant performance degradation and even system instability in real-time systems.

**Solutions:**

- **Priority Inheritance:** Temporarily raise the priority of the low-priority task holding the resource to the priority of the blocked high-priority task. This ensures that the low-priority task can quickly release the resource and avoid delaying the high-priority task.

# Task Synchronization

# Task Synchronization

Consider a system with two tasks: a data acquisition task and a data processing task. The data acquisition task collects sensor data and stores it in a shared buffer. The data processing task reads data from the buffer and performs calculations. Synchronization mechanisms (like semaphores) are crucial to:

- Prevent the data acquisition task from overwriting data while the processing task is reading it.
- Ensure that the data processing task only reads valid data.

**Key Considerations:**

- **Choosing the Right Synchronization Mechanism:** The choice of synchronization mechanism depends on the specific requirements of the application, such as the number of tasks, the nature of resource sharing, and performance constraints.
- **Avoiding Deadlocks:** Careful design and implementation of synchronization mechanisms are essential to avoid deadlocks, where tasks become permanently blocked waiting for resources held by other blocked tasks.
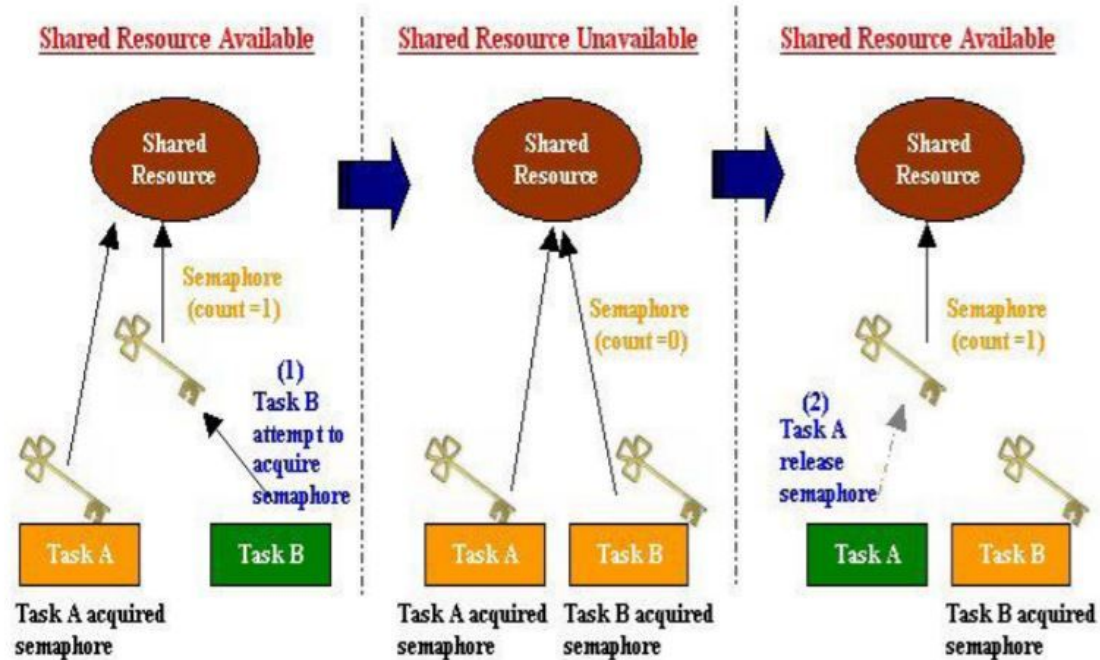
# Task Synchronization

Task synchronization ensures that tasks coordinate their execution to avoid conflicts and ensure correct operation. Here are some common mechanisms:

1. **Semaphores**: Used to control access to shared resources. A semaphore can be binary (0 or 1) or counting (non-negative integer).
2. **Mutexes**: Similar to semaphores but with ownership. Only the task that locked the mutex can unlock it.
3. **Event Flags**: Used to signal events between tasks. A task can wait for a specific flag to be set before proceeding.
4. **Message Queues**: Allow tasks to send and receive messages. This is useful for passing data between tasks without shared memory.
5. **Mailboxes**: Similar to message queues but typically used for sending messages to a single task.
6. **Task Notifications**: Used to signal a task to perform a specific action. This can be more efficient than using semaphores or event flags.

# Use of Semaphores



**Counting Semaphore**

Shared Resource Available | Shared Resource Unavailable | Shared Resource Available

Shared Resource

Semaphore (count =1)

(1) Task B attempt to acquire semaphore

Task A | Task B

Task A acquired semaphore

Shared Resource

Semaphore (count =0)

Task A | Task B

Task A acquired semaphore | Task B acquired semaphore

Shared Resource

Semaphore (count =1)

(2) Task A release semaphore

Task A | Task B

Task B acquired semaphore

# Semaphore

- Semaphores are a powerful synchronization mechanism used in operating systems to control access to shared resources by multiple processes or threads.

**Key Concepts**

- **Semaphore:** Essentially, a semaphore is an integer variable that is used to control access to a shared resource.

**Two Atomic Operations:**

- **wait() (also known as P()):** If the value of the semaphore is greater than 0, it is decremented. If the value is 0, the process that invoked wait() is blocked (suspended) until the semaphore value becomes greater than 0.
- **signal() (also known as V()):** Increments the value of the semaphore. If any processes are blocked waiting for the semaphore, one of them is unblocked and allowed to proceed.

# Types of Semaphores

**Binary Semaphore:**

- Can take on only two values: 0 or 1.
- Used for mutual exclusion, ensuring that only one process can access a shared resource at a time.
- Acts like a lock:If the semaphore is 1 (available), a process can acquire the lock (wait()), use the resource, and then release the lock (signal()).
- If the semaphore is 0 (unavailable), the process must wait until another process releases the lock.

# Types of Semaphores

**Counting Semaphore:**

- Can take on any non-negative integer value.
- Represents the number of available resources.
- For example, if a semaphore has a value of 3, it means that three processes can access the resource concurrently.
- Used to control access to a pool of resources, such as available memory blocks or communication buffers.

# Key Uses of Semaphores

- **Mutual Exclusion:** Ensuring that only one process or thread accesses a shared resource at a time.
- **Synchronization:** Coordinating the execution of multiple processes or threads, such as in producer-consumer problems.
- **Resource Allocation:** Controlling access to limited resources, such as memory, disk drives, and network connections.

# Example

- Imagine a shared printer. A binary semaphore can be used to control access to the printer.
- Initially, the semaphore's value is 1, indicating that the printer is available.
- When a process wants to print, it must first acquire the semaphore (wait()).
- If the semaphore value is 1, the process can proceed to print and then release the semaphore (signal()), making it available for other processes.
- If the semaphore value is 0 (printer is in use), the process is blocked until the current printing job is finished and the semaphore is released.

# Mutual Exclusion through Busy Waiting/Spin Lock

- A simple technique to enforce mutual exclusion, where a thread continuously checks a shared lock variable until it becomes available.
- The thread "spins" in a loop, constantly checking the lock status, hence the name "spinlock."

**Implementation:**

- A shared variable (e.g., lock) is used.
  - lock = 0 indicates the resource is available.
  - lock = 1 indicates the resource is currently being used.

**Acquiring the Lock:**

- A thread repeatedly checks the value of lock.
- If lock is 1, the thread continues to loop (busy wait).
- If lock is 0, the thread sets lock to 1, acquiring the lock.

# Mutual Exclusion through Busy Waiting/Spin Lock

**Advantages:**

- Simple to implement.
- Low overhead when the lock is frequently acquired and released quickly.

**Disadvantages:**

- **Wastes CPU cycles:** While waiting for the lock, the thread consumes CPU time without performing any useful work. This can significantly degrade system performance, especially with multiple threads competing for the lock.
- **Starvation:** If a high-priority thread is constantly preempted by lower-priority threads that are busy waiting, it may never get a chance to acquire the lock, leading to starvation.
- **Inefficient for long waits:** Busy waiting is highly inefficient for long wait times as it wastes significant CPU resources.

# Mutual Exclusion through sleep and wake up

- Mutual exclusion using sleep and wakeup provides a more efficient alternative to busy waiting by allowing threads to relinquish the CPU while waiting for a resource.
- This approach reduces CPU waste and improves overall system performance.
- However, it introduces some overhead due to the context switching involved in sleeping and waking up threads.

# Core Idea

- **Initial State:**
  - A shared variable (e.g., lock) is initialized to 0, indicating that the resource is available.
  - A mechanism to track waiting threads is introduced (e.g., a linked list or a queue).
- **Acquiring the Lock:**
  - If the lock is 1 (resource is in use), the thread calls a system function like sleep() to voluntarily relinquish the CPU.
  - The operating system schedules another thread to run.
  - The blocked thread remains inactive until it is woken up by the thread currently holding the lock.

# Core Idea

- **Critical Section:**

  - If the lock is 0, the thread sets lock to 1, acquiring the lock.
  - The thread then enters the critical section and accesses the shared resource.
- **Releasing the Lock:**

  - After accessing the shared resource, the thread sets lock back to 0.
  - If any threads are waiting for the lock, the thread wakes up one of the waiting threads using a system call like wakeup().

# Advantage

**Advantages:**

- **Reduced CPU Usage:** By sleeping, the thread avoids wasting CPU cycles while waiting for the lock, allowing other threads to execute.
- **Improved Fairness:** Reduces the risk of starvation, as threads are not constantly competing for the CPU.