

<h3>Section 3 : Concurrency/Synchronization</h3> <h4>Synchronization Constructs</h4> <p>Mutex (Mutual Exclusion Lock)</p> <p>A mutex is a simple lock that makes sure only one thread can enter the critical section at a time.</p> <p>How does it work?</p> <ul style="list-style-type: none"> • A thread that wants to enter a critical section locks the mutex • If another thread tries to lock it while its already taken, it has to wait • Once the thread is done, it unlocks the mutex, letting another one in. <p>Ensures mutual exclusion. May cause busy waiting or put threads to sleep depending on implementation</p> <h4>Semaphores</h4> <p>Flexible synchronization tool based on an integer counter. Allows multiple threads to access shared resources up to a limit.</p> <p>How does it work?</p> <ul style="list-style-type: none"> • <code>wait()</code> (P) decreases semaphore value, if value is already 0, thread blocks or spins (waits) • <code>signal()</code> (also called V) increases the value, possibly waking up a waiting thread. <p>Two types</p> <ol style="list-style-type: none"> 1. Counting Semaphore <ul style="list-style-type: none"> • Value can be greater than 1, used when you'd have multiple identical resources 2. Binary Semaphore <ul style="list-style-type: none"> • Value is either 0 or 1, acts exactly like a mutex <p>Analogy: Think of a parking lot</p> <ul style="list-style-type: none"> • <code>wait()</code> means you enter and use a spot (counter goes down) • <code>signal()</code> means you leave, freeing a spot (counter goes up) • If its full (<code>value = 0</code>), you have to wait until someone leaves 	<p>Solutions to the critical section problem</p> <p>a. Software-only (Peterson's Algorithm)</p> <p>Peterson's algorithm is a class software based method that provides a solution to the critical section problem for two threads</p> <p>Key Ideas</p> <ul style="list-style-type: none"> • Each thread has a flag to say "I'm interested in the critical section" • A shared turn variable decides whose turn it is to enter the critical section • Both threads cooperate and check these two variables before proceeding <p>Downside</p> <ul style="list-style-type: none"> • It relies on busy waiting, meaning the thread constantly checks a condition in a loop (also called a spinlock). This wastes CPU cycles while waiting <p>Def - Spinlock - A type of lock here a thread waits in a loop (spins), constantly checking if it can enter the critical section.</p> <ul style="list-style-type: none"> • It keeps checking a condition until the lock becomes available • No sleeping, just checking over and over, which wastes CPU cycles, but can be fast if the wait time is short. <p>b. Hardware Support</p> <p>a. Test and set</p> <p>This is an atomic instruction that test the value of a variable and sets it to true (locked)</p> <ul style="list-style-type: none"> • If the var was already true (locked), the thread waits, if false, locks and proceeds • Useful for implementing basic locks, but also uses busy waiting <p>b. Compare and Swap (CAS)</p> <p>Also atomic, checks whether a value matches an expected value, and only then updates</p> <ul style="list-style-type: none"> • This is more flexible and powerful than test and set • Commonly used in lock free data structure, and modern mutex implementations <p>c. Disabling Interrupts (Kernel Only)</p> <p>Sometimes need to disable interrupts before entering critical section so no other thread can preempt current one.</p> <ul style="list-style-type: none"> • This works only in single CPU systems or inside the kernel • Not recommended in user space or multicore environments 	<p>Blocking vs Busy Waiting</p> <p>Busy Waiting</p> <ul style="list-style-type: none"> • The thread spins checking a condition repeatedly • Wastes CPU cycles - bad for performance <p>Blocking</p> <ul style="list-style-type: none"> • The OS suspends the thread until it can enter the critical section • More efficient, especially with many threads <p>What is Synchronization</p> <p>Synchronization is the coordination of multiple threads or processes so shared resources are accessed safely & program behaves correctly even when many things are happening at once.</p> <p>Requirements for a correct solution</p> <p>Mutual Exclusion - Only one process/thread in the critical section (CS) at a time</p> <p>Why its important - If two threads modify the same resource at once, you get data corruption or race conditions.</p> <p>Progress - If no process is in critical section, and some threads want to enter, the system must allow one of them to go in without unnecessary delay</p> <p>Why its important - We want the system to be fair and efficient. If the critical section is free, we shouldn't block everyone just because the program is indecisive or stuck</p> <p>Bounded waiting - Every thread should get a fair turn. If a thread wants to enter critical section, there should be a limit on how many times other threads cut in line before it gets a chance</p> <p>Why its important - Without bounded waiting, threads can be starved, never get their turn</p>
<p>Classic Synchronization Problems</p> <p>a. Producer-Consumer (Bounded Buffer)</p> <ul style="list-style-type: none"> • Shared buffer between a producer (writes data) and consumer (reads it) • Use semaphores to track <ul style="list-style-type: none"> ◦ Empty slots, Filled slots, Mutual exclusion for buffer access <p>b. Readers-Writers</p> <ul style="list-style-type: none"> • Multiple readers can access data simultaneously • But only one write can update at a time - and no readers during writing <p>Variations</p> <ul style="list-style-type: none"> • First readers' problem (writer starvation possible) • First writers' problem (reader starvation possible) <p>c. Dining Philosophers</p> <ul style="list-style-type: none"> • Five philosophers sharing five forks • Must grab two forks (shared resources) to eat • Can lead to: deadlock, starvation, or livelock depending on implementation 	<p>Critical Section</p> <p>A critical section is any code that accesses shared resources that must not be concurrently modified. They key rule is : only one thread/process should access the critical section at a time.</p> <p>The problem</p> <p>If multiple threads/processes access shared resources w/out coordination → race conditions</p> <p>Def - Race Condition - Happens when two or more threads race to read and write the same data at the same time. It can lead to unpredictable results.</p> <p>Def - Atomic Instruction - A low level CPU operation that runs completely without interruption - meaning no other thread or process can interfere while its happening/</p> <ul style="list-style-type: none"> • It is indivisible - either it finishes entirely or it doesn't happen at all • This is important for synchronization, because it avoids race conditions. 	
<p>Monitors (High level Abstraction)</p> <p>Built in synchronization feature provided by some programming languages that combines</p> <ul style="list-style-type: none"> • A mutex to control who can enter the critical section • Condition variables for threads to wait and signal each other inside the monitor <p>Features</p> <ul style="list-style-type: none"> • Monitors abstract away the complexity - you don't need to manually lock/unlock • They bundle data and synchronization logic together • Internally, they use mutexes and condition variables <p>Def - Condition Variable - Allows threads to wait for certain conditions to become true and lets other threads signal when those conditions are met</p> <p>Def - trylock() - Non blocking version of lock().</p> <ul style="list-style-type: none"> • It tries to acquire the lock, if the lock is already held by another thread, it doesn't wait - it just returns an error (or false), if it succeeds, the thread enters the critical section. 	<p>What is Concurrency</p> <p>Concurrency is the ability of a system to allow multiple tasks (processes or threads) to make progress at the same <i>logical time</i>. It doesn't require multiple CPUs.</p> <p>Logical concurrency can be achieved even on a single core processor by interleaving execution.</p> <p>Def: Interleaving - Alternating the execution of instructions from multiple tasks on a single CPU, so that it appears as if tasks are running at the same time.</p> <ul style="list-style-type: none"> • A single core processor can only execute one single instruction at a time, the OS rapidly switches between tasks - this is called context switching. Doing this frequently creates the illusion of tasks progressing simultaneously, even though only one is running. <p>Types of Concurrency</p> <ul style="list-style-type: none"> • Logical Concurrency - Tasks appear to run simultaneously • Physical Concurrency - Tasks actually run at the same time <p>Why do we care?</p> <ul style="list-style-type: none"> • OS's must handle multiple active entities at once : processes, IO, user interaction, etc • It enables parallelism, responsiveness, and efficiency. <ul style="list-style-type: none"> ◦ Def - Parallelism - When multiple tasks are literally executed at the same time using multiple processing units (multiple cores for ex). Different than interleaving, when tasks just appear to run at the same time but actually take turns. 	