

Operating Systems: Coordination and Modern Systems

Week 1-12

SDB

Autumn 2025

Week 1

Agenda

- ① The Need for Synchronization
- ② Race Conditions and Critical Sections
- ③ Synchronization Requirements
- ④ Software & Hardware Approaches
- ⑤ Busy Waiting and Spinlocks
- ⑥ Real-World Examples
- ⑦ Summary & Key Takeaways
- ⑧ Next Week's Preview

Why Synchronization?

The Problem of Shared Resources

- Concurrent processes and threads often share resources: memory, files, and devices.
- Uncoordinated access to these shared resources can lead to data inconsistency and corruption.
- Synchronization mechanisms ensure a consistent, predictable, and correct outcome.

Why Synchronization?

The Problem of Shared Resources

- Concurrent processes and threads often share resources: memory, files, and devices.
- Uncoordinated access to these shared resources can lead to data inconsistency and corruption.
- Synchronization mechanisms ensure a consistent, predictable, and correct outcome.

Analogy: Imagine two people updating a shared bank account balance at the same time. If one reads the balance, then the other reads it before the first one has written the new balance, the final result will be wrong.

Why Synchronization?

The Problem of Shared Resources

- Concurrent processes and threads often share resources: memory, files, and devices.
- Uncoordinated access to these shared resources can lead to data inconsistency and corruption.
- Synchronization mechanisms ensure a consistent, predictable, and correct outcome.

Analogy: Imagine two people updating a shared bank account balance at the same time. If one reads the balance, then the other reads it before the first one has written the new balance, the final result will be wrong.

Goal: The objective of synchronization is to maintain correctness without sacrificing the benefits of concurrency.

What is a Race Condition?

Definition

A **race condition** occurs when multiple threads or processes access shared data concurrently, and the final outcome depends on the specific, unpredictable order of their execution.

What is a Race Condition?

Definition

A **race condition** occurs when multiple threads or processes access shared data concurrently, and the final outcome depends on the specific, unpredictable order of their execution.

Example: Two threads incrementing a shared counter.

```
int counter = 0; // Shared resource
```

```
// Thread A
```

```
int temp = counter; // Reads 0
temp = temp + 1;     // Becomes 1
counter = temp;      // Writes 1
```

```
// Thread B (interleaves)
```

```
// ...reads 0...
// ...writes 1...
```


What is a Race Condition?

Definition

A **race condition** occurs when multiple threads or processes access shared data concurrently, and the final outcome depends on the specific, unpredictable order of their execution.

Example: Two threads incrementing a shared counter.

```
int counter = 0; // Shared resource
```

```
// Thread A
```

```
int temp = counter; // Reads 0
temp = temp + 1;     // Becomes 1
counter = temp;      // Writes 1
```

```
// Thread B (interleaves)
```

```
// ...reads 0...
// ...writes 1...
```

The final value of 'counter' could be 1 instead of the expected 2. Race conditions are difficult to debug because they are non-deterministic.

The Critical Section Problem

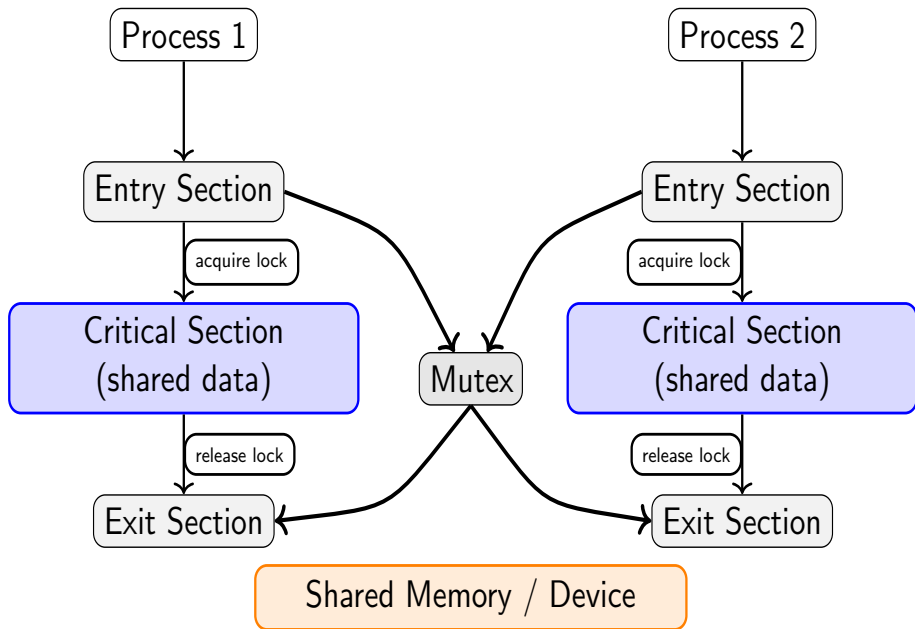
Definition: A **critical section** is a portion of code that accesses shared resources (like a shared variable or file) and must be executed by only one process or thread at a time.

The Critical Section Problem

Definition: A **critical section** is a portion of code that accesses shared resources (like a shared variable or file) and must be executed by only one process or thread at a time. **Three Essential Conditions for a Correct Solution:**

- **Mutual Exclusion:** If one process is executing in its critical section, no other process may be executing in its critical section.
- **Progress:** If no process is in its critical section and some processes wish to enter, then only those processes not in the remainder section can participate in the decision of which will enter next, and this decision cannot be postponed indefinitely.
- **Bounded Waiting:** There must be a limit on the number of times that other processes are allowed to enter their critical sections after a process has made a request to enter its critical section and before that request is granted.

Critical Section Visualization



Busy Waiting and Its Drawbacks

- Many simple locking mechanisms use a loop that continuously checks a condition (a **spinlock**).
- **Busy waiting** (or "spinning") is a technique where a thread sits in a tight loop, constantly consuming CPU cycles while waiting for a lock to be released.

Busy Waiting and Its Drawbacks

- Many simple locking mechanisms use a loop that continuously checks a condition (a **spinlock**).
- **Busy waiting** (or "spinning") is a technique where a thread sits in a tight loop, constantly consuming CPU cycles while waiting for a lock to be released.

Drawbacks:

- Wastes CPU time that could be used by other processes.
- Can lead to priority inversion on single-core systems.

Busy Waiting and Its Drawbacks

- Many simple locking mechanisms use a loop that continuously checks a condition (a **spinlock**).
- **Busy waiting** (or "spinning") is a technique where a thread sits in a tight loop, constantly consuming CPU cycles while waiting for a lock to be released.

Drawbacks:

- Wastes CPU time that could be used by other processes.
- Can lead to priority inversion on single-core systems.

When is it acceptable?

- On multiprocessor systems where the wait time is expected to be very short.
- As the basis for implementing higher-level blocking synchronization primitives.

Hardware & Software Approaches

Approach	Description	Examples
Software	Algorithms that rely on shared variables and logical checks to satisfy the critical section requirements.	Peterson's Algorithm, Bakery Algorithm
Hardware	Special, low-level atomic CPU instructions that perform read-modify-write operations in a single, indivisible step.	'test-and-set', 'compare-and-swap'

Hardware & Software Approaches

Approach	Description	Examples
Software	Algorithms that rely on shared variables and logical checks to satisfy the critical section requirements.	Peterson's Algorithm, Bakery Algorithm
Hardware	Special, low-level atomic CPU instructions that perform read-modify-write operations in a single, indivisible step.	'test-and-set', 'compare-and-swap'

Key Point: Modern operating systems almost exclusively rely on hardware-based atomic instructions as the foundation for all higher-level synchronization tools (like mutexes and semaphores) because they are faster and more reliable than complex software algorithms.

Example: Fixing the Race with a Mutex

```
#include <pthread.h>

int counter = 0;
pthread_mutex_t lock;

void* increment(void* arg) {
    for (int i = 0; i < 100000; ++i) {
        pthread_mutex_lock(&lock); // Acquire the lock
        counter++;
        pthread_mutex_unlock(&lock); // Release the lock
    }
    return NULL;
}
```

Example: Fixing the Race with a Mutex

```
#include <pthread.h>

int counter = 0;
pthread_mutex_t lock;

void* increment(void* arg) {
    for (int i = 0; i < 100000; ++i) {
        pthread_mutex_lock(&lock); // Acquire the lock
        counter++;
        pthread_mutex_unlock(&lock); // Release the lock
    }
    return NULL;
}
```

Analysis: The mutex ensures that only one thread can access the 'counter++' statement at any given time, thus preventing a race condition and guaranteeing the correct final value.

Key Takeaways

- **Synchronization** is essential in concurrent systems to prevent data corruption from **race conditions**.
- The **critical section problem** requires that solutions satisfy **mutual exclusion**, **progress**, and **bounded waiting**.
- Simple locking can lead to **busy waiting**, which is inefficient on single-core systems.
- Modern OSes use **atomic hardware instructions** as the efficient, reliable foundation for building higher-level synchronization primitives.

Next Week Preview: Semaphores and Classical Problems

Building on our understanding of critical sections and hardware support, we will delve into the powerful tools used for synchronization.

- **Semaphores:** The core concept of a semaphore and the 'wait()' / 'signal()' operations.
- **Mutexes vs. Semaphores:** Understanding the differences and when to use each.
- **Classical Synchronization Problems:** Solving famous synchronization challenges like the Producer-Consumer Problem and the Dining Philosophers Problem.

Appendix: Quiz

Conceptual Questions

- ① What is the difference between a race condition and a deadlock?
- ② In the context of the critical section problem, explain the difference between the 'Progress' and 'Bounded Waiting' conditions.
- ③ Why is busy waiting a significant problem on a single-core CPU but less of an issue on a multi-core CPU?

Appendix: Exercises

Coding Exercise

- ① Write a C program that creates two threads. Each thread should attempt to increment a shared global integer variable 100,000 times.
- ② Run the program without any synchronization and observe the final value. It should be less than 200,000.
- ③ Add a mutex to the code to protect the shared variable. Re-run the program and confirm that the final value is exactly 200,000.

Appendix: Advanced Topics to Explore

Further Reading

- **Peterson's Algorithm:** A classic software-only solution for the critical section problem for two processes.
- **The Bakery Algorithm:** An elegant software solution that extends to 'n' processes.
- **Memory Models:** The relationship between hardware synchronization instructions and the memory consistency model of a multi-core processor.
- **Transactional Memory:** An advanced research topic that allows for atomic execution of a block of code, offering a higher-level abstraction for concurrency.

Week 2

Agenda

- ① Semaphores: Concept and Usage
- ② wait() and signal() Semantics
- ③ Binary vs Counting Semaphores
- ④ Classical Problems:
 - ▶ Bounded Buffer
 - ▶ Dining Philosophers
 - ▶ Readers–Writers
- ⑤ Summary & Key Takeaways
- ⑥ Next Week's Preview

What is a Semaphore?

A Signaling Mechanism

A semaphore is an integer variable used for signaling between processes, originally proposed by Edsger W. Dijkstra.

What is a Semaphore?

A Signaling Mechanism

A semaphore is an integer variable used for signaling between processes, originally proposed by Edsger W. Dijkstra.

The Two Atomic Operations:

- **wait(S):** Decrements the semaphore value 'S'. If 'S' becomes negative, the process blocks. Also known as 'P()' or 'down()'.
- **signal(S):** Increments the semaphore value 'S'. If 'S' was negative before the increment, it unblocks a process waiting on 'S'. Also known as 'V()' or 'up()'.

What is a Semaphore?

A Signaling Mechanism

A semaphore is an integer variable used for signaling between processes, originally proposed by Edsger W. Dijkstra.

The Two Atomic Operations:

- **wait(S):** Decrements the semaphore value 'S'. If 'S' becomes negative, the process blocks. Also known as 'P()' or 'down()'.
- **signal(S):** Increments the semaphore value 'S'. If 'S' was negative before the increment, it unblocks a process waiting on 'S'. Also known as 'V()' or 'up()'.

Analogy: Think of a semaphore as a key counter. 'wait()' decrements the available keys (access to a resource), and 'signal()' returns a key. If no keys are available, a process must wait.

Binary vs Counting Semaphores

Feature	Binary Semaphore (Mutex)	Counting Semaphore
Value Range	0 or 1	Any non-negative integer
Purpose	Enforces mutual exclusion, protecting a critical section	Controls access to a limited pool of resources
Use Case	Protecting a shared variable from a race condition	Synchronizing a producer-consumer buffer with N slots

Binary vs Counting Semaphores

Feature	Binary Semaphore (Mutex)	Counting Semaphore
Value Range	0 or 1	Any non-negative integer
Purpose	Enforces mutual exclusion, protecting a critical section	Controls access to a limited pool of resources
Use Case	Protecting a shared variable from a race condition	Synchronizing a producer-consumer buffer with N slots

Relationship: A binary semaphore can be seen as a special case of a counting semaphore, initialized to 1. Many systems implement mutexes as optimized binary semaphores.

The Bounded Buffer Problem

The Problem: A producer and consumer share a fixed-size buffer. The producer generates data and puts it in the buffer; the consumer takes data out.

The Bounded Buffer Problem

The Problem: A producer and consumer share a fixed-size buffer. The producer generates data and puts it in the buffer; the consumer takes data out. **Synchronization with Semaphores:**

- `mutex` (binary, initialized to 1): For mutual exclusion, ensuring only one process can access the buffer at a time.
- `empty` (counting, initialized to N): Counts the number of empty slots in the buffer.
- `full` (counting, initialized to 0): Counts the number of filled slots in the buffer.

The Bounded Buffer Problem

The Problem: A producer and consumer share a fixed-size buffer. The producer generates data and puts it in the buffer; the consumer takes data out. **Synchronization with Semaphores:**

- **mutex** (binary, initialized to 1): For mutual exclusion, ensuring only one process can access the buffer at a time.
- **empty** (counting, initialized to N): Counts the number of empty slots in the buffer.
- **full** (counting, initialized to 0): Counts the number of filled slots in the buffer.

Producer Code Fragment:

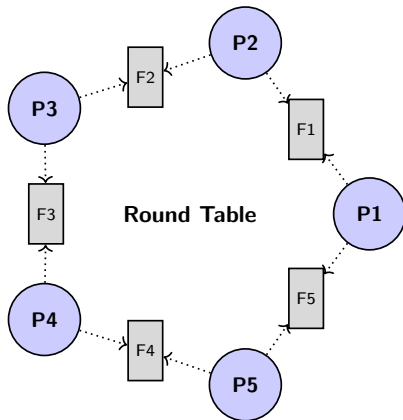
```
wait(empty);      // Wait for an empty slot
wait(mutex);      // Acquire exclusive access to buffer
    // Insert an item into the buffer
signal(mutex);    // Release exclusive access
signal(full);      // Signal that a slot is full
```

The Dining Philosophers Problem

The Problem: Five philosophers sit at a round table, with a single chopstick between each pair. Each philosopher must pick up both adjacent chopsticks to eat.

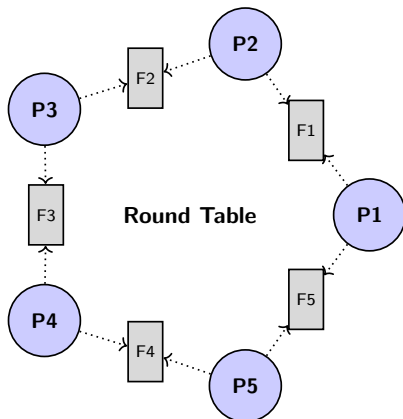
The Dining Philosophers Problem

The Problem: Five philosophers sit at a round table, with a single chopstick between each pair. Each philosopher must pick up both adjacent chopsticks to eat.



The Dining Philosophers Problem

The Problem: Five philosophers sit at a round table, with a single chopstick between each pair. Each philosopher must pick up both adjacent chopsticks to eat.



The Risk: This problem is a classic example of **deadlock**. If all five philosophers pick up their right chopstick simultaneously, they will all be waiting indefinitely for their left chopstick.

The Readers–Writers Problem

The Problem: A shared data object is accessed by multiple processes.

- **Readers:** Can read the data concurrently with other readers.
- **Writers:** Must have exclusive access to the data.

The Readers–Writers Problem

The Problem: A shared data object is accessed by multiple processes.

- **Readers:** Can read the data concurrently with other readers.
- **Writers:** Must have exclusive access to the data.

Semaphores Used:

- `mutex` (binary, initialized to 1): Used to protect the shared 'readcount' variable from race conditions.
- `rw_mutex` (binary, initialized to 1): Enforces mutual exclusion for writers and the first/last reader.
- 'readcount' (shared integer): Tracks how many readers are currently in the critical section.

Reader Code Fragment

```
wait(mutex);
readcount++;
if (readcount == 1)
    wait(rw_mutex); // First reader locks out writers
signal(mutex);

read(); // Reading is the critical section

wait(mutex);
readcount--;
if (readcount == 0)
    signal(rw_mutex); // Last reader allows writers
signal(mutex);
```


Key Takeaways

- **Semaphores** are a powerful and simple primitive for synchronization, providing both mutual exclusion and signaling.
- **Binary semaphores** are used for mutexes; **counting semaphores** are for resource pools.
- The **classical problems** (Bounded Buffer, Dining Philosophers, Readers–Writers) are excellent illustrations of how semaphores can be used to solve complex synchronization challenges and prevent issues like deadlock.

Next Week Preview: Deadlock Modeling and Conditions

Building on our understanding of synchronization challenges, we will formally define and model deadlocks.

- **Resource Allocation Graphs:** A visual model for detecting deadlocks.
- **Coffman's Conditions:** The four necessary conditions that must be met for a deadlock to occur.
- **Deadlock Examples:** We'll revisit the Dining Philosophers problem and analyze its deadlock potential in detail.

Appendix: Quiz

Conceptual Questions

- ① Explain the key difference in purpose between a mutex and a semaphore.
- ② How does a semaphore's 'wait()' operation differ from a busy-waiting loop?
- ③ In the Readers–Writers problem, why is a separate 'mutex' needed to protect the 'readcount' variable?

Appendix: Exercises

Coding Exercise

- ① Extend the Bounded Buffer producer code to include the consumer's logic. Write a full program that simulates producers and consumers passing data through a shared buffer.
- ② Use counting semaphores to manage the 'empty' and 'full' slots and a binary semaphore to enforce mutual exclusion on the buffer.
- ③ Experiment with changing the buffer size and the number of producers and consumers.

Appendix: Advanced Topics to Explore

Further Reading

- **Monitors:** A higher-level synchronization primitive that combines a mutex and condition variables into a single object-oriented construct.
- **Futexes (Fast User-space Mutex):** A low-level synchronization tool in Linux that avoids kernel intervention for simple cases, improving performance.
- **Condition Variables:** Primitives used with mutexes to allow threads to wait for a specific condition to become true without busy-waiting.

Week 3

Agenda

- ① Introduction to Deadlocks
- ② System Model and Resources
- ③ Coffman's Necessary Conditions
- ④ Resource Allocation Graph (RAG)
- ⑤ Examples and Scenarios
- ⑥ Summary & Key Takeaways
- ⑦ Next Week's Preview

What is a Deadlock?

The Core Problem

A set of processes is in a **deadlock** if each process in the set is waiting for a resource that is currently held by another process in the set.

What is a Deadlock?

The Core Problem

A set of processes is in a **deadlock** if each process in the set is waiting for a resource that is currently held by another process in the set.

Analogy: A classic example is a traffic gridlock at a four-way intersection where each car is waiting for the one in front of it to move. No car can proceed, and the entire system is stuck.

What is a Deadlock?

The Core Problem

A set of processes is in a **deadlock** if each process in the set is waiting for a resource that is currently held by another process in the set.

Analogy: A classic example is a traffic gridlock at a four-way intersection where each car is waiting for the one in front of it to move. No car can proceed, and the entire system is stuck. **Examples in Computing:**

- Process P1 holds Mutex A and is waiting for Mutex B, which is held by P2.
- Process P2 holds Mutex B and is waiting for Mutex A, which is held by P1.

System Model and Resource Types

- Processes request and use various types of **resources**, which can include CPU cycles, memory blocks, I/O devices, file locks, or synchronization primitives like semaphores.
- Each resource type can have one or more identical **instances**.
- The OS must allocate these resources in a way that avoids deadlock.

System Model and Resource Types

- Processes request and use various types of **resources**, which can include CPU cycles, memory blocks, I/O devices, file locks, or synchronization primitives like semaphores.
- Each resource type can have one or more identical **instances**.
- The OS must allocate these resources in a way that avoids deadlock.

Process Resource Protocol

Each process follows a simple protocol when using a resource:

- ① **Request:** The process requests the resource.
- ② **Use:** The process uses the resource.
- ③ **Release:** The process releases the resource.

Coffman's Necessary Conditions

For a deadlock to occur, all four of these conditions must hold simultaneously in the system.

Condition	Description
Mutual Exclusion	At least one resource must be held in a non-sharable mode. Only one process can use it at a time.
Hold and Wait	A process must be holding at least one resource and waiting to acquire additional resources that are currently held by other processes.
No Preemption	A resource cannot be forcibly taken away from a process that is holding it. The process must release it voluntarily.
Circular Wait	There must exist a set of waiting processes (P_0, P_1, \dots, P_n) such that P_0 is waiting for a resource held by P_1 , P_1 is waiting for a resource held by P_2 , and so on, with P_n waiting for a resource held by P_0 .

Resource Allocation Graph (RAG)

Modeling Deadlocks Visually

The Resource Allocation Graph is a directed graph used to model and detect potential deadlocks.

Resource Allocation Graph (RAG)

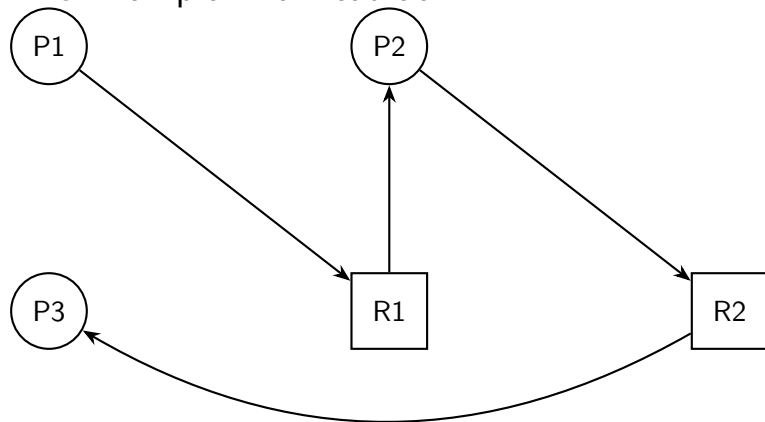
Modeling Deadlocks Visually

The Resource Allocation Graph is a directed graph used to model and detect potential deadlocks.

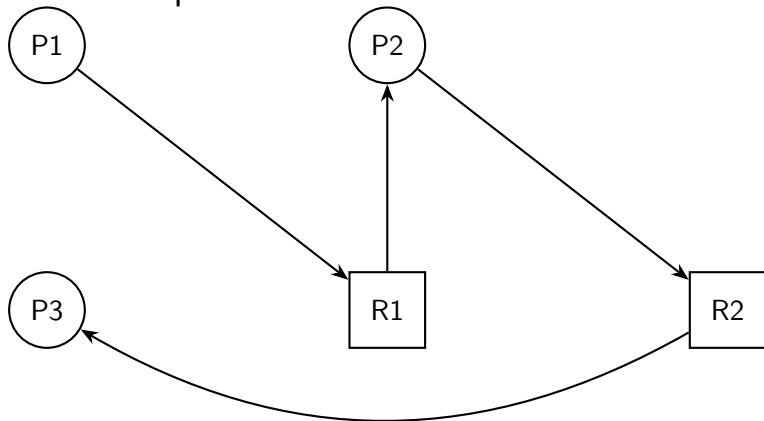
Components:

- **Process Nodes:** Represented by circles (\bigcirc).
- **Resource Nodes:** Represented by squares (\square). Dots within the square denote instances of the resource type.
- **Request Edge:** A directed edge from a process to a resource ($P_i \rightarrow R_j$).
- **Assignment Edge:** A directed edge from a resource instance to a process ($R_j \rightarrow P_i$).

RAG Example: No Deadlock



RAG Example: No Deadlock

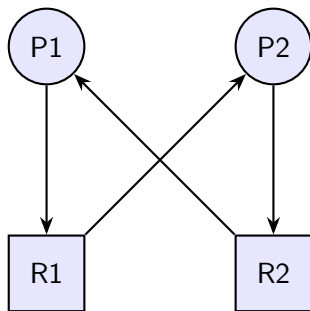


Analysis:

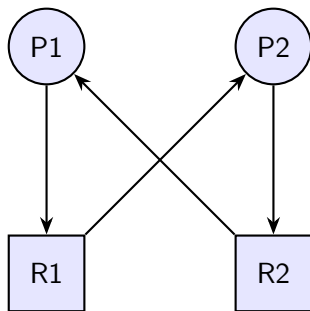
- There are no cycles in this graph.
- The path from P1 to R1 to P2 to R2 to P3 is linear.

Conclusion: No cycle in the RAG implies that the system is not in a deadlock state.

RAG Example: Deadlock



RAG Example: Deadlock



Deadlock Analysis

- **Cycle:** A clear cycle exists: $P_1 \rightarrow R_1 \rightarrow P_2 \rightarrow R_2 \rightarrow P_1$.
- **Conclusion:** Because each resource has only a single instance, a cycle in the RAG is a necessary and sufficient condition for a deadlock to exist.

Key Takeaways

- Deadlocks are a state where processes are permanently blocked, waiting for resources held by each other.
- **Coffman's four conditions** (Mutual Exclusion, Hold and Wait, No Preemption, Circular Wait) are necessary for a deadlock to occur.
- The **Resource Allocation Graph (RAG)** is a valuable tool for modeling resource requests and assignments.
- A cycle in the RAG indicates the **possibility** of a deadlock; for single-instance resources, a cycle is a definitive sign of deadlock.

Next Week Preview: Deadlock Prevention and Avoidance

Building on our understanding of how deadlocks are formed, we will explore strategies for managing them.

- **Deadlock Prevention:** Strategies to ensure at least one of Coffman's conditions is never met.
- **Deadlock Avoidance:** The Banker's Algorithm and how to dynamically check for safe states.
- **Deadlock Detection and Recovery:** How to identify a deadlock after it has occurred and recover the system.

Appendix: Quiz

Conceptual Questions

- ① Explain how violating the "No Preemption" condition can prevent a deadlock.
- ② If a resource type has multiple instances, is a cycle in the Resource Allocation Graph a guarantee of a deadlock? Explain why or why not.
- ③ Consider a system with a single resource and two processes. Can a deadlock occur? Justify your answer using Coffman's conditions.

Appendix: Exercises

RAG Analysis Exercise

- ① Draw a Resource Allocation Graph for the following state:
 - ▶ Process P1 holds R1, requests R2.
 - ▶ Process P2 holds R2, requests R3.
 - ▶ Process P3 holds R3, requests R1.
- ② Does a deadlock exist in this scenario? If so, identify the processes and resources involved in the cycle.

Appendix: Advanced Topics to Explore

Further Reading

- **Livelock:** A state similar to deadlock where processes are not blocked but are continuously changing their state in response to each other, preventing progress.
- **Starvation:** A situation where a process is repeatedly denied access to a resource, even though it's available, because other processes are always given priority.
- **Wait-for Graph:** A simplification of the RAG used for deadlock detection when all resources have a single instance.

Week 4

Agenda

- ① Deadlock Prevention Overview
- ② Violating Coffman's Conditions
- ③ Safe vs. Unsafe States
- ④ The Banker's Algorithm
- ⑤ Limitations and Real-World OS Practices
- ⑥ Summary & Key Takeaways
- ⑦ Next Week's Preview

Deadlock Prevention: The Strategy

Proactive Elimination

Deadlock prevention is a set of design-time policies and protocols that ensure a system never enters a deadlock state by eliminating one of the four necessary **Coffman's conditions**.

Deadlock Prevention: The Strategy

Proactive Elimination

Deadlock prevention is a set of design-time policies and protocols that ensure a system never enters a deadlock state by eliminating one of the four necessary **Coffman's conditions**.

Key Characteristics:

- It's a proactive approach, built into the system's design.
- It's often simple to implement but can lead to poor resource utilization.
- The goal is to make deadlock impossible by violating a fundamental condition.

Violating Coffman's Conditions

Condition to Break	Strategy
Mutual Exclusion	Allow resources to be shared (e.g., read-only files). This is often not possible for non-sharable resources like a printer or mutex.
Hold and Wait	Force processes to request all required resources at once, or release all held resources before requesting new ones. This can lead to resource starvation .
No Preemption	Allow the OS to preempt (forcefully take) a resource from a process that is holding it and waiting for others. This is only feasible for resources where state can be easily saved and restored (e.g., CPU registers).
Circular Wait	Impose a total global ordering on all resource types. A process can only request resources in an increasing order of enumeration.

Violating Coffman's Conditions

Condition to Break	Strategy
Mutual Exclusion	Allow resources to be shared (e.g., read-only files). This is often not possible for non-sharable resources like a printer or mutex.
Hold and Wait	Force processes to request all required resources at once, or release all held resources before requesting new ones. This can lead to resource starvation .
No Preemption	Allow the OS to preempt (forcefully take) a resource from a process that is holding it and waiting for others. This is only feasible for resources where state can be easily saved and restored (e.g., CPU registers).
Circular Wait	Impose a total global ordering on all resource types. A process can only request resources in an increasing order of enumeration.

Most Practical Strategies: In general-purpose operating systems, breaking the **Hold and Wait** or **Circular Wait** conditions are the most practical approaches.

Deadlock Avoidance: Safe and Unsafe States

The Safe State Concept

A system is in a **safe state** if there exists a **safe sequence** of all processes, where each process can acquire all its resources and run to completion without causing a deadlock.

Deadlock Avoidance: Safe and Unsafe States

The Safe State Concept

A system is in a **safe state** if there exists a **safe sequence** of all processes, where each process can acquire all its resources and run to completion without causing a deadlock.

- An **unsafe state** is a state where the OS cannot guarantee a safe sequence exists.
- It's important to note that an unsafe state is **not** a deadlock, but it has the potential to lead to one.
- **Deadlock avoidance** is the strategy of ensuring the system never enters an unsafe state by carefully managing resource allocation.

The Banker's Algorithm: Intuition

Concept: The Banker's Algorithm is a classic deadlock avoidance algorithm designed by Dijkstra. It's an analogy to a banker who only grants loans if they can guarantee that all customers can eventually pay back their loans.

The Banker's Algorithm: Intuition

Concept: The Banker's Algorithm is a classic deadlock avoidance algorithm designed by Dijkstra. It's an analogy to a banker who only grants loans if they can guarantee that all customers can eventually pay back their loans. **Assumptions:**

- Each process must declare its **maximum possible resource needs** upfront.
- The system tracks the current **allocated** resources, the **maximum** needs, and the current **available** resources.

The Banker's Algorithm: Intuition

Concept: The Banker's Algorithm is a classic deadlock avoidance algorithm designed by Dijkstra. It's an analogy to a banker who only grants loans if they can guarantee that all customers can eventually pay back their loans. **Assumptions:**

- Each process must declare its **maximum possible resource needs** upfront.
- The system tracks the current **allocated** resources, the **maximum** needs, and the current **available** resources.

How it Works: When a process requests resources, the Banker's Algorithm performs a **safety check**. The request is only granted if doing so leaves the system in a safe state. Otherwise, the process must wait.

Banker's Algorithm – Allocation Example

Scenario: A system with 10 instances of a single resource type.

Process	Max Need	Allocated	Need
P0	7	0	7
P1	5	2	3
P2	3	2	1

$$\text{Available} = 10 - (0+2+2) = 6$$

Banker's Algorithm – Allocation Example

Scenario: A system with 10 instances of a single resource type.

Process	Max Need	Allocated	Need
P0	7	0	7
P1	5	2	3
P2	3	2	1

Available = $10 - (0+2+2) = 6$ Safety Analysis:

- We have 6 available instances.
- Can we satisfy P2's need of 1? Yes. Run P2 to completion. Available becomes $6 + 2 = 8$.
- Now can we satisfy P1's need of 3? Yes. Run P1 to completion. Available becomes $8 + 2 = 10$.
- Finally, can we satisfy P0's need of 7? Yes. Run P0. Available becomes $10 + 0 = 10$.

Banker's Algorithm – Allocation Example

Scenario: A system with 10 instances of a single resource type.

Process	Max Need	Allocated	Need
P0	7	0	7
P1	5	2	3
P2	3	2	1

Available = $10 - (0+2+2) = 6$ Safety Analysis:

- We have 6 available instances.
- Can we satisfy P2's need of 1? Yes. Run P2 to completion. Available becomes $6 + 2 = 8$.
- Now can we satisfy P1's need of 3? Yes. Run P1 to completion. Available becomes $8 + 2 = 10$.
- Finally, can we satisfy P0's need of 7? Yes. Run P0. Available becomes $10 + 0 = 10$.

Conclusion: The sequence **(P2, P1, P0)** is a safe sequence. The system is in a safe state.

Limitations of the Banker's Algorithm

- **Requires Prior Knowledge:** Processes must declare their maximum resource needs in advance, which is often not possible.
- **High Overhead:** The safety algorithm must be run every time a process requests resources, which is computationally expensive.
- **Static Process Set:** Assumes a fixed number of processes that don't change dynamically.
- **Infrequent Use:** Due to these limitations, the Banker's Algorithm is rarely used in its full form in modern, general-purpose operating systems.

How Real OSes Handle Deadlocks

General-purpose operating systems typically do not implement full deadlock avoidance or prevention. Instead, they use a mix of practical strategies:

- **Resource Ordering:** For critical shared resources like mutexes, a strict, global ordering is often enforced to break the circular wait condition.
- **Timeout-based Locking:** A process requests a lock with a timeout. If it cannot acquire the lock within the time limit, it releases all held locks and tries again later.
- **Hierarchical Locking:** Locks are grouped into a hierarchy, and processes must acquire locks in a specific order (e.g., always acquire a higher-level lock before a lower-level one).
- **Deadlock Detection and Recovery:** A background process periodically checks for deadlocks and, if one is found, takes action to resolve it (e.g., terminating a process). This is what most OSes do.

Key Takeaways

- **Deadlock Prevention** aims to make deadlock impossible by violating one of the four necessary conditions.
- **Deadlock Avoidance** is a more dynamic approach that uses information about future requests to ensure the system remains in a **safe state**.
- The **Banker's Algorithm** is a classic example of an avoidance algorithm, but its high overhead and strict requirements limit its practical use.
- **Real-world OSes** often use a combination of partial prevention strategies and deadlock detection rather than full prevention or avoidance.

Next Week Preview: Deadlock Detection and Recovery

In our final week on deadlocks, we will examine how to deal with them when they are not prevented or avoided.

- **Detection Algorithms:** The logic behind how the OS can identify a deadlock after it has occurred.
- **Recovery Techniques:** Strategies for breaking a deadlock, including process termination, resource preemption, and rollback.
- **Practical Implementation:** A look at how detection and recovery are used in database systems.

Appendix: Quiz

Conceptual Questions

- ① Explain the key difference between deadlock prevention and deadlock avoidance.
- ② A system is in an unsafe state. Does this mean a deadlock has occurred? Explain.
- ③ Describe a real-world scenario where a strict global resource ordering could be applied to prevent a deadlock.

Appendix: Exercises

Banker's Algorithm Practice

- ① A system has 12 instances of a resource type. The following allocation is made:
 - ▶ P0: Allocated 5, Max 10
 - ▶ P1: Allocated 2, Max 4
 - ▶ P2: Allocated 2, Max 9
- ② Is the system in a safe state? If so, provide a safe sequence.
- ③ If P1 requests one more resource, will the Banker's Algorithm grant the request? Justify your answer.

Appendix: Advanced Topics to Explore

Further Reading

- **Wait-for Graphs:** A simplified tool for deadlock detection in single-instance resource systems.
- **Graph Reduction:** The core algorithm used to find deadlocks in Resource Allocation Graphs.
- **Deadlock in Distributed Systems:** The unique challenges of detecting and resolving deadlocks when resources and processes are spread across multiple machines.

Week 5

Agenda

- ① Deadlock Detection Strategy
- ② Detection with Resource Allocation Graph (RAG)
- ③ Detection Algorithm for Multiple Instances
- ④ Recovery Strategies
- ⑤ Evaluation and Real-World Examples
- ⑥ Summary & Key Takeaways
- ⑦ Next Week's Preview

Deadlock Detection Strategy

Why and When to Use Detection

A deadlock detection and recovery strategy is used in systems where deadlocks are not a frequent occurrence. Instead of preventing deadlocks with strict policies, the system allows them to happen, then identifies and resolves them.

Deadlock Detection Strategy

Why and When to Use Detection

A deadlock detection and recovery strategy is used in systems where deadlocks are not a frequent occurrence. Instead of preventing deadlocks with strict policies, the system allows them to happen, then identifies and resolves them.

Trade-offs:

- **Advantages:** Maximizes resource utilization and concurrency by not imposing strict resource allocation policies.
- **Disadvantages:** The system must be able to periodically check for deadlocks and incur the overhead of recovery, which can be complex and may require aborting processes.

Detection with RAG (Single-Instance Resources)

- In a system where each resource type has only a **single instance**, a deadlock is equivalent to a **cycle in the Resource Allocation Graph (RAG)**.
- The detection strategy involves building the RAG and periodically running a cycle-detection algorithm on it.

Detection with RAG (Single-Instance Resources)

- In a system where each resource type has only a **single instance**, a deadlock is equivalent to a **cycle in the Resource Allocation Graph (RAG)**.
- The detection strategy involves building the RAG and periodically running a cycle-detection algorithm on it.

Algorithm:

- ① Construct the RAG based on current resource allocations and requests.
- ② Use a standard graph traversal algorithm, such as **Depth-First Search (DFS)**, starting from each process node.
- ③ If a back edge is found during the traversal, a cycle exists, and the system is in a deadlock.

Detection with RAG (Single-Instance Resources)

- In a system where each resource type has only a **single instance**, a deadlock is equivalent to a **cycle in the Resource Allocation Graph (RAG)**.
- The detection strategy involves building the RAG and periodically running a cycle-detection algorithm on it.

Algorithm:

- ① Construct the RAG based on current resource allocations and requests.
- ② Use a standard graph traversal algorithm, such as **Depth-First Search (DFS)**, starting from each process node.
- ③ If a back edge is found during the traversal, a cycle exists, and the system is in a deadlock.

Example: The Dining Philosophers problem RAG will form a cycle, signifying a deadlock.

Detection Algorithm – Multiple Instances

For systems with multiple instances of a resource type, a cycle in the RAG is a necessary but not sufficient condition for deadlock. A more general algorithm is needed.

Detection Algorithm – Multiple Instances

For systems with multiple instances of a resource type, a cycle in the RAG is a necessary but not sufficient condition for deadlock. A more general algorithm is needed. **Input Data Structures:**

- **Available Vector (V):** A vector of size ' m ' (number of resource types) showing the number of available instances of each type.
- **Allocation Matrix (A):** An ' $n \times m$ ' matrix where ' $A[i,j]$ ' is the number of instances of resource ' j ' allocated to process ' i '.
- **Request Matrix (R):** An ' $n \times m$ ' matrix where ' $R[i,j]$ ' is the number of instances of resource ' j ' that process ' i ' is currently requesting.

Detection Algorithm – Multiple Instances

For systems with multiple instances of a resource type, a cycle in the RAG is a necessary but not sufficient condition for deadlock. A more general algorithm is needed. **Input Data Structures:**

- **Available Vector (V):** A vector of size ' m ' (number of resource types) showing the number of available instances of each type.
- **Allocation Matrix (A):** An ' $n \times m$ ' matrix where ' $A[i,j]$ ' is the number of instances of resource ' j ' allocated to process ' i '.
- **Request Matrix (R):** An ' $n \times m$ ' matrix where ' $R[i,j]$ ' is the number of instances of resource ' j ' that process ' i ' is currently requesting.

Algorithm Logic: The algorithm is similar to the safety algorithm in Banker's. It tries to find a sequence of processes that can finish. Any process that cannot finish is deadlocked.

Example: Deadlock Detection Matrix

Scenario: A system with 3 resource types (A, B, C).

Available = (3, 3, 2)

Process	Allocation			Request			Finish?
	A	B	C	A	B	C	
P0	0	1	0	0	0	0	Yes
P1	2	0	0	2	0	2	No
P2	3	0	2	0	0	0	Yes
P3	2	1	1	1	0	0	No

Example: Deadlock Detection Matrix

Scenario: A system with 3 resource types (A, B, C).

Available = (3, 3, 2)

Process	Allocation			Request			Finish?
	A	B	C	A	B	C	
P0	0	1	0	0	0	0	Yes
P1	2	0	0	2	0	2	No
P2	3	0	2	0	0	0	Yes
P3	2	1	1	1	0	0	No

Analysis:

- P0 and P2 can finish initially since their requests are 0.
- 'Work' starts at (3, 3, 2) and increases to (3, 4, 2) after P0 finishes, then to (6, 4, 4) after P2 finishes.
- Neither P1 nor P3 can be satisfied with 'Work' = (6, 4, 4). They are deadlocked.

Deadlock Recovery Techniques

Once a deadlock is detected, a recovery mechanism is needed to break the cycle and restore the system to a non-deadlocked state.

Deadlock Recovery Techniques

Once a deadlock is detected, a recovery mechanism is needed to break the cycle and restore the system to a non-deadlocked state. **1. Process**

Termination

- **Kill all deadlocked processes:** The most drastic option. Breaks the cycle but can be very costly.
- **Kill one by one:** Terminate one process in the cycle at a time and re-run the detection algorithm. This is more efficient but takes longer.

Deadlock Recovery Techniques

Once a deadlock is detected, a recovery mechanism is needed to break the cycle and restore the system to a non-deadlocked state. **1. Process Termination**

1. Process Termination

- **Kill all deadlocked processes:** The most drastic option. Breaks the cycle but can be very costly.
- **Kill one by one:** Terminate one process in the cycle at a time and re-run the detection algorithm. This is more efficient but takes longer.

2. Resource Preemption

- **Preempt a resource:** Forcibly take a resource from one process and give it to another.
- This requires the ability to **roll back** the victim process to a previous "safe" state. This is often achieved through **checkpointing**.

Recovery Policy Decisions

When choosing a victim for termination or preemption, the OS must make a policy decision based on cost.

Recovery Policy Decisions

When choosing a victim for termination or preemption, the OS must make a policy decision based on cost. **Factors to consider:**

- The **priority** of the process.
- How much **CPU time** a process has consumed so far.
- The number of **resources held** by the process.
- The number of **resources needed** to complete the process.
- How many other processes will need to be rolled back.

Recovery Policy Decisions

When choosing a victim for termination or preemption, the OS must make a policy decision based on cost. **Factors to consider:**

- The **priority** of the process.
- How much **CPU time** a process has consumed so far.
- The number of **resources held** by the process.
- The number of **resources needed** to complete the process.
- How many other processes will need to be rolled back.

Important: Choosing a victim wisely minimizes the cost of recovery and ensures the system can continue operating efficiently.

Deadlock Handling in Real Systems

Most modern, general-purpose OSes take a pragmatic approach to deadlocks, often relying on **detection and recovery** in specific contexts rather than a global strategy.

- **Linux:** The kernel itself does not implement a global deadlock detection algorithm. Instead, it relies on strict **lock ordering** protocols in its core code to prevent deadlocks and uses a timeout mechanism for specific resource locks.
- **Windows:** Similar to Linux, it relies on lock ordering and hierarchical locking to avoid deadlocks. Components like the **NTFS file system** and **SQL Server** implement their own deadlock detection and recovery mechanisms (often via transaction rollback).

Key Takeaways

- Deadlock **detection** is a strategy that identifies a deadlock after it has occurred, trading simplicity for the risk of stalled processes.
- For single-instance resources, a cycle in the RAG is sufficient for deadlock detection.
- For multiple-instance resources, a more complex **matrix-based algorithm** is needed to determine if a safe state can be reached.
- Recovery involves making tough policy decisions, typically through **process termination** or **resource preemption**.

Next Week Preview: Virtual Machines and Hypervisors

Now that we have a solid understanding of how a single operating system manages its resources, we will explore how multiple operating systems can coexist on a single hardware platform.

- **Types of Virtualization:** Full, Para, and Hardware-Assisted.
- **Hypervisor Architecture:** The role of the Hypervisor in managing VMs.
- **VM Isolation and Resource Sharing:** How VMs are isolated from each other while sharing the same physical resources.

Appendix: Quiz

Conceptual Questions

- ① Under what conditions is a cycle in a Resource Allocation Graph a necessary and sufficient condition for a deadlock?
- ② What is the main trade-off of using a detection and recovery strategy over a prevention strategy?
- ③ How can the "deadlock detection algorithm for multiple instances" distinguish between a deadlock and a safe state?

Appendix: Exercises

Detection Algorithm Practice

- ① Given a system with 3 processes and 3 resources, and the following state:
 - ▶ Available: $A=2, B=1, C=0$
 - ▶ P1: Allocated $(1,0,1)$, Request $(0,1,0)$
 - ▶ P2: Allocated $(1,1,0)$, Request $(0,0,1)$
 - ▶ P3: Allocated $(0,1,0)$, Request $(2,0,0)$
- ② Is the system in a deadlock state? Show your work by applying the detection algorithm.

Appendix: Advanced Topics to Explore

Further Reading

- **Wait-for Graphs:** A simplified tool for deadlock detection in single-instance resource systems.
- **Deadlock Detection in Distributed Systems:** The unique challenges of detecting and resolving deadlocks when resources and processes are spread across multiple machines.
- **Deadlock in Databases:** How database management systems (DBMS) use deadlock detection and transaction rollback for concurrency control.

Week 6

Agenda

- ① Introduction to Virtualization
- ② Types of Virtualization Techniques
- ③ Hypervisors: Type 1 and Type 2 Architectures
- ④ VMM Architecture and Responsibilities
- ⑤ Use Cases and OS Examples
- ⑥ Summary & Key Takeaways
- ⑦ Next Week's Preview

What is Virtualization?

The Core Concept

Virtualization is the process of creating a software-based, logical representation of physical hardware. A **Virtual Machine (VM)** is the logical entity that encapsulates a complete operating system and its applications.

What is Virtualization?

The Core Concept

Virtualization is the process of creating a software-based, logical representation of physical hardware. A **Virtual Machine (VM)** is the logical entity that encapsulates a complete operating system and its applications.

Analogy: Think of a large physical server as an apartment building. Virtualization allows you to partition that building into multiple isolated apartments (VMs), each with its own tenants (guest OS) and furnishings (applications), all sharing the building's core resources (CPU, memory, storage).

What is Virtualization?

The Core Concept

Virtualization is the process of creating a software-based, logical representation of physical hardware. A **Virtual Machine (VM)** is the logical entity that encapsulates a complete operating system and its applications.

Analogy: Think of a large physical server as an apartment building. Virtualization allows you to partition that building into multiple isolated apartments (VMs), each with its own tenants (guest OS) and furnishings (applications), all sharing the building's core resources (CPU, memory, storage).

The **Virtual Machine Monitor (VMM)** or **Hypervisor** is the software layer that manages the creation and execution of these virtual machines.

Why Virtualization? Key Benefits

- **Resource Utilization:** By consolidating multiple workloads onto a single physical machine, virtualization dramatically increases resource efficiency and reduces hardware costs.
- **Isolation and Security:** VMs are isolated from each other. A failure or security breach in one VM does not affect the others on the same host.
- **Portability:** A VM is just a file or a set of files. It can be easily moved, copied, or backed up, making it highly portable.
- **Simplified Management:** Virtualization simplifies system administration, testing, and disaster recovery.

Why Virtualization? Key Benefits

- **Resource Utilization:** By consolidating multiple workloads onto a single physical machine, virtualization dramatically increases resource efficiency and reduces hardware costs.
- **Isolation and Security:** VMs are isolated from each other. A failure or security breach in one VM does not affect the others on the same host.
- **Portability:** A VM is just a file or a set of files. It can be easily moved, copied, or backed up, making it highly portable.
- **Simplified Management:** Virtualization simplifies system administration, testing, and disaster recovery.

Impact: Virtualization is the foundational technology for modern cloud computing, enterprise data centers, and development/testing environments.

Virtualization Techniques

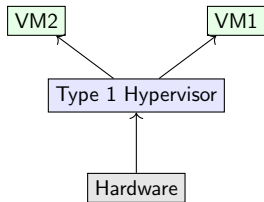
Technique	Guest OS Awareness	Performance	Mechanism
Full Virtualization	Unaware	Moderate	Binary translation of privileged instructions
Paravirtualization	Modified	High	Guest OS uses "hypercalls" to communicate with hypervisor
Hardware-Assisted	Unaware	Highest	Hardware extensions (Intel VT-x, AMD-V) trap privileged instructions

Hypervisor Types: Type 1 vs. Type 2

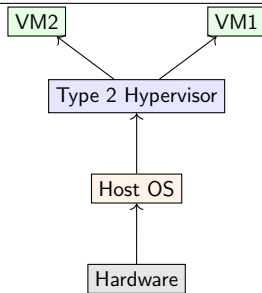
Type	Type 1 (Bare-metal)	Type 2 (Hosted)
Architecture	Runs directly on the host hardware	Runs as an application on top of a host OS
Performance	Very high, low overhead	Lower, since it must pass through the host OS's layers
Use Case	Data centers, enterprise servers, cloud computing	Desktop use, development, running legacy applications
Examples	Xen, VMware ESXi, Microsoft Hyper-V	VirtualBox, VMware Workstation, QEMU

Hypervisor Types: Type 1 vs. Type 2

Type	Type 1 (Bare-metal)	Type 2 (Hosted)
Architecture	Runs directly on the host hardware	Runs as an application on top of a host OS
Performance	Very high, low overhead	Lower, since it must pass through the host OS's layers
Use Case	Data centers, enterprise servers, cloud computing	Desktop use, development, running legacy applications
Examples	Xen, VMware ESXi, Microsoft Hyper-V	VirtualBox, VMware Workstation, QEMU

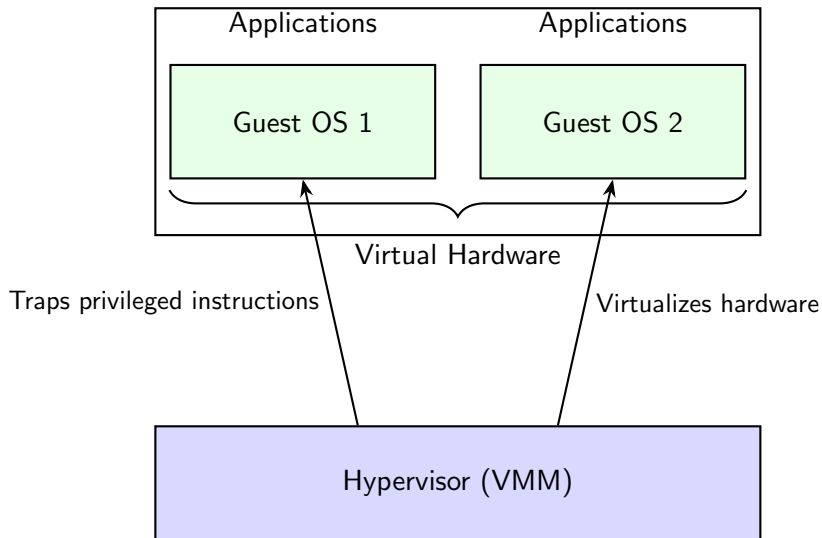


Type 1 Hypervisor



Type 2 Hypervisor

VMM Architecture Overview



VMM Responsibilities

Key Functions of a Hypervisor

- **Resource Scheduling:** Manages and allocates CPU time, memory, and I/O bandwidth to each VM.
- **Hardware Emulation:** Provides virtual devices (e.g., virtual network cards, virtual disks) to each VM.
- **Privilege Virtualization:** Intercepts and handles privileged instructions from the guest OS, ensuring isolation and preventing direct hardware access.
- **Memory Isolation:** Manages the translation of guest virtual addresses to physical host addresses, ensuring VMs cannot access each other's memory.

OS Virtualization Examples

- **KVM (Kernel-based Virtual Machine):** A Type 1 hypervisor integrated into the Linux kernel. It leverages hardware-assisted virtualization (VT-x/AMD-V) to run VMs efficiently.
- **Xen:** One of the first hypervisors, supporting both paravirtualization and hardware-assisted modes. It's a key component in many cloud platforms.
- **VMware ESXi:** A bare-metal Type 1 hypervisor known for its robust features and high performance in enterprise data centers.
- **VirtualBox:** A popular, easy-to-use Type 2 hypervisor for desktop use, enabling users to run different operating systems as applications.

Key Takeaways

- **Virtualization** abstracts hardware to run multiple isolated OSes on a single machine, improving resource utilization and security.
- **Hypervisors (VMMs)** are the core software that manages VMs, handling scheduling, resource allocation, and isolation.
- There are two main types of hypervisors: **Type 1 (Bare-metal)** for performance-critical environments and **Type 2 (Hosted)** for general desktop use.
- Virtualization is the foundation for modern cloud computing and is enabled by techniques like **paravirtualization** and **hardware-assisted virtualization**.

Next Week Preview: OS-Level Virtualization and Containers

Building on our understanding of virtual machines, we'll explore a different form of virtualization that is even more lightweight and efficient.

- **Containers vs. Virtual Machines:** A key comparison of their differences in architecture and performance.
- **Linux Namespaces and Cgroups:** The core OS primitives that make containerization possible.
- **Container Runtimes:** The role of tools like Docker and 'runc' in managing containers.

Appendix: Quiz

Conceptual Questions

- ① Describe the main architectural difference between a Type 1 and a Type 2 hypervisor.
- ② Why is hardware-assisted virtualization generally considered more efficient than full virtualization using binary translation?
- ③ What is the purpose of a "hypercall" in paravirtualization?

Appendix: Exercises

Research & Discussion

- ① Research and compare the performance overhead of running an application in a VM on a Type 1 vs. a Type 2 hypervisor.
- ② Investigate a specific hardware-assisted virtualization extension (e.g., Intel VT-x) and explain how it helps the hypervisor trap privileged instructions.

Appendix: Advanced Topics to Explore

Further Reading

- **Nested Virtualization:** The ability to run a hypervisor inside another hypervisor.
- **Memory Ballooning:** A memory management technique used by hypervisors to reclaim unused memory from VMs.
- **I/O Virtualization (PCI Passthrough):** Giving a VM direct access to a physical hardware device to bypass the hypervisor for maximum performance.

Week 7

Agenda

- ① VMs vs. Containers: The Core Differences
- ② OS-Level Virtualization: The Foundational Concept
- ③ Linux Namespaces and Control Groups (cgroups)
- ④ The Container Runtime Stack (Docker, runc)
- ⑤ Isolation, Performance, and Security
- ⑥ Summary & Key Takeaways
- ⑦ Next Week's Preview

VMs vs. Containers: The Core Differences

Architectural Comparison

Feature	Virtual Machines	Containers
Kernel	Each VM has its own guest OS kernel	Share the host OS kernel
Isolation	Hardware-level (Hypervisor)	OS-level (Namespaces, cgroups)
Size	Large (GBs), full OS image	Small (MBs), application + dependencies
Startup Time	Slow (minutes)	Fast (seconds or milliseconds)
Resource Use	High, each VM has dedicated resources	Low, share host resources

VMs vs. Containers: The Core Differences

Architectural Comparison

Feature	Virtual Machines	Containers
Kernel	Each VM has its own guest OS kernel	Share the host OS kernel
Isolation	Hardware-level (Hypervisor)	OS-level (Namespaces, cgroups)
Size	Large (GBs), full OS image	Small (MBs), application + dependencies
Startup Time	Slow (minutes)	Fast (seconds or milliseconds)
Resource Use	High, each VM has dedicated resources	Low, share host resources

Analogy: A VM is like a full house with its own utilities and kitchen, whereas a container is like a single apartment within a building, sharing the building's infrastructure.

OS-Level Virtualization

The Container Foundation

OS-level virtualization is a technique where the operating system kernel allows for the existence of multiple isolated user-space instances. Instead of emulating hardware, it partitions the OS's resources.

OS-Level Virtualization

The Container Foundation

OS-level virtualization is a technique where the operating system kernel allows for the existence of multiple isolated user-space instances. Instead of emulating hardware, it partitions the OS's resources.

The Linux kernel provides two key primitives that enable this:

- **Linux Namespaces:** Isolate processes so they have their own view of system resources (e.g., process IDs, network interfaces, filesystems).
- **Control Groups (cgroups):** Enforce resource limits, ensuring a container doesn't consume all of the host's CPU, memory, or I/O.

A container is essentially a process (or a group of processes) wrapped in a set of namespaces and restricted by cgroups.

Linux Namespaces: The Isolation Mechanism

Namespaces provide a per-process view of system resources, giving each container the illusion that it is running on a clean, isolated system.

Namespace Type	Resource Isolated
PID	Process IDs
NET	Network interfaces, IP addresses, routing tables
MNT	Filesystem mounts
UTS	Hostname and domain name
IPC	Interprocess communication mechanisms
USER	User and group IDs

Linux Namespaces: The Isolation Mechanism

Namespaces provide a per-process view of system resources, giving each container the illusion that it is running on a clean, isolated system.

Namespace Type	Resource Isolated
PID	Process IDs
NET	Network interfaces, IP addresses, routing tables
MNT	Filesystem mounts
UTS	Hostname and domain name
IPC	Interprocess communication mechanisms
USER	User and group IDs

Example: The 'MNT' namespace allows a container to have its own root filesystem, separate from the host, using a feature like 'chroot' or 'overlayfs'.

Linux Control Groups (cgroups): The Resource Manager

Resource Enforcement

Control Groups are a kernel feature that hierarchically organizes processes and controls their access to system resources.

Linux Control Groups (cgroups): The Resource Manager

Resource Enforcement

Control Groups are a kernel feature that hierarchically organizes processes and controls their access to system resources.

cgroups allow the OS to:

- **Limit:** Set hard limits on a container's CPU, memory, and I/O.
- **Monitor:** Track resource usage of a container.
- **Prioritize:** Allocate different levels of access to resources.

This is critical for a multi-tenant environment, as it prevents a single misbehaving container from monopolizing resources and impacting the performance of others.

The Container Runtime Stack

Docker is a high-level tool that simplifies the entire container lifecycle, from building to running. It's built on a lower-level standard.

The Container Runtime Stack

Docker is a high-level tool that simplifies the entire container lifecycle, from building to running. It's built on a lower-level standard. **Key**

Components:

- **Docker CLI:** The command-line interface that users interact with.
- **Docker Daemon:** A background service that manages container lifecycle, images, volumes, and networks.
- **Container Runtime:** A lower-level component that is responsible for actually running the container.

The Container Runtime Stack

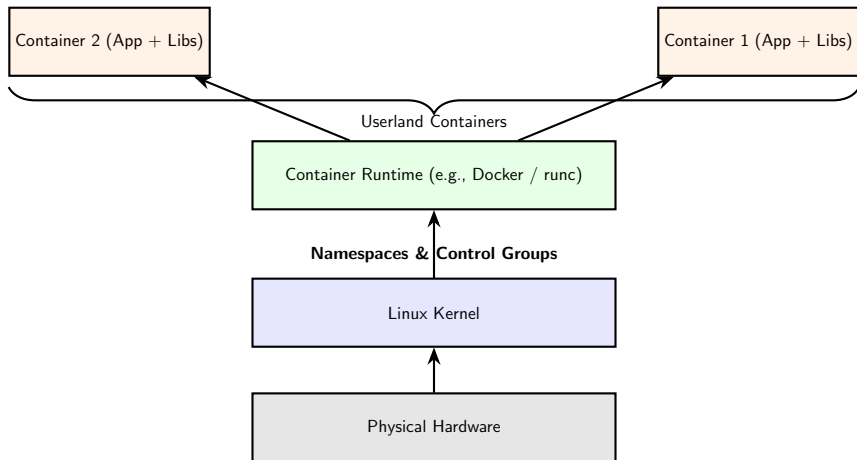
Docker is a high-level tool that simplifies the entire container lifecycle, from building to running. It's built on a lower-level standard. **Key**

Components:

- **Docker CLI:** The command-line interface that users interact with.
- **Docker Daemon:** A background service that manages container lifecycle, images, volumes, and networks.
- **Container Runtime:** A lower-level component that is responsible for actually running the container.

runc: A lightweight, portable container runtime that implements the **Open Container Initiative (OCI)** specification. Docker uses 'runc' under the hood to handle the core task of creating and running containers based on namespaces and cgroups.

Container Stack (Linux Host)



Container Security and Isolation

While containers share a kernel, their security is managed through multiple layers of defense.

- **Namespaces & cgroups:** The core isolation mechanisms.
- **Seccomp Filters:** A mechanism to restrict the system calls a container can make to the kernel, reducing the attack surface.
- **Linux Capabilities:** Instead of giving a container 'root' access, specific, fine-grained privileges can be granted (e.g., 'CAP_NET_ADMIN' for network configuration).
- **AppArmor / SELinux:** Mandatory Access Control (MAC) systems that enforce policies on what processes can access, based on labels rather than user IDs.

Key Takeaways

- Containers offer a lightweight and fast alternative to VMs by sharing the host OS kernel and providing process-level isolation.
- The core technologies enabling containers are **Linux Namespaces** (for isolation) and **Control Groups** (for resource management).
- Tools like Docker simplify the entire workflow, while a low-level runtime like 'runc' handles the actual container creation and execution.
- Containers are the foundation of modern **microservices architectures** and **CI/CD** pipelines.

Next Week Preview: Isolation Internals – A Deep Dive

To truly understand how containers work, we will peel back the layers and look at the underlying mechanisms in more detail.

- **Namespaces in Detail:** A technical breakdown of how each namespace type is implemented and what it isolates.
- **cgroups V1 vs. V2:** Understanding the evolution of control groups and their differences.
- **Container Runtime Lifecycle:** A step-by-step look at how a command like 'docker run' actually creates a container.

Appendix: Quiz

Conceptual Questions

- ① What is the main security implication of a container sharing the host OS kernel, compared to a VM?
- ② A user wants to limit a container's memory to 512MB and its CPU usage to 50%. Which Linux kernel primitive would they use?
- ③ Explain the role of a union filesystem (like 'overlayfs') in container images.

Appendix: Exercises

Hands-on Practice

- ① Use the 'unshare' command to create a new shell with its own PID namespace. Verify that the new shell has a PID of 1.
- ② Use the 'docker inspect' command on a running container to view the cgroups path and the applied resource limits.

Appendix: Advanced Topics to Explore

Further Reading

- **Container Networking:** The various ways containers are given network connectivity (e.g., bridge mode, overlay networks).
- **Container Orchestration:** Tools like Kubernetes and Swarm that automate the deployment, scaling, and management of containers.
- **Container Sandboxes:** Alternative container runtimes (e.g., Kata Containers) that provide VM-like isolation while maintaining container-like speed.

Week 8

Agenda

- ① Linux Namespaces: The Isolation Primitives
- ② Control Groups (cgroups) v1 vs. v2
- ③ The Container Lifecycle from an OS Perspective
- ④ Coordinating Containers via OS Interfaces
- ⑤ Summary & Key Takeaways
- ⑥ Next Week's Preview

Linux Namespaces – The Isolation Primitives

Recap: The Building Blocks

Linux Namespaces provide process-level isolation by giving a process its own separate view of a specific system resource. A process can be part of multiple namespaces simultaneously, each one providing a different layer of isolation.

Linux Namespaces – The Isolation Primitives

Recap: The Building Blocks

Linux Namespaces provide process-level isolation by giving a process its own separate view of a specific system resource. A process can be part of multiple namespaces simultaneously, each one providing a different layer of isolation.

How it works:

- A process running on a host is part of the host's default set of namespaces.
- When creating a container, the runtime uses a system call like 'clone()' or 'unshare()' to create new namespaces for the child process.
- The child process (and its children) will then have their own isolated view of that specific resource.

Key Namespaces: PID and USER

PID Namespace: Process ID Isolation

- Isolates the process ID (PID) space. The first process in a new PID namespace gets PID 1 and acts as the 'init' process for that namespace.
- The container's PID 1 is not the same as the host's PID 1, preventing a container from killing critical host processes.
- Enables clean process trees within the container, allowing for clean shutdown.

Key Namespaces: PID and USER

PID Namespace: Process ID Isolation

- Isolates the process ID (PID) space. The first process in a new PID namespace gets PID 1 and acts as the 'init' process for that namespace.
- The container's PID 1 is not the same as the host's PID 1, preventing a container from killing critical host processes.
- Enables clean process trees within the container, allowing for clean shutdown.

USER Namespace: Privilege Isolation

- Isolates user and group ID space. A user inside a container (e.g., 'root') can be mapped to an unprivileged user on the host.
- This is a critical security feature, as it allows containers to run as 'root' without having root privileges on the host system.

Key Namespaces: Network and Mount

Network (NET) Namespace

- Each network namespace has its own network stack, including network interfaces, IP addresses, routing tables, and firewall rules.
- The host can use virtual Ethernet ('veth') pairs to connect a container's network namespace to the host's network or to other containers.

Key Namespaces: Network and Mount

Network (NET) Namespace

- Each network namespace has its own network stack, including network interfaces, IP addresses, routing tables, and firewall rules.
- The host can use virtual Ethernet ('veth') pairs to connect a container's network namespace to the host's network or to other containers.

Mount (MNT) Namespace

- Isolates the filesystem's mount points.
- This allows each container to have its own filesystem hierarchy, independent of the host.
- This is often implemented using 'chroot' and a layered filesystem like 'overlayfs' to present a full root filesystem to the container.

Other Namespaces: UTS and IPC

UTS Namespace

- Isolates the system's hostname and domain name.
- This allows each container to have its own unique hostname.

Other Namespaces: UTS and IPC

UTS Namespace

- Isolates the system's hostname and domain name.
- This allows each container to have its own unique hostname.

IPC Namespace

- Isolates inter-process communication resources.
- This prevents processes in one container from interfering with System V IPC or POSIX message queues used by processes in other containers.

Control Groups (cgroups): v1 vs. v2

cgroups v1

- The original version, in use for many years.
- Has a separate hierarchy for each controller (e.g., CPU, Memory, I/O). This can be complex to manage.
- A single process can be in multiple cgroups, one for each resource controller.

Control Groups (cgroups): v1 vs. v2

cgroups v1

- The original version, in use for many years.
- Has a separate hierarchy for each controller (e.g., CPU, Memory, I/O). This can be complex to manage.
- A single process can be in multiple cgroups, one for each resource controller.

cgroups v2

- The modern, unified version of cgroups.
- All resource controllers share a single, unified hierarchy.
- Enforces a "single writer" rule, where a process can only be in a single cgroup. This simplifies management and provides better consistency.
- Supports better resource delegation and hierarchical management.

The Container Lifecycle (from an OS View)

Step-by-Step Container Creation

When a user runs a command like 'docker run', the following low-level steps occur:

- ① The container runtime calls 'clone()' with flags to create new namespaces (e.g., 'CLONE_NEWPID', 'CLONE_NEWNET').
- ② The runtime creates a new cgroup and adds the container's process to it. Resource limits are configured.
- ③ The container's root filesystem is mounted using 'pivot_root' or 'chroot'.
- ④ The hostname is set, and user/group ID mappings are configured.
- ⑤ The specified command (e.g., '/bin/bash') is executed as PID 1 within the new namespace.

Coordinating Containers via OS Interfaces

- Container orchestration systems like **Kubernetes** are built on these core OS primitives.
- Kubernetes' 'kubelet' process on each node is responsible for telling the container runtime to create and manage containers.
- It uses cgroups to enforce a Pod's resource requests and limits, ensuring predictable performance.
- It leverages namespaces to provide network and process isolation between different Pods on the same node.

Coordinating Containers via OS Interfaces

- Container orchestration systems like **Kubernetes** are built on these core OS primitives.
- Kubernetes' 'kubelet' process on each node is responsible for telling the container runtime to create and manage containers.
- It uses cgroups to enforce a Pod's resource requests and limits, ensuring predictable performance.
- It leverages namespaces to provide network and process isolation between different Pods on the same node.

Conclusion: The Linux kernel acts as the ultimate orchestrator, providing the necessary low-level hooks for higher-level tools to build complex, multi-tenant systems.

Key Takeaways

- **Linux Namespaces** are the kernel's core mechanism for isolating system resources, enabling a process to have its own view of PIDs, network stack, filesystem, etc.
- **cgroups** are the kernel's mechanism for managing and enforcing resource usage (CPU, memory, I/O) to ensure fair sharing.
- The combination of these two features allows for lightweight and secure **OS-level virtualization**.
- The entire container lifecycle, from creation to destruction, is managed through these fundamental kernel interfaces.

Next Week Preview: Coordination in Distributed OS

Having covered single-host OS virtualization, we will now turn our attention to the challenges of managing resources and state across multiple machines.

- **Clock Synchronization:** The problem of keeping time consistent across a distributed system.
- **Remote Procedure Calls (RPC):** How processes on different machines can communicate.
- **Consensus Algorithms:** How distributed systems agree on a single state (e.g., Raft, Paxos).

Appendix: Quiz

Conceptual Questions

- ① Explain why a container running as 'root' is not as dangerous as a process running as 'root' on the host, using the concept of namespaces.
- ② What is the primary benefit of the unified hierarchy in cgroups v2 compared to cgroups v1?
- ③ Which combination of namespaces is essential for a container to have its own isolated filesystem and process tree?

Appendix: Exercises

Hands-on Practice

- ① Using the 'unshare' command, create a shell in a new mount namespace. Create a file inside that shell and then exit. Check the host's filesystem to confirm the file is not visible.
- ② Create a new cgroup, configure a memory limit for it, and then run a process inside it that tries to exceed that limit. Observe the kernel's reaction.

Appendix: Advanced Topics to Explore

Further Reading

- **eBPF (extended Berkeley Packet Filter):** A powerful kernel technology that allows for dynamic, safe, and efficient tracing and networking, which is heavily used in modern container networking and security.
- **OverlayFS:** A type of union filesystem used extensively in containers to create a writable layer on top of a read-only base image, making images efficient to share.
- **The OCI (Open Container Initiative) Specification:** The open standard that defines how container images and runtimes should operate, ensuring interoperability between tools.

Week 9

Agenda

- ① Coordination in Distributed Systems
- ② Time and Clock Synchronization
- ③ Distributed Mutual Exclusion
- ④ Naming, Resource Sharing, and Transparency
- ⑤ Examples of Distributed OS Features
- ⑥ Summary & Key Takeaways
- ⑦ Next Week's Preview

Why Is Coordination Challenging?

The Fundamental Challenges

Distributed systems are composed of multiple independent nodes that communicate over a network. This presents fundamental challenges that don't exist in a single-machine OS.

Why Is Coordination Challenging?

The Fundamental Challenges

Distributed systems are composed of multiple independent nodes that communicate over a network. This presents fundamental challenges that don't exist in a single-machine OS.

Key Issues:

- **No Shared Memory:** Processes can't directly inspect each other's state; all communication must be through messages.
- **No Global Clock:** Each node has its own clock that can drift, making it difficult to establish a global ordering of events.
- **Message Delays and Failures:** Messages can be lost, corrupted, or delayed, and nodes can fail independently.

Why Is Coordination Challenging?

The Fundamental Challenges

Distributed systems are composed of multiple independent nodes that communicate over a network. This presents fundamental challenges that don't exist in a single-machine OS.

Key Issues:

- **No Shared Memory:** Processes can't directly inspect each other's state; all communication must be through messages.
- **No Global Clock:** Each node has its own clock that can drift, making it difficult to establish a global ordering of events.
- **Message Delays and Failures:** Messages can be lost, corrupted, or delayed, and nodes can fail independently.

The goal of a distributed OS is to hide these challenges from the user, providing an illusion of a single, consistent system.

Transparency Goals in a Distributed OS

Transparency is a core design principle of distributed systems, aiming to hide the complexity of the underlying architecture.

Transparency Type	Description
Access	Use the same interface for local and remote resources. (e.g., 'read()' works on a local file or a remote file.)
Location	Hide the physical location of a resource. The user doesn't need to know which node a file or service is on.
Migration	The ability to move processes or data without affecting the user or application.
Replication	Hide the fact that resources might be replicated for fault tolerance or performance.
Concurrency	Allow multiple users to safely access a shared resource at the same time.
Failure	Mask failures of nodes or communication links so the system continues to operate.

Time and Clock Synchronization

The Problem of Time

Each computer has its own physical clock, which can drift and become unsynchronized. Establishing a globally consistent time is crucial for logging, event ordering, and mutual exclusion.

Time and Clock Synchronization

The Problem of Time

Each computer has its own physical clock, which can drift and become unsynchronized. Establishing a globally consistent time is crucial for logging, event ordering, and mutual exclusion.

Algorithms for Clock Synchronization:

- **Cristian's Algorithm:** A simple client-server model. A client sends a request to a time server and adjusts its clock based on the response time, factoring in network latency.
- **Berkeley Algorithm:** An active server polls all clients for their time, calculates an average, and tells each client how to adjust its clock (no sudden jumps).
- **NTP (Network Time Protocol):** The standard for time synchronization on the internet. It uses a hierarchical, fault-tolerant structure of time servers to achieve high accuracy.

Logical Clocks (Lamport Timestamps)

The "Happens-Before" Relationship

Logical clocks don't synchronize physical time. Instead, they provide a way to establish a **causal order** of events, answering the question: "Did event A happen before event B?"

Logical Clocks (Lamport Timestamps)

The "Happens-Before" Relationship

Logical clocks don't synchronize physical time. Instead, they provide a way to establish a **causal order** of events, answering the question: "Did event A happen before event B?"

Lamport's Algorithm:

- Each process maintains a local counter (the logical clock).
- The clock is incremented before each event.
- When a process sends a message, it includes its current timestamp.
- When a process receives a message, it updates its own clock to be the maximum of its current value and the message's timestamp, plus one.

Distributed Mutual Exclusion

Sharing a Critical Section

The goal of distributed mutual exclusion is to ensure that only one process at a time can access a shared resource, even though it's distributed across multiple nodes.

Distributed Mutual Exclusion

Sharing a Critical Section

The goal of distributed mutual exclusion is to ensure that only one process at a time can access a shared resource, even though it's distributed across multiple nodes.

Algorithm	How it Works	Trade-offs
Centralized	A single coordinator grants access requests.	Simple, but a single point of failure and bottleneck.
Ricart-Agrawala	Processes request access from all others, and enter if they get 'OK' replies.	No single point of failure, but high message overhead.
Token-based	A special token circulates. The process holding the token can enter the critical section.	Low message count, but complex to recover if the token is lost.

Naming and Resource Sharing

Providing a Unified View

A distributed OS must provide a global naming service that allows processes to locate and access resources (files, devices, services) without knowing their physical location.

Naming and Resource Sharing

Providing a Unified View

A distributed OS must provide a global naming service that allows processes to locate and access resources (files, devices, services) without knowing their physical location.

Examples:

- **Remote File Systems:** Tools like **NFS** (Network File System) and **AFS** (Andrew File System) allow clients to access and manage files on a remote server as if they were local.
- **Distributed Shared Memory (DSM):** Provides the illusion of a shared memory space across multiple nodes, abstracting away the need for explicit message passing.
- **Naming Services:** Systems like **DNS** and directory services (e.g., LDAP) provide a hierarchical, location-independent way to find resources.

Distributed OS Examples

- **Amoeba:** A microkernel-based distributed OS from the 1980s. It focused on an object-based naming system, where all resources were treated as objects with capabilities.
- **Sprite:** A research OS from the 1990s that provided a single-system image, allowing processes to migrate between nodes and share a global file system.
- **Plan 9 from Bell Labs:** A research OS that pioneered the concept of "everything is a file," providing a uniform, hierarchical namespace for all system resources.
- **Google Fuchsia:** A modern distributed-ready OS with a microkernel architecture and a component-based model, designed for a wide range of devices from phones to large-scale data centers.

Key Takeaways

- The fundamental challenges of distributed systems are the absence of shared memory and a global clock.
- **Transparency** is the core design goal to hide this complexity from the user.
- **Clock synchronization** (physical and logical) and **distributed mutual exclusion** are essential for coordinating events and resource access.
- **Global naming services** and **distributed shared memory** are key features that provide a unified, single-system view.

Next Week Preview: Distributed Synchronization and Consensus

Building on our understanding of distributed coordination, we will delve into the problem of achieving consensus and electing a leader in a fault-prone distributed system.

- **Election Algorithms:** The process of choosing a coordinator or leader from a set of processes.
- **Leader-based Coordination:** How a leader can simplify resource management and message passing.
- **Paxos/Raft Consensus Overview:** The core concepts of famous consensus algorithms that ensure all nodes agree on a single value, even with failures.

Appendix: Quiz

Conceptual Questions

- ① What is the key difference in purpose between a physical clock synchronization algorithm (like NTP) and a logical clock (like Lamport timestamps)?
- ② Describe a scenario where a system with **concurrency transparency** would be beneficial.
- ③ What is the main drawback of using a centralized approach for distributed mutual exclusion?

Appendix: Exercises

Design & Discussion

- ① You are designing a distributed file system. Which types of transparency from our list would be the most important to implement, and why?
- ② Consider a three-node system where each node has a Lamport clock. Trace the logical clock values for a series of events: P1 sends a message to P2, then P2 sends a message to P3, then P3 sends a message to P1.

Appendix: Advanced Topics to Explore

Further Reading

- **Vector Clocks:** A more powerful extension of logical clocks that can determine causal relationships between any two events.
- **The CAP Theorem:** A fundamental theorem in distributed systems that describes the trade-off between Consistency, Availability, and Partition Tolerance.
- **Two-Phase Commit:** A protocol used to ensure that all nodes in a distributed transaction either commit or abort the transaction.

Week 10

Agenda

- ① The Need for Distributed Synchronization
- ② Leader Election Algorithms (Bully, Ring)
- ③ Distributed Consensus: The Core Problem
- ④ Paxos and Raft: A Comparison
- ⑤ Fault Tolerance in Distributed Systems
- ⑥ Summary & Key Takeaways
- ⑦ Next Week's Preview

The Need for Distributed Synchronization

Why do we need it?

In a distributed system, nodes must coordinate their actions and maintain a consistent view of the system's state without a central authority or shared memory. This is essential for reliability and correctness.

The Need for Distributed Synchronization

Why do we need it?

In a distributed system, nodes must coordinate their actions and maintain a consistent view of the system's state without a central authority or shared memory. This is essential for reliability and correctness.

Key Use Cases:

- **Leader Election:** Choosing a single node to act as a coordinator for a specific task.
- **Atomic Operations:** Ensuring that a series of operations either all succeed or all fail across all nodes (e.g., a distributed database transaction).
- **State Replication:** Keeping the state of a service consistent across multiple machines.

Goal: Achieve a state of agreement and consistency despite unpredictable network delays and node failures.

Leader Election Problem

The Problem of Coordination

The leader election problem is a classic challenge in distributed systems: how do a group of nodes, with no central coordinator, collectively choose one node to be the leader?

Leader Election Problem

The Problem of Coordination

The leader election problem is a classic challenge in distributed systems: how do a group of nodes, with no central coordinator, collectively choose one node to be the leader?

Requirements for a Leader Election Algorithm:

- **Uniqueness:** At any given time, only one leader is elected.
- **Termination:** The process must eventually terminate, with a leader being chosen.
- **Robustness:** The algorithm should be able to handle node failures and new nodes joining the system.

Bully Algorithm

The "Bully" Approach

In the Bully Algorithm, a process detects the leader has failed and starts an election. It sends 'ELECTION' messages to all nodes with higher IDs.

Bully Algorithm

The "Bully" Approach

In the Bully Algorithm, a process detects the leader has failed and starts an election. It sends 'ELECTION' messages to all nodes with higher IDs.

How it works:

- If a higher-ID node responds with an 'OK' message, the bully stands down and waits for the new leader's 'COORDINATOR' message.
- If no response is received after a timeout, the bully assumes it's the highest-ID node and declares itself the new leader by sending a 'COORDINATOR' message to all nodes.

Bully Algorithm

The "Bully" Approach

In the Bully Algorithm, a process detects the leader has failed and starts an election. It sends 'ELECTION' messages to all nodes with higher IDs.

How it works:

- If a higher-ID node responds with an 'OK' message, the bully stands down and waits for the new leader's 'COORDINATOR' message.
- If no response is received after a timeout, the bully assumes it's the highest-ID node and declares itself the new leader by sending a 'COORDINATOR' message to all nodes.

Trade-offs:

- **Pros:** Simple to understand; the highest-ID node always wins.
- **Cons:** Can generate a high volume of messages; assumes reliable message passing for failure detection.

Ring-Based Election Algorithm

Passing a Token

In the Ring Algorithm, all nodes are arranged in a logical ring. A failed leader is detected, and the next node in the ring starts the election by sending an 'ELECTION' message.

Ring-Based Election Algorithm

Passing a Token

In the Ring Algorithm, all nodes are arranged in a logical ring. A failed leader is detected, and the next node in the ring starts the election by sending an 'ELECTION' message.

How it works:

- The 'ELECTION' message contains the initiator's ID.
- As the message circulates around the ring, each node adds its own ID if it's higher than the current ID in the message.
- The message returns to the initiator, which then sends a 'COORDINATOR' message containing the highest ID it received.

Ring-Based Election Algorithm

Passing a Token

In the Ring Algorithm, all nodes are arranged in a logical ring. A failed leader is detected, and the next node in the ring starts the election by sending an 'ELECTION' message.

How it works:

- The 'ELECTION' message contains the initiator's ID.
- As the message circulates around the ring, each node adds its own ID if it's higher than the current ID in the message.
- The message returns to the initiator, which then sends a 'COORDINATOR' message containing the highest ID it received.

Trade-offs:

- **Pros:** Message-efficient (at most $2N$ messages for N nodes).
- **Cons:** Slow to elect a leader in a large ring; failure of a single node can break the ring.

Distributed Consensus: The Core Problem

Agreement on a Single Value

Distributed Consensus is the problem of getting a group of nodes to agree on a single, shared value. It is a fundamental problem in distributed systems and is a core component of replicated state machines.

Distributed Consensus: The Core Problem

Agreement on a Single Value

Distributed Consensus is the problem of getting a group of nodes to agree on a single, shared value. It is a fundamental problem in distributed systems and is a core component of replicated state machines.

Key Properties of Consensus Algorithms:

- **Safety:** Nothing bad happens.
 - ▶ **Agreement:** All nodes that decide on a value, decide on the same value.
 - ▶ **Validity:** The decided value must be one of the proposed values.
- **Liveness:** Something good eventually happens.
 - ▶ **Termination:** All non-faulty processes eventually decide on a value.

Paxos Overview (Simplified)

The First Consensus Algorithm

Paxos is a family of protocols for solving the consensus problem. It's notoriously difficult to understand and implement correctly.

Paxos Overview (Simplified)

The First Consensus Algorithm

Paxos is a family of protocols for solving the consensus problem. It's notoriously difficult to understand and implement correctly.

Roles:

- **Proposer:** Proposes a value to be agreed upon.
- **Acceptor:** Votes on proposed values. A majority of acceptors constitutes a quorum.
- **Learner:** Learns the final decided value.

Paxos Overview (Simplified)

The First Consensus Algorithm

Paxos is a family of protocols for solving the consensus problem. It's notoriously difficult to understand and implement correctly.

Roles:

- **Proposer:** Proposes a value to be agreed upon.
- **Acceptor:** Votes on proposed values. A majority of acceptors constitutes a quorum.
- **Learner:** Learns the final decided value.

Two Phases:

- ① **Phase 1 (Prepare/Promise):** A proposer asks acceptors if they're willing to accept a value.
- ② **Phase 2 (Accept/Commit):** If a majority agrees, the proposer sends the value to be committed.

Raft Consensus Algorithm

A Modern, Understandable Alternative

Raft is a leader-based consensus algorithm designed to be more understandable than Paxos. Its key idea is to first elect a stable leader, and then have the leader manage the replicated log.

Raft Consensus Algorithm

A Modern, Understandable Alternative

Raft is a leader-based consensus algorithm designed to be more understandable than Paxos. Its key idea is to first elect a stable leader, and then have the leader manage the replicated log.

Roles:

- **Leader:** The single node that handles all client requests and manages log replication.
- **Follower:** Passive nodes that only respond to leader requests.
- **Candidate:** A follower that has timed out and is attempting to become the new leader.

Raft Consensus Algorithm

A Modern, Understandable Alternative

Raft is a leader-based consensus algorithm designed to be more understandable than Paxos. Its key idea is to first elect a stable leader, and then have the leader manage the replicated log.

Roles:

- **Leader:** The single node that handles all client requests and manages log replication.
- **Follower:** Passive nodes that only respond to leader requests.
- **Candidate:** A follower that has timed out and is attempting to become the new leader.

Key Mechanisms:

- **Leader Election:** Followers use randomized timeouts to trigger elections.
- **Log Replication:** The leader replicates a sequence of commands (the log) to all followers. A command is committed once a majority of followers have stored it.

Use Cases: etcd, Consul, and the Kubernetes API Server.

Handling Faults in Coordination

Ensuring Reliability

A core goal of these protocols is to function correctly despite node crashes, network partitions, and message loss.

Handling Faults in Coordination

Ensuring Reliability

A core goal of these protocols is to function correctly despite node crashes, network partitions, and message loss.

Common Fault Tolerance Techniques:

- **Quorum-based Protocols:** A decision is only made after a majority of nodes have agreed. This ensures the system can continue operating as long as a majority of nodes are available (e.g., $2f + 1$ nodes can tolerate f failures).
- **Timeouts and Retries:** Used to detect presumed failures. If a message isn't acknowledged within a certain time, the sender assumes the receiver is down and retries the action.
- **Persistent State:** All critical state (e.g., the log in Raft) is written to a persistent disk before being acknowledged. This allows nodes to recover their state after a crash.

Key Takeaways

- **Distributed synchronization** is essential for creating reliable systems that operate without a central point of control.
- **Leader election algorithms** like Bully and Ring provide ways for nodes to choose a coordinator.
- **Consensus protocols** like Paxos and Raft are the gold standard for ensuring a consistent, replicated state across a cluster.
- All these protocols rely on a **majority-based quorum** to ensure safety and liveness in the face of faults.

Next Week Preview: OS Virtualization & Coordination Case Studies

In our final week, we'll connect the concepts we've learned to real-world systems.

- **Linux Containers & Kubernetes:** How container orchestration leverages namespaces, cgroups, and leader election.
- **VMware & vMotion:** How a virtualized environment coordinates live migration of VMs.
- **Zookeeper & etcd:** An overview of these production-grade coordination services and the algorithms they use.

Appendix: Quiz

Conceptual Questions

- ① What is the main difference between the 'Coordinator' message in the Bully algorithm and the one in the Ring algorithm?
- ② What is the key advantage of Raft's log replication model over Paxos?
- ③ In a 5-node cluster, how many nodes must be non-faulty to achieve consensus?

Appendix: Exercises

Thought Experiment

- ① Design a simple leader election algorithm for a network where nodes are organized in a 2D grid.
- ② Research the role of the 'etcd' component in a Kubernetes cluster and explain how it uses a consensus algorithm.

Appendix: Advanced Topics to Explore

Further Reading

- **Consensus in the BFT (Byzantine Fault Tolerance) Model:** A more complex problem where nodes can be malicious.
- **The FLP Impossibility Result:** A famous theorem in distributed computing that shows that it's impossible to design a perfect consensus algorithm in an asynchronous system with even a single crash failure.
- **The Raft paper:** The original paper is a great example of a well-written technical paper, designed for clarity.

Week 11

Agenda

- ① Real-World OS Coordination Scenarios
- ② Linux + Containers: Kernel-Level Coordination
- ③ Kubernetes: A Distributed OS for Containers
- ④ Hypervisor-Based Systems: VMware vs. KVM
- ⑤ Integrated Reflection and Architectural Insights
- ⑥ Summary & Key Takeaways
- ⑦ Next Week's Preview

Linux Kernel Coordination: Containers

The Local Orchestrator

The Linux kernel itself acts as a local coordinator for processes. For containers, it provides the low-level primitives needed for isolation and resource management.

Linux Kernel Coordination: Containers

The Local Orchestrator

The Linux kernel itself acts as a local coordinator for processes. For containers, it provides the low-level primitives needed for isolation and resource management.

Core Kernel Primitives for Coordination:

- **Namespaces:** Provide isolated views of system resources like PIDs, network interfaces, and the filesystem, ensuring a container cannot see or interact with others on the same host.
- **cgroups:** Enforce resource limits, ensuring a container can't consume all the CPU, memory, or I/O, providing a form of Quality-of-Service (QoS).
- **Seccomp and Capabilities:** Fine-tune the security, limiting a container's ability to make certain system calls or access privileged operations.

These primitives are the foundation of all container engines, including Docker.

Docker Coordination Overview

Single-Node Management

The Docker Daemon is the primary coordination component on a single host. It manages the entire lifecycle of a container, from image management to execution.

Docker Coordination Overview

Single-Node Management

The Docker Daemon is the primary coordination component on a single host. It manages the entire lifecycle of a container, from image management to execution.

Key Coordination Tasks:

- **Image Layering:** Uses copy-on-write filesystems (like 'overlays') to manage container images efficiently, sharing a read-only base layer across multiple containers.
- **Restart Policies:** Automatically restarts containers that fail or exit, ensuring high availability at the single-node level.
- **Event Management:** Provides hooks for external tools to respond to container events (e.g., creation, stop, restart).

It's important to note that the Docker Daemon does not inherently manage multi-node coordination, which is the role of tools like Kubernetes.

Kubernetes Case Study

A Distributed OS for Containers

Kubernetes is an orchestrator that acts like a distributed operating system for a cluster of machines. It abstracts away the individual nodes and provides a single control plane for managing containerized workloads.

Kubernetes Case Study

A Distributed OS for Containers

Kubernetes is an orchestrator that acts like a distributed operating system for a cluster of machines. It abstracts away the individual nodes and provides a single control plane for managing containerized workloads.

Key Components and their Roles:

- **Kube-Scheduler:** Decides which node a new Pod should run on, based on resource requirements and other constraints.
- **Kubelet:** The agent on each node that communicates with the scheduler and the container runtime to start and manage Pods.
- **Controllers:** A set of control loops (e.g., ReplicationController) that ensure the actual state of the cluster matches the desired state.
- **etcd:** A distributed key-value store that serves as the cluster's "source of truth." It uses the **Raft consensus algorithm** to ensure high availability and consistency.

Kubernetes Coordination Techniques

Scaling OS Concepts

Kubernetes applies classic OS coordination principles to a distributed environment.

Kubernetes Coordination Techniques

Scaling OS Concepts

Kubernetes applies classic OS coordination principles to a distributed environment.

Key Concepts:

- **Resource Management:** Kubernetes uses 'cgroups' and 'namespaces' under the hood. It defines Pod-level resource 'requests' (for scheduling) and 'limits' (for enforcement).
- **Scheduling:** The scheduler uses a two-phase process (filtering and ranking) to find the best node for a Pod, extending the concept of a local process scheduler.
- **Self-Healing:** Kubernetes' core loop constantly checks the state of Pods. If a node or container fails, the controller automatically recreates the Pod on a healthy node, providing robust **failure transparency**.

VMware vSphere – Coordination via Hypervisor

Virtualization-Based Coordination

Unlike the container model, VMware's vSphere coordinates full virtual machines. The hypervisor is the central point of control, and it's built on a bare-metal architecture.

VMware vSphere – Coordination via Hypervisor

Virtualization-Based Coordination

Unlike the container model, VMware's vSphere coordinates full virtual machines. The hypervisor is the central point of control, and it's built on a bare-metal architecture.

Key Coordination Features:

- **vMotion:** Allows for the live migration of a running virtual machine from one physical server to another with zero downtime, providing seamless **migration transparency**.
- **DRS (Dynamic Resource Scheduling):** A feature that continuously monitors resource usage across a cluster of servers and automatically migrates VMs to balance the load.
- **HA (High Availability):** If a physical server fails, the hypervisor automatically restarts its VMs on another server in the cluster, providing **failure transparency**.

KVM + libvirt Stack

The Open-Source Alternative

KVM (Kernel-based Virtual Machine) is a Linux kernel feature that turns the kernel into a Type 1 hypervisor. The 'libvirt' toolkit provides the management layer for KVM.

KVM + libvirt Stack

The Open-Source Alternative

KVM (Kernel-based Virtual Machine) is a Linux kernel feature that turns the kernel into a Type 1 hypervisor. The 'libvirt' toolkit provides the management layer for KVM.

Roles of the Stack:

- **KVM:** Provides the core hardware virtualization support, allowing guest OSes to run with near-native performance.
- **QEMU:** Emulates the hardware devices (e.g., virtual network cards, disks) that the guest OS sees.
- **libvirt:** A daemon that provides a consistent API to manage VMs, storage, and networking. It can perform coordination tasks like VM migration and resource control.

Summary: Coordination Models Comparison

System	Primary Isolation	Scheduling	
Kubernetes	Kernel Namespaces	Distributed Scheduler	Self-
VMware vSphere	Hypervisor Abstraction	DRS Load Balancing	HA (

Integrated Reflection

The Big Picture

The principles of operating systems extend from a single machine to entire data centers.

Integrated Reflection

The Big Picture

The principles of operating systems extend from a single machine to entire data centers.

- **Resource Management:** Single-OS schedulers and memory managers become distributed schedulers and load balancers in a cluster.
- **Synchronization:** Semaphores and locks evolve into consensus algorithms like Raft and Paxos to manage shared state across many nodes.
- **Fault Tolerance:** Simple process restart policies become complex self-healing and failover mechanisms.

The goal remains the same: provide a consistent, reliable, and efficient computing environment.

Key Takeaways

- Modern systems use a combination of kernel primitives and high-level tools for coordination.
- **Kubernetes** acts as a distributed OS, managing containerized workloads using concepts like scheduling and self-healing.
- **VMware vSphere** provides similar coordination for full VMs, leveraging the hypervisor as the central control point.
- The fundamental challenges of OS design, such as resource management and fault tolerance, are directly applicable to distributed systems.

Next Week Preview: Final Integration + Capstone Projects

This is our final topic for the course. Next week, we will synthesize everything we've learned into a comprehensive capstone.

- **Course-Wide Integration:** Tying together the concepts of processes, memory, file systems, and distributed coordination.
- **Project Synthesis:** Discussing how these modules can be combined into a single, cohesive project.
- **Final Q&A:** An open forum for questions and review.

Appendix: Quiz

Conceptual Questions

- ① How does Kubernetes' use of 'etcd' and the Raft algorithm relate to the concepts of "distributed consensus" we discussed?
- ② What is the main difference between how Kubernetes and VMware's vSphere achieve high availability (HA)?
- ③ What is the primary role of the 'libvirt' daemon in a KVM environment?

Appendix: Exercises

Design & Analysis

- ① Design a simple, multi-node container orchestration system that uses a centralized coordinator. What are its major drawbacks compared to a decentralized system like Kubernetes?
- ② Research the concept of a "sidecar container" in a Kubernetes Pod. How does it leverage OS-level coordination to solve a common problem?

Appendix: Advanced Topics to Explore

Further Reading

- **Service Meshes (Istio, Linkerd):** A network layer that handles communication, security, and observability between services in a distributed system.
- **Kubernetes Operators:** A method of packaging, deploying, and managing a Kubernetes application.
- **The 'containerd' Project:** The core container runtime that provides the container creation and management functionality used by Docker and Kubernetes.

Week 12

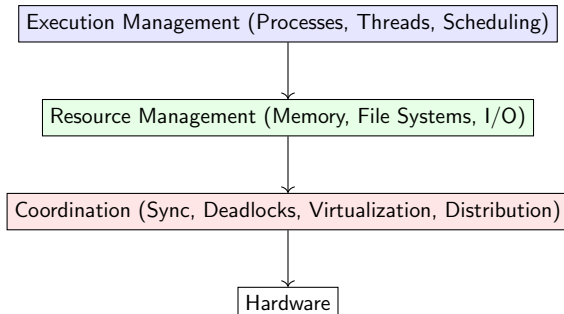
Agenda

- ① Course-Wide Integration: A Unified View of OS Concepts
- ② Revisit of Key OS Themes: Execution, Resources, Coordination
- ③ Capstone Project Opportunities and Scoping
- ④ Emerging Trends and Real-World Challenges
- ⑤ Final Thoughts and Takeaways

Operating System Model: A Layered Revisit

The Unified OS Stack

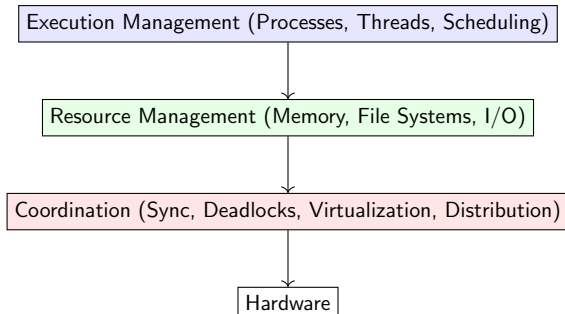
We've explored the operating system in distinct modules, but in reality, all components are deeply interconnected. This layered model shows how each part builds upon the last.



Operating System Model: A Layered Revisit

The Unified OS Stack

We've explored the operating system in distinct modules, but in reality, all components are deeply interconnected. This layered model shows how each part builds upon the last.



The system is a cohesive unit, where changes in one layer have cascading effects on the others.

Integrated View of Concepts

Interconnected Themes

The boundaries between our course modules are artificial. Many topics we discussed are intertwined, requiring a holistic understanding.

Integrated View of Concepts

Interconnected Themes

The boundaries between our course modules are artificial. Many topics we discussed are intertwined, requiring a holistic understanding.

- **Execution & Scheduling** directly impact **Resource Demand** (e.g., a CPU scheduler must consider memory locality).
- **Memory & I/O** management rely on **Coordination** mechanisms (e.g., a file system requires locks to prevent race conditions).
- **Synchronization** primitives extend naturally to **Distributed Coordination** (e.g., a mutex on a single machine becomes a consensus protocol in a cluster).
- **Virtualization** is a synthesis of all layers, providing isolated **Execution** and **Resource** management with advanced **Coordination** and security.

Capstone Project Ideas (Conceptual Scopes)

Putting Theory into Practice

Your capstone project is an opportunity to apply a deeper understanding of one or more of the course modules.

Capstone Project Ideas (Conceptual Scopes)

Putting Theory into Practice

Your capstone project is an opportunity to apply a deeper understanding of one or more of the course modules.

Execution & Scheduling Focused:

- Build a custom CPU scheduler simulator that visually demonstrates how different algorithms (e.g., Round Robin, CFS) affect process performance.
- Develop a thread scheduling performance benchmarking tool to compare kernel vs. user-level threading models.

Capstone Project Ideas (Conceptual Scopes)

Putting Theory into Practice

Your capstone project is an opportunity to apply a deeper understanding of one or more of the course modules.

Execution & Scheduling Focused:

- Build a custom CPU scheduler simulator that visually demonstrates how different algorithms (e.g., Round Robin, CFS) affect process performance.
- Develop a thread scheduling performance benchmarking tool to compare kernel vs. user-level threading models.

Resource & File System Focused:

- Create a page replacement policy visualizer (e.g., LRU, FIFO) to show hit/miss ratios in real-time.
- Implement a simple file system simulator with a custom block allocator and inode management.

Capstone Project Ideas (Conceptual Scopes)

Putting Theory into Practice

Your capstone project is an opportunity to apply a deeper understanding of one or more of the course modules.

Execution & Scheduling Focused:

- Build a custom CPU scheduler simulator that visually demonstrates how different algorithms (e.g., Round Robin, CFS) affect process performance.
- Develop a thread scheduling performance benchmarking tool to compare kernel vs. user-level threading models.

Resource & File System Focused:

- Create a page replacement policy visualizer (e.g., LRU, FIFO) to show hit/miss ratios in real-time.
- Implement a simple file system simulator with a custom block allocator and inode management.

Coordination & Distributed Systems Focused:

- Build a simplified container runtime that manages 'namespaces' and 'cgroups'.
- Implement a Paxos or Raft-based replicated key-value store to demonstrate distributed consensus.

End-to-End Project Suggestions

Synthesizing Across Modules

These ideas combine multiple course concepts into a single, comprehensive project.

End-to-End Project Suggestions

Synthesizing Across Modules

These ideas combine multiple course concepts into a single, comprehensive project.

- **Build a Container-Aware Micro-OS:**
 - ▶ **Execution:** Implement a basic process scheduler for your micro-OS.
 - ▶ **Resources:** Create memory sandboxing for each container.
 - ▶ **Coordination:** Use 'namespaces' and 'cgroups' to manage and isolate container processes and resources.
- **Simulate a Distributed Resource Manager:**
 - ▶ **Coordination:** Implement a leader election algorithm (e.g., Bully).
 - ▶ **Resources:** Have the leader manage resource tracking and dynamic scheduling.
 - ▶ **Execution:** Simulate job failures and have the system recover using the elected leader.

Open Challenges and Emerging Trends

The Future of OS Design

The field of operating systems is constantly evolving to address new hardware and security needs.

Open Challenges and Emerging Trends

The Future of OS Design

The field of operating systems is constantly evolving to address new hardware and security needs.

- **Kernel-Level Security:** Using technologies like 'eBPF' and 'seccomp' to create advanced security policies for containers without sacrificing performance.
- **Real-Time Scheduling:** Ensuring predictable performance for time-sensitive applications in multi-tenant environments.
- **Confidential Computing:** Using hardware-enforced memory encryption (e.g., Intel TDX, AMD SEV) to protect data even from the cloud provider or hypervisor.
- **Micro-Virtualization:** Lightweight VMs like 'Firecracker' offer the speed of containers with the strong isolation of traditional virtualization, a key technology for serverless computing.

Final Takeaways

- **Operating systems** are not a static field; they are living systems that evolve with the underlying hardware and the applications they support.
- **Coordination** is the fundamental challenge of all OS design, from managing concurrent threads on a single core to a globally distributed cluster.
- To be an effective systems thinker, you must understand the interplay between **execution, resource management, and coordination**.
- The capstone project is your opportunity to synthesize these concepts and build something impactful.