

## 18. Operating System Concepts

by Abraham Silberschatz, Peter B. Galvin, Greg Gagne

Audio (MP3) version: [https://books.kim/mp3/book/www.books.kim\\_747\\_summary-18\\_\\_Operating\\_System.mp3](https://books.kim/mp3/book/www.books.kim_747_summary-18__Operating_System.mp3)

### Summary:

Book 18 of Operating System Concepts by Abraham Silberschatz, Peter B. Galvin, Greg Gagne is an introduction to the fundamentals of operating systems. It covers topics such as process management, memory management, file-system implementation and protection mechanisms. The book also provides a comprehensive overview of distributed systems and networking concepts.

The first chapter introduces the basic concepts of operating systems including processes, threads, CPU scheduling algorithms and deadlocks. It also discusses virtual machines and their role in modern computing environments. The second chapter focuses on memory management techniques such as segmentation and paging along with virtual memory support for multiprogramming environments.

Chapter three examines file system implementation issues such as directory structures, access control lists (ACLs) and disk scheduling algorithms. Chapter four looks at I/O subsystems including device drivers, interrupt handling strategies and DMA controllers. Chapters five through seven cover various aspects of distributed systems including client-server architectures, remote procedure calls (RPCs), network protocols and security measures.

The eighth chapter explores protection mechanisms used to ensure secure operation within a computer system while chapters nine through eleven discuss different types of communication networks from local area networks (LANs) to wide area networks (WANs). Finally chapters twelve through fourteen provide an overview of advanced topics such as real-time operating systems (RTOSes), embedded systems programming languages like C++ or Java.

### Main ideas:

**#1. *Process Management: Process management is the management of the various processes that are running on a computer system. It involves the creation, scheduling, and termination of processes, as well as the management of resources used by the processes.***

Process management is an important part of any operating system. It involves the creation, scheduling, and termination of processes that are running on a computer system. Processes can be created by either the user or the operating system itself, depending on what type of task needs to be accomplished. The process scheduler then determines which processes should run at any given time based on their priority level and other factors such as memory availability.

Once a process has been scheduled for execution, it must be managed in order to ensure that it runs efficiently and does not interfere with other processes running on the same machine. This includes managing resources such as CPU time, memory usage, disk accesses, network connections etc., so that each process gets its fair share of resources without overloading the system or causing conflicts between different programs.

Finally, when a process is no longer needed it must be terminated properly in order to free up resources for other tasks. This requires careful coordination between all parts of the operating system in order to make sure that everything is cleaned up correctly before terminating a process.

**#2. *Memory Management: Memory management is the process of allocating and deallocating memory to processes in order to ensure that the system runs efficiently. It involves the use of virtual memory, segmentation, and paging to manage memory.***

Memory management is an essential part of any operating system. It involves the allocation and deallocation of memory to processes in order to ensure that the system runs efficiently. Memory management techniques such as virtual memory, segmentation, and paging are used to manage memory effectively. Virtual memory allows a process to access more physical memory than what is available on the computer by using disk space as additional RAM. Segmentation divides a program into logical segments which can be stored separately in different areas of main memory or even swapped out onto secondary storage devices when not needed. Paging breaks up large chunks of data into smaller pages for easier retrieval from main memory.

The goal of effective memory management is to provide enough resources for all running programs while minimizing wastage due to fragmentation or inefficient use of available resources. To achieve this, modern operating systems employ sophisticated algorithms and techniques such as garbage collection, compaction, caching, prefetching etc., which help optimize resource utilization.

**#3. *Storage Management: Storage management is the process of managing the storage of data on a computer system. It involves the use of file systems, disk scheduling algorithms, and RAID systems to ensure that data is stored efficiently and securely.***

Storage management is an important part of any computer system. It involves the use of file systems, disk scheduling algorithms, and RAID systems to ensure that data is stored efficiently and securely. File systems are used to organize data into logical units such as files and directories. Disk scheduling algorithms determine how requests for data are handled by the storage device in order to optimize performance. RAID (Redundant Array of Independent Disks) systems provide redundancy so that if one disk fails, the other disks can be used to recover lost data.

In addition to these components, storage management also includes backup strategies which help protect against accidental or malicious loss of data. Backup strategies involve regularly copying critical files onto a separate medium such as tape or another hard drive in case something happens to the original copy on the primary storage device. This ensures that there is always a recent version available should anything happen.

Finally, security measures must be taken when dealing with sensitive information stored on a computer system. Encryption techniques can be used to scramble confidential information so it cannot be read without authorization from someone who knows the encryption key.

**#4. *Protection and Security: Protection and security are essential for ensuring the integrity of a computer system. It involves the use of access control mechanisms, authentication protocols, and encryption algorithms to protect data from unauthorized access.***

Protection and security are essential for ensuring the integrity of a computer system. It involves the use of access control mechanisms, authentication protocols, and encryption algorithms to protect data from unauthorized access. Access control mechanisms allow only authorized users to gain access to resources within a system. Authentication protocols verify that users attempting to gain access are who they claim to be. Encryption algorithms scramble data so that it is unreadable by anyone without the proper key or password.

In addition, firewalls can be used as an additional layer of protection against malicious attacks on a network or system. Firewalls act as gatekeepers between networks, allowing only certain types of traffic through while blocking others. They also monitor incoming and outgoing traffic for suspicious activity.

Finally, antivirus software can help protect systems from viruses and other malware threats by scanning files for known patterns associated with malicious code before they are allowed into the system.

**#5. *Networking: Networking is the process of connecting computers together in order to share resources and information. It involves the use of protocols, such as TCP/IP, to enable communication between computers.***

Networking is an essential part of modern computing. It allows computers to communicate with each other and share resources, such as files, printers, and applications. Networking also enables users to access the Internet from any location in the world. The process of networking involves connecting computers together using a variety of protocols, such as TCP/IP. These protocols enable communication between different types of computers and networks.

The most common type of network is a local area network (LAN). A LAN consists of two or more connected devices that are located within close proximity to one another. This type of network can be used for sharing files, printers, and other resources among multiple users on the same physical network segment. Other types of networks include wide area networks (WANs), which connect geographically dispersed locations; metropolitan area networks (MANs), which connect multiple LANs within a city; and virtual private networks (VPNs), which provide secure remote access over public infrastructure.

Network security is an important consideration when setting up any kind of computer network. Security measures should be taken to protect data from unauthorized access or malicious attacks by hackers or viruses. Firewalls can be used to restrict incoming traffic while encryption techniques can help ensure that data remains confidential even if it is intercepted by third parties.

**#6. *Distributed Systems: Distributed systems are computer systems that are composed of multiple computers that are connected together. It involves the use of distributed algorithms, such as distributed mutual exclusion, to ensure that the system runs efficiently.***

Distributed systems are computer systems that are composed of multiple computers connected together. These distributed systems allow for the sharing of resources and data between different nodes in the system, allowing for greater scalability and flexibility than traditional centralized computing architectures. Distributed algorithms such as distributed mutual exclusion ensure that the system runs efficiently by ensuring that only one node can access a resource at any given time. This prevents conflicts from occurring when two or more nodes attempt to access the same resource simultaneously.

In addition to providing scalability and flexibility, distributed systems also provide fault tolerance. If one node fails, other nodes can take over its tasks until it is restored or replaced. This ensures that services remain available even if individual components fail.

Distributed systems have become increasingly popular due to their ability to scale easily with increasing demand and their ability to provide reliable service even in cases of component failure.

**#7. *Deadlocks: Deadlocks are a situation in which two or more processes are waiting for resources that are held by the other processes. It involves the use of deadlock detection and prevention algorithms to ensure that the system does not enter a deadlock state.***

Deadlocks are a common problem in computer systems, and can cause serious issues if not managed properly. A deadlock occurs when two or more processes are waiting for resources that are held by the other processes. This situation is known as a circular wait, where each process is waiting on another to release its resource before it can proceed. In order for the system to continue functioning normally, these deadlocks must be avoided or resolved quickly.

To prevent deadlocks from occurring, there are several techniques that can be used. One of the most popular methods is called "deadlock detection" which involves monitoring the state of all active processes and detecting any potential conflicts between them. If a conflict is detected then an algorithm will take action to resolve it before it becomes a full-blown deadlock.

Another technique used to avoid deadlocks is called "deadlock prevention" which involves making sure that certain conditions cannot occur in order for a deadlock to happen in the first place. For example, one condition might require

that no process holds multiple resources at once; this would ensure that no single process could block access to those resources and create a circular wait situation.

Finally, there are also algorithms available for resolving existing deadlocks should they occur despite preventive measures being taken. These algorithms typically involve identifying one or more processes involved in the circular wait and releasing their resources so that others may use them instead.

**#8. *Virtual Machines: Virtual machines are software-based systems that emulate the behavior of a physical computer. It involves the use of virtualization technology to create multiple virtual machines on a single physical machine.***

Virtual machines are a powerful tool for creating and managing multiple computing environments on a single physical machine. By using virtualization technology, it is possible to create multiple virtual machines that can run different operating systems and applications independently of each other. This allows users to have the flexibility to switch between different operating systems or applications without having to reboot their computer.

The use of virtual machines also provides an efficient way of running multiple programs simultaneously on one machine. It eliminates the need for additional hardware resources such as memory, storage space, and processing power which would otherwise be required if all programs were running directly on the host system.

In addition, virtual machines provide enhanced security by isolating each program from others in its own environment. This prevents malicious code from affecting other programs or data stored on the same physical machine.

Overall, virtual machines offer many advantages over traditional computing methods including increased efficiency, cost savings, improved security and flexibility. As such they are becoming increasingly popular among businesses looking for ways to maximize their IT infrastructure while minimizing costs.</p></div>

**#9. *Operating System Structures: Operating system structures are the components of an operating system that are responsible for managing the resources of a computer system. It involves the use of kernel, device drivers, and system libraries to manage the resources of the system.***

Operating system structures are the components of an operating system that are responsible for managing the resources of a computer system. The kernel is the core component of an operating system and is responsible for providing basic services such as memory management, process scheduling, and device drivers. Device drivers provide access to hardware devices such as printers, scanners, and disk drives. System libraries provide functions that allow applications to interact with the operating system in order to perform tasks such as file I/O or network communication. These components work together to manage all aspects of a computers resources.

The kernel provides low-level services which can be used by other parts of the operating system or by user programs. It also manages physical memory and virtual memory so that processes have enough space to run without interfering with each other. The device drivers enable applications to communicate with hardware devices in order to read from or write data onto them. System libraries contain code which allows user programs access certain features provided by the OS like networking capabilities or graphical interfaces.

In addition, there are various tools available which help administrators configure their systems according to their needs. These include utilities for setting up users accounts, configuring security settings, monitoring performance metrics etcetera.

Overall, these components form the basis upon which modern day computing systems operate efficiently and securely while providing users with powerful features they need in order complete their tasks.</p></div>

**#10. Process Synchronization: Process synchronization is the process of ensuring that multiple processes are able to access shared resources without interfering with each other. It involves the use of semaphores, monitors, and message passing to ensure that processes are synchronized.**

Process synchronization is an important concept in operating systems, as it allows multiple processes to access shared resources without interfering with each other. It involves the use of semaphores, monitors, and message passing to ensure that processes are synchronized. Semaphores are used to control access to a shared resource by allowing only one process at a time to have access. Monitors provide mutual exclusion for critical sections of code so that two or more processes cannot execute them simultaneously. Message passing is used when two or more processes need to communicate with each other in order for synchronization to occur.

The main goal of process synchronization is ensuring that all operations on the shared resource take place in the correct order and no data corruption occurs due to concurrent access from multiple processes. This can be achieved through various techniques such as locking mechanisms, atomic transactions, and deadlock avoidance algorithms.

Process synchronization plays an important role in distributed computing systems where multiple computers must coordinate their activities over a network connection. In this case, message passing protocols such as TCP/IP are often used for communication between nodes.

**#11. Deadlock Handling: Deadlock handling is the process of dealing with deadlocks when they occur. It involves the use of deadlock detection and recovery algorithms to ensure that the system is able to recover from a deadlock state.**

Deadlock handling is an important part of operating system design. A deadlock occurs when two or more processes are waiting for resources that cannot be allocated to them due to a circular wait condition. This can lead to the system becoming unresponsive and unable to continue processing requests. To prevent this, deadlock detection and recovery algorithms must be implemented in order to detect and resolve any potential deadlocks.

Deadlock detection algorithms typically involve monitoring the state of each process in the system, as well as their resource requirements. If a circular wait condition is detected, then it can be resolved by either pre-empting one of the processes involved or by rolling back some of its operations so that it no longer requires the resources it was waiting for. Once a deadlock has been identified and resolved, recovery algorithms can then be used to restore normal operation.

In addition, prevention techniques such as resource ordering schemes may also be employed in order to reduce the likelihood of a deadlock occurring in the first place. These schemes ensure that all resources are requested in an ordered fashion which prevents any circular waits from forming.

**#12. Memory Allocation: Memory allocation is the process of allocating memory to processes in order to ensure that the system runs efficiently. It involves the use of segmentation, paging, and virtual memory to manage memory efficiently.**

Memory allocation is an important part of operating system design. It involves the use of segmentation, paging, and virtual memory to manage memory efficiently. Segmentation divides a program into logical segments that can be loaded separately into main memory when needed. Paging allows programs to access more physical memory than is available in main memory by swapping out portions of the program not currently being used. Virtual memory provides a larger address space for programs than what is physically available in main memory by using disk storage as an extension of RAM.

The goal of efficient memory allocation is to ensure that all processes have enough resources to run without causing any conflicts or bottlenecks. This requires careful management and monitoring of both physical and virtual memories so that each process has sufficient resources while still allowing other processes to run smoothly. Memory allocation also helps reduce fragmentation, which occurs when there are many small chunks of free space scattered throughout the systems

address space.

In order for effective resource management, it is important for the operating system to keep track of how much physical and virtual memories are allocated at any given time. The operating system must also be able to determine which processes need more or less resources depending on their current state and priority level.

**#13. File Systems: File systems are the components of an operating system that are responsible for managing the storage of data on a computer system. It involves the use of file systems, such as FAT and NTFS, to ensure that data is stored efficiently and securely.**

A file system is a set of rules and procedures that an operating system uses to manage the storage, retrieval, and updating of data on a computer. It provides the means for organizing files and directories, as well as keeping track of which areas of memory are used by which files. File systems also provide security measures such as access control lists (ACLs) to ensure that only authorized users can access certain files or directories.

The most common type of file system is the hierarchical file system (HFS), which organizes data into folders and subfolders in a tree-like structure. This allows users to easily navigate through their data without having to remember exact locations or filenames. Other popular types include networked file systems (NFS), distributed file systems (DFS), object-based storage devices (OBSDs), and virtualized file systems.

File systems play an important role in ensuring that data is stored securely and efficiently on computers. They allow users to quickly locate specific information when needed, while also providing safeguards against unauthorized access or accidental deletion. As technology advances, new types of file systems are being developed with improved features such as better scalability, faster performance, increased reliability, enhanced security measures, etc.

**#14. I/O Systems: I/O systems are the components of an operating system that are responsible for managing the input and output of data on a computer system. It involves the use of device drivers, interrupt handlers, and DMA controllers to ensure that data is transferred efficiently.**

I/O systems are an essential part of any operating system. They provide the means for a computer to interact with its environment, allowing it to receive input from users and other devices, as well as send output back out. I/O systems involve the use of device drivers, interrupt handlers, and DMA controllers in order to ensure that data is transferred efficiently between the CPU and external devices.

Device drivers act as intermediaries between hardware components and software applications. They allow software programs to access hardware resources without needing direct knowledge of how those resources work. Interrupt handlers are responsible for responding to interrupts generated by external devices or internal processes within the system. Finally, DMA controllers manage Direct Memory Access (DMA) operations which enable faster data transfers than would be possible using only processor cycles.

The combination of these components allows computers to communicate with their environment in a reliable manner while also ensuring efficient utilization of available resources. Without them, modern computing would not be possible.

**#15. Security: Security is the process of protecting data from unauthorized access. It involves the use of access control mechanisms, authentication protocols, and encryption algorithms to ensure that data is secure.**

Security is an essential component of any operating system. It involves the use of access control mechanisms, authentication protocols, and encryption algorithms to ensure that data is secure from unauthorized access. Access control mechanisms are used to restrict user access to certain resources or operations within a system. Authentication protocols verify the identity of users attempting to gain access to a system, while encryption algorithms protect data by scrambling it so that only authorized users can read it.



The security measures implemented in an operating system must be tailored for its specific environment and purpose. For example, if the OS is designed for use in a corporate setting with multiple users accessing sensitive information, then more stringent security measures may need to be put into place than if it were being used on a personal computer at home.

In addition, security measures should also take into account potential threats such as malware or hackers who might attempt to gain unauthorized access. Security policies should be regularly reviewed and updated as needed in order to keep up with changing technology and new threats.

**#16. *Distributed File Systems: Distributed file systems are file systems that are composed of multiple computers that are connected together. It involves the use of distributed algorithms, such as distributed mutual exclusion, to ensure that the system runs efficiently.***

Distributed file systems are a type of computer system that allows multiple computers to access and share files over a network. This type of system is composed of multiple computers connected together, each with its own set of resources. The distributed file system uses distributed algorithms such as distributed mutual exclusion to ensure the efficient operation of the entire system.

The main advantage of using a distributed file system is that it provides users with access to data from any location on the network. It also enables users to store large amounts of data in an organized manner, making it easier for them to find what they need quickly and efficiently. Additionally, since all nodes in the network can access shared files, this makes collaboration between different users much simpler.

In order for a distributed file system to work properly, there must be some form of coordination among all nodes involved in the process. This coordination ensures that no two processes try to modify or delete the same file at once and helps maintain consistency across all nodes in the network. Furthermore, security measures must be taken into account when setting up a distributed file system so as not to compromise sensitive information.

**#17. *Operating System Design: Operating system design is the process of designing an operating system that is able to meet the needs of the user. It involves the use of modular design, layered design, and microkernel design to ensure that the system is efficient and reliable.***

Operating system design is a complex process that requires careful consideration of the users needs. Modular design involves breaking down the operating system into smaller, more manageable components that can be easily modified and updated as needed. Layered design allows for different levels of abstraction to be used in order to simplify the development process and make it easier to maintain. Finally, microkernel design focuses on creating a minimal set of core functions that are necessary for an operating system to function properly.

The goal of any operating system designer is to create an efficient and reliable product that meets all the requirements of its users. This means taking into account factors such as performance, scalability, security, reliability, portability, extensibility and usability when designing an OS. It also involves making sure that all components work together seamlessly so they can provide a consistent experience across multiple platforms.

Designing an effective operating system requires knowledge in many areas including computer architecture, software engineering principles and algorithms. It also requires creativity in order to come up with innovative solutions for problems encountered during development.

**#18. *Real-Time Systems: Real-time systems are computer systems that are designed to respond to events within a certain time frame. It involves the use of real-time scheduling algorithms, such as Earliest Deadline First, to ensure that the system is able to meet its deadlines.***

Real-time systems are computer systems that are designed to respond to events within a certain time frame. These

systems must be able to process data quickly and accurately in order for them to meet their deadlines. To ensure this, real-time scheduling algorithms such as Earliest Deadline First (EDF) are used. EDF is an algorithm which assigns priorities based on the deadline of each task; tasks with earlier deadlines will have higher priority than those with later ones. This ensures that all tasks can be completed before their respective deadlines.

In addition, real-time systems also employ other techniques such as preemption and rate monotonic analysis (RMA). Preemption allows the system to interrupt a running task if it needs more resources or if another task has a higher priority. RMA is an algorithm which assigns priorities based on how often each task needs to run; tasks that need to run more frequently will have higher priority than those that don't need to run as often.

Real-time systems are essential for many applications where accuracy and speed of response is critical, such as medical equipment, industrial automation, military operations and air traffic control. By using these scheduling algorithms and techniques, real-time systems can guarantee timely responses while still ensuring accuracy.

**#19. *Distributed Operating Systems: Distributed operating systems are operating systems that are composed of multiple computers that are connected together. It involves the use of distributed algorithms, such as distributed mutual exclusion, to ensure that the system runs efficiently.***

Distributed operating systems are a type of computer system that is composed of multiple computers connected together. These computers can be located in different physical locations, or they can even be virtual machines running on the same physical machine. The distributed operating system allows for the sharing of resources and data between all nodes in the network, as well as providing fault tolerance and scalability.

In order to ensure efficient operation, distributed algorithms such as distributed mutual exclusion must be used. This algorithm ensures that only one process at a time has access to shared resources, preventing any conflicts from occurring. Additionally, it also helps with load balancing by ensuring that each node receives an equal amount of work.

The use of distributed operating systems provides many advantages over traditional single-node systems. For example, since multiple nodes are involved in processing tasks, there is increased redundancy which makes them more reliable than single-node systems. Furthermore, due to their ability to scale easily across multiple nodes they are able to handle larger workloads than traditional single-node systems.

**#20. *System Performance: System performance is the process of measuring the performance of a computer system. It involves the use of performance metrics, such as response time and throughput, to measure the performance of the system.***

System performance is an important factor in the success of any computer system. It involves measuring how well a system performs its intended tasks, and can be used to identify areas for improvement. Performance metrics such as response time and throughput are commonly used to measure the performance of a system. Response time measures how quickly a task is completed, while throughput measures the amount of work that can be done within a given period of time.

Performance tuning is an important part of improving system performance. This involves making changes to hardware or software components in order to improve their efficiency and reduce resource consumption. For example, increasing memory size or adding faster processors may help improve response times by reducing wait times for data retrieval from storage devices.

Monitoring tools are also useful for tracking system performance over time. These tools allow administrators to track key metrics such as CPU utilization, disk I/O rates, network traffic levels, and more. By monitoring these metrics regularly, administrators can detect potential problems before they become serious issues.