

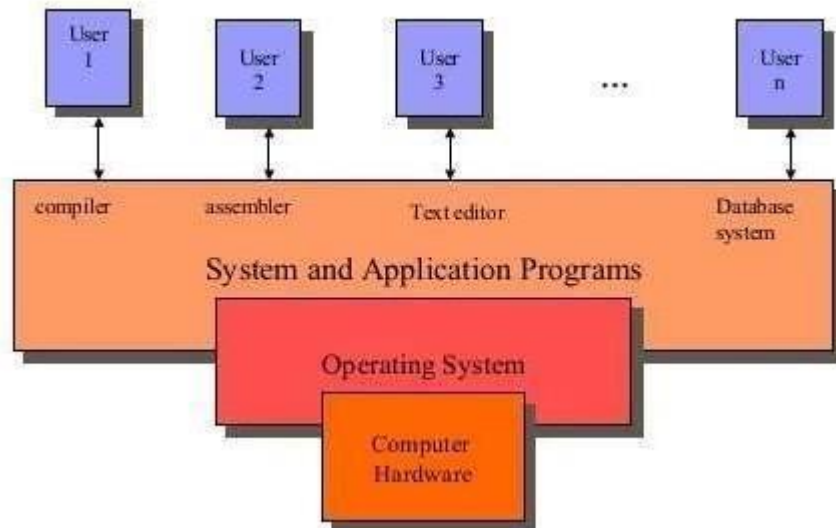
UNIT-1

Operating System Introduction-Structure-Simple batch, Multiprogramming, Time-shared, Personal Computer, Parallel, Distributed System, Real time System, System Components, OS services, System calls.

Process-Process concepts and scheduling, Operations on process, Cooperating process, Threads

Introduction of Operating System

- An OS act as an interface between user and system hardware.
- Computer consists of the hardware, Operating System, system programs, application programs.
- The hardware consists of memory, CPU, ALU, I/O device, storage device and peripheral device.
- System program consists of compilers, loaders, editors, OS etc.
- Application program consists of database programs, business programs.
- Every computer must have an OS to run other programs.
- The OS controls & coordinates the use of the hard ware among the various system programs and application programs for various tasks.
- Its imply provides an environment with in which other programs can do useful work.



OPERATINGSYSTEM

Definition

- In the 1960's one might have defined OS as **“The software that controls the hardware”**.
- Operating System performs all the basic tasks like managing files, processes,

OPERATING SYSTEM(23CS403)

And memory. Thus operating system acts as the manager of all the resources, i.e.

Resource manager.

- Operating system becomes an interface between the user and the machine. It is one of the most required software that is present in the device.
- Operating System is a type of software that works as an interface between the system program and the hardware.

Concept of OS

- The OS is a set of special programs that run on a computer system that allow it to work properly.
- It performs basic task as recognizing input from the keyboard, keeping track of files and directories on the disk, sending output to the display screen and controlling a peripheral device.
- The OS must support the following tasks. They are,
 - Provides the facilities to create, modification of program and data file using an editor.
 - Access to the compiler for translating the user program from high level language to machine language.
 - Provide a loader program to move the compiled program code to the computer memory for execution.

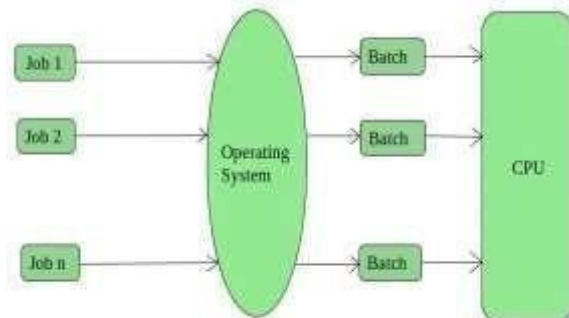
Types of Operating Systems

There are several types of Operating Systems which are mentioned below.

- Batch Operating System
- Multi-Programming System
- Time-Sharing Operating System
- Personal Computers
- Parallel Operating System
- Distributed Operating System
- Real-Time Operating System

1. Batch Operating System

This type of operating system does not interact with the computer directly. There is an operator which takes similar job shaving the same requirement and groups them into batches. It is the responsibility of the operator to sort jobs with similar needs.



Advantages

OPERATING SYSTEM (23CS403)

- Processors of the batch systems know how long the job would be when it is in the queue.
- Multiple users can share the batch systems.
- The idle time for the batch system is very less.
- It is easy to manage large work repeatedly in batch systems.

Disadvantages

- The computer operators should be well known with batch systems.
- Batch systems are hard to debug.
- It is sometimes costly.
- The other jobs will have to wait for an unknown time if any job fails.
- It is very difficult to guess or know the time required for any job to complete.

Examples

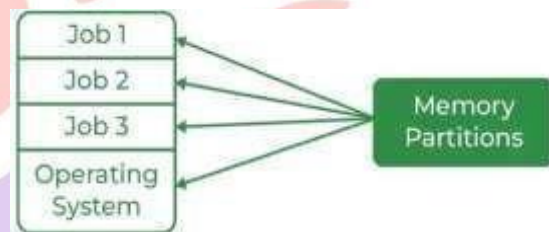
Payroll Systems, Bank Statements, etc.

2. Multi-Programming Operating System

Multi programming Operating Systems can be simply illustrated as more than one program is present in the main memory and any one of them can be kept in execution.

This is basically used for better

Execution of resources.



Advantages of Multi-Programming Operating System

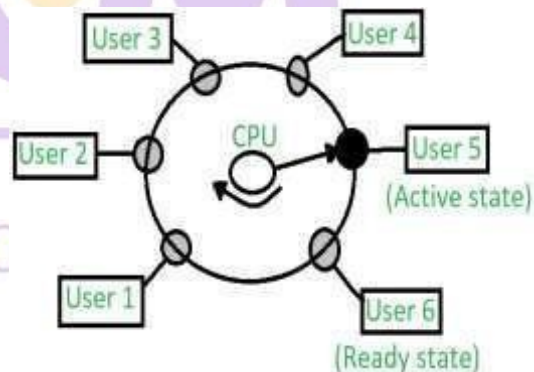
- Multi Programming increases the Through put of the System.
- It helps in reducing there sponse time.

Disadvantages of Multi-Programming Operating System

- There is not any facility for user interaction of system resources with the system.

3. Time-Sharing Operating Systems

Each task is given some time to execute so that all the tasks work smoothly. Each user gets the time of the CPU as they use a single system. These systems are also known as Multitasking Systems. The task can be from a single user or different users also. The time that each task gets to execute is called quantum. After this time interval is over OS switches over to the next task.



Advantages

- Each task gets an equal opportunity.
- Fewer chances of duplication of software.
- CPU idle time can be reduced.
- **Resource Sharing:** Time-sharing systems allow multiple users to share hardware resources such as the CPU, memory, and peripherals, reducing the cost of hardware and increasing efficiency.
- **Improved Productivity:** Time-sharing allows users to work concurrently, thereby reducing the waiting time for their turn to use the computer. This increased productivity translates to more work getting done in less time.
- **Improved User Experience:** Time-sharing provides an interactive environment that allows users to communicate with the computer in real time, providing a better user experience than batch processing.

Disadvantages

- Reliability problem.
- One must have to take care of the security and integrity of user programs and data.
- Data communication problem.
- **High Overhead:** Time-sharing systems have a higher overhead than other operating systems due to the need for scheduling, context switching, and other overheads that come with supporting multiple users.
- **Complexity:** Time-sharing systems are complex and require advanced software to manage multiple users simultaneously. This complexity increases the chance of bugs and errors.
- **Security Risks:** With multiple users sharing resources, the risk of security breaches increases. Time-sharing systems require careful management of user access, authentication, and authorization to ensure the security of data and software.
- - sharing operating system that allows multiple users to access a Windows server remotely. Users can run their own applications and access shared resources, such as printers and network storage, in real-time.

4. Personal Computer

A personal computer (PC) is a microcomputer designed for use by one person at a time.

Prior to the PC, computers were designed for -- and only affordable for -- companies that attached terminals for multiple users to a single large mainframe computer whose resources were shared among all users. By the 1980s, technological advances made it feasible to build a small computer that an individual could own and use as a word processor and for other computing functions.

OPERATING SYSTEM(23CS403)

Whether they are home computers or business ones, PCs can be used to store, retrieve and process data of all kinds. A PC runs firmware that supports an operating system (OS), which supports a spectrum of other software. This software lets consumers and business users perform a range of general-purpose tasks, such as the following:

- word processing
- spreadsheets
- email
- instant messaging
- accounting
- database management
- internet access
- listening to music
- network-attached storage
- graphic design
- music composition
- video gaming
- software development
- network reconnaissance
- multimedia servers
- wireless network access hot spots
- video conferencing

Types

Personal computers fall into various categories, such as the following:

- **Desktop computers** usually have a tower, monitor, keyboard and mouse.
- **Tablets** are mobile devices with a touch screen display.
- **Smart phones** are phones with computing capabilities.
- **Wearables** are devices users wear, such as smartwatches and various types of smart clothing.
- **Laptop computers** are portable personal computers that usually come with an attached keyboard and trackpad.
- **Note book computers** are light weight laptops.
- **Handheld computers** include advanced calculators and various gaming devices.

5. Parallel Operating System

Parallel Systems are designed to speed up the execution of programs by dividing the programs into multiple fragments and processing these fragments at the same time.

Advantages

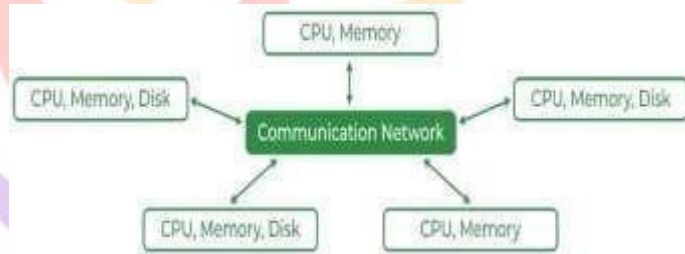
- **High Performance:** Parallel systems can execute computationally intensive tasks more quickly compared to single processor systems.
- **Cost Effective:** Parallel systems can be more cost-effective compared to distributed systems, as they do not require additional hardware for communication.

Disadvantages

- **Limited Scalability:** Parallel systems have limited scalability as the number of processors or cores in a single computer is finite.
- **Complexity:** Parallel systems are more complex to program and debug compared to single processor systems.
- **Synchronization Overhead:** Synchronization between processors in a parallel system can add overhead and impact performance.

6. Distributed Operating System

These types of operating systems are a recent advancement in the world of computer technology and are being widely accepted all over the world and, that too, at a great pace.



Various autonomous interconnected computers communicate with each other using a shared communication network. Independent systems possess their own memory unit and CPU. These are referred to as loosely coupled systems or distributed systems. These systems' processors differ in size and function.

The major benefit of working with these types of the operating system is that it is always possible that one user can access the files or software which are not actually present on his system but some other system connected within this network i.e., remote access is enabled within the devices connected in that network.

Types of Distributed Systems

The nodes in the distributed systems can be arranged in the form of client/server systems or peer to peer systems. Details about these are as follows –

Client/Server Systems

In client server systems, the client requests a resource and the server provides that resource. A server may serve multiple clients at the same time while a client is in contact with only one server. Both the client and server usually communicate via a computer network and so they are a part of distributed systems.

Peer to Peer Systems

The peer to peer systems contains nodes that are equal participants in data sharing. All the tasks are equally divided between all the nodes. The nodes interact with each other as required as share resources. This is done with the help of a network.

Advantages

- Failure of one will not affect the other network communication, as all systems are independent of each other.
- Electronic mail increases the data exchange speed.
- Sincere resources are being shared, computation is highly fast and durable.
- Load on host computer reduces.
- These systems are easily scalable as many systems can be easily added to the network.
- Delay in data processing reduces.

Disadvantages

- Failure of the main network will stop the entire communication.
- To establish distributed systems the language is used not well-defined yet.
- These types of systems are not readily available as they are very expensive. Not only that the underlying software is highly complex and not understood well yet.

Example: LOCUS

7. Real-Time Operating System

These types of OSs serve real-time systems. The time interval required to process and respond to inputs is very small. This time interval is called **response time**.

Real-time systems are used when there are time requirements that are very strict like missile systems, air traffic control systems, robots, etc.

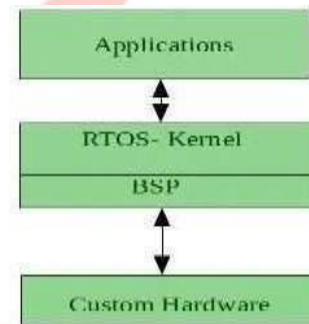
Types:

1. Hard Real-Time Systems

Hard Real-Time OSs are meant for applications where time constraints are very strict and even the shortest possible delay is not acceptable. These systems are built for saving life like automatic parachutes or airbags which are required to be readily available in case of an accident. Virtual memory is rarely found in these systems.

2. Soft Real-Time Systems

These OSs are for applications where time-constraint is less strict.



Advantages

- **Maximum Consumption:** Maximum utilization of devices and systems, thus more output from all the resources.
- **Task Shifting:** The time assigned for shifting tasks in these systems is very less. For example, in older systems, it takes about 10 microseconds in shifting from one task to another, and in the latest systems, it takes 3 microseconds.

OPERATING SYSTEM(23CS403)

- **Focus on Application:** Focus on running applications and less importance on applications that are in the queue.
- **Real-time operating system in the embedded system:** Since the size of programs is small, RTOS can also be used in embedded systems like in transport and others.
- **Error Free:** These types of systems are error-free.
- **Memory Allocation:** Memory allocation is best managed in these types of systems.

Disadvantages

- **Limited Tasks:** Very few tasks run at the same time and their concentration is very less on a few applications to avoid errors.
- **Use heavy system resources:** Sometimes the system resources are not so good and they are expensive as well.
- **Complex Algorithms:** The algorithms are very complex and difficult for the designer to write on.
- **Device driver and interrupt signals:** It needs specific device drivers and interrupt signals to respond earliest to interrupts.
- **Thread Priority:** It is not good to set thread priority as these systems are very less prone to switching tasks.

Examples

Scientific experiments, medical imaging systems, industrial control systems, weapon systems, robots, air traffic control systems, etc.

Operating System Services

- **User Interface** - User interface is essential and all operating systems provide it. Users either interface with the operating system through command-line interface (CUI) or graphical user interface (GUI). Command interpreter executes next user-specified command. A GUI offers the user a mouse-based window and menu system as an interface.
- **Program execution** - The system must be able to load a program into memory and to run that program, end execution, either normally or abnormally (indicating error)
- **I/O operations** - A running program may require I/O, which may involve a file or an I/O device.
- **File-system manipulation** - The file system is of particular interest. Obviously, programs need to read and write files and directories, create and delete them, search them, list file information, permission management.
- **Communications** - Processes may exchange information, on the same

OPERATING SYSTEM(23CS403)

computer or between computers over a network. Communications may be via shared memory or through message passing (packets moved by the OS)

- **Error detection** – OS needs to be constantly aware of possible errors may occur in the CPU and memory hardware, in I/O devices, in user program. For each type of error, OS should take the appropriate action to ensure correct and consistent computing. Debugging facilities can greatly enhance the user's and programmer's abilities to efficiently use the system.

Another set of OS functions exists for ensuring the efficient operation of the system itself via resource sharing

- **Resource allocation** - When multiple users or multiple jobs running concurrently, resources must be allocated to each of them. Many types of resources such as CPU cycles, main memory, and file storage may have special allocation code, others such as I/O devices may have general request and release code.
- **Accounting** - To keep track of which users use how much and what kinds of computer resources
- **Protection and security** - The owners of information stored in a multiuser or networked computer system may want to control use of that information, concurrent processes should not interfere with each other. **Protection** involves ensuring that all access to system resources is controlled. **Security** of the system from outsiders requires user authentication, extends to defending external I/O devices from invalid access attempts. If a system is to be protected and secure, precautions must be instituted throughout it. A chain is only as strong as its weakest link.

System Calls

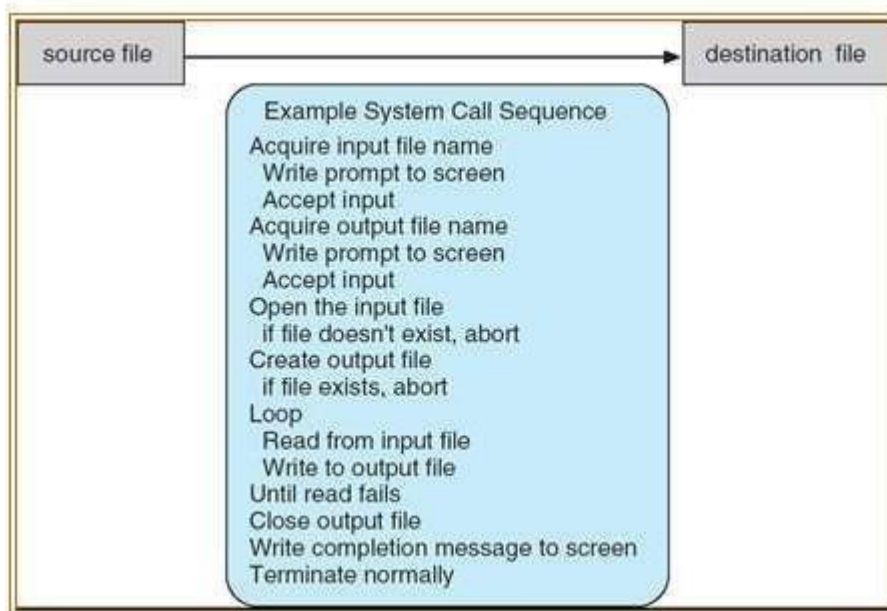
- A system call is a way for a user program to interface with the operating system. The program requests several services, and the OS responds by invoking a series of system calls to satisfy the request.
- A system call can be written in assembly language or a high-level language like **C, C++ or Pascal**.
- System calls are predefined functions that the operating system may directly invoke if a high-level language is used.
- A system call is a method for a computer program to request a service from the kernel of the operating system on which it is running.
- A system call is a method of interacting with the operating system via programs.
- A system call is a request from computer software to an operating system's kernel.

OPERATINGSYSTEM(23CS403)

- A simple system call may take few nanoseconds to provide the result, like retrieving the system date and time. A more complicated system call, such as connecting to a network device, may take a few seconds. Most operating systems launch a distinct kernel thread for each system call to avoid bottlenecks. Modern operating systems are multi-threaded, which means they can handle various system calls at the same time.
- The **Application Program Interface (API)** connects the operating system's functions to user programs. It acts as a link between the operating system and a process, allowing user-level programs to request operating system services. The kernel system can only be accessed using system calls. System calls are required for any programs that use resources.
- When computer software needs to access the operating system's kernel, it makes a system call. The system call uses an API to expose the operating system's services to user programs. It is the only method to access the kernel system. All programs or processes that require resources for execution must use system calls, as they serve as an interface between the operating system and user programs.

Example of System Calls

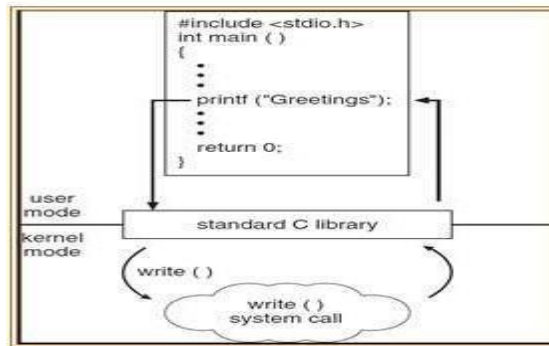
System call sequence to copy the contents of one file to another file



your roots to success...

Standard C Library Example

C program invoking printf() library call, which call write() system call



There are various situations where we must require system calls in the operating system. Following of the situations are as follows:

1. It must require when a file system wants to create or delete a file.
2. Network connections require the system call to send and receive data packets.
3. If you want to read or write a file, you need system calls.
4. If you want to access hardware devices, including a printer, scanner, you need a system call.
5. System calls are used to create and manage new processes.

Types of System Calls

There are commonly five types of system calls. These are as follows:

1. **Process Control**
2. **File Management**
3. **Device Management**
4. **Information Maintenance**
5. **Communication**

Process Control

Process control is the system call that is used to direct the processes. Some process control examples include creating, load, abort, end, execute, process, terminate the process, etc.

File Management

File management is a system call that is used to handle the files. Some file management examples include creating files, delete files, open, close, read, write, etc.

Device Management

Device management is a system call that is used to deal with devices. Some examples of device management include read, device, write, get device attributes,

OPERATINGSYSTEM(23CS403)

releasedevice,etc.

InformationMaintenance

Information maintenance is a system call that is used to maintain information. There are some examples of information maintenance, including getting system data, set time or date, get time or date, set system data, etc.

Communication

Communication is a system call that is used for communication. There are some examples of communication, including create, delete communication connections, send, receive messages, etc.

Examples of Windows and Unix system calls

Process	Windows	Unix
Process Control	CreateProcess() ExitProcess() WaitForSingleObject()	Fork() Exit() Wait()
File Manipulation	CreateFile() ReadFile() WriteFile() CloseHandle()	Open() Read() Write() Close()
Device Management	SetConsoleMode() ReadConsole() WriteConsole()	Ioctl() Read() Write()
Information Maintenance	GetCurrentProcessID() SetTimer() Sleep()	Getpid() Alarm() Sleep()
Communication	CreatePipe() CreateFileMapping() MapViewOfFile()	Pipe() Shmget() Mmap()
Protection	SetFileSecurity() InitializeSecurityDescriptor() SetSecurityDescriptorGroup()	Chmod() Umask() Chown()

open()

The **open()** system call allows you to access a file on a file system. It allocates resources to the file and provides a handle that the process may refer to. Many processes can open a file at once or by a single process only. It's all based on the file

system and structure.

read()

It is used to obtain data from a file on the file system. It accepts three arguments in general:

- A file descriptor.
- A buffer to store read data.
- The number of bytes to read from the file.

The file descriptor of the file to be read could be used to identify it and open it using **open()** before reading.

wait()

In some systems, a process may have to wait for another process to complete its execution before proceeding. When a parent process makes a child process, the parent process execution is suspended until the child process is finished. The **wait()** system call is used to suspend the parent process. Once the child process has completed its execution, control is returned to the parent process.

write()

It is used to write data from a user buffer to a device like a file. This system call is one way for a program to generate data. It takes three arguments in general:

- A file descriptor.
- A pointer to the buffer in which data is saved.
- The number of bytes to be written from the buffer.

fork()

Processes generate clones of themselves using the **fork()** system call. It is one of the most common ways to create processes in operating systems. When a parent process spawns a child process, execution of the parent process is interrupted until the child process completes. Once the child process has completed its execution, control is returned to the parent process.

close()

It is used to end file system access. When this system call is invoked, it signifies that the program no longer requires the file, and the buffers are flushed, the file information is altered, and the file resources are de-allocated as a result.

exec()

When an executable file replaces an earlier executable file in an already executing process, this system function is invoked. As a new process is not built, the old process identification stays, but the new process replaces data, stack, data, head, etc.

exit()

The **exit()** is a system call that is used to end program execution. This call indicates that the thread execution is complete, which is especially useful in multi-

OPERATING SYSTEM(23CS403)

threaded environments. The operating system reclaims resources spent by the process following the use of the **exit()** system function.

System components in OS:-

An operating system is a large and complex system that can only be created by partitioning into small parts. These pieces should be a well-defined part of the system, carefully defining inputs, outputs, and functions.

Although Windows, Mac, UNIX, Linux, and other OS do not have the same structure, most operating systems share similar OS system components, such as file, memory, process, I/O device management.

The components of an operating system play a key role to make a variety of computer system parts work together. There are the following components of an operating system, such as:

1. Process Management
2. File Management
3. Network Management
4. Main Memory Management
5. Secondary Storage Management
6. I/O Device Management
7. Security Management
8. Command Interpreter System

Operating system components help you get the correct computing by detecting CPU and memory hardware errors.



1. Process Management

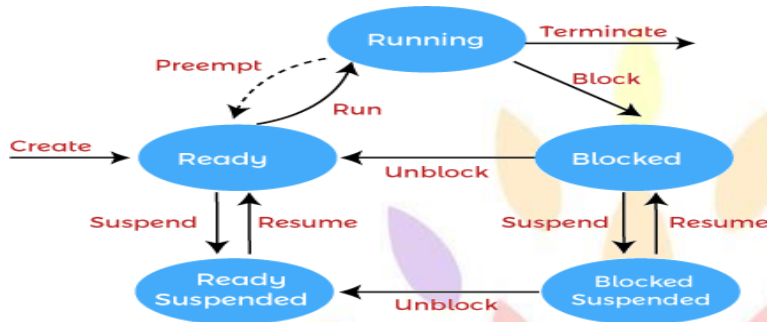
The process management component is a procedure for managing many processes running simultaneously on the operating system. Every running software application program has one or more processes associated with them.

For example, when you use a search engine like Chrome, there is a process running for that browser program.

OPERATING SYSTEM(23CS403)

Process management keeps processes running efficiently. It also uses memory allocated to them and shutting them down when needed.

The execution of a process must be sequential so, at least one instruction should be executed on behalf of the process.



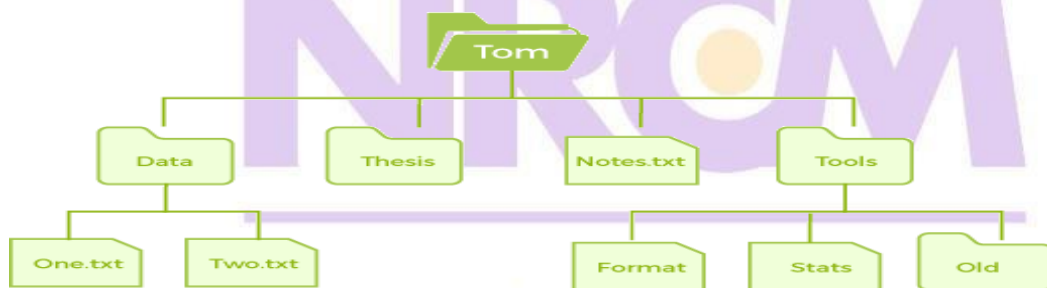
Functions of process management

Here are the following functions of process management in the operating system, such as:

- Process creation and deletion.
- Suspension and resumption.
- Synchronization process
- Communication process

2. File Management

A file is a set of related information defined by its creator. It commonly represents programs (both source and object forms) and data. Data files can be alphabetic, numeric, or alphanumeric.



Function of file management

The operating system has the following important activities in connection with file management:

- File and directory creation and deletion.
- For manipulating files and directories.
- Mapping files onto secondary storage.

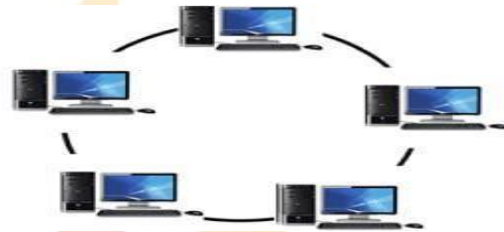
- Backup files on stable storage media.

3. Network Management

Network management is the process of administering and managing computer networks. It includes performance management, provisioning of networks, fault analysis, and maintaining the quality of service.

Computer Networks

When we hook up computers together using data communication facilities, we call this a computer network.



A distributed system is a collection of computers or processors that never share their memory and clock. In this type of system, all the processors have their local memory, and the processors communicate with each other using different communication cables, such as fibre optics or telephoned lines.

The computers in the network are connected through a communication network, which can configure in many different ways. The network can fully or partially connect in network management, which helps users design routing and connection strategies that overcome connection and security issues.

Functions of Network Management

Network management provides the following functions, such as:

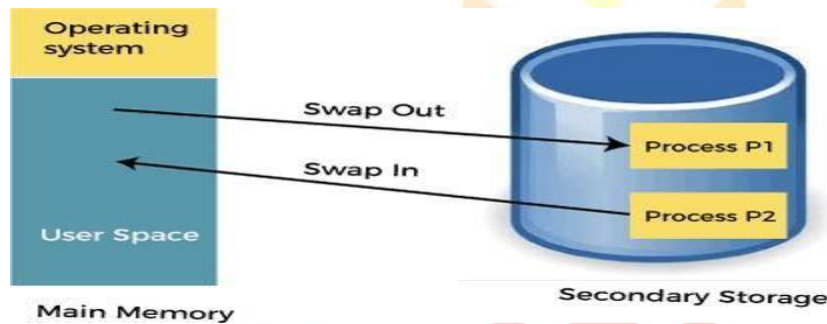
- Distributed systems help you to various computing resources in size and function. They may involve minicomputers, microprocessors, and many general-purpose computer systems.
- A distributed system also offers the user access to the various resources the network shares.
- It helps to access shared resources that help computation to speed up or offers data availability and reliability.

4. Main Memory Management

Main memory is a large array of storage or bytes, which has an address. The memory management process is conducted by using a sequence of reads or writes of specific memory addresses.

It should be mapped to absolute addresses and loaded inside the memory to execute a program. The selection of a memory management method depends on several factors.

However, it is mainly based on the hardware design of the system. Each algorithm requires corresponding hardware support. Main memory offers fast storage that can be accessed directly by the CPU. It is costly and hence has a lower storage capacity. However, for a program to be executed, it must be in the main memory.



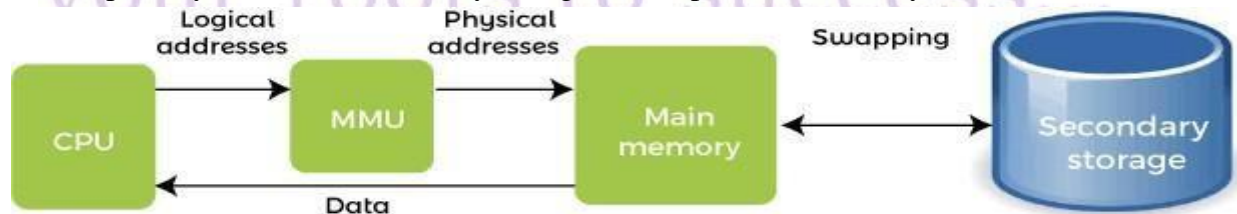
Functions of Memory management

An Operating System performs the following functions for Memory Management in the operating system:

- It helps you to keep track of primary memory.
- Determine what part of it are in use by whom, what part is not in use.
- In a multi-programming system, the OS decides which process will get memory and how much.
- Allocate the memory when a process requests.
- It also de-allocates the memory when a process no longer requires or has been terminated.

5. Secondary-Storage Management

The most important task of a computer system is to execute programs. These programs help you to access the data from the main memory during execution. This memory of the computer is very small to store all data and programs permanently. The computer system offers secondary storage to backup the main memory.



OPERATING SYSTEM (23CS403)

Today modern computers use hard drives/SSD as the primary storage of both programs and data. However, the secondary storage management also works with storage devices, such as USB flash drives and CD/DVD drives. Programs like assemblers and compilers are stored on the disk until they are loaded into memory, and then the disk is used as a source and destination for processing.

Functions of Secondary Storage Management

Here are some major functions of secondary storage management in the operating system:

- Storage allocation
- Free space management
- Disk scheduling
-

6. I/O Device Management

One of the important uses of an operating system that helps to hide the variations of specific hardware devices from the user.



Functions of I/O Management

The I/O management system offers the following functions, such as:

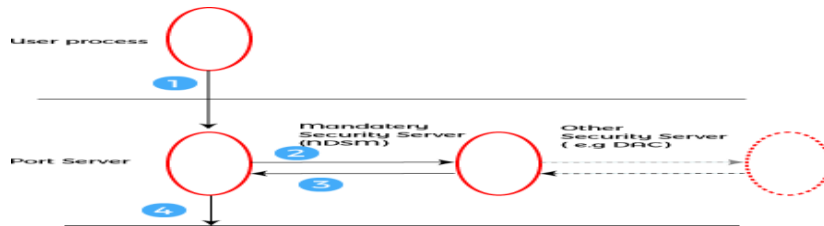
- It offers a buffer caching system
- It provides general device driver code
- It provides drivers for particular hardware devices.
- I/O helps you to know the individualities of a specific device.

7. Security Management

The various processes in an operating system need to be secured from other activities. Therefore, various mechanisms can ensure those processes that want to operate files, memory CPU, and other hardware resources should have proper authorization from the operating system.

Security refers to a mechanism for controlling the access of programs, processes, or

users to the resources defined by computer control to be imposed, together with some means of enforcement.

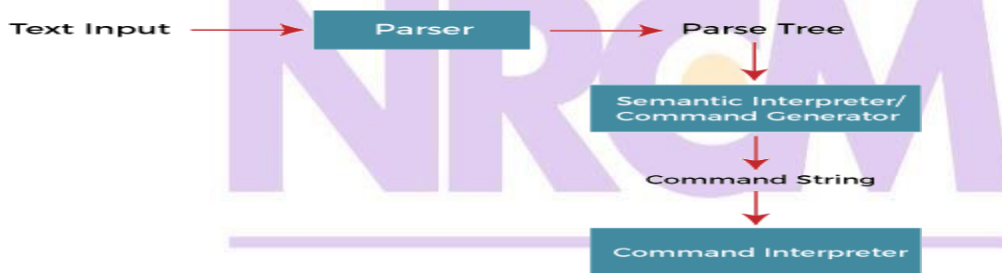


For example, memory addressing hardware helps to confirm that a process can be executed within its own address space. The time ensures that no process has control of the CPU without renouncing it. Lastly, no process is allowed to do its own I/O to protect, which helps you to keep the integrity of the various peripheral devices.

Security can improve reliability by detecting latent errors at the interfaces between components/subsystems. Early detection of interface errors can prevent the foulness of a healthy subsystem by a malfunctioning subsystem. An unprotected resource cannot be misused by an unauthorized or incompetent user.

9. Command Interpreter System

One of the most important components of an operating system is its command interpreter. The command interpreter is the primary interface between the user and the rest of the system.



Many commands are given to the operating system by control statements. A program that reads and interprets control statements is automatically executed when a new job is started in a batch system or a user logs into a time-shared system. This program is variously called.

- The control card interpreter,
- The command-line interpreter,
- The shell (in UNIX), and soon.

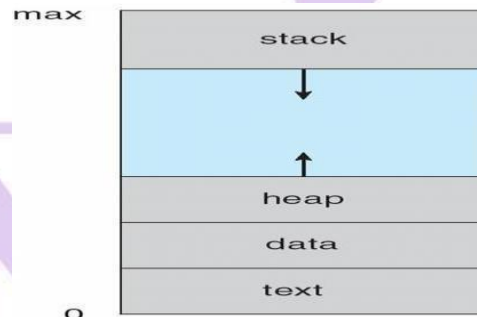
Its function is quite simple, get the next command statement, and execute it. The command statements deal with process management, I/O handling, secondary storage management, main memory management, file system access, protection, and networking.

PROCESS: A process can be thought of as a program in execution. A process is the unit of work in most systems.

A process will need certain resources—such as CPU time, memory, files, and I/O devices to accomplish its task. These resources are allocated to the process either when it is created or while it is executing.

Structure of a Process in Memory

- A process is more than the program code, which is sometimes known as the **text section**.
- It also includes the current activity, as represented by the value of the **program counter** and the contents of the processor's registers.
- A process generally also includes the process **stack**, which contains temporary data (such as function parameters, return addresses, and local variables).
- A **data section**, which contains global variables.
- A process may also include a **heap**, which is memory that is dynamically allocated during process run time.



When a Program becomes Process?

A program is a *passive* entity, such as a file containing a list of instructions stored on disk (Often called as **executable file**). In contrast, a process is an *active* entity, with a program counter specifying the next instruction to execute and a set of associated resources. A program becomes a process when an executable file is loaded into memory.

Two common techniques for loading executable files are double-clicking an icon representing the executable file and entering the name of the executable file on the command line (as in prog.exe or a.out).

If two processes are associated with the same program, are they the same or different? (Or) Explain if you run same program twice, what section would be shared in memory?

OPERATINGSYSTEM(23CS403)

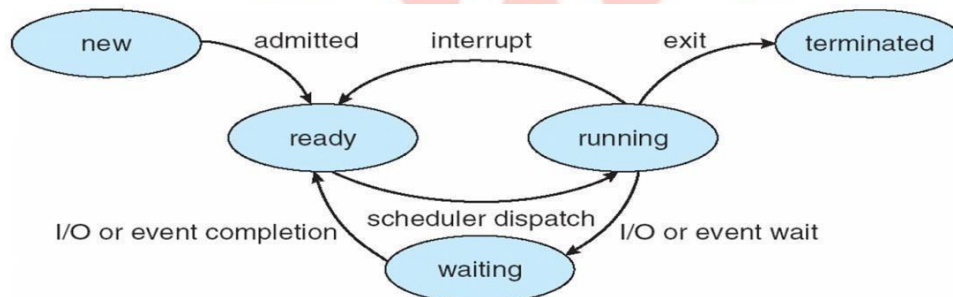
Although two processes may be associated with the same program, they are nevertheless considered two separate execution sequences. For instance, several users may be running different copies of the mail program, or the same user may invoke many copies of the web browser program. Each of these is a separate process; and although the text sections are equivalent, the data, heap, and stack sections vary. It is also common to have a process that spawns many processes as it runs.

1. Process State

As a process executes, it changes **state**. The state of a process is defined in part by the current activity of that process.

A process may be in one of the following states:

- **New:** The process is being created.
- **Running:** Instructions are being executed.
- **Waiting:** The process is waiting for some event to occur (such as an I/O completion or reception of a signal).
- **Ready:** The process is waiting to be assigned to a processor.
- **Terminated:** The process has finished execution.



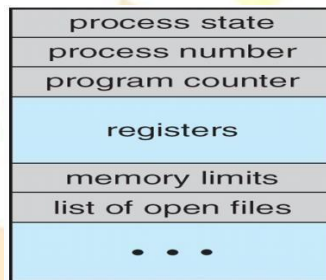
2. Process Control Block

Each process is represented in the operating system by a **Process Control Block (PCB)** or **Task Control Block**. It contains many pieces of information associated with a specific process, including these:

- **Process state:** The state may be new, ready, running, and waiting, halted, and soon.
- **Program counter.** The counter indicates the address of the next instruction to be executed for this process.
- **CPU registers.** The registers vary in number and type, depending on the computer architecture. They include accumulators, index registers, stack pointers, and general-purpose registers, plus any condition-code information. Along with the program counter, this state information must be saved when an interrupt occurs, to allow the process to be continued correctly afterward.
- **CPU-scheduling information.** This information includes a process priority, pointers to scheduling queues, and any other scheduling parameters.

OPERATINGSYSTEM(23CS403)

- **Memory-management information.** This information may include such items as the value of the base and limit registers and the page tables, or the segment tables, depending on the memory system used by the operating system.
- **Accounting information.** This information includes the amount of CPU and real time used, time limits, account numbers, job or process numbers, and so on.
- **I/O status information.** This information includes the list of I/O devices allocated to the process, a list of open files, and so on.



Process Scheduling

The objective of multiprogramming is to have some process running at all times, to maximize CPU utilization.

The objective of time sharing is to switch the CPU among processes so frequently that users can interact with each program while it is running.

To meet these objectives, the **process scheduler** selects an available process (possibly from a set of several available processes) for program execution on the CPU.

1. Scheduling Queues

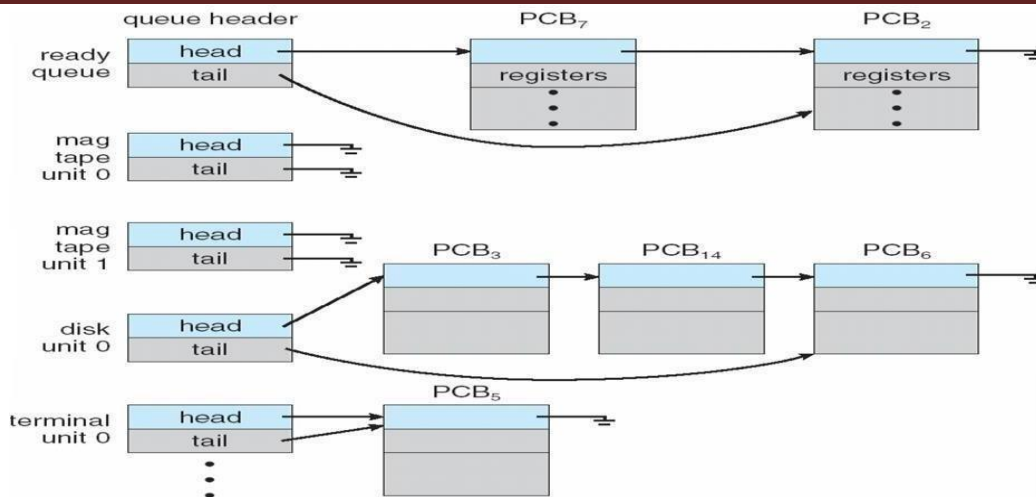
The following are the different queues available,

a. Job Queue

- As processes enter the system, they are put into a **job queue**, which consists of all processes in the system.

b. Ready Queue

- The processes that are residing in main memory and are ready and waiting to execute are kept on a list called the **ready queue**.
- This queue is generally stored as a linked list. A ready-queue header contains pointers to the first and final PCBs in the list. Each PCB includes a pointer field that points to the next PCB in the ready queue.



c. Device Queue

- The list of processes waiting for a particular I/O device is called a **device queue**.
- Each device has its own device queue.

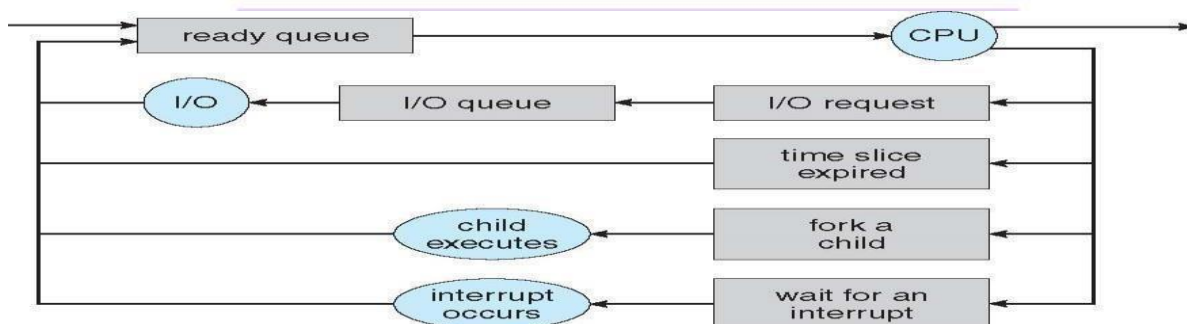
Queuing-diagram representation of process scheduling

A common representation of process scheduling is a **queuing diagram**. Each rectangular box represents a queue. Two types of queues are present: the ready queue and a set of device queues. The circles represent the resources that serve the queues, and the arrows indicate the flow of processes in the system.

A new process is initially put in the ready queue. It waits there until it is selected for execution, or **dispatched**. Once the process is allocated the CPU and is executing, one of several events could occur:

- The process could issue an I/O request and then be placed in an I/O queue.
- The process could create a new child process and wait for the child's termination.
- The process could be removed forcibly from the CPU, as a result of an interrupt, and be put back in the ready queue.

In the first two cases, the process eventually switches from the waiting state to the ready state and is then put back in the ready queue. A process continues this cycle until it terminates, at which time it is removed from all queues and has its PCB and resources deallocated.



2. Schedulers

Definition: A process migrates among the various scheduling queues throughout its lifetime. The operating system must select, for scheduling purposes, processes from these queues in some fashion. The selection process is carried out by the appropriate **scheduler**.

Types of Schedulers

a. Long-Term Scheduler or Job Scheduler

- Often, in a batch system, more processes are submitted than can be executed immediately. These processes are spooled to a mass-storage device (typically a disk), where they are kept for later execution.
- The **long-term scheduler**, or **job scheduler**, selects processes from this pool and loads them into memory for execution.
- The long-term scheduler executes much less frequently; minutes may separate the creation of one new process and the next.
- The long-term scheduler controls the **degree of multiprogramming** (the number of processes in memory).
- If the degree of multiprogramming is stable, then the average rate of process creation must be equal to the average departure rate of processes leaving the system. Thus, the long-term scheduler may need to be invoked only when a process leaves the system.
- Because of the longer interval between executions, the long-term scheduler can afford to take more time to decide which process should be selected for execution. It is important that the long-term scheduler select a good *process mix* of I/O-bound and CPU-bound processes.
- On some systems, the long-term scheduler may be absent or minimal.

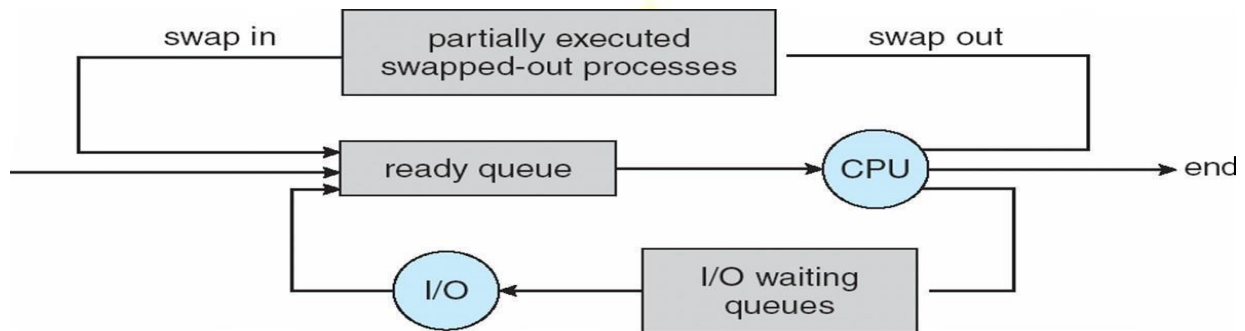
b. Short-Term Scheduler, Or CPU Scheduler

- The **short-term scheduler**, or **CPU scheduler**, selects from among the processes that are ready to execute and allocates the CPU to one of them.
- The short-term scheduler must select a new process for the CPU frequently.
- A process may execute for only a few milliseconds before waiting for an I/O request. Often, the short-term scheduler executes at least once every 100 milliseconds.
- Because of the short time between executions, the short-term scheduler must be fast.

c. Medium-Term Scheduler

- Some operating systems, such as time-sharing systems, may introduce an additional, intermediate level of scheduling.
- The key idea behind a medium-term scheduler is that sometimes it can be advantageous to remove a process from memory (and from active contention for the CPU) and thus reduce the degree of multiprogramming.
- Later, the process can be reintroduced into memory, and its execution can be continued where it left off. This scheme is called **swapping**.

- The process is swapped out, and is later swapped in, by the medium-term scheduler. Swapping may be necessary to improve the process mix or because a change in memory requirements has overcommitted available memory, requiring memory to be freed up.



3. Context Switch

Definition: Switching the CPU to another process requires performing a state save of the current process and a state restore of a different process. This task is known as a **context switch**.

When a context switch occurs, the kernel saves the context of the old process in its PCB and loads the saved context of the new process scheduled to run.

Overhead: Context-switch time is pure overhead, because the system does no useful work while switching.

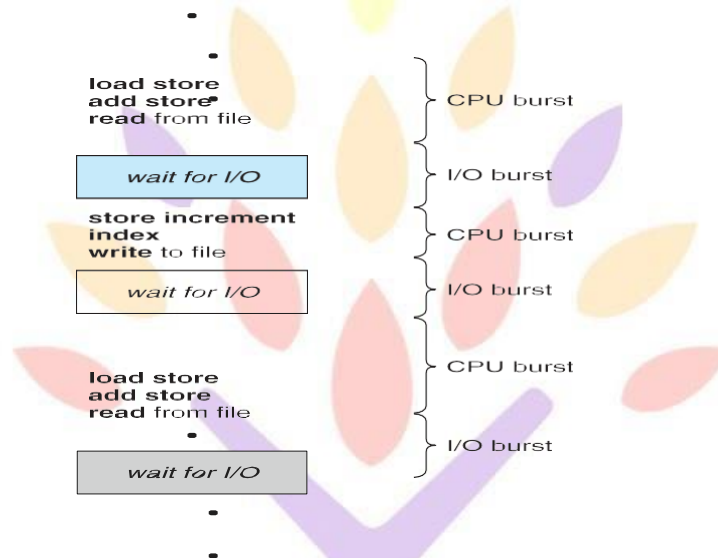
Switching Speed: Switching speed varies from machine to machine, depending on the memory speed, the number of registers that must be copied, and the existence of special instructions (such as a single instruction to load or store all registers). A typical speed is a few milliseconds.

Hardware Support: Context-switch times are highly dependent on hardware support. A context switch here simply requires changing the pointer to the current register set. Of course, if there are more active processes than there are register sets, the system resorts to copying register data to and from memory, as before. Also, the more complex the operating system, the greater the amount of work that must be done during a context switch.

your roots to success...

4. CPU-I/O Burst Cycle

The success of CPU scheduling depends on an observed property of processes: process execution consists of a **cycle** of CPU execution and I/O wait. Processes alternate between these two states. Process execution begins with a **CPU burst**. That is followed by an **I/O burst**, which is followed by another CPU burst, then another I/O burst, and so on. Eventually, the final CPU burst ends with a system request to terminate execution.



Definition of Non Preemptive Scheduling

Under nonpreemptive scheduling, once the CPU has been allocated to a process, the process keeps the CPU until it releases the CPU either by terminating or by switching to the waiting state. This scheduling method was used by Microsoft Windows 3.x.

Definition of Preemptive Scheduling

Under this, a running process may be replaced by higher priority process at any time. Used from Windows 95 to till now. Incurs the cost associated with access to shared data. It also affects the design of OS.

Dispatcher

Another component involved in the CPU-scheduling function is the **dispatcher**. The dispatcher is the module that gives control of the CPU to the process selected by the short-term scheduler. This function involves the following:

- Switching context
- Switching to user mode
- Jumping to the proper location in the user program to restart that program

The dispatcher should be as fast as possible, since it is invoked during every process switch.

OPERATINGSYSTEM(23CS403)

Dispatch Latency: The time it takes for the dispatcher to stop one process and start another running is known as the **dispatch latency**.

Operations on processes (OR) System call interface for process management - fork, exit, wait, waitpid, exec

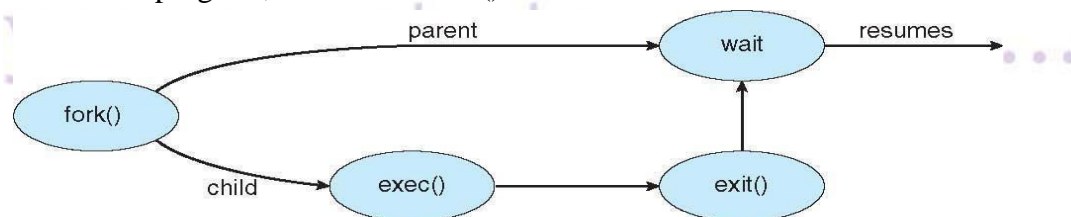
The processes in most systems can execute concurrently, and they may be created and deleted dynamically. Thus, these systems must provide a mechanism for process creation and termination.

a. Process Creation

During the course of execution, a process may create several new processes. The creating process is called a parent process, and the new processes are called the children of that process. Each of these new processes may in turn create other processes, forming a **tree** of processes.

System Calls

- **fork()**
 - Most operating systems (including UNIX, Linux, and Windows) identify processes according to a unique **process identifier** (or **pid**), which is typically an integer number.
 - A new process is created by the `fork()` system call. The new process consists of a copy of the address space of the original process.
 - This mechanism allows the parent process to communicate easily with its child process. Both processes (the parent and the child) continue execution at the instruction after the `fork()`, with one difference: the return code for the `fork()` is zero for the new (child) process, whereas the (nonzero) process identifier of the child is returned to the parent.
- **exec()**
 - After a `fork()` system call, one of the two processes typically uses the `exec()` system call to replace the process's memory space with a new program.
 - The `exec()` system call loads a binary file into memory and starts its execution. In this manner, the two processes are able to communicate and then go their separate ways.
- **wait()**
 - The parent can then create more children; or, if it has nothing else to do while the child runs, it can issue a `wait()` system call to move itself off the ready queue until the termination of the child. Because the call to `exec()` overlays the process's address space with a new program, the call to `exec()` does not return control unless an error occurs.



b. Process Termination

A process terminates when it finishes executing its final statement and asks the operating system to delete it by using the `exit ()` system call. At that point, the process may return a status value (typically an integer) to its parent process (via the `wait()` system call). All the resources of the process—including physical and virtual memory, open files and I/O buffers—are deallocated by the operating system.

Termination can occur in other circumstances as well. A process can cause the termination of another process via an appropriate system call (for example, `TerminateProcess()` in Windows). Usually, such a system call can be invoked only by the parent of the process that is to be terminated. Otherwise, users could arbitrarily kill each other's jobs.

COOPERATING PROCESSES

- The concurrent processes executing in the OS may be either independent process or cooperating process.
- Independent process **cannot** affect or be affected by the execution of another process.
- Cooperating process **can** affect or be affected by the execution of another process.

Advantages of process cooperation

1. **Information sharing:** several users may be interested in the same piece of information.
2. **Computations speed-up:** If we want a particular task to run faster, we must break it into subtasks and run in parallel.
3. **Modularity:** Constructing the system in modular fashion, dividing the system functions into separate process.
4. **Convenience:** User will have many tasks to work in parallel (Editing, compiling, printing).

Processes can communicate with each other through both:

- Shared Memory
- Message passing

The following figure shows a basic structure of communication between processes via the shared memory method and via the message passing method.



Figure 1 - Shared Memory and Message Passing

(i) SharedMemory

Communication between processes using shared memory requires processes to share some variable, and it completely depends on how the programmer will implement it.

One way of communication using shared memory can be imagined like this: Suppose process1 and process2 are executing simultaneously, and they share some resources or use some information from another process. Process1 generates information about certain computations or resources being used and keeps it as a record in shared memory. When process2 needs to use the shared information, it will check in the record stored in shared memory and take note of the information generated by process1 and act accordingly.

Processes can use shared memory for extracting information as a record from another process as well as for delivering any specific information to other processes.

Ex: Producer-Consumer problem

A producer process produces information that is consumed by a consumer process. For example, a print program produces characters that are consumed by the printer driver.

A producer can produce one item while the consumer is consuming another item. The Producer and Consumer must be synchronized. The consumer does not try to consume an item, the consumer must wait until an item is produced.

Unbounded-Buffer

- no practical limit on the size of the buffer.
- Producer can produce any number of items.
- Consumer may have to wait

Bounded-Buffer

- assume that there is a fixed buffer size.

Bounded-Buffer–Shared-Memory Solution:

Shared data

```
#define BUFFER_SIZE 10
typedef struct
{
  ...
} item;
item buffer[BUFFER_SIZE];
int in = 0;
int out = 0;
```

Bounded-Buffer–Producer Process:

```
itemnextProduced;
while (1)
{
while(((in+1)% BUFFER_SIZE)==out);          /*donothing
*/buffer[in]= nextProduced;
in=(in+1)%BUFFER_SIZE;
}
```

Bounded-Buffer-ConsumerProcess:

```
itemnextConsumed;
while (1)
{
while(in==out);          /*donothing
*/next Consumed = buffer[out];
out=(out+1)%BUFFER_SIZE;
}
```

(ii) MessagingPassingMethod

In this method, processes communicate with each other without using any kind of shared memory. If two processes want to communicate with each other, they proceed as follows



- Establish a communication link (if a link already exists, no need to establish it again.)
- Start exchanging messages using basic primitives.
- The message size can be of fixed size or of variable size. If it is of fixed size, it is easy for an OS designer but complicated for a programmer and if it is of variable size then it is easy for a programmer but complicated for the OS designer.
- Cooperating process to communicate with each other via an inter process communication (IPC).
- IPC provides a Mechanism to allow processes to communicate and to synchronize their actions.
- If two processes want to communicate, a communication link exists between them and

OPERATING SYSTEM(23CS403)

exchange messages via send/receive. OS provides this facility.

- IPC facility provides two operations:
 - **Send** (*message*) – message size fixed or variable.
 - **Receive** (*message*)

- Implementation of communication link by following.
 - physical (e.g., shared memory, hardware bus)
 - logical (e.g., logical properties)

Methods for logical implementation of a link

- i. Direct communication.
- ii. Indirect communication.

Direct Communication

- Each process must name each other explicitly:
 - **Send** (*P, message*) – send a message to process P.
 - **Receive** (*Q, message*) – receive a message from process Q.
- Links are established automatically.
- A link is associated with exactly one pair of communicating processes.
- Between each pair there exists exactly one link.
- The link may be unidirectional, but is usually bi-directional.
- This exhibits both symmetry and asymmetry in addressing

Symmetry:

Both the sender and the receiver processes must name the other to communicate.

Asymmetry:

Only the sender names the recipient, the recipient is not required to name the sender. The send and receive primitives are as follows.

- **Send** (*P, message*) – send a message to process P.
- **Receive** (*id, message*) – receive a message from any process.

Disadvantage of direct communication

Changing a name of the process creates problems.

Indirect Communication

- The messages are sent and received from mailboxes (also referred to as ports).
- A mailbox is an object
- Process can place messages.
- Process can remove messages.

OPERATING SYSTEM (23CS403)

- Two processes can communicate only if they have a shared mailbox.
- Primitives are defined as:
 - **send**(A, message) – send a message to mailbox A
 - **receive**(A, message) – receive a message from mailbox A.
- A mailbox may be owned either by a process or by the OS.
- If the mailbox is owned by a process, then we distinguish b/w the owner (who can only receive msg through this mailbox) and the user (who can only send msg to the mailbox).
- A mailbox may be owned by the OS independently and provide a mechanism,
 - create a mailbox
 - receive messages through mailbox
 - destroy a mailbox.

Mailbox sharing problem

The processes $P_1, P_2,$ and P_3 all share mailbox A. Processes P_1 sends; P_2 and P_3 receive the message from A. Who gets a message?

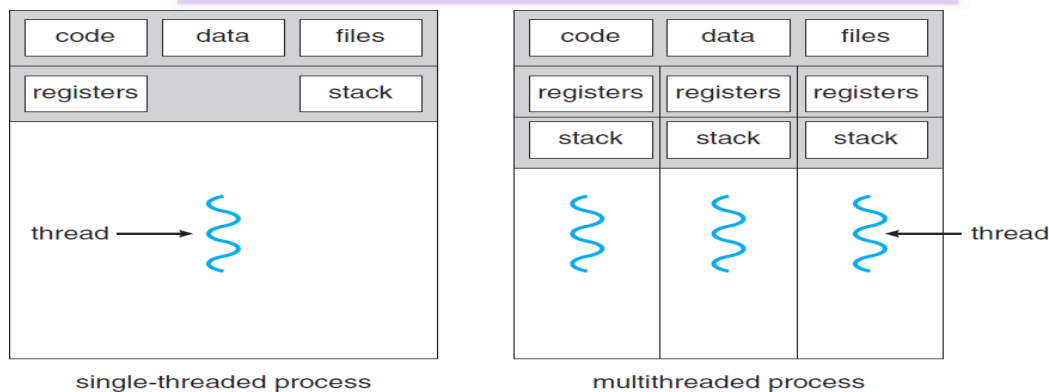
Solutions

- Allow a link to be associated with at most two processes.
- Allow only one process at a time to execute a receive operation.
- Allow the system to select arbitrarily the receiver. The system may identify the receiver to the sender.

Defining Thread

A thread is a lightweight process. A thread is a flow of control execution of the program. A thread is a basic unit of CPU utilization; it comprises a thread ID, a program counter, a register set, and a stack. It shares with other threads belonging to the same process its code section, data section, and other operating-system resources, such as open files and signals.

A traditional (or **heavyweight**) process has a single thread of control. If a process has multiple threads of control, it can perform more than one task at a time.



SingleThread

- A process is a program that performs a single thread of execution.
 - For example, when a process is running a word-processor program, a single thread of instructions is being executed.
 - This single thread of control allows the process to perform only one task at a time. The user cannot simultaneously type in characters and run the spell checker within the same process, for example.
- ### Multi Thread
- Most modern operating systems have extended the process concept to allow a process to have multiple threads of execution and thus to perform more than one task at a time.
 - This feature is especially beneficial on multicore systems, where multiple threads can run in parallel.
 - On a system that supports threads, the PCB is expanded to include information for each thread. Other changes throughout the system are also needed to support threads.

Multithreading Models

Support for threads may be provided either at the user level, for **user threads**, or by the kernel, for **kernel threads**. User threads are supported above the kernel and are managed without kernel support, whereas kernel threads are supported and managed directly by the operating system. Virtually all contemporary operating systems—including Windows, Linux, Mac OS X, and Solaris support kernel threads.

Ultimately, a relationship must exist between user threads and kernel threads. The following are the three common ways of establishing such a relationship: the many-to-one model, the one-to-one model, and the many-to-many models.

1. Many-to-One Model

- The many-to-one model maps many user-level threads to one kernel thread.

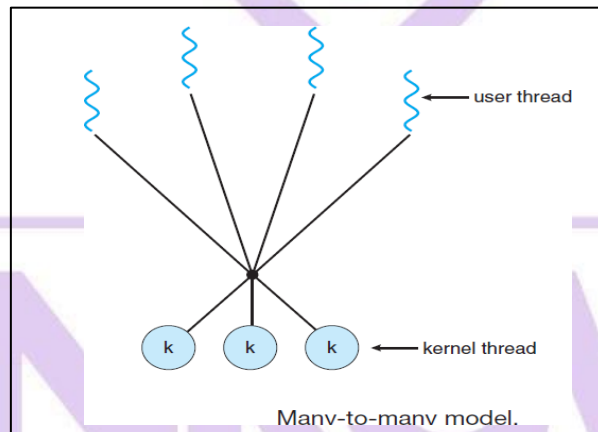
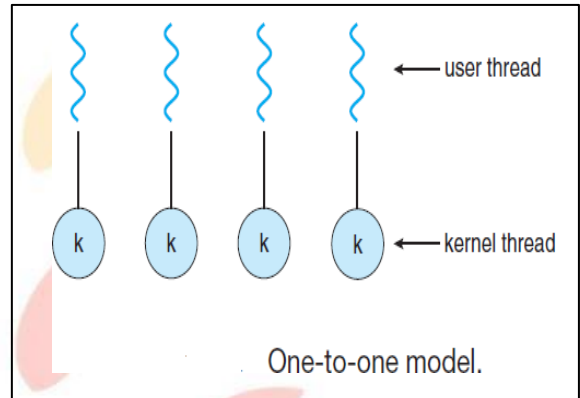
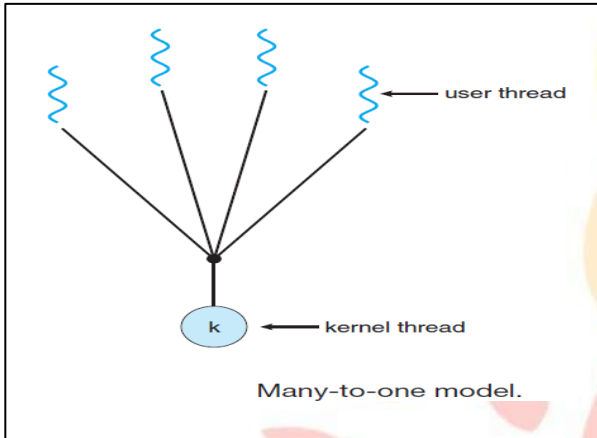
2. One-to-One Model

- The one-to-one model maps each user thread to a kernel thread.

3. Many-to-Many Model

- It multiplexes many user-level threads to a smaller or equal number of kernel threads.

your roots to success...



your roots to success...

CPU scheduling- Scheduling criteria, Scheduling Algorithms, Multiple Processor Scheduling, System call interface for process management-fork, exit, wait, waitpid, exec.

Deadlocks- system Model, Deadlock Characterization, Methods for handling deadlocks, Deadlock Prevention, Deadlock Avoidance, deadlock Detection, and recovery from deadlock

CPU scheduling

CPU scheduling is the process of deciding which process will own the CPU to use while another process is suspended. The main function of the CPU scheduling is to ensure that whenever the CPU remains idle, the OS has at least selected one of the processes available in the ready-to-use line.

In Multiprogramming, if the long-term scheduler selects multiple I / O binding processes then most of the time, the CPU remains an idle. The function of an effective program is to improve resource utilization.

If most operating systems change their status from performance to waiting then there may always be a chance of failure in the system. So in order to minimize this excess, the OS needs to schedule tasks in order to make full use of the CPU and avoid the possibility of deadlock.

Objectives of Process Scheduling Algorithm

- Utilization of CPU at maximum level. Keep CPU as busy as possible.
- Allocation of CPU should be fair.
- Throughput should be Maximum. i.e. Number of processes that complete their execution per time unit should be maximized.
- Minimum turnaround time, i.e. time taken by a process to finish execution should be the least.
- There should be a minimum waiting time and the process should not starve in the ready queue.
- Minimum response time. It means that the time when a process produces the first response should be as less as possible.

Terminologies

- **Arrival Time:** Time at which the process arrives in the ready queue.
- **Completion Time:** Time at which process completes its execution.
- **Burst Time:** Time required by a process for CPU execution.
- **Turn Around Time:** Time Difference between completion time and arrival time.
 $Turn\ Around\ Time = Completion\ Time - Arrival\ Time$
- **Waiting Time (W.T):** Time Difference between turnaround time and burst

time.

$$\text{Waiting Time} = \text{Turn Around Time} - \text{Burst Time}$$

THE SCHEDULING CRITERIA

CPU Utilization:

The main purpose of any CPU algorithm is to keep the CPU as busy as possible. Theoretically, CPU usage can range from 0 to 100 but in a real-time system, it varies from 40 to 90 percent depending on the system load.

Throughput:

The average CPU performance is the number of processes performed and completed during each unit. This is called throughput. The output may vary depending on the length or duration of the processes.

Turnaround Time:

For a particular process, the important conditions are how long it takes to perform that process. The time elapsed from the time of process delivery to the time of completion is known as the conversion time. Conversion time is the amount of time spent waiting for memory access, waiting in line, using CPU, and waiting for I/O.

Waiting Time:

The Scheduling algorithm does not affect the time required to complete the process once it has started performing. It only affects the waiting time of the process i.e. the time spent in the waiting process in the ready queue.

Response Time:

In a collaborative system, turn around time is not the best option. The process may produce something early and continue to compute the new results while the previous results are released to the user. Therefore another method is the time taken in the submission of the application process until the first response is issued. This measure is called response time.

Types of CPU Scheduling Algorithms

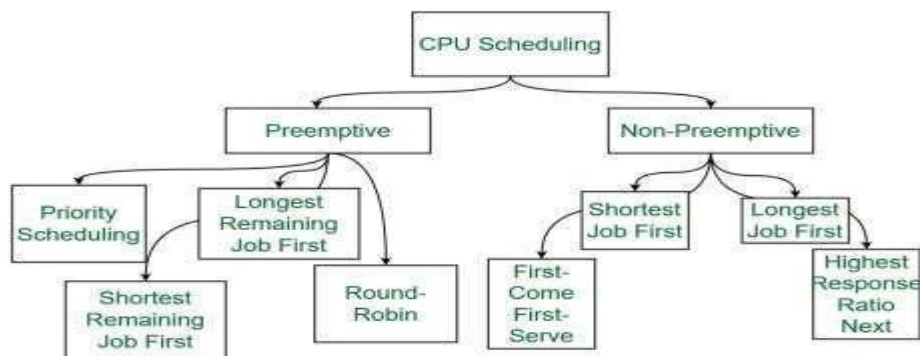
There are mainly two types of scheduling methods:

Preemptive Scheduling:

Preemptive scheduling is used when a process switches from running state to ready state or from the waiting state to the ready state.

Non-Preemptive Scheduling:

Non-Preemptive scheduling is used when a process terminates, or when a process switches from running state to waiting state.



1. FirstComeFirstServeScheduling:

FCFS considered to be the simplest of all operating system scheduling algorithms. First come first serve scheduling algorithm states that the process that requests the CPU first is allocated the CPU first and is implemented by using FIFO queue.

Characteristics:

- FCFS supports non-preemptive and preemptive CPU scheduling algorithms.
- Tasks are always executed on a First-come, First-serve concept.
- FCFS is easy to implement and use.
- This algorithm is not much efficient in performance, and the wait time is quite high.

Advantages:

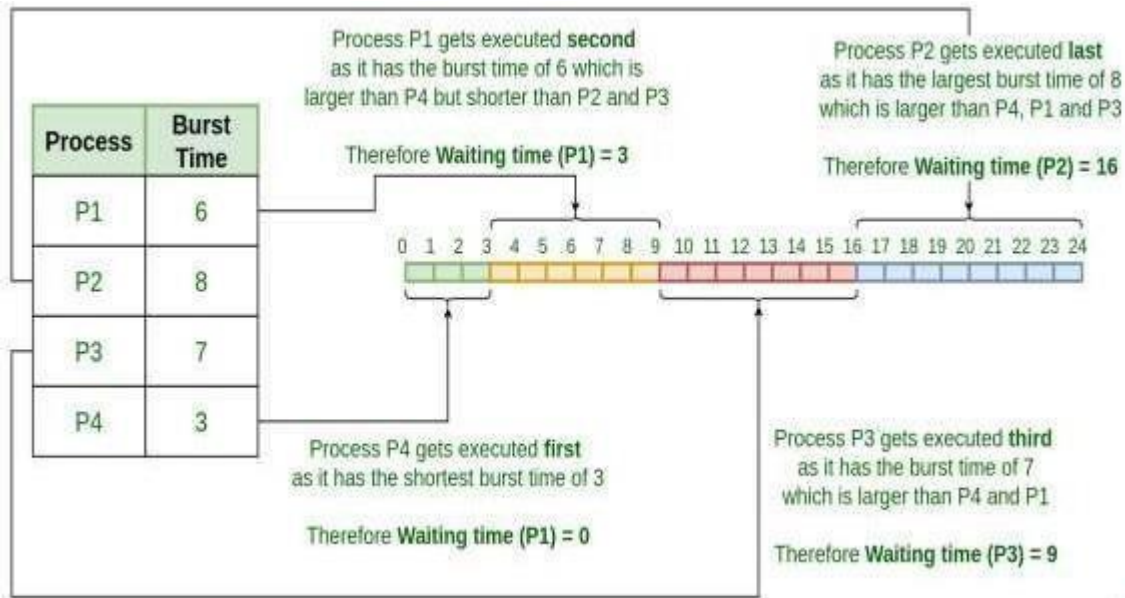
- Easy to implement
- First come, first serve method

Disadvantages:

- FCFS suffers from **Convoy effect**.
- The average waiting time is much higher than the other algorithms.
- FCFS is very simple and easy to implement and hence not much efficient.

2. ShortestJobFirst (SJF)Scheduling:

Shortest job first (SJF) is a scheduling process that selects the waiting process with the smallest execution time to execute next. This scheduling method may or may not be preemptive. Significantly reduces the average waiting time for other processes waiting to be executed. The full form of SJF is Shortest Job First.



Characteristics:

- Shortest Job first has the advantage of having a minimum average waiting time among all operating system scheduling algorithms.
- It is associated with each task as a unit of time to complete.
- It may cause starvation if shorter processes keep coming. This problem can be solved using the concept of ageing.

Advantages:

- As SJF reduces the average waiting time, thus, it is better than the first come first serve scheduling algorithm.
- SJF is generally used for long term scheduling.

Disadvantages:

- One of the demerits of SJF is starvation.
- Many times it becomes complicated to predict the length of the upcoming CPU request.

3. Longest Job First (LJF) Scheduling:

This is just opposite of shortest job first (SJF), as the name suggests this algorithm is based upon the fact that the process with the largest burst time is processed first. Longest Job First is non-preemptive in nature.

Characteristics:

- Among all the processes waiting in a waiting queue, CPU is always assigned to the process having largest burst time.
- If two processes have the same burst time then tie is broken using FCFS, i.e. the process that arrived first is processed first.

- LJFCPU Scheduling can be of both preemptive and non-preemptive types.

Advantages:

- No other task can schedule until the longest job or process executes completely.
- All the jobs or processes finish at the same time approximately.

Disadvantages:

- Generally, the LJF algorithm gives a very high average waiting time and average turn-around time for a given set of processes.
- This may lead to convoy effect.

4. Priority Scheduling:

Preemptive Priority CPU Scheduling Algorithm is a pre-emptive method of CPU scheduling algorithm that works **based on the priority** of a process. In this algorithm, the editor sets the functions to be as important, meaning that the most important process must be done first. In the case of any conflict, that is, where there are more than one processor with equal value, then the most important CPU planning algorithm works on the basis of the FCFS **Characteristics:**

- Schedule tasks based on priority.
- When the higher priority work arrives while a task with less priority is executed, the higher priority work takes the place of the less priority one and
- The latter is suspended until the execution is complete.
- Lower is the number assigned, higher is the priority level of a process.

Advantages:

- The average waiting time is less than FCFS
- Less complex

Disadvantages:

- One of the most common demerits of the Preemptive priority CPU scheduling algorithm is the Starvation Problem. This is the problem in which a process has to wait for a longer amount of time to get scheduled into the CPU. This condition is called the starvation problem.

5. Round Robin Scheduling:

Round Robin is a CPU scheduling algorithm where each process is cyclically assigned a fixed time slot. It is the preemptive version of First come First Serve CPU Scheduling algorithm. Round Robin CPU Algorithm generally focuses on Time Sharing technique.

Characteristics:

- It's simple, easy to use, and starvation-free as all processes get the balanced CPU allocation.
- One of the most widely used methods in CPU scheduling as a core.

- It is considered preemptive as the processes are given to the CPU for a very limited time.

Advantages:

- Round robin seems to be fair as every process gets an equal share of CPU.
- The newly created process is added to the end of the ready queue.

6. Shortest Remaining Time First Scheduling (SRTF):

SRTF is the preemptive version of the Shortest job first which we have discussed earlier where the processor is allocated to the job closest to completion. In SRTF the process with the smallest amount of time remaining until completion is selected to execute.

Characteristics:

- SRTF algorithm makes the processing of the jobs faster than SJF algorithm, given its overhead charges are not counted.
- The context switch is done a lot more times in SRTF than in SJF and consumes the CPU's valuable time for processing. This adds up to its processing time and diminishes its advantage of fast processing.

Advantages:

- In SRTF the short processes are handled very fast.
- The system also requires very little overhead since it only makes a decision when a process completes or a new process is added.

Disadvantages:

- Like the shortest job first, it also has the potential for process starvation.
- Long processes may be held off indefinitely if short processes are continually added.

7. Longest Remaining Time First:

The longest remaining time first is a preemptive version of the longest job first scheduling algorithm. This scheduling algorithm is used by the operating system to program incoming processes for use in a systematic way. This algorithm schedules those processes first which have the longest processing time remaining for completion.

Characteristics:

- Among all the processes waiting in a waiting queue, the CPU is always assigned to the process having the largest burst time.
- If two processes have the same burst time then the tie is broken using FCFS i.e. the process that arrived first is processed first.
- LJFCPU scheduling can be of both preemptive and non-preemptive types.

Advantages:

- No other process can execute until the longest task executes completely.

- All the jobs or processes finish at the same time approximately.

Disadvantages:

- This algorithm gives a very high average waiting time and average turn-around time for a given set of processes.
- This may lead to a convoy effect.

8. Highest Response Ratio Next:

Highest Response Ratio Next is a non-preemptive CPU Scheduling algorithm and it is considered as one of the most optimal scheduling algorithms. The name itself states that we need to find the response ratio of all available processes and select the one with the highest Response Ratio. A process once selected will run till completion.

Characteristics:

- The criteria for HRRN is **Response Ratio** and the **mode** is **NonPreemptive**.
- HRRN is considered as the modification of Shortest Job First to reduce the problem of starvation.
- In comparison with SJF, during the HRRN scheduling algorithm, the CPU is allotted to the next process which has the **highest response ratio** and not to the process having less burst time.

$$\text{Response Ratio} = (W + S) / S$$

Here, **W** - Waiting time of the process
S - Burst time of the process.

Advantages:

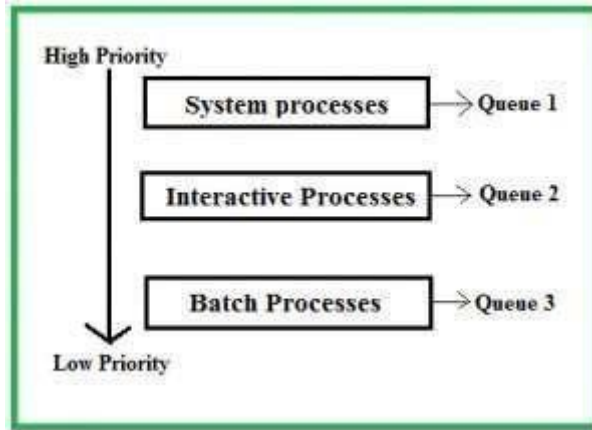
- HRRN Scheduling algorithm generally gives better performance than the shortest job first Scheduling.
- There is a reduction in waiting time for longer jobs and also it encourages shorter jobs.

Disadvantages:

- The implementation of HRRN scheduling is not possible as it is not possible to know the burst time of every job in advance.
- In this scheduling, there may occur an overload on the CPU.

9. Multiple Queue Scheduling:

Processes in the ready queue can be divided into different classes where each class has its own scheduling needs. For example, a common division is a **foreground (interactive)** process and a **background (batch)** process. These two classes have different scheduling needs. For this kind of situation **Multilevel Queue Scheduling** is used.



The description of the processes in the above diagram is as follows:

- **System Processes:** The CPU itself has its process to run, generally termed as System Process.
- **Interactive Processes:** An Interactive Process is a type of process in which there should be the same type of interaction.
- **Batch Processes:** Batch processing is generally a technique in the Operating system that collects the programs and data together in the form of a **batch** before the **processing** starts.

Advantages:

- The main merit of the multilevel queue is that it has a low scheduling overhead.

Disadvantages:

- Starvation problem
- It is inflexible in nature

10. Multilevel Feedback Queue Scheduling:

Multilevel Feedback Queue Scheduling (MLFQ) CPU Scheduling is like **Multilevel Queue Scheduling** but in this process can move between the queues. And thus, much more efficient than multilevel queue scheduling.

Characteristics:

- In a multilevel queue-scheduling algorithm, processes are permanently assigned to a queue on entry to the system, and processes are not allowed to move between queues.
- As the processes are permanently assigned to the queue, this setup has the advantage of low scheduling overhead,
- But on the other hand, a disadvantage of being inflexible.

Advantages:

- It is more flexible

OPERATING SYSTEM (23CS403)

- It allows different processes to move between different queues

Disadvantages:

- It also produces CPU overheads
- It is the most complex algorithm.

Comparison between various CPU Scheduling algorithms

Here is a brief comparison between different CPU scheduling algorithms:

Algorithm	Allocation	Complexity	Average waiting time (AWT)	Preemption	Starvation	Performance
FCFS	According to the arrival time of the processes, the CPU is allocated.	Simple and easy to implement	Large.	No	No	Slow
SJF	Based on the lowest CPU burst time (BT).	More complex than FCFS	Smaller than FCFS	No	Yes	Good
SRTF	Same as SJF the allocation of the CPU is based on the lowest CPU burst time (BT). But it is preemptive.	More complex than FCFS	Depending on arrival time, process size	Yes	Yes	Good
RR	According to the order of the process arrives with fixed time quantum (TQ)	The complexity depends on TQ	Large than SJF and Priority scheduling.	Yes	No	Fair
Priority Pre-emptive	According to the priority. The bigger task executes first	Less complex	Smaller than FCFS	Yes	Yes	Well

OPERATING SYSTEM (23CS403)

Priority non-preemptive	According to the priority monitoring the new incoming priority jobs with the higher	Less complex than Priority preemptive	Smaller than No	Yes	Most beneficial with batch systems
--------------------------------	---	---------------------------------------	-----------------	-----	------------------------------------

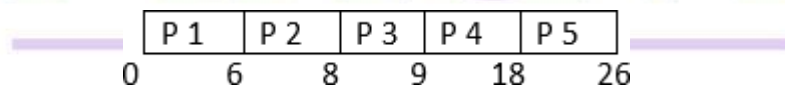
Algorithm	Allocation	Complexity	Average waiting time (AWT)	Preemption	Starvation	Performance
MLQ	According to the process that resides in the bigger queue priority	More complex than the priority	Smaller than FCFS	No	Yes	Good
MLFQ	According to the process of a bigger priority queue.	It is the most Complex	Smaller than all scheduling	No	No	Good

Example 1 (FCFS)

1. Process ID Process Name Burst Time (ms)

P1	A	6
P2	B	2
P3	C	1
P4	D	9
P5	E	8

Gantt Chart



Process ID	Arrival Time (ms)	Burst Time (ms)	Completion Time (ms)	Turn Around Time (ms)	Waiting Time (ms)
P1	0	6	6	6	0
P2	2	2	8	8	6

OPERATINGSYSTEM(23CS403)

P3	3	1	9	9	8
P4	4	9	18	18	9
P5	5	8	26	26	18

AverageTurnAroundTime=(6+8 +9+18 +26)/5=67 /5=13.4ms

AverageWaitingTime =(0 +6 +8+9 +18)/ 5 =41 / 5=8.2ms

Example2(FCFS)

ProcessID	ProcessName	Burst Time
P1	A	79
P2	B	2
P3	C	3
P4	D	1
P5	E	25
P6	F	3

ProcessId	BurstTime (BT)	CompletionTime (CT)	Turn Around Time(TAT)	Waiting Time(WT)
P1	79	79	79	0
P2	2	81	81	79
P3	3	84	84	81
P4	1	85	85	84
P5	25	110	110	85
P6	3	113	113	110

AvgWaitingTime=(0 +79+81 +84 +85+110)/6 =73.17ms

AvgTurnAroundTime= (79 +81+84+85+110+113)/ 6=92 ms

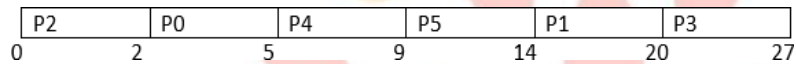
OPERATINGSYSTEM(23CS403)

Example3(SJF)

ProcessID	Arrival Time	BurstTime
P0	1	3
P1	2	6
P2	1	2
P3	3	7
P4	2	4
P5	5	5

NonPre-Emptive ShortestJobFirstCPUScheduling

GanttChart:



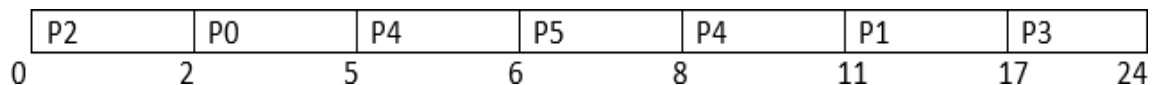
ProcessID	Arrival Time	Burst Time	Completion Time	TurnAround Time TAT=CT-AT	Waiting Time WT=CT-BT
P0	1	3	5	4	1
P1	2	6	20	18	12
P2	0	2	2	2	0
P3	3	7	27	24	17
P4	2	4	9	7	4
P5	5	5	14	10	5

AverageWaitingTime=(1+12 +17 +0+5+4)/6=39/ 6=6.5 ms

AverageTurnAround Time=(4+18+2+24+7+10)/6=65/6=10.83ms Pre Emptive...

Shortest Job First CPU Scheduling

Ganttchart:



OPERATING SYSTEM (23CS403)

Process ID	Arrival Time	Burst Time	Completion Time	Turn Around Time TAT=CT-AT	Waiting Time WT=CT-BT
P0	1	3	5	4	1
P1	2	6	17	15	9
P2	0	2	2	2	0
P3	3	7	24	21	14
P4	2	4	11	9	5
P5	6	2	8	2	0

Average Turn Around Time = $(4 + 15 + 2 + 21 + 9 + 2) / 6 = 53 / 6 = 8.83$ ms

Average Waiting Time = $(1 + 9 + 0 + 14 + 5 + 0) / 6 = 29 / 6 = 4.83$ ms

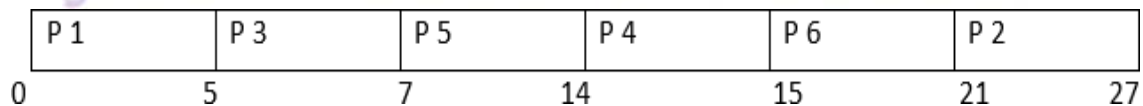
Example 4 (PRIORITY)

S.No	Process ID	Arrival Time	Burst Time	Priority
1	P1	0	5	5
2	P2	1	6	4
3	P3	2	2	0
4	P4	3	1	2
5	P5	4	7	1
6	P6	4	6	3

(5 has the least priority and 0 has the highest priority)

Solution:

Gantt Chart:



OPERATINGSYSTEM(23CS403)

Process Id	Arrival Time	Burst Time	Priority	Completion Time	TurnAround Time TAT=CT-AT	Waiting Time WT=TAT-BT
P1	0	5	5	5	5	0
P2	1	6	4	27	26	20
P3	2	2	0	7	5	3
P4	3	1	2	15	12	11
P5	4	7	1	14	10	3
P6	4	6	3	21	17	11

AvgWaiting Time = $(0+20+3 +11+3 +11) / 6 = 48 / 6 = 8$ ms

AvgTurnAroundTime = $(5 +26+5 +11+ 10 +17) / 6 = 74 / 6 = 12.33$ ms

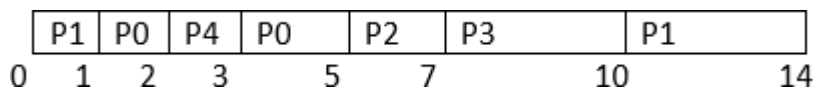
Example5(RoundRobin)

TimeQuantum=1 ms

ProcessID	Arrival Time	Burst Time
P0	1	3
P1	0	5
P2	3	2
P3	4	3
P4	2	1

Solution:

Ganttchart:



ProcessID	Arrival Time	Burst Time	Completion Time	Turn AroundTime	Waiting Time
P0	1	3	5	4	1
P1	0	5	14	14	9
P2	3	2	7	4	2
P3	4	3	10	6	3
P4	2	1	3	1	0

AvgTurnAroundTime=(4+14+4+6+1)/5=5.8 ms
AvgWaitingTime =(1+9+2+3+0)/5=3 ms

DEADLOCK

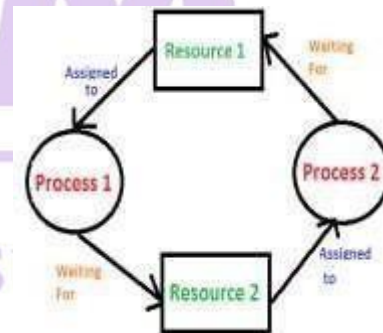
A process in operating system uses resources in the following way.

- (i) Requests a resource
- (ii) Uses the resource
- (iii) Releases the resource

A **deadlock** is a situation where a set of processes are blocked because each process is holding a resource and waiting for another resource acquired by some other process.

Consider an example when two trains are coming toward each other on the same track and there is only one track, none of the trains can move once they are in front of each other.

A similar situation occurs in operating systems when there are two or more processes that hold some resources and wait for resources held by other(s). For example, in the below diagram, Process 1 is holding Resource 1 and waiting for Resource 2 which is acquired by Process 2, and Process 2 is waiting for Resource 1.



Examples of Deadlock

1. The system has 2 tape drives. P1 and P2 each hold one tape drive and each needs

OPERATING SYSTEM (23CS403)

another one.

2. Semaphores A and B, initialized to 1, P0, and P1 are in deadlock as follows:

P0 executes wait(A) and preempts. P1 executes wait(B).

Now P0 and P1 enter in deadlock.

P0	P1
wait(A);	wait(B)
wait(B); and the following	wait(A)

3. Assume the space is available for allocation of 200K byte sequence of events occurs.

P0	P1
Request 80KB;	Request 70KB;
Request 60KB;	Request 80KB;

System model:

A system consists of a finite number of resources to be distributed among a number of competing processes. The resources are partitioned into several types, each consisting of some number of identical instances. Memory space, CPU cycles, files, I/O devices are examples of resource types. If a system has 2 CPUs, then the resource type CPU has 2 instances.

A process must request a resource before using it and must release the resource after using it. A process may request as many resources as it requires to carry out its task. The number of resources as it requires to carry out its task. The number of resources requested may not exceed the total number of resources available in the system. A process cannot request 3 printers if the system has only two.

A process may utilize a resource in the following sequence:

- (I) **REQUEST:** The process requests the resource. If the request cannot be granted immediately (if the resource is being used by another process), then the requesting process must wait until it can acquire the resource.
- (II) **USE:** The process can operate on the resource. If the resource is a printer, the process can print on the printer.
- (III) **RELEASE:** The process releases the resource.

For each use of a kernel managed by a process the operating system checks that the process has requested and has been allocated the resource. A system table records whether each resource is free (or) allocated. For each resource that is allocated, the table also records the process to which it is allocated. If a process requests a resource that is currently allocated to another process, it can be added to a queue of processes waiting for this resource.

OPERATING SYSTEM(23CS403)

To illustrate a deadlocked state, consider a system with 3 CD RW drives. Each of 3

OPERATING SYSTEM (23CS403)

processes hold one of these CDRW drives. If each process now requests another drive, the 3 processes will be in a deadlocked state. Each is waiting for the event "CDRW is released" which can be caused only by one of the other waiting processes. This example illustrates a deadlock involving the same resource type.

Deadlocks may also involve different resource types. Consider a system with one printer and one DVD drive. The process P_i is holding the DVD and process P_j is holding the printer. If P_i requests the printer and P_j requests the DVD drive, a deadlock occurs.

NECESSARY CONDITIONS FOR DEADLOCK

- **Mutual Exclusion**
Two or more resources are non-shareable (Only one process can use at a time)
- **Hold and Wait**
A process is holding at least one resource and waiting for resources.
- **No Pre-emption**
A resource cannot be taken from a process unless the process releases the resource.
- **Circular Wait**
A set of processes waiting for each other in circular form.

Resource Allocation Graph

The resource allocation graph is the pictorial representation of the state of a system. As its name suggests, the resource allocation graph is the complete information about all the processes which are holding some resources or waiting for some resources.

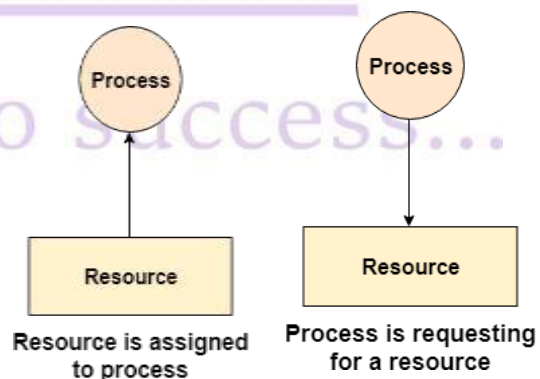
It also contains the information about all the instances of all the resources whether they are available or being used by the processes.

In Resource allocation graph, the process is represented by a Circle while the Resource is represented by a rectangle.

Vertices are mainly of two types, Resource and Process. Each of them will be represented by a different shape. Circle represents process while rectangle represents resource. A resource can have more than one instance. Each instance will be represented by a dot inside the rectangle.

Edges in RAG are also of two types, one represents **Assignment Edge** and other represents the wait of a process for a resource i.e. **Request Edge**.

A resource is shown as assigned to a process if the tail of the arrow is attached to an instance of the resource and the head is attached to a process.



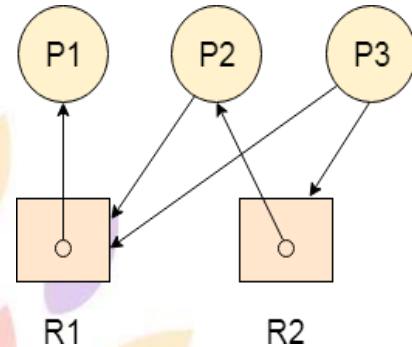
A process is shown as waiting for a resource if the tail of an arrow is attached to the process while the head is pointing towards the resource.

Example

Consider 3 processes P1, P2 and P3 and two types of resources R1 and R2. The resources are having 1 instance each.

According to the graph, R1 is being used by P1, P2 is holding R2 and waiting for R1, P3 is waiting for R1 as well as R2.

The graph is deadlock free since no cycle is being formed in the graph.



Using Resource Allocation Graph, it can be easily detected whether system is in a Deadlock state or not. The rules are

Rule-01: In a Resource Allocation Graph where all the resources are single instance,

- If a cycle is being formed, then system is in a deadlock state.
- If no cycle is being formed, then system is not in a deadlock state.

Rule-02: In a Resource Allocation Graph where all the resources are **NOT** single instance,

- If a cycle is being formed, then system may be in a deadlock state.
- **Banker's Algorithm** is applied to confirm whether system is in a deadlock state or not.
- If no cycle is being formed, then system is not in a deadlock state.
- Presence of a cycle is a necessary but not a sufficient condition for the occurrence of deadlock.

METHODS FOR HANDLING DEADLOCK

There are three ways to handle deadlock

1) **Deadlock prevention or avoidance** **PREVENTION**

The idea is to not let the system into a deadlock state. This system will make sure that above mentioned four conditions will not arise. These techniques are very costly so we use this in cases where our priority is making a system deadlock-free.

One can zoom into each category individually, Prevention is done by negating one of the four necessary conditions for deadlock.

Eliminate mutual exclusion

It is not possible to dis-satisfy the mutual exclusion because some resources, such as the tape drive and printer, are inherently non-shareable.

Solve hold and wait

Allocate all required resources to the process before the start of its execution, this way hold and wait condition is eliminated but it will lead to low device utilization. For example, if a process requires a printer at a later time and we have allocated a printer before the start of its execution printer will remain blocked till it has completed its execution. The process will make a new request for resources after releasing the current set of resources. This solution may lead to starvation.

Allow pre-emption

Preempt resources from the process when resources are required by other high-priority processes.

Circular wait solution

Each resource will be assigned a numerical number. A process can request the resources to increase/decrease. order of numbering. For Example, if the P1 process is allocated R5 resources, now next time if P1 asks for R4, R3 lesser than R5 such a request will not be granted, only a request for resources more than R5 will be granted.

AVOIDANCE

Avoidance is kind of futuristic. By using the strategy of "Avoidance", we have to make an assumption. We need to ensure that all information about resources that the process will need is known to us before the execution of the process.

Resource Allocation Graph

The resource allocation graph (RAG) is used to visualize the system's current state as a graph. The Graph includes all processes, the resources that are assigned to them, as well as the resources that each Process requests. Sometimes, if there are fewer processes, we can quickly spot a deadlock in the system by looking at the graph rather than the tables we use in Banker's algorithm.

Banker's Algorithm

Banker's Algorithm is a resource allocation and deadlock avoidance algorithm which tests all the request made by processes for resources, it checks for the safe state, and after granting a request system remains in the safe state it allows the request, and if there is no safe state it doesn't allow the request made by the process.

In prevention and avoidance, we get the correctness of data but performance

decreases.

2) Deadlock detection and recovery

If deadlock prevention or avoidance is not applied to the software then we can handle this by deadlock detection and recovery, which consist of two phases.

In the first phase, we examine the state of the process and check whether there is a deadlock or not in the system.

If found deadlock in the first phase then we apply the algorithm for recovery of the deadlock.

3) Deadlock ignorance:

If a deadlock is very rare, then let it happen and reboot the system. This is the approach that both Windows and UNIX take. We use the ostrich algorithm for deadlock ignorance.

In Deadlock, ignorance performance is better than the above two methods but not the correctness of data.

SAFESTATE

A safe state can be defined as a state in which there is no deadlock. It is achievable if:

- If a process needs an unavailable resource, it may wait until the same has been released by a process to which it has already been allocated. If such a sequence does not exist, it is an unsafe state.
- All the requested resources are allocated to the process.

BANKER'S ALGORITHM

It is a banker algorithm used to **avoid deadlock** and **allocate resources** safely to each process in the computer system. The '**S-State**' examines all possible tests or activities before deciding whether the allocation should be allowed to each process. It also helps the operating system to successfully share the resources between all the processes.

The banker's algorithm is named because it checks whether a person should be sanctioned a loan amount or not to help the bank system safely simulate allocation resources.

Suppose the number of account holders in a particular bank is 'n', and the total money in a bank is 'T'. If an account holder applies for a loan; first, the bank subtracts the loan amount from full cash and then estimates the cash difference is greater than T to approve the loan amount. These steps are taken because if another person applies for a loan or withdraws some amount from the bank, it helps the bank manage and operate all things without any restriction in the functionality of the banking system.

Similarly, it works in an **operating system**. When a new process is created in a computer system, the process must provide all types of information to the

OPERATING SYSTEM(23CS403)

operating system like upcoming processes, requests for their resources, counting them, and delays.

Based on these criteria, the operating system decides which process sequence should be executed or waited so that no deadlock occurs in a system. Therefore, it is also known as **deadlock avoidance algorithm** or **deadlock detection** in the operating system.

When working with a banker's algorithm, it requests to know about three things:

1. How much each process can request for each resource in the system. It is denoted by the **[MAX]** request.
2. How much each process is currently holding each resource in a system. It is denoted by the **[ALLOCATED]** resource.
3. It represents the number of each resource currently available in the system. It is denoted by the **[AVAILABLE]** resource.

Following are the important data structures/terms applied in the banker's algorithm as follows:

Suppose n is the number of processes, and m is the number of each type of resource used in a computer system.

1. **Available:** It is an array of length ' m ' that defines each type of resource available in the system. When $Available[j] = K$, means that ' K ' instances of Resources type $R[j]$ are available in the system.
2. **Max:** It is a $[n \times m]$ matrix that indicates each process $P[i]$ can store the maximum number of resources $R[j]$ (each type) in a system.
3. **Allocation:** It is a matrix of $m \times n$ orders that indicates the type of resources currently allocated to each process in the system. When $Allocation[i, j] = K$, it means that process $P[i]$ is currently allocated K instances of Resources type $R[j]$ in the system.
4. **Need:** It is an $M \times N$ matrix sequence representing the number of remaining resources for each process. When the $Need[i][j] = k$, then process $P[i]$ may require K more instances of resources type R_j to complete the assigned work.
 $Need[i][j] = Max[i][j] - Allocation[i][j]$.
5. **Finish:** It is the vector of the order m . It includes a Boolean value (true/false) indicating whether the process has been allocated to the requested resources, and all resources have been released after finishing its task.

OPERATINGSYSTEM(23CS403)

The Banker's Algorithm is the combination of the safety algorithm and the resource request algorithm to control the processes and avoid deadlock.

Safety Algorithm

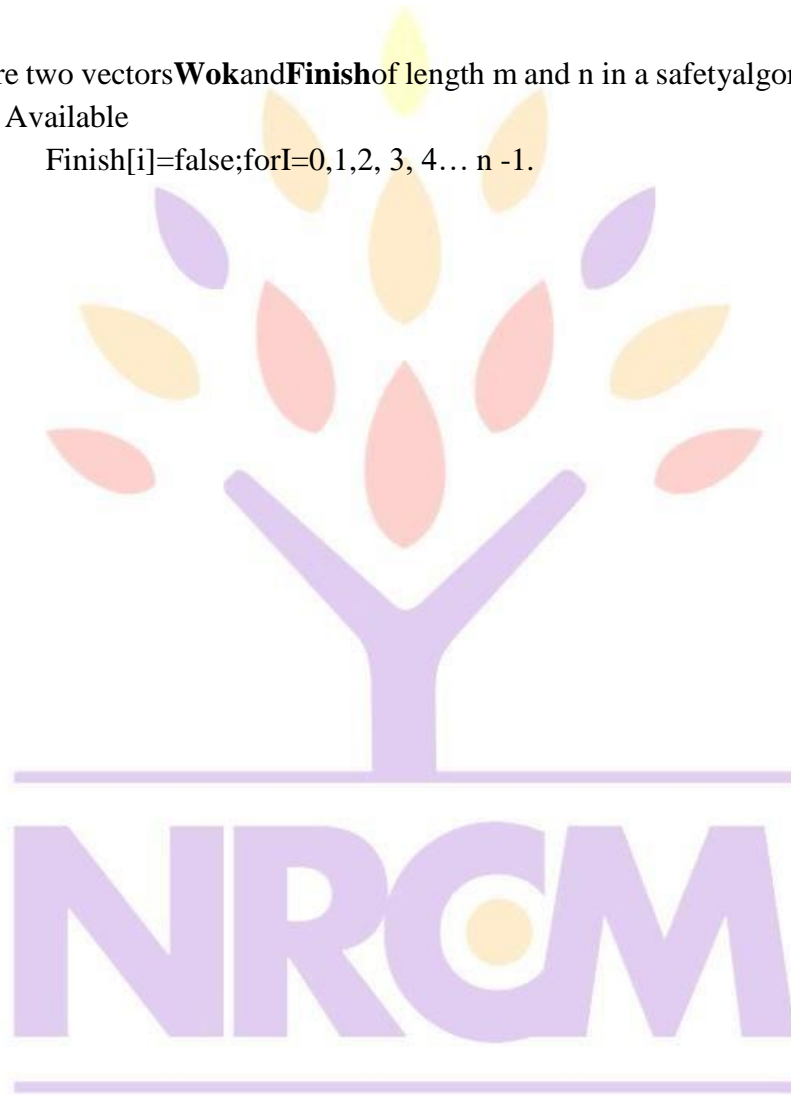
It is a safety algorithm used to check whether or not a system is in a safe state or follows the safe sequence in a banker's algorithm:

Step1:

There are two vectors **Work** and **Finish** of length m and n in a safety algorithm. Initialize:

Work = Available

Finish[i] = false; for $i = 0, 1, 2, 3, 4 \dots n - 1$.



your roots to success...

OPERATING SYSTEM(23CS403)

Step2:

Check the availability status for each type of resources [i], such as: $Need[i] \leq Work$
 $Finish[i] == false$
If it does not exist, go to step 4. Step3:

Step4:

$Work = Work + Allocation(i)$ //to get new resource allocation
 $Finish[i] = true$
Go to step 2 to check the status of resource availability for the next process. If $Finish[i] == true$;

it means that the system is safe for all processes.

Resource Request Algorithm

Let create a resource request array $R[i]$ for each process $P[i]$.

Step1:

When the number of **requested resources** of each type is less than the **Need** resources, go to step 2 and if the condition fails, which means that the process $P[i]$ exceeds its maximum claim for the resource. As the expressions suggests:

If $Request(i) \leq Need$, then go to step 2, Else raise an error message.

Step2:

And when the number of requested resources of each type is less than the available resource for each process, go to step (3). As the expressions suggests: If $Request(i) \leq Available$, then go to step 3.

Else Process $P[i]$ must wait for the resource.

Step3:

When the requested resource is allocated to the process by changing state: $Available = Available - Request$
 $Allocation(i) = Allocation(i) + Request(i)$ $Need_i = Need_i - Request_i$

When the resource allocation state is safe, its resources are allocated to the process $P(i)$. And if the new state is unsafe, the Process $P(i)$ has to wait for each type of Request $R(i)$ and restore the old resource-allocation state.

OPERATING SYSTEM(23CS403)

Example:

Consider a system that contains five processes P1, P2, P3, P4, P5 and the three resource types A, B and C. Following are the resources types: A has 10, B has 5 and the resource type C has 7 instances.

Process	Allocation			A	Max B	C	Available		
	A	B	C				A	B	C
P1	0	1	0	7	5	3	3	3	2
P2	2	0	0	3	2	2			
P3	3	0	2	9	0	2			
P4	2	1	1	2	2	2			
P5	0	0	2	4	3	3			

Answer the following questions using the banker's algorithm:

1. What is the reference of the need matrix?
2. Determine if the system is safe or not.
3. What will happen if the resource request (1,0,2) for process P1 can the system accept this request immediately?
4. What will happen if the resource request (3,3,0) for process P5?
5. What will happen if the resource request (0,2,0) for process P1?

Ans.1:

Context of the need matrix is as $Need[i] = Max[i] - Allocation[i]$

Need for P1: $(7,5,3) - (0,1,0) = 7, 4, 3$

Need for P2: $(3,2, 2) - (2,0,0) = 1, 2, 2$

Need for P3: $(9,0, 2) - (3,0,2) = 6, 0, 0$

Need for P4: $(2,2, 2) - (2,1,1) = 0, 1, 1$

Need for P5: $(4,3, 3) - (0,0,2) = 4, 3, 1$

Process	Need		
	A	B	C
P1	7	4	3
P2	1	2	2
P3	6	0	0
P4	0	1	1
P5	4	3	1

OPERATINGSYSTEM(23CS403)

Ans.2:ApplytheBanker'sAlgorithm:

AvailableResourcesofA,BandCare3, 3,and 2.

Nowwecheckifeachtypeofresourcerequestisavailableforeachprocess.

Step1:

ForProcessP1:

Need \leq Available

7, 4, 3 \leq 3, 3, 2 condition is **false**.

So,we examine another process, P2.

Step2:

ForProcessP2:

Need \leq Available

1,2, 2 \leq 3,3, 2condition**true**

New available = available + Allocation(3,
3, 2) + (2, 0, 0) \Rightarrow 5, 3, 2

Similarly,weexamineanotherprocessP3.

Step3:

ForProcessP3:

P3Need \leq Available

6, 0, 0 \leq 5, 3, 2 condition is **false**.

Similarly,we examine another process, P4.

Step4:

ForProcessP4:

P4Need \leq Available

0,1, 1 \leq 5,3, 2conditionis **true**

NewAvailableresource=Available+Allocation5,3, 2 +
2, 1, 1 \Rightarrow 7, 4, 3

Similarly,weexamineanotherprocessP5.

Step5:

ForProcessP5:

P5Need \leq Available

4,3, 1 \leq 7,4, 3conditionis **true**

New available resource= Available+Allocation7, 4, 3
+ 0, 0, 2 \Rightarrow 7, 4, 5

Now,weagainexamineeachtypeofresourcerequestforprocesses P1 and P3.

OPERATING SYSTEM(23CS403)

Step6:

For Process P1:

$P1_{Need} \leq Available$

$7, 4, 3 \leq 7, 4, 5$ condition is **true**

$New Available Resource = Available + Allocation$ $7, 4, 5 +$

$0, 1, 0 \Rightarrow 7, 5, 5$

So, we examine another process P2.

Step7:

For Process P3:

$P3_{Need} \leq Available$

$6, 0, 0 \leq 7, 5, 5$ condition is **true**

$New Available Resource = Available + Allocation$ $7, 5, 5 +$

$3, 0, 2 \Rightarrow 10, 5, 7$

Hence, we execute the banker's algorithm to find the safe state and the safe sequence like P2, P4, P5, P1 and P3.

Ans.3:

For granting the Request (1,0,2), first we have to check that

Request \leq Available, that is $(1, 0, 2) \leq (3, 3, 2)$, Since

the condition is true, the process P2 may get

the request immediately.

Allocation for P2 is (3,0,2) and new Available is (2, 3, 0)

Context of the need matrix is as follows: Need [i]

$= Max [i] - Allocation [i]$

Need for P1: $(7, 5, 3) - (0, 1, 0) = 7, 4, 3$

Need for P2: $(3, 2, 2) - (3, 0, 2) = 0, 2, 0$

Need for P3: $(9, 0, 2) - (3, 0, 2) = 6, 0, 0$

Need for P4: $(2, 2, 2) - (2, 1, 1) = 0, 1, 1$

Need for P5: $(4, 3, 3) - (0, 0, 2) = 4, 3, 1$

Process	Need		
	A	B	C
P1	7	4	3
P2	0	2	0
P3	6	0	0
P4	0	1	1
P5	4	3	1

Apply the Banker's Algorithm:

Available Resources of A, B and C are 2, 3, and 0.

Now we check if each type of resource request is available for each process.

your roots to success...

OPERATINGSYSTEM(23CS403)

Step1:

ForProcessP1:

Need \leq Available

7, 4, 3 \leq 2, 3, 0 condition is **false**.

So, we examine another process, P2.

Step2:

ForProcessP2:

Need \leq Available

1, 2, 2 \leq 2, 3, 0 condition **true**

New available = available + Allocation(2,

3, 0) + (3, 0, 2) \Rightarrow 5, 3, 2

Similarly, we examine another process P3.

Step3:

ForProcessP3:

P3Need \leq Available

6, 0, 0 \leq 5, 3, 2 condition is **false**.

Similarly, we examine another process, P4.

Step4:

ForProcessP4:

P4Need \leq Available

0, 1, 1 \leq 5, 3, 2 condition is **true**

New Available resource = Available + Allocation 5, 3, 2 +

2, 1, 1 \Rightarrow 7, 4, 3

Similarly, we examine another process P5.

Step5:

ForProcessP5:

P5Need \leq Available

4, 3, 1 \leq 7, 4, 3 condition is **true**

New available resource = Available + Allocation 7, 4, 3 +

0, 0, 2 \Rightarrow 7, 4, 5

Now, we again examine for processes P1 and P3.

Step6:

ForProcessP1:

P1Need \leq Available

7, 4, 3 \leq 7, 4, 5 condition is **true**

New Available Resource = Available + Allocation 7, 4, 5 +

0, 1, 0 \Rightarrow 7, 5, 5

So, we examine another process P2.

Step7:

ForProcessP3:

$P3Need \leq Available$

$6, 0, 0 \leq 7, 5, 5$ condition is true

$NewAvailableResource = Available + Allocation$ $7, 5, 5 +$

$3, 0, 2 \Rightarrow 10, 5, 7$

Hence, P2 granted immediately and the safe sequence like P2, P4, P5, P1 and P3.

Ans.4:

For granting the Request (3,3, 0) by P5, first we have to check that **Request** \leq **Available**, that is $(3, 3, 0) \leq (2, 3, 0)$, Since the condition is false. So the request for (3,3,0) by process P5 cannot be granted.

Ans.5:

For granting the Request (0,2, 0) by P1, first we have to check that **Request** \leq **Available**, that is $(0, 2, 0) \leq (2, 3, 0)$, Since the condition is true. So the request for (0,2,0) by process P1 may be granted.

Allocation for P1 is (0,3,0)

Context of the need matrix is as follows: Need [i]

= Max [i] - Allocation [i]

Need for P1: $(7, 5, 3) - (0, 3, 0) = 7, 2, 3$

Process	Need		
	A	B	C
P1	7	2	3
P2	0	2	0
P3	6	0	0
P4	0	1	1
P5	4	3	1

Apply the Banker's Algorithm:

Available Resources of A, B and C are 2, 1, and 0.

For Process P1: $7, 2, 3 \leq 2, 1, 0$ condition is **false**.

For Process P2: $0, 2, 0 \leq 2, 1, 0$ condition is **false**.

For Process P3: $6, 0, 0 \leq 2, 1, 0$ condition is **false**.

For Process P4: $0, 1, 1 \leq 2, 1, 0$ condition is **false**.

For Process P5: $4, 3, 1 \leq 2, 1, 0$ condition is **false**.

Hence, the state is unsafe, P1 cannot be granted immediately.

DEADLOCKDETECTION

If a system does not employ either a deadlock prevention or deadlock avoidance algorithm then a deadlock situation may occur. In this case-

- Apply an algorithm to examine the system's state to determine whether a deadlock has occurred.
- Apply an algorithm to recover from the deadlock.

A deadlock detection algorithm is a technique used by an operating system to identify deadlocks in the system. This algorithm checks the status of processes and resources to determine whether any deadlock has occurred and takes appropriate actions to recover from the deadlock.

The algorithm employs several times varying data structures:

Available – A vector of length m indicates the number of available resources of each type.

Allocation – An $n \times m$ matrix defines the number of resources of each type currently allocated to a process. The column represents resource and rows represent a process.

Request – An $n \times m$ matrix indicates the current request of each process. If $request[i][j]$ equals k then process P_i is requesting k more instances of resource type R_j .

The Bankers algorithm includes a **Safety Algorithm / Deadlock Detection Algorithm**. The algorithm for finding out whether a system is in a safe state can be described as follows:

Steps of Algorithm:

1. Let *Work* and *Finish* be vectors of length m and n respectively. Initialize $Work = Available$. For $i = 0, 1, \dots, n-1$, if $Request_i = 0$, then $Finish[i] = true$; otherwise, $Finish[i] = false$.
2. Find an index i such that both
 - a) $Finish[i] = false$
 - b) $Request_i \leq Work$If no such i exists, go to step 4.
3. $Work = Work + Allocation_i$
 $Finish[i] = true$
Goto Step 2.
4. If $Finish[i] = false$ for some $i, 0 \leq i < n$, then the system is in a deadlocked state. Moreover, if $Finish[i] = false$ the process P_i is deadlocked.

Forexample,

	Allocation			Request			Available		
	A	B	C	A	B	C	A	B	C
P0	0	1	0	0	0	0	0	0	0
P1	2	0	0	2	0	2			
P2	3	0	3	0	0	0			
P3	2	1	1	1	0	0			
P4	0	0	2	0	0	2			

1. In this, $Work = [0, 0, 0]$ & $Finish = [false, false, false, false, false]$
2. $i=0$ is selected as both $Finish[0]=false$ and $[0,0,0] \leq [0,0,0]$.
 $Work = [0, 0, 0] + [0, 1, 0] \Rightarrow [0, 1, 0]$ & $Finish = [true, false, false, false, false]$.
3. $i=1$ is selected as both $Finish[1]=false$ and $[0,1,0] \leq [0,1,0]$.
 $Work = [0, 1, 0] + [2, 0, 0] \Rightarrow [2, 1, 0]$ & $Finish = [true, true, false, false, false]$.
4. $i=2$ is selected as both $Finish[2]=false$ and $[2,1,0] \leq [2,1,0]$.
 $Work = [2, 1, 0] + [3, 0, 3] \Rightarrow [5, 1, 3]$ & $Finish = [true, true, true, false, false]$.
5. $i=3$ is selected as both $Finish[3]=false$ and $[5,1,3] \leq [5,1,3]$.
 $Work = [5, 1, 3] + [2, 1, 1] \Rightarrow [7, 2, 4]$ & $Finish = [true, true, true, true, false]$.
6. $i=4$ is selected as both $Finish[4]=false$ and $[7,2,4] \leq [7,2,4]$.
 $Work = [7, 2, 4] + [0, 0, 2] \Rightarrow [7, 2, 6]$ & $Finish = [true, true, true, true, true]$.
7. Since $Finish$ is a vector of all true it means **there is no dead lock** in this example.

your roots to success...

There are several algorithms for detecting deadlocks in an operating system, including:

1. Wait-For Graph:

A graphical representation of the system's processes and resources. A directed edge is created from a process to a resource if the process is waiting for that resource. A cycle in the graph indicates a deadlock.

2. Banker's Algorithm:

A resource allocation algorithm that ensures that the system is always in a safe state, where deadlocks cannot occur.

3. Resource Allocation Graph:

A graphical representation of processes and resources, where a directed edge from a process to a resource means that the process is currently holding that resource. Deadlocks can be detected by looking for cycles in the graph.

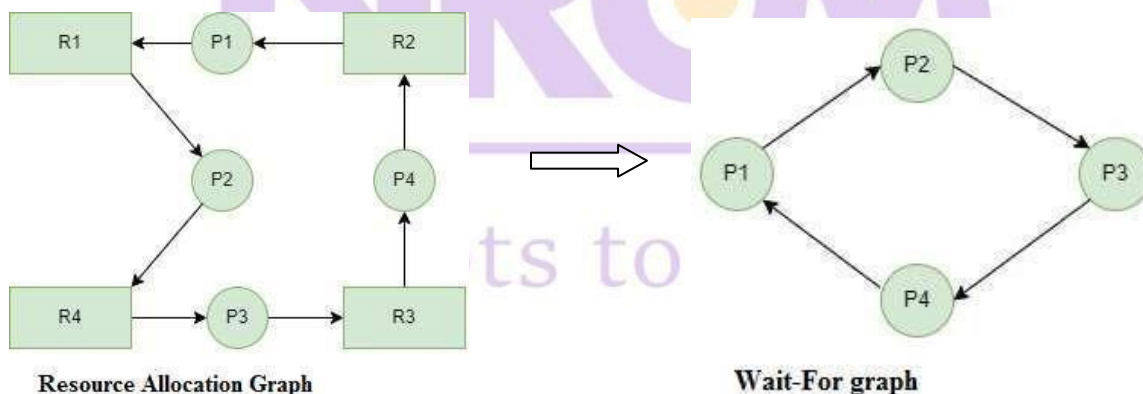
4. Detection by System Modeling:

A mathematical model of the system is created, and deadlocks can be detected by finding a state in the model where no process can continue to make progress.

5. Timestamping:

Each process is assigned a timestamp, and the system checks to see if any process is waiting for a resource that is held by a process with a lower timestamp.

These algorithms are used in different operating systems and systems with different resource allocation and synchronization requirements. The choice of algorithm depends on the specific requirements of the system and the trade-offs between performance, complexity and accuracy.



RECOVERY FROM DEADLOCK

The OS will use various recovery techniques to restore the system if it encounters any deadlocks. When a Deadlock Detection Algorithm determines that a deadlock has occurred in the system, the system must recover from that deadlock.

Approaches to Breaking a Deadlock

(a) Process Termination

To eliminate the deadlock, we can simply kill one or more processes. For this, we use two methods:

1. Abort all the Deadlocked Processes:

Aborting all the processes will certainly break the deadlock but at a great expense. The deadlocked processes may have been computed for a long time, and the result of those partial computations must be discarded and there is a probability of recalculating them later.

2. Abort one process at a time until the deadlock is eliminated:

Abort one deadlocked process at a time, until the deadlock cycle is eliminated from the system. Due to this method, there may be considerable overhead, because, after aborting each process, we have to run a deadlock detection algorithm to check whether any processes are still deadlocked.

(b) Resource Preemption

To eliminate deadlocks using resource preemption, we preempt some resources from processes and give those resources to other processes. This method will raise three issues –

1. Selecting a victim:

We must determine which resources and which processes are to be preempted and also in order to minimize the cost.

2. Rollback:

We must determine what should be done with the process from which resources are preempted. One simple idea is total rollback. That means aborting the process and restarting it.

3. Starvation:

In a system, it may happen that the same process is always picked as a victim. As a result, that process will never complete its designated task. This situation is called **Starvation** and must be avoided. One solution is that a process must be picked as a victim only a finite number of times.

UNIT-3

Process Management and Synchronization-

The critical section problems, Synchronization hardware, Semaphore, and Classical problems of Synchronization, Critical region, Monitor.

Inter process communication Mechanism- IPC between process on a single computer system, IPC between process on different system, Using Pipes, FIFOs, Message Queue, Shared memory

SYNCHRONIZATION

Process Synchronization is the coordination of execution of multiple processes in a multi-process system to ensure that they access shared resources in a controlled and predictable manner. It aims to resolve the problem of race conditions and other synchronization issues in a concurrent system.

The main objective of process synchronization is to ensure that multiple processes access shared resources without interfering with each other and to prevent the possibility of inconsistent data due to concurrent access. To achieve this, various synchronization techniques such as semaphores, monitors and critical sections are used.

On the basis of synchronization, processes are categorized as one of the following two types:

- **Independent Process:** The execution of one process does not affect the execution of other processes.
- **Cooperative Process:** A process that can affect or be affected by other processes executing in the system.

Process synchronization problem arises in the case of Cooperative processes also because resources are shared in Cooperative processes.

Race Condition

A race condition is a condition when there are many processes and every process shares the data with each other and accessing the data concurrently and the output of execution depends on a particular sequence in which they share the data and access.

(OR)

When more than one process is executing the same code or accessing the same memory or a shared variable in that condition there is a possibility that the output or the value of the shared variable is wrong so for that all the processes doing the race to say that my output is correct. This condition is known as **race condition**.

Several processes access and process the manipulations over the same data concurrently, then the outcome depends on the particular order in which the access takes place.

Example:

Let's say there are two processes P1 and P2 which share common variable (shared=10), both processes are present in ready – queue and waiting for its turn to be executed.

Suppose, Process P1 first comes under execution, initialized as X=10 and increments it by 1 (ie. X=11), after then when CPU reads line sleep(1), it switches from current process P1 to process P2 present in ready-queue. The process P1 goes in waiting state for 1 second.

Process1	Process2
intX=shared	intY=shared
X++	Y--
sleep(1)	sleep(1)
shared= X	shared= Y

Now CPU executes the Process P2, initializes Y=10 and decrements Y by 1 (ie. Y=9), after then when CPU reads sleep(1), the current process P2 goes in waiting state and CPU remains idle for some time as there is no process in ready-queue.

After completion of 1 second of process P1 when it comes in ready-queue, CPU takes the process P1 under execution and executes the remaining line of code and shared=11.

After completion of 1 second of Process P2, when process P2 comes in ready-queue, CPU starts executing the further remaining line of Process P2 and shared=9.

Note:

We are assuming the final value of common variable (shared) after execution of Process P1 and Process P2 is 10 (as Process P1 increments variable by 1 and Process P2 decrements variable by 1 and finally it becomes shared=10). But we are getting undesired value due to lack of proper synchronization.

Actual meaning of race-condition

- If the order of execution of process (first P1 -> then P2) then we will get the value of common variable (shared) = 9.
- If the order of execution of process (first P2 -> then P1) then we will get the final value of common variable (shared) = 11.

Basically, here the (value1 = 9) and (value2=11) are racing, if we execute these two processes in our computer system then sometime we will get 9 and sometime we will get 10 as final value of common variable (shared). This phenomenon is called **Race-Condition**.

CRITICALSECTIONPROBLEM

A critical section is a code segment that can be accessed by only one process at a time. The critical section contains shared variables that need to be synchronised to maintain the consistency of data variables. So the critical section problem means designing a way for cooperative processes to access shared resources without creating data inconsistencies.

In the entry section, the process requests for entry in the **Critical Section**.

Any solution to the critical section problem must satisfy three requirements:

- **Mutual Exclusion:** If a process is executing in its critical section, then no other process is allowed to execute in the critical section.
- **Progress:** If no process is executing in the critical section and other processes are waiting outside the critical section, then only those processes that are not executing in their remainder section can participate in deciding which will enter in the critical section next, and the selection can't be postponed indefinitely.
- **Bounded Waiting:** A bound must exist 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.

PETERSON'S SOLUTION

Peterson's solution is a classical software-based solution to the critical section problem. In Peterson's solution, we have two shared variables:

- **boolean flag[i]:** Initialized to FALSE, initially no one is interested in entering the critical section
- **int turn:** The process whose turn is to enter the critical section.

```
//code for producer i
do
{
flag[i]=true;turn
=j;
while(flag[j]==true&&turn==j);
criticalsection
flag[i]=false;
remindersection
}while(TRUE);
```

```
//codeforconsumerj
do
{
flag[j]=true;turn
=i;
while(flag[i]==true&&turn==i);
criticalsection
flag[i]=false;
remindersection
}while(TRUE);
```

In the solution, i represents the Producer and j represents the Consumer. Initially, the flags are false. When a process wants to execute its critical section, it sets its flag to true and turn into the index of the other process. This means that the process wants to execute but it will allow the other process to run first. The process performs busy waiting until the other process has finished its own critical section. After this, the current process enters its critical section and adds or removes a random number from the shared buffer. After completing the critical section, it sets its own flag to false, indicating it does not wish to execute anymore.

Peterson's Solution preserves all three conditions:

- Mutual Exclusion is assured as only one process can access the critical section at any time.
- Progress is also assured, as a process outside the critical section does not block other processes from entering the critical section.
- Bounded Waiting is preserved as every process gets a fair chance.

Disadvantages of Peterson's Solution

- It involves busy waiting.
- It is limited to 2 processes.
- Peterson's solution cannot be used in modern CPU architectures.

Synchronization Hardware

- Problems of Critical Section are also solvable by hardware.
- Uniprocessor systems disable interrupts while a Process P_i is using the CS but it is a great disadvantage in multiprocessor systems

- Some systems provide a lock functionality where a Process acquires a lock while entering the CS and releases the lock after leaving it. Thus another process trying to enter CS cannot enter as the entry is locked. It can only do so if it is free by acquiring the lock itself
- Another advanced approach is the **Atomic Instructions** (Non-Interruptible instructions).

MUTEX LOCKS

- As the synchronization hardware solution is not easy to implement from everywhere, a strict software approach called Mutex Locks was introduced. In this approach, in the entry section of code, a LOCK is acquired over the critical resources modified and used inside the critical section, and in the exit section that LOCK is released. As the resource is locked while a process executes its critical section hence no other process can access

SEMAPHORES

Semaphore is a Hardware Solution. This Hardware solution is written or given to critical section problem. The Semaphore is just a normal integer. The Semaphore cannot be negative. The least value for a Semaphore is zero (0). The Maximum value of a Semaphore can be anything. The Semaphores usually have two operations. The two operations have the capability to decide the values of the semaphores.

The two Semaphore Operations are:

1. Wait()
2. Signal()

Wait Semaphore Operation

The Wait operation works on the basis of Semaphore or Mutex Value. If the Semaphore value is greater than zero, then the Process can enter the Critical Section Area. If the Semaphore value is equal to zero then the Process has to wait. If the process exits the Critical Section, then have to reduce the value of Semaphore.

Definition of wait()

```
wait(Semaphore S)
{
    while (S <= 0); //no operation
    S--;
}
```

Signal Semaphore Operation

The most important part is that this Signal Operation or V Function is executed

OPERATINGSYSTEM(23CS403)

only when the process comes out of the critical section. The value of semaphore cannot be incremented before the exit of process from the critical section.

```
Definition of signal()
signal(S)
{
    S++;
}
```

There are two types of semaphores:

➤ **Binary Semaphores:**

They can only be either 0 or 1. They are also known as mutex locks, as the locks can provide mutual exclusion. All the processes can share the same mutex semaphore that is initialized to 1. Then, a process has to wait until the lock becomes 0. Then, the process can make the mutex semaphore 1 and start its critical section. When it completes its critical section, it can reset the value of the mutex semaphore to 0 and some other process can enter its critical section.

➤ **Counting Semaphores:**

They can have any value and are not restricted over a certain domain. They can be used to control access to a resource that has a limitation on the number of simultaneous accesses. The semaphore can be initialized to the number of instances of the resource. Whenever a process wants to use that resource, it checks if the number of remaining instances is more than zero, i.e., the process has an instance available. Then, the process can enter its critical section thereby decreasing the value of the counting semaphore by 1. After the process is over with the use of the instance of the resource, it can leave the critical section thereby adding 1 to the number of available instances of the resource.

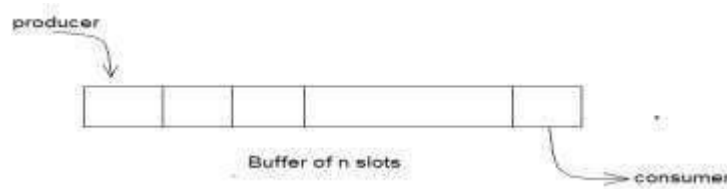
CLASSICAL PROBLEMS OF SYNCHRONIZATION

The following problems of synchronization are considered as classical problems:

1. Bounded-buffer (or Producer-Consumer) Problem,
2. Dining-Philosophers Problem,
3. Readers and Writers Problem,

Bounded-buffer (or Producer-Consumer) Problem

Bounded Buffer problem is also called **producer consumer problem** and it is one of the classic problems of synchronization. This problem is generalized in terms of the



Producer-Consumer problem. Solution to this problem is, creating two counting semaphores “full” and “empty” to keep track of the current number of full and empty buffers respectively. Producers produce a product and consumers consume the product, but both use one of the containers each time.

A producer tries to insert data into an empty slot of the buffer. A consumer tries to remove data from a filled slot in the buffer. There needs to be a way to make the producer and consumer work in an independent manner.

One solution of this problem is to use semaphores. The semaphores which will be used here are:

- **m, a binary semaphore** which is used to acquire and release the lock.
- **empty, a counting semaphore** whose initial value is the number of slots in the buffer, since, initially all slots are empty.
- **full, a counting semaphore** whose initial value is 0.

At any instant, the current value of empty represents the number of empty slots in the buffer and full represents the number of occupied slots in the buffer.

The Producer Operation

```
do
{
    wait(empty);           //wait until empty > 0 and then decrement 'empty'
    wait(mutex);         // acquire lock

    /*perform the insert operation in a slot*/

    signal(mutex);       // release lock
    signal(full);        //increment 'full'

} while(TRUE);
```

- Looking at the above code for a producer, we can see that a producer first waits until there is at least one empty slot.
- Then it decrements the **empty** semaphore because, there will now be one less empty slot, since the producer is going to insert data in one of those slots.

OPERATINGSYSTEM(23CS403)

- Then, it acquires lock on the buffer, so that the consumer cannot access the buffer

until producer completes its operation.

- After performing the insert operation, the lock is released and the value of **full** is incremented because the producer has just filled a slot in the buffer.

The Consumer Operation

```
do
{
    wait(full);
                                //wait until full > 0 and then decrement 'full'
    ex);                          // acquire the lock

    /*perform the remove operation in a slot*/

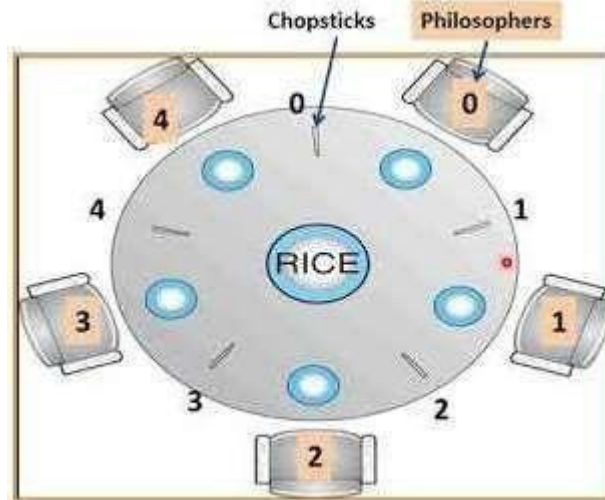
    signal(mutex);               //release the lock
    signal(empty);               // increment 'empty'

} while(TRUE);
```

- The consumer waits until there is at least one full slot in the buffer.
- Then it decrements the **full** semaphore because the number of occupied slots will be decreased by one, after the consumer completes its operation.
- After that, the consumer acquires lock on the buffer.
- Following that, the consumer completes the removal operation so that the data from one of the full slots is removed.
- Then, the consumer releases the lock.
- Finally, the **empty** semaphore is incremented by 1, because the consumer has just removed data from an occupied slot, thus making it empty.

Dining-Philosophers Problem

The Dining Philosopher Problem states that K philosophers seated around a circular table with one chopstick between each pair of philosophers. There is one chopstick between each philosopher. A philosopher may eat if he can pick up the two chopsticks adjacent to him. One chopstick may be picked up by anyone of its adjacent followers but not both. This problem involves the allocation of limited resources to a group of processes in a deadlock-free and starvation-free manner.



The design of the problem was to illustrate the challenges of avoiding deadlock, a deadlock state of a system is a state in which no progress of system is possible. Consider a proposal where each philosopher is instructed to behave as follows:

- The philosopher is instructed to think till the left fork is available, when it is available, hold it.
- The philosopher is instructed to think till the right fork is available, when it is available, hold it.
- The philosopher is instructed to eat when both forks are available.
- then, put the right fork down first
- then, put the left fork down next
- repeat from the beginning.

The structure of Philosopher i is as follows. do

```
{  
  Wait(take_chopstick[i]);  
  Wait(take_chopstick[(i+1)%5]);  
  ...  
  EAT  
  ...  
  Signal(put_chopstick[i]);  
  Signal(put_chopstick[(i+1)%5]);  
  ...  
  THINK  
} while(TRUE);
```

In the above code, first wait operation is performed on `take_chopstick[i]` and `take_chopstick[(i+1)%5]`. This shows philosopher i has picked up the chopsticks

from its left and right. The eating function is performed after that.

On completion of eating by philosopher i the, signal operation is performed on `take_chopstick[i]` and `take_chopstick[(i+1) % 5]`. This shows that the philosopher i have eaten and put down both the left and right chopsticks. Finally, the philosopher starts thinking again.

Let value of $i = 0$ (initial value), Suppose Philosopher P_0 wants to eat, it will enter in `Philosopher()` function, and execute `Wait(take_chopstick[i])`; by doing this it holds **C0 chopstick** and reduces semaphore C_0 to 0, after that it execute `Wait(take_chopstick[(i+1) % 5])`; by doing this it holds **C1 chopstick** (since $i = 0$, therefore $(0 + 1) \% 5 = 1$) and reduces semaphore C_1 to 0.

Similarly, suppose now Philosopher P_1 wants to eat, it will enter in `Philosopher()` function, and execute `Wait(take_chopstick[i])`; by doing this it will try to hold **C1 chopstick** but will not be able to do that, since the value of semaphore C_1 has already been set to 0 by philosopher P_0 , therefore it will enter into an infinite loop because of which philosopher P_1 will not be able to pick chopstick C_1 whereas if Philosopher P_2 wants to eat, it will enter in `Philosopher()` function, and execute `Wait(take_chopstickC[i])`; by doing this it holds **C2 chopstick** and reduces semaphore C_2 to 0, after that, it executes `Wait(take_chopstickC[(i+1) % 5])`; by doing this it holds **C3 chopstick** (since $i = 2$, therefore $(2 + 1) \% 5 = 3$) and reduces semaphore C_3 to 0.

Hence the above code is providing a solution to the dining philosopher problem, A philosopher can only eat if both immediate left and right chopsticks of the philosopher are available else philosopher needs to wait. Also at one go two independent philosophers can eat simultaneously (i.e., philosopher **P0 and P2, P1 and P3 & P2 and P4** can eat simultaneously as all are the independent processes and they are following the above constraint of dining philosopher problem)

The drawback of the above solution of the dining philosopher problem

- Not two neighbouring philosophers can eat at the same point in time.
- This solution can lead to a deadlock condition. This situation happens if all the philosophers pick their left chopstick at the same time, which leads to the condition of deadlock and none of the philosophers can eat.

To avoid deadlock, some of the solutions are as follows:

- Maximum number of philosophers on the table should not be more than four, in this case, chopstick C_4 will be available for philosopher P_3 , so P_3 will start eating and after the finish of his eating procedure, he will put down his both the chopstick C_3 and C_4 , i.e. semaphore C_3 and C_4 will now be incremented to 1. Now philosopher P_2 which was holding chopstick C_2 will also have chopstick C_3 available, hence similarly, he will put down his chopstick after eating and enable other philosophers to eat.
- A philosopher at an even position should pick the right chopstick and then the left

OPERATINGSYSTEM(23CS403)

chopstick while a philosopher at an odd position should pick the left chopstick and then the right chopstick.

- Only in case if both the chopsticks (left and right) are available at the same time, only then a philosopher should be allowed to pick their chopsticks
- All the four starting philosophers (P0, P1, P2, and P3) should pick the left chopstick and then the right chopstick, whereas the last philosopher P4 should pick the right chopstick and then the left chopstick. This will force P4 to hold his right chopstick first since the right chopstick of P4 is C0, which is already held by philosopher P0 and its value is set to 0, i.e C0 is already 0, because of which P4 will get trapped into an infinite loop and chopstick C4 remains vacant. Hence philosopher P3 has both left C3 and right C4 chopstick available, therefore it will start eating and will put down its both chopsticks once finishes and let others eat which removes the problem of deadlock.

Readers and Writers Problem

Suppose that a database is to be shared among several concurrent processes. Some of these processes may want only to read the database, whereas others may want to update (that is, to read and write) the database. We distinguish between these two types of processes by referring to the former as readers and to the latter as writers. Precisely in OS we call this situation as the readers-writers problem. Problem parameters:

- One set of data is shared among a number of processes.
- Once a writer is ready, it performs its write. Only one writer may write at a time.
- If a process is writing, no other process can read it.
- If at least one reader is reading, no other process can write.
- Readers may not write and only read.

There are four types of cases that could happen here.

Case	Process1	Process2	Allowed/Not Allowed
Case1	Writing	Writing	Not Allowed
Case2	Writing	Reading	Not Allowed
Case3	Reading	Writing	Not Allowed
Case4	Reading	Reading	Allowed

Three variables are used: **mutex, wrt, readcnt**

OPERATINGSYSTEM(23CS403)

1. Semaphore **mutex** is used to ensure mutual exclusion when **readcnt** is updated i.e. when any reader enters or exits from the critical section.
2. Semaphore **wrt** is used by both readers and writers.
3. **readcnt** tells the number of processes performing read in the critical section, initially 0 and it is integer variable.

Functions for semaphore

wait(): decrements the semaphore value.

signal(): increments the semaphore value.

Reader process

- Reader requests the entry to critical section.
- If allowed:
 - ❖ it increments the count of number of readers inside the critical section. If this reader is the first reader entering, it locks the **wrt** semaphore to restrict the entry of writers if any reader is inside.
 - ❖ It then, signals **mutex** as any other reader is allowed to enter while others are already reading.
 - ❖ After performing reading, it exits the critical section. When exiting, it checks if no more reader is inside, it signals the semaphore “**wrt**” as now, writer can enter the critical section.
- If not allowed, it keeps on waiting.

```
do
{
    wait(mutex);           // Reader wants to enter the critical section
    readcnt++;             // The number of readers has now increased by 1

    if (readcnt==1)       //there is at least one reader in the critical section wait(wrt);
                          // no writer can enter if there is even one reader

    signal(mutex);        //other readers can enter where otherer is inside

    ..... perform READING

    wait(mutex);          //a reader wants to leave
    readcnt--;

    if (readcnt == 0)     //no reader is left in the critical section,
        signal(wrt);     // writers can enter
```

```
signal(mutex); //readerleaves
```

```
}while(true);
```

Writerprocess

1. Writerrequeststheentrytocriticalsection.
2. Ifallowed i.e. wait() gives a true value, it enters and performs the write. If not allowed, it keeps on waiting.
3. It exitsthecriticalsection.

```
do
{
    wait(wrt); //writerrequestsforcriticalsection
    ...performWRITING
    signal(wrt); //leavesthecriticalsection
}while(true);
```

Thus, the semaphore „wrt, is queued on both readers and writers in a manner such that preference is given to readers if writers are also there. Thus, no reader is waiting simply because a writer has requested to enter the critical section.

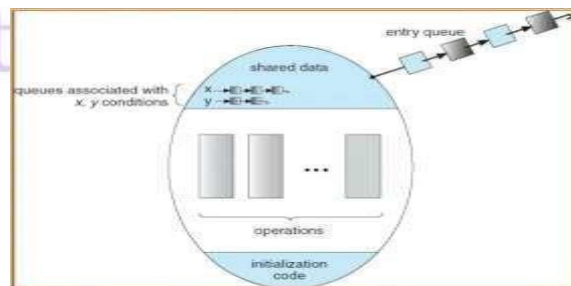
MONITOR

It is a synchronization technique that enables threads to mutual exclusion and the wait() for a given condition to become true. It is an abstract data type. It has a shared variable and a collection of procedures executing on the shared variable. A process may not directly access the shared data variables, and procedures are required to allow several processes to access the shared data variables simultaneously.

At any particular time, only one process may be active in a monitor. Other processes that require access to the shared variables must queue and are only granted access after the previous process releases the shared variables.

Syntax:

```
monitor
{
    //sharedvariabledeclarations
```




```
data variables;  
ProcedureP1(){...}  
ProcedureP2(){...}  
.  
.  
.  
Procedure Pn() { ... }  
Initialization Code() { ... }  
}
```

Advantages

- Mutualexclusionisautomaticin monitors.
- Monitorsarelessdifficulttoimplementthansemaphores.
- Monitors may overcome the timing errors that occur when semaphores areused.
- Monitorsareacollectionofproceduresandconditionvariablesthatarecombined in a special type of module.

Disadvantages

- Monitorsmustbeimplementedintotheprogramminglanguage. The
- compiler should generate code for them.
- It gives the compiler the additional burden of knowing what operating system features is available for controlling access to crucial sections in concurrent processes.



NRCM

your roots to success...

OPERATING SYSTEM (23CS403)

Comparison between the Semaphore and Monitor

Features	Semaphore	Monitor
Definition	A semaphore is an integer variable that allows many processes in a parallel system to manage access to a common resource like a multitasking OS.	It is a synchronization process that enables threads to have mutual exclusion and the wait() for a given condition to become true.
Syntax	<pre>// Wait Operation wait(Semaphore S) { while(S <= 0); S--; } // Signal Operation signal(Semaphore S) { S++; }</pre>	<pre>Monitor { // shared variable declarations Procedure P1(){...} Procedure P2(){...} . . . Procedure Pn(){...} InitializationCode(){...} }</pre>
Basic	Integer variable	Abstract data type
Access	When a process uses shared resources, it calls the wait() method on S, and when it releases them, it uses the signal() method on S.	When a process uses shared resources in the monitor, it has to access them via procedures.
Action	The semaphore's value shows the number of shared resources available in the system.	The Monitor type includes shared variables as well as a set of procedures that operate on them.
Condition Variable	No condition variables.	It has condition variables.

What is Inter-Process Communication

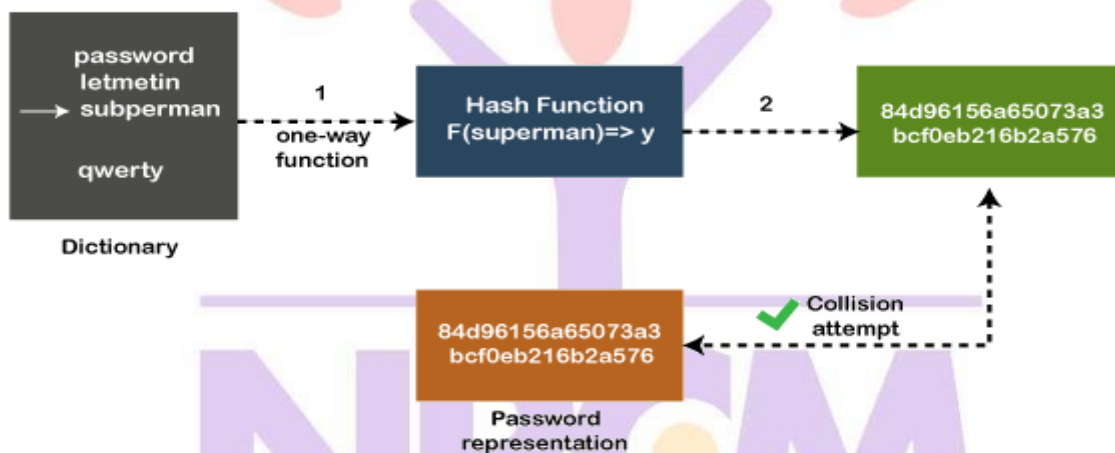
In general, Inter Process Communication is a type of mechanism usually provided by the operating system (or OS). The main aim or goal of this mechanism is to provide communications in between several processes. In short, the intercommunication allows a process letting another process know that some event has occurred.

Let us now look at the general definition of inter-process communication, which will explain the same thing that we have discussed above.

Definition

"Inter-process communication is used for exchanging useful information between numerous threads in one or more processes (or programs)."

To understand interprocess communication, you can consider the following given diagram that illustrates the importance of inter-process communication:



Role of Synchronization in Inter-Process Communication

It is one of the essential parts of inter process communication. Typically, this is provided by inter process communication control mechanisms, but sometimes it can also be controlled by communication processes.

These are the following methods that used to provide the synchronization:

1. **Mutual Exclusion**
2. **Semaphore**
3. **Barrier**
4. **Spinlock**

Mutual Exclusion:-

It is generally required that only one process thread can enter the critical section at a time. This also helps in synchronization and creates a stable state to avoid the race condition.

Semaphore:-

Semaphore is a type of variable that usually controls the access to the shared resources by several processes. Semaphore is further divided into two types which are as follows:

1. Binary Semaphore
2. Counting Semaphore

Barrier:-

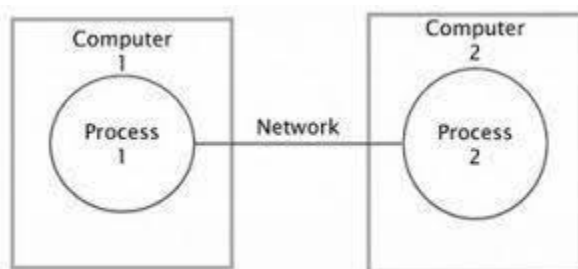
A barrier typically not allows an individual process to proceed unless all the processes does not reach it. It is used by many parallel languages, and collective routines impose barriers.

Spinlock:-

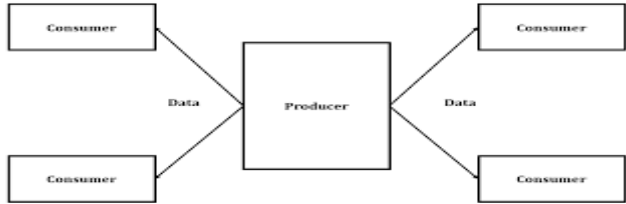
Spinlock is a type of lock as its name implies. The processes are trying to acquire the spinlock waits or stays in a loop while checking that the lock is available or not. It is known as busy waiting because even though the process active, the process does not perform any functional operation (or task).

IPC between processes on a single computer system:-

IPC refers to the mechanisms and techniques that operating systems use to facilitate communication between different processes. In a multitasking environment, numerous processes are running concurrently, and IPC serves as the bridge that allows them to exchange information and coordinate their actions.

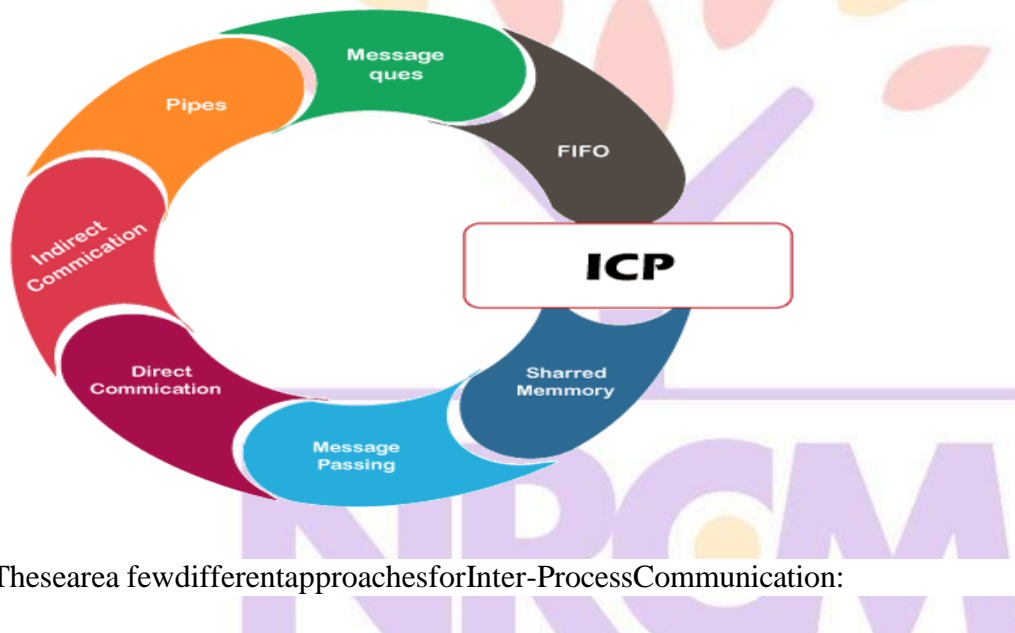


IPCbetweenprocessesondifferent system:-



Approachesto Interprocess Communication

We will now discuss some different approaches to inter-process communication which are as follows:



These are a few different approaches for Inter-Process Communication:

1. **Pipes**
2. **Shared Memory**
3. **Message Queue**
4. **Direct Communication**
5. **Indirect communication**
6. **Message Passing**
7. **FIFO**

To understand them in more detail, we will discuss each of them individually.

Pipe:-

The pipe is a type of data channel that is unidirectional in nature. It means that the data in this type of data channel can be moved in only a single direction at a time. Still, one can use two-channel of this type, so that he can able to send and receive data in two processes. Typically, it uses the standard methods for input and output. These pipes are used in all types of POSIX systems and in different versions of window operating systems as well.

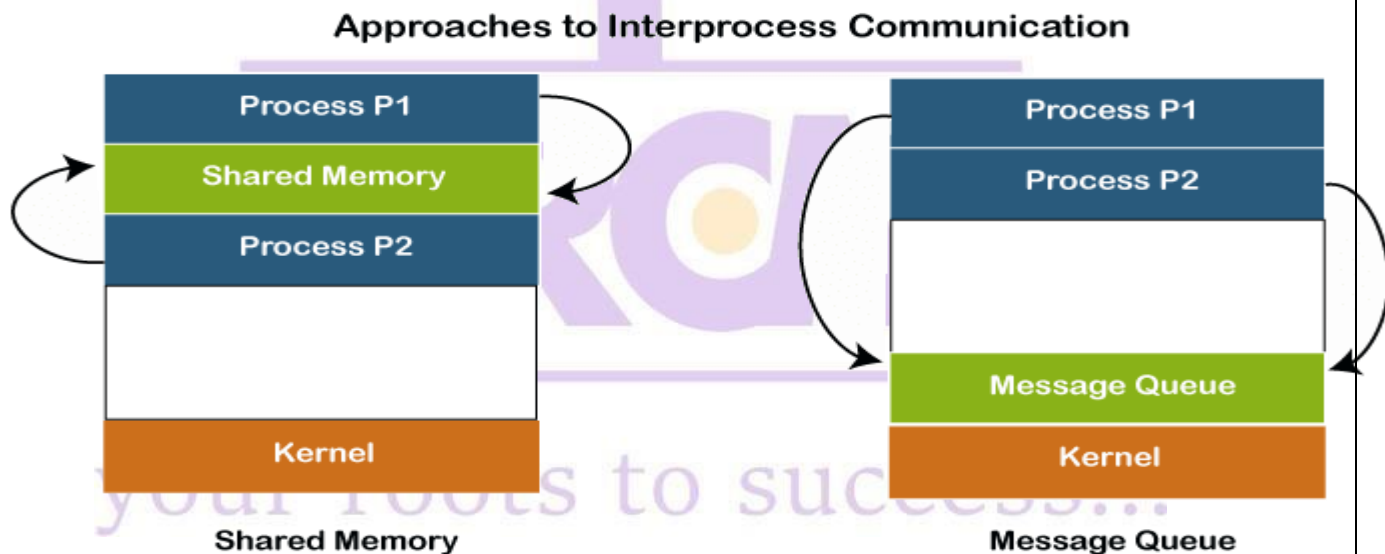
Shared Memory:-

It can be referred to as a type of memory that can be used or accessed by multiple processes simultaneously. It is primarily used so that the processes can communicate with each other. Therefore the shared memory is used by almost all POSIX and Windows operating systems as well.

Message Queue:-

In general, several different messages are allowed to read and write the data to the message queue. In the message queue, the messages are stored or stay in the queue unless their recipients retrieve them. In short, we can also say that the message queue is very helpful in inter-process communication and used by all operating systems.

To understand the concept of Message queue and Shared memory in more detail, let's take a look at its diagram given below:



Message Passing:-

It is a type of mechanism that allows processes to synchronize and communicate with each other. However, by using the message passing, the processes can communicate with each other without restoring the shared variables.

OPERATING SYSTEM(23CS403)

Usually, the inter-process communication mechanism provides two operations that are as follows:

- send(message)
- received(message)

Direct Communication:-

In this type of communication process, usually, a link is created or established between two communicating processes. However, in every pair of communicating processes, only one link can exist.

Indirect Communication

Indirect communication can only exist or be established when processes share a common mailbox, and each pair of these processes shares multiple communication links. These shared links can be unidirectional or bi-directional.

FIFO:-

It is a type of general communication between two unrelated processes. It can also be considered as full-duplex, which means that one process can communicate with another process and vice versa.

Some other different approaches

- **Socket:-**

It acts as a type of endpoint for receiving or sending the data in a network. It is correct for data sent between processes on the same computer or data sent between different computers on the same network. Hence, it is used by several types of operating systems.

- **File:-**

file server. Another most important thing is that several processes can access that file as required or needed.

- **Signal:-**

As its name implies, they are a type of signal used in inter process communication in a minimal way. Typically, they are the messages of systems that are sent by one process to another. Therefore, they are not used for sending data but for remote commands between multiple processes.

Usually, they are not used to send the data but to remote commands in between several processes.

Why we need interprocess communication?

There are numerous reasons to use inter-process communication for sharing the data. Here are some of the most important reasons that are given below:

- It helps to speed up modularity
- Computational
- Privilege separation
- Convenience
- Helps operating system to communicate with each other and synchronize their actions as well.



your roots to success...

UNIT-4

Memory Management and virtual memory-Logical versus physical address space, Swapping, Contiguous allocations, Paging, segmentation, segmentation with paging, Demand paging, Page replacement, Page Replacement algorithms.

Memory Management:-Memory is central to the operation of a modern computer system. Memory consists of a large array of bytes, each with its own address.

A typical instruction-execution cycle, for example, first fetches an instruction from memory. The instruction is then decoded and may cause operands to be fetched from memory. After the instruction has been executed on the operands, results may be stored back in memory.

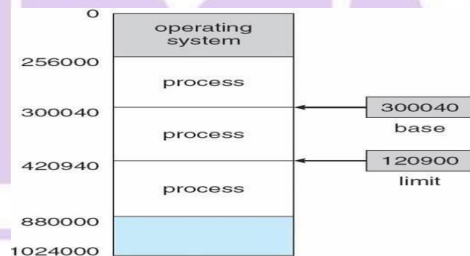
1. Basic Hardware

Main memory and the registers built into the processor itself are the only general-purpose storage that the CPU can access directly. Therefore, any instructions in execution, and any data being used by the instructions, must be in one of these direct-access storage devices. If the data are not in memory, they must be moved there before the CPU can operate on them.

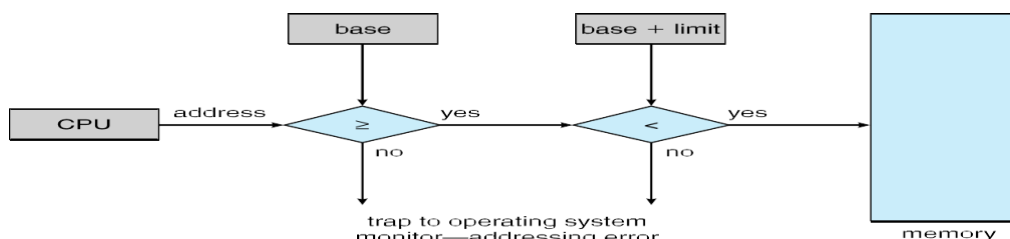
Protecting user processes from one another:

We first need to make sure that each process has a separate memory space. Separate per-process memory space protects the processes from each other and is fundamental to having multiple processes loaded in memory for concurrent execution. To separate memory spaces, we need the ability to determine the range of legal addresses that the process may access and to ensure that the process can access only these legal addresses. We can provide this protection by using two registers, usually a base and a limit

The **base register** holds the smallest legal physical memory address; the **limit register** specifies the size of the range. For example, if the base register holds 300040 and the limit register is 120900, and then the program can legally access all addresses from 300040 through 420939 (inclusive).



Any attempt by a program executing in user mode to access operating-system memory or other users' memory results in a trap to the operating system.



The base and limit registers can be loaded only by the operating system, which uses a

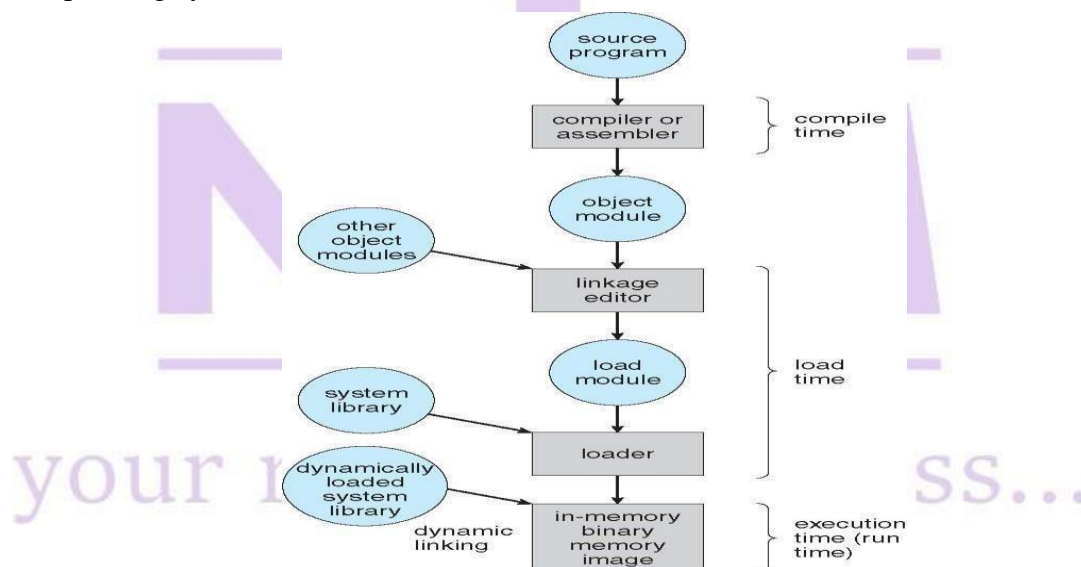
special privileged instruction. Since privileged instructions can be executed only in kernel mode, and since only the operating system executes in kernel mode, only the operating system can load the base and limit registers.

2. Address Binding

Usually, a program resides on a disk as a binary executable file. To be executed, the program must be brought into memory and placed within a process. Depending on the memory management in use, the process may be moved between disk and memory during its execution. The processes on the disk that are waiting to be brought into memory for execution form the **input queue**.

Classically, the binding of instructions and data to memory addresses can be done at any step along the way:

- **Compile time.** If you know at compile time where the process will reside in memory, then **absolute code** can be generated. For example, if you know that a user process will reside starting at location R , then the generated compiler code will start at that location and extend up from there. If, at some later time, the starting location changes, then it will be necessary to recompile this code. The MS-DOS .COM-format programs are bound at compile time.
- **Load time.** If it is not known at compile time where the process will reside in memory, then the compiler must generate **relocatable code**. In this case, final binding is delayed until load time. If the starting addresses change, we need only reload the user code to incorporate this changed value.
- **Execution time.** If the process can be moved during its execution from one memory segment to another, then binding must be delayed until run time. Most general-purpose operating systems use this method.



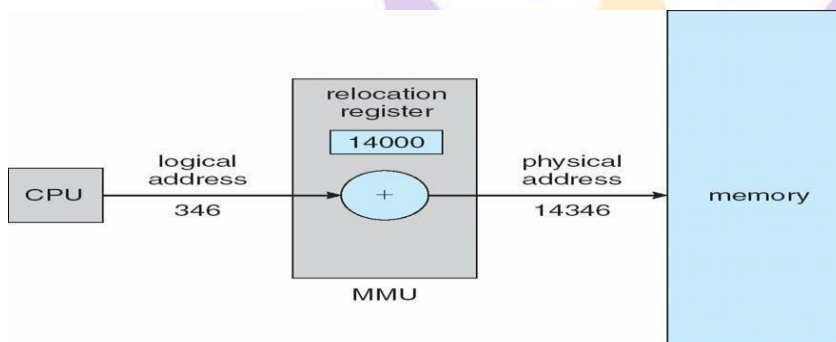
3. Logical Versus Physical Address Space

- An address generated by the CPU is commonly referred to as **logical address** or **virtual address**.
- An address seen by the memory unit—that is, the one loaded into the **memory-address register** of the memory—is commonly referred to as a **physical address**.
- These set of all logical addresses generated by a program is a **logical address space**.

- This set of all physical addresses corresponding to these logical addresses is a **physical address space**.

Memory-Management Unit (MMU)

- The run-time mapping from virtual to physical addresses is done by a hardware device called the **memory-management unit (MMU)**.
- The base register is now called a **relocation register**. The value in the relocation register is added to every address generated by a user process at the time the address is sent to memory.
- For example, if the base is at 14000, then an attempt by the user to address location 0 is dynamically relocated to location 14000; an access to location 346 is mapped to location 14346.



What is Swapping?

A process must be in memory to be executed. A process, however, can be **swapped** temporarily out of memory to a **backing store** and then brought back into memory for continued execution. Swapping makes it possible for the total physical address space of all processes to exceed the real physical memory of the system, thus increasing the degree of multiprogramming in a system.

1. Standard Swapping

Standard swapping involves moving processes between main memory and a backing store. The backing store is commonly a fast disk. It must be large enough to accommodate copies of all memory images for all users, and it must provide direct access to these memory images.

Ready Queue: The system maintains a **ready queue** consisting of all processes whose memory images are on the backing store or in memory and are ready to run.

Dispatcher: Whenever the CPU scheduler decides to execute a process, it calls the dispatcher. The dispatcher checks to see whether the next process in the queue is in memory. If it is not, and if there is no free memory region, the dispatcher swaps out a process currently in memory and swaps in the desired process. It then reloads registers and transfers control to the selected process.

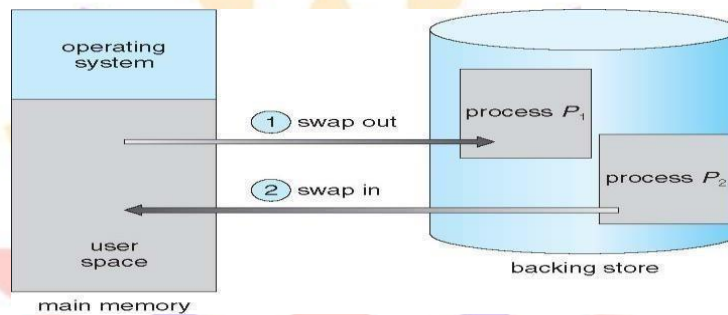
Factors

- The context-switch time in such a swapping system is fairly high.
- The total transfer time is directly proportional to the amount of memory swapped.

- If we want to swap a process, we must be sure that it is completely idle.

Standard swapping in modern operating systems

- Standard swapping is not used in modern operating systems. It requires too much swapping time and provides too little execution time to be a reasonable memory-management solution.
- Modified versions of swapping, however, are found on many systems, including UNIX, Linux, and Windows.
- In one common variation, swapping is normally disabled but will start if the amount of free memory falls below a threshold amount. Swapping is halted when the amount of free memory increases.
- Another variation involves swapping portions of processes—rather than entire processes—to decrease swap time.



2. Swapping on Mobile Systems

Mobile systems typically do not support swapping in any form.

Reasons

- Mobile devices generally use flash memory rather than hard disks. The resulting space constraints avoid swapping.
- The limited number of writes that flash memory can tolerate before it becomes unreliable.
- The poor throughput between main memory and flash memory in these devices.

Mechanisms instead of Swapping

- **Apple's iOS** asks applications to voluntarily relinquish allocated memory. Any applications that fail to free up sufficient memory may be terminated by the operating system.
- **Android** does not support swapping and adopts a strategy similar to that used by iOS. It may terminate a process if insufficient free memory is available. However, before terminating a process, Android writes its **application state** to flash memory so that it can be quickly restarted.

Contiguous Memory Allocation

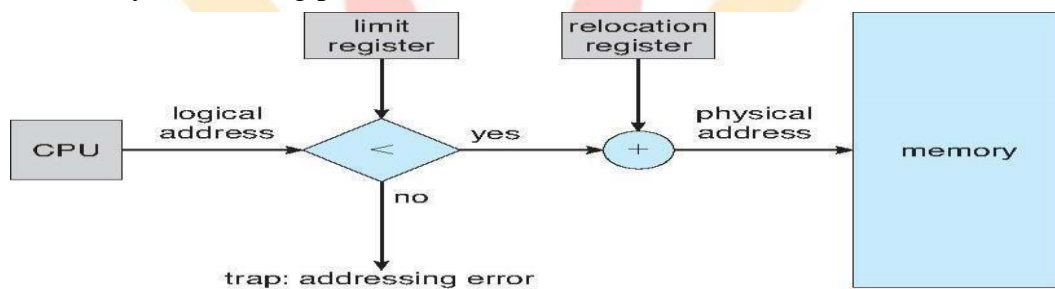
We usually want several user processes to reside in memory at the same time. We therefore need to consider how to allocate available memory to the processes that are in the input queue waiting to be brought into memory. In **contiguous memory allocation**, each process is contained in a single section of memory that is contiguous to the section containing the next process.

1. Memory Protection

We can prevent a process from accessing memory it does not own by combining two ideas. If we have a system with a relocation register, together with a limit register, we accomplish our goal.

Process

- The relocation register contains the value of the smallest physical address; the limit register contains the range of logical addresses (for example, relocation = 100040 and limit = 74600).
- Each logical address must fall within the range specified by the limit register.
- The MMU maps the logical address dynamically by adding the value in the relocation register. This mapped address is sent to memory.
- When the CPU scheduler selects a process for execution, the dispatcher loads the relocation and limit registers with the correct values as part of the context switch.
- Because every address generated by a CPU is checked against these registers, we can protect both the operating system and the other users programs and data from being modified by this running process.



2. Memory allocation methods for memory allocation

a. Fixed-Sized Partitions

- One of the simplest methods for allocating memory is to divide memory into several fixed-sized **partitions**.
- Each partition may contain exactly one process.
- In this **multiple partition method**, when a partition is free, a process is selected from the input queue and is loaded into the free partition.
- When the process terminates, the partition becomes available for another process.
- This method was originally used by the IBM OS/360 operating system (called MFT) but is no longer in use.

b. Variable Sized -Partition

- In the **variable-partition** scheme, the operating system keeps a table indicating which parts of memory are available and which are occupied.
- Initially, all memory is available for user processes and is considered one large block of available memory, a **hole**.
- When a process arrives and needs memory, the system searches the set for a hole that is large enough for this process.

- If the hole is too large, it is split into two parts. One part is allocated to the arriving process; the other is returned to the set of holes.
- When a process terminates, it releases its block of memory, which is then placed back in the set of holes.
- If the new hole is adjacent to other holes, these adjacent holes are merged to form one larger hole.
- At this point, the system may need to check whether there are processes waiting for memory and whether this newly freed and recombined memory could satisfy the demands of any of these waiting processes.

Dynamic Storage Allocation Problem (Memory Allocation Techniques)

This concerns how to satisfy a request of size n from a list of free holes. There are many solutions to this problem. The **first-fit**, **best-fit**, and **worst-fit** strategies are the ones most commonly used to select a free hole from the set of available holes.

- **First fit.** Allocate the first hole that is big enough. Searching can start either at the beginning of the set of holes or at the location where the previous first-fit search ended. We can stop searching as soon as we find a free hole that is large enough.
- **Best fit.** Allocate the smallest hole that is big enough. We must search the entire list, unless the list is ordered by size. This strategy produces the smallest leftover hole.
- **Worst fit.** Allocate the largest hole. Again, we must search the entire list, unless it is sorted by size. This strategy produces the largest leftover hole, which may be more useful than the smaller leftover hole from a best-fit approach.

Comparison:

- First fit and Best fit are better than Worst fit in terms of decreasing time and storage utilization.
- Neither first fit nor Best fit is clearly better than the other in terms of storage utilization, but First fit is generally faster.

3. Fragmentation

Memory fragmentation can be internal as well as external.

a. Internal Fragmentation

- The overhead to keep track of this hole will be substantially larger than the hole itself. The general approach to avoiding this problem is to break the physical memory into fixed-sized blocks and allocate memory in units based on block size.
- With this approach, the memory allocated to a process may be slightly larger than the requested memory.
- The difference between these two numbers is **internal fragmentation**—unused memory that is internal to a partition.

b. External Fragmentation

- Both the first-fit and best-fit strategies for memory allocation suffer from **external fragmentation**. As processes are reloaded and removed from memory, the free memory space is broken into little pieces.

- External fragmentation exists when there is enough total memory space to satisfy a request but the available spaces are not contiguous: storage is fragmented into a large number of small holes.

50-percent rule: Depending on the total amount of memory storage and the average process size, external fragmentation may be a minor or a major problem. Statistical analysis of firstfit, for instance, reveals that, even with some optimization, given N allocated blocks, another $0.5 N$ blocks will be lost to fragmentation. That is, one-third of memory may be unusable! This property is known as the **50-percent rule**.

Solution to External Fragmentation

a. Compaction

- The goal is to shuffle the memory contents so as to place all free memory together in one large block.
- Compaction is not always possible, however. If relocation is static and is done at assembly or load time, compaction cannot be done. It is possible only if relocation is dynamic and is done at execution time.
- The simplest compaction algorithm is to move all processes toward one end of memory; all holes move in the other direction, producing one large hole of available memory. This scheme can be expensive.

b. Noncontiguous logical address space

- This permits the logical address space of the processes to be noncontiguous, thus allowing a process to be allocated physical memory wherever such memory is available.
- Two complementary techniques achieve this solution: segmentation and paging.

Segmentation

Dealing with memory in terms of its physical properties is inconvenient to both the operating system and the programmer. What if the hardware could provide a memory mechanism that mapped the programmer's view to the actual physical memory? The system would have more freedom to manage memory, while the programmer would have a more natural programming environment. Segmentation provides such a mechanism.

1. Basic Method

Segmentation is a memory-management scheme that supports the programmer view of memory. A logical address space is a collection of variable sized segments. Each segment has a name and a length. The addresses specify both the segment name and the offset within the segment. The programmer therefore specifies each address by two quantities: a segment name and an offset.

segments are numbered and are referred to by a segment number, rather than by a segment name. Thus, a logical address consists of a *two tuple*:

$\langle \text{segment-number, offset} \rangle$.

Example of Segments

When a program is compiled, the compiler automatically constructs segments reflecting the input program.

A C compiler might create separate segments for the following:

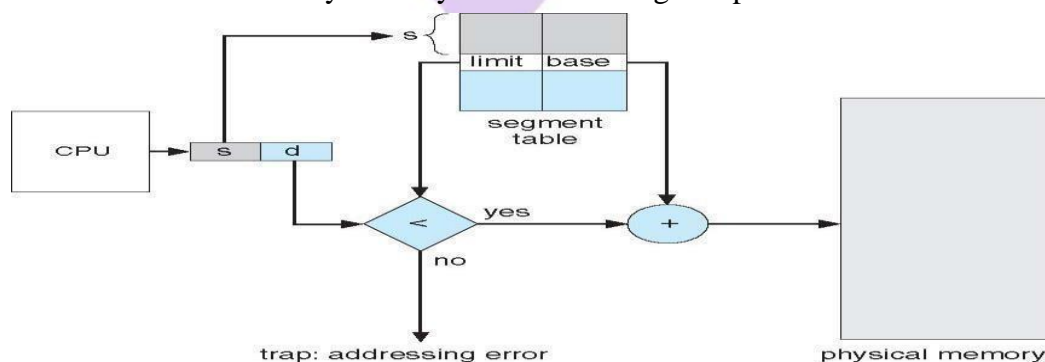
1. The code
2. Global variables
3. The heap, from which memory is allocated
4. The stacks used by each thread
5. The standard C library

2. Segmentation Hardware

Although the programmer can now refer to objects in the program by a two-dimensional address, the actual physical memory is still, of course, a one-dimensional sequence of bytes. Thus, we must define an implementation to map two-dimensional user-defined addresses into one-dimensional physical addresses. This mapping is affected by a **segment table**. Each entry in the segment table has a **segment base** and a **segment limit**.

- **Segment base:** The segment base contains the starting physical address where the segment resides in memory.
- **Segment limit:** These segment limits specify the length of the segment.

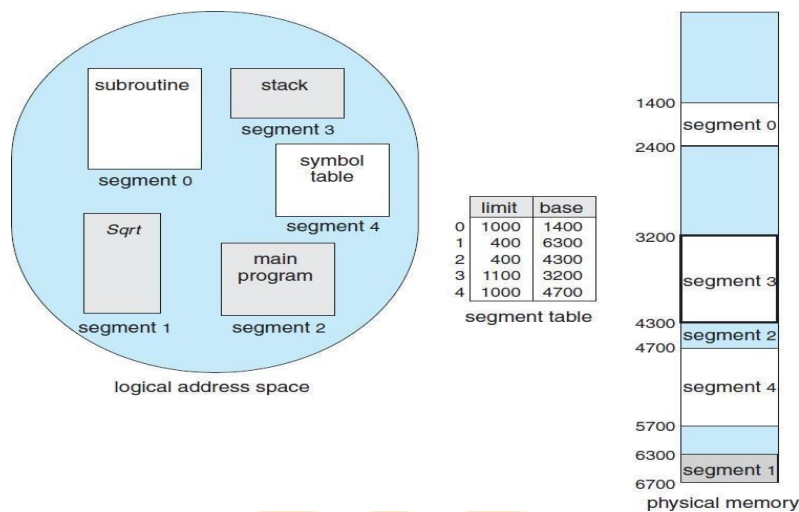
A logical address consists of two parts: a segment number, s , and an offset into that segment, d . The segment number are used as an index to the segment table. The offset d of the logical address must be between 0 and the segment limit. If it is not, we trap to the operating system (logical addressing attempt beyond end of segment). When an offset is legal, it is added to the segment base to produce the address in physical memory of the desired byte. The segment table is thus essentially an array of base–limit register pairs.



Example:

We have five segments numbered from 0 through 4. The segments are stored in physical memory. The segment table has a separate entry for each segment, giving the beginning address of the segment in physical memory (or base) and the length of that segment (or limit).

Consider, segment 2 is 400 bytes long and begins at location 4300. Thus, a reference to byte 53 of segment 2 is mapped onto location $4300 + 53 = 4353$. A reference to segment 3, byte 852, is mapped to 3200 (the base of segment 3) + 852 = 4052. A reference to byte 1222 of segment 0 would result in a trap to the operating system, as this segment is only 1,000 bytes long.



Paging

Paging is another memory-management scheme that offers physical address space of a process to be non-contiguous. Paging also avoids external fragmentation and the need for compaction, whereas segmentation does not. Because of its advantages, paging in its various forms is used in most operating systems, from mainframes to smart phones.

1. Basic Method

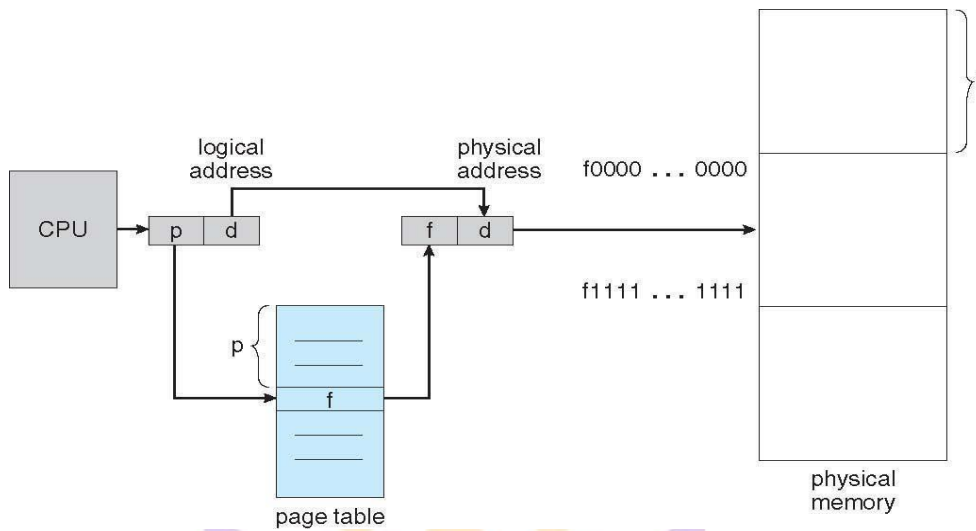
- **Frames:** Paging involves breaking physical memory into fixed-sized blocks called **frames**.
- **Pages:** Breaking logical memory into blocks of the same size called **pages**.

When a process is to be executed, its pages are loaded into any available memory frames from their source (a file system or the backing store).

Hardware Support for Paging

Every address generated by the CPU is divided into two parts: a **page number (p)** and a **page offset (d)**.

- **Page Table:** The page number is used as an index into a **page table**. The page table contains the base address of each page in physical memory. This base address is combined with the page offset to define the physical memory address that is sent to the memory unit.
- **Frame Table:** Since the operating system is managing physical memory, it must be aware of the allocation details of physical memory—which frames are allocated, which frames are available, how many total frames there are, and so on? This information is generally kept in a data structure called a **frame table**. The frame table has one entry for each physical page frame, indicating whether the latter is free or allocated and, if it is allocated, to which page of which process or processes.



Defining of Page Size

The page size (like the frame size) is defined by the hardware. The size of a page is a power of 2, varying between 512 bytes and 1 GB per page, depending on the computer architecture. The selection of a power of 2 as a page size makes the translation of a logical address into a page number and page offset particularly easy.

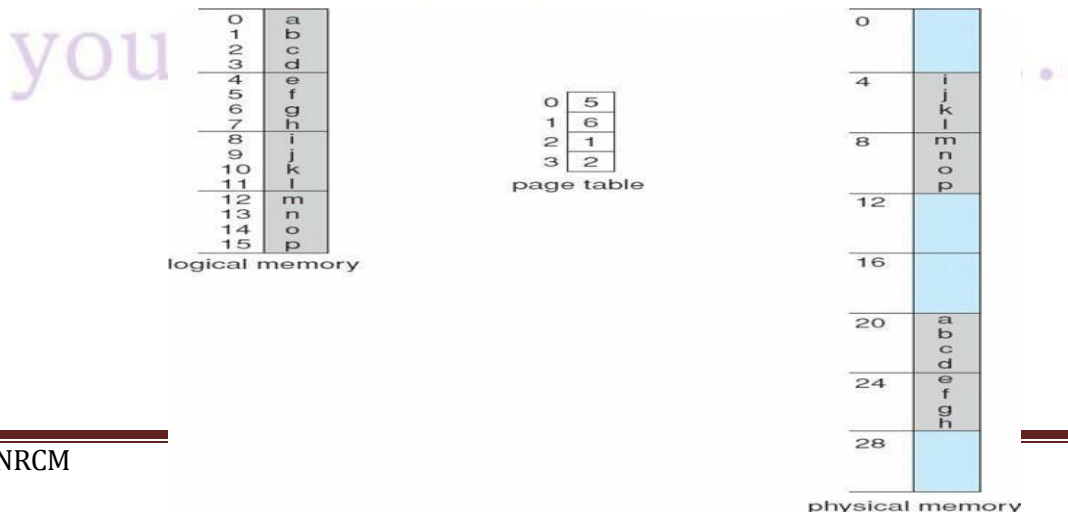
If the size of the logical address space is 2^m , and a page size is 2^n bytes, then the high-order $m - n$ bits of a logical address designate the page number, and the n low-order bits designate the page offset. Thus, the logical address is as follows:



where p is an index into the page table and d is the displacement within the page.

Example

Here, in the logical address, $n = 2$ and $m = 4$. Using a page size of 4 bytes and a physical memory of 32 bytes (8 pages). Logical address 0 is page 0, offset 0. Indexing into the page table, we find that page 0 is in frame 5. Thus, logical address 0 maps to physical address 20 $[(5 \times 4) + 0]$. Logical address 3 (page 0, offset 3) maps to physical address 23 $[(5 \times 4) + 3]$. Logical address 4 is page 1, offset 0; according to the page table, page 1 is mapped to frame 6. Thus, logical address 4 maps to physical address 24 $[(6 \times 4) + 0]$. Logical address 13 maps to physical address 9.



2. Hardware Support

Methods for storing page table: Each operating system has its own methods for storing page tables.

- a) Some allocate a page table for each process. A pointer to the page table is stored with the other register values (like the instruction counter) in the process control block. When the dispatcher is told to start a process, it must reload the user registers and define the correct hardware page-table values from the stored user page table.
- b) Other operating systems provide one or at most a few page tables, which decreases the overhead involved when processes are context-switched.

Hardware Implementation of the Page Table

Registers

- In the simplest case, the page table is implemented as a set of dedicated **registers**. These registers should be built with very high-speed logic to make the paging-address translation efficient.
- Every access to memory must go through the paging map, so efficiency is a major consideration.
- The CPU dispatcher reloads these registers, just as it reloads the other registers. Instructions to load or modify the page-table registers are, of course, privileged, so that only the operating system can change the memory map. The use of registers for the page table is satisfactory if the page table is reasonably small.
- The DECPDP-11 is an example.

Page-Table Base Register (PTBR)

Most contemporary computers, allow the page table to be very large (for example, 1 million entries). For these machines, the use of fast registers to implement the page table is not feasible. Rather, the page table is kept in main memory, and a **page-table base register (PTBR)** points to the page table. Changing page tables requires changing only this one register, substantially reducing context-switch time.

Problem

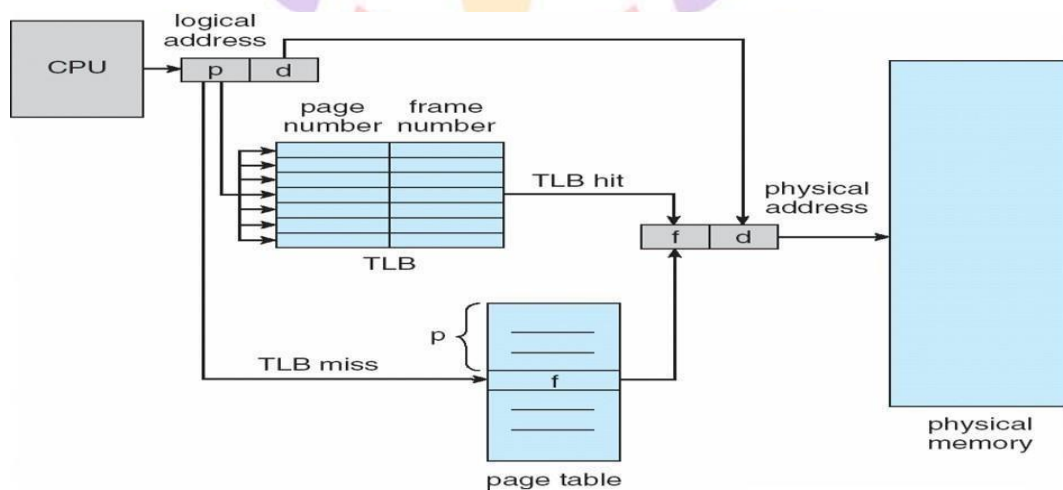
The problem with this approach is the time required to access a user memory location. If we want to access location i , we must first index into the page table, using the value in the PTBR offset by the page number for i . This task requires a memory access. It provides us with the frame number, which is combined with the page offset to produce the actual address. We can then access the desired place in memory. With this scheme, *two* memory accesses are needed to access a byte (one for the page-table entry, one for the byte). Thus, memory access is slowed by a factor of 2.

Solution: Translation Look-Aside Buffer (TLB).

The standard solution to this problem is to use a special, small, fast lookup hardware cache called a **translation look-aside buffer (TLB)**. The TLB is associative, high-speed memory. Each entry in the TLB consists of two parts: a key (or tag) and a value.

Working of translation look-aside buffer (TLB):

- The TLB is used with page tables in the following way. The TLB contains only a few of the page-table entries.
- When a logical address is generated by the CPU, its page number is presented to the TLB. If the page number is found, its frame number is immediately available and is used to access memory.
- If the page number is not in the TLB (known as a **TLB miss**), a memory reference to the page table must be made.
- Depending on the CPU, this may be done automatically in hardware or via an interrupt to the operating system.
- When the frame number is obtained, we can use it to access memory. In addition, we add the page number and frame number to the TLB, so that they will be found quickly on the next reference. If the TLB is already full of entries, an existing entry must be selected for replacement.

**Address-Space Identifiers (ASIDs)**

- Some TLBs store **address-space identifiers (ASIDs)** in each TLB entry. An ASID uniquely identifies each process and is used to provide address-space protection for that process.
- When the TLB attempts to resolve virtual page numbers, it ensures that the ASID for the currently running process matches the ASID associated with the virtual page. If the ASIDs do not match, the attempt is treated as a TLB miss.

Hit ratio: The percentage of times that the page number of interest is found in the TLB is called the **hit ratio**

3. Protection

Memory protection in a paged environment is accomplished by protection bits associated with each frame. Normally, these bits are kept in the page table.

Read-Write or Read-Only Bit

- One bit can define a page to be read-write or read-only. Every reference to memory goes through the page table to find the correct frame number.
- At the same time that the physical address is being computed, the protection bits can be checked to verify that no writes are being made to a read-only page.

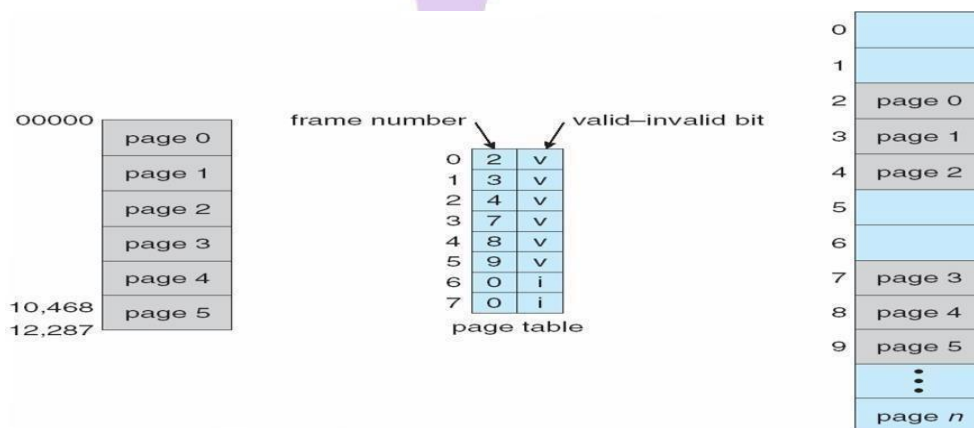
- An attempt to write to a read-only page causes a hardware trap to the operating system (or memory-protection violation).
- We can easily expand this approach to provide a finer level of protection.
- We can create hardware to provide read-only, read-write, or execute-only protection; or, by providing separate protection bits for each kind of access, we can allow any combination of these accesses. Illegal attempts will be trapped to the operating system.

Valid-Invalid Bit

- One additional bit is generally attached to each entry in the page table: a **valid-invalid** bit.
- When this bit is set to *valid*, the associated page is in the process's logical address space and is thus a legal (or valid) page.
- When the bit is set to *invalid*, the page is not in the process's logical address space. Illegal addresses are trapped by use of the valid-invalid bit.
- The operating system sets this bit for each page to allow or disallow access to the page.

Example

- Suppose, for example, that in a system with a 14-bit address space (0 to 16383), we have a program that should use only addresses 0 to 10468.
- Given a page size of 2 KB, addresses in pages 0, 1, 2, 3, 4, and 5 are mapped normally through the page table.
- Any attempt to generate an address in pages 6 or 7, however, will find that the valid-invalid bit is set to invalid, and the computer will trap to the operating system (invalid page reference).



Hardware for Protection: Page-Table Length Register (PTLR)

Some systems provide hardware, in the form of a **page-table length register (PTLR)**, to indicate the size of the page table. This value is checked against every logical address to verify that the address is in the valid range for the process. Failure of this test causes an error trap to the operating system.

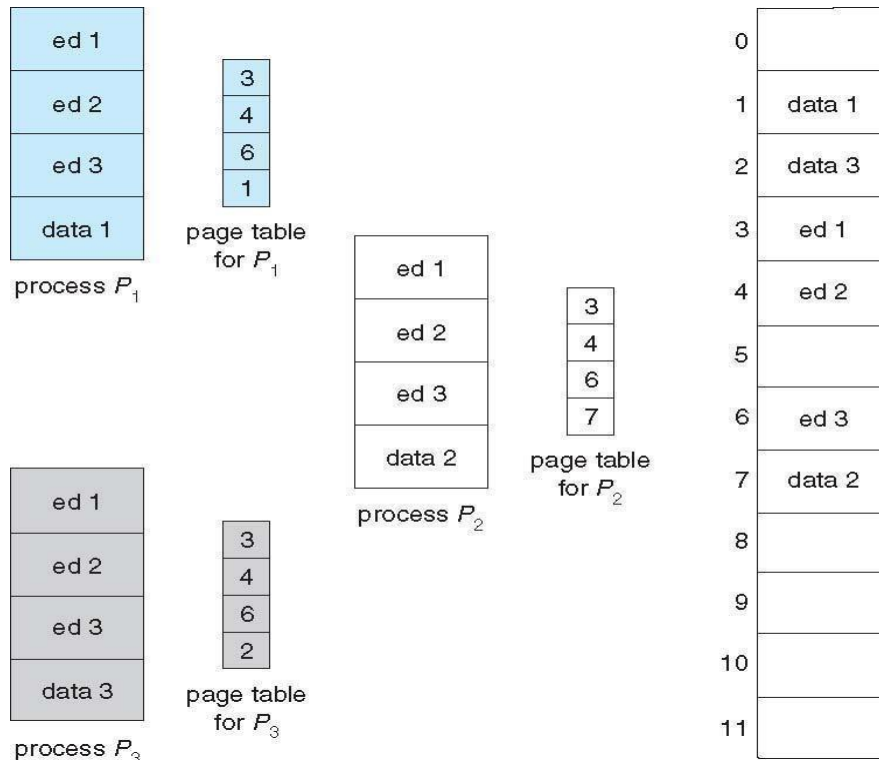
4. Shared Pages

An advantage of paging is the possibility of *sharing* common code. This consideration is particularly important in a time-sharing environment.

Example:

This

Consider a system that supports 40 users, each of whom executes a text editor. If the text editor consists of 150 KB of code and 50 KB of data space, we need 8,000 KB to support the 40 users. If the code is **reentrant code** or **pure code** (Reentrant code is non-self-modifying code: it never changes during execution. Thus, two or more processes can execute the same code at the same time. However, it can be shared.



Each process has its own copy of registers and data storage to hold the data for the process's execution. The data for two different processes will, of course, be different. Only one copy of the editor need be kept in physical memory. Each user's page table maps onto the same physical copy of the editor, but data pages are mapped onto different frames. Thus, to support 40 users, we need only one copy of the editor (150 KB), plus 40 copies of the 50 KB of data space per user. The total space required is now 2,150 KB instead of 8,000 KB—a significant savings.

Other heavily used programs can also be shared—compilers, window systems, run-time libraries, database systems, and so on.

Segmentation with Paging

Pure segmentation is not very popular and not being used in many of the operating systems. However, Segmentation can be combined with Paging to get the best features out of both the techniques.

In Segmented Paging, the main memory is divided into variable size segments which are further divided into fixed size pages. Pages are smaller than segments. Each Segment has a page table which means every program has multiple page tables.

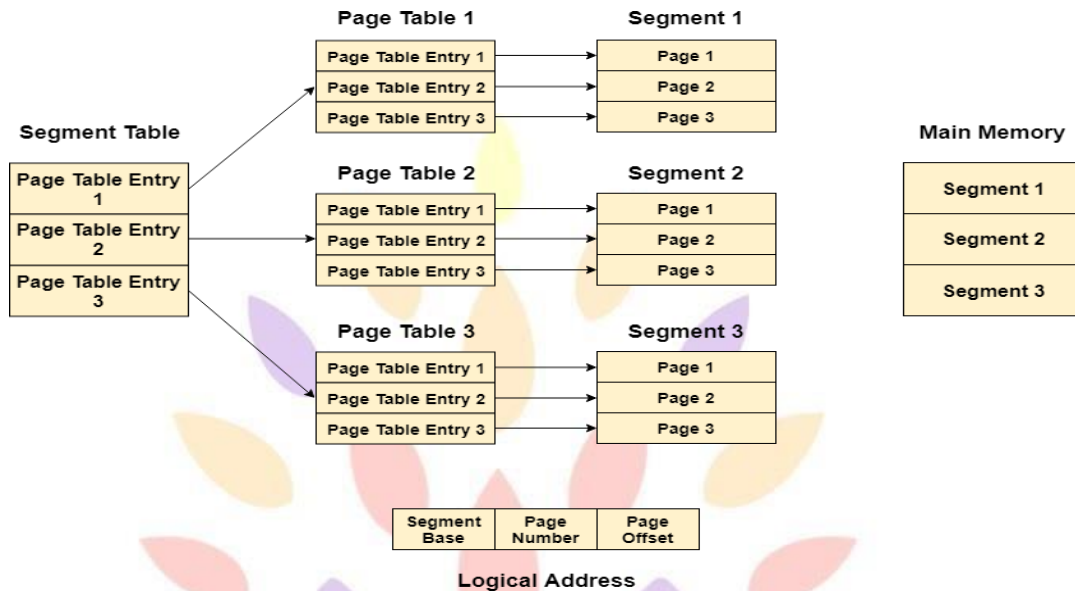
The logical address is represented as Segment Number (base address), Page number and page offset.

Segment Number → It points to the appropriate Segment Number. Page

Number → It Points to the exact page within the segment

Page Offset → Used as an offset within the page frame

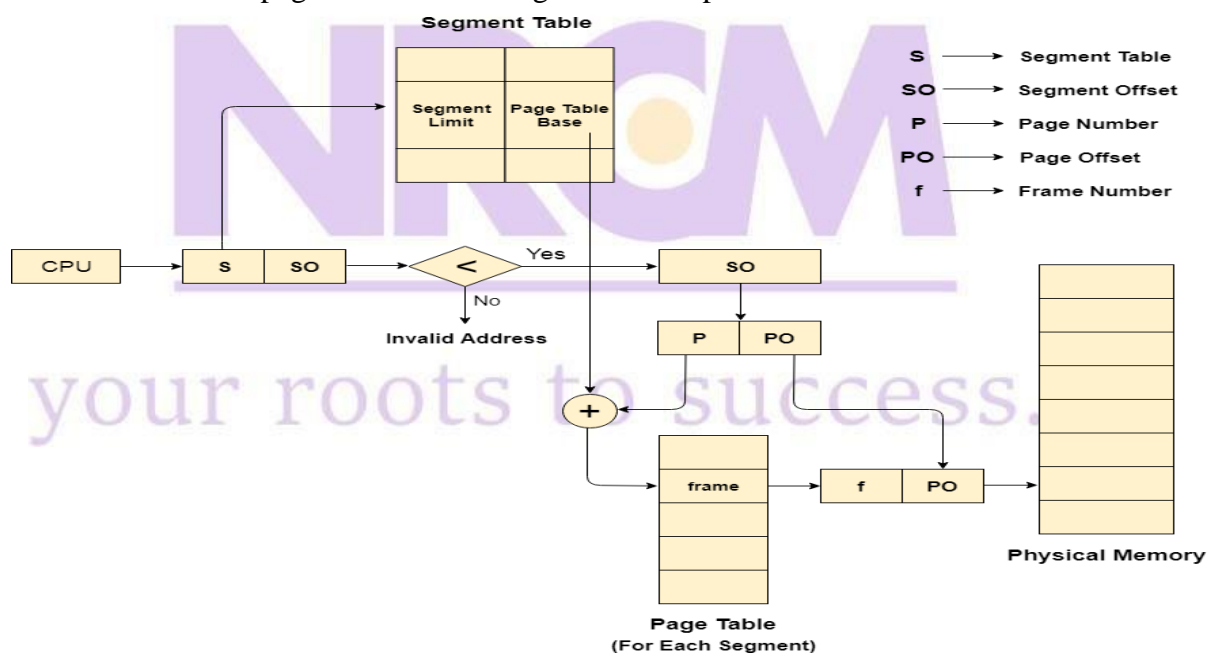
Each Page table contains the various information about every page of the segment. The Segment Table contains the information about every segment. Each segment table entry points to a page table entry and every page table entry is mapped to one of the page within a segment.



Translation of logical address to physical address

The CPU generates a logical address which is divided into two parts: Segment Number and Segment Offset. The Segment Offset must be less than the segment limit. Offset is further divided into Page number and Page Offset. To map the exact page number in the page table, the page number is added into the page table base.

The actual frame number with the page offset is mapped to the main memory to get the desired word in the page of the certain segment of the process.



Advantages of Segmented Paging

1. It reduces memory usage.
2. Page table size is limited by the segment size.
3. Segment table has only one entry corresponding to one actual segment.
4. External fragmentation is not there.
5. It simplifies memory allocation.

Disadvantages of Segmented Paging

1. Internal fragmentation will be there.
2. The complexity level will be much higher as compared to paging.
3. Page tables need to be contiguously stored in the memory.

Demand Paging

Definition



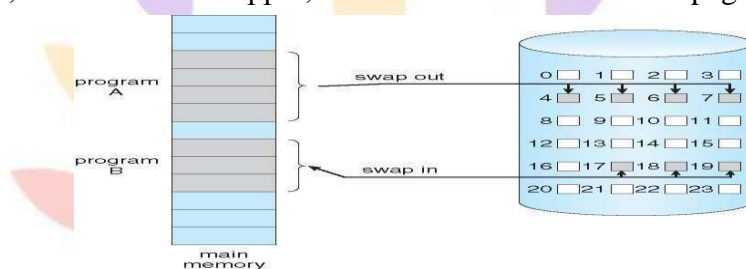
your roots to success...

Loading the entire program into memory results in loading the executable code for *all* options, regardless of whether or not an option is ultimately selected by the user. An alternative strategy is to load pages only as they are needed. This technique is known as **demand paging** and is commonly used in virtual memory systems.

With demand-paged virtual memory, pages are loaded only when they are demanded during program execution. Pages that are never accessed are thus never loaded into physical memory.

Lazy Swapper

A demand-paging system is similar to a paging system with swapping where processes reside in secondary memory (usually a disk). When we want to execute a process, we swap it into memory. Rather than swapping the entire process into memory, though, we use a **lazy swapper**. A lazy swapper never swaps a page into memory unless that page will be needed. In the context of a demand-paging system, use of the term “swapper” is technically incorrect. We thus use “pager,” rather than “swapper,” in connection with demand paging.



Transfer of paged memory to contiguous disk space.

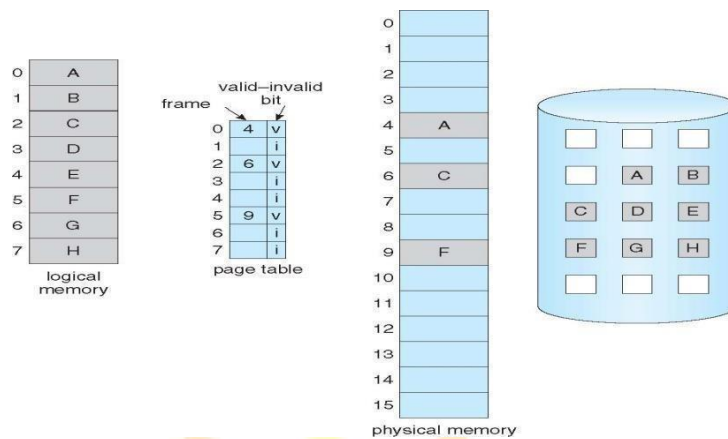
1. Basic Concepts

When a process is to be swapped in, the pager guesses which pages will be used before the process is swapped out again. Instead of swapping in a whole process, the pager brings only those pages into memory. Thus, it avoids reading into memory pages that will not be used anyway, decreasing the swap time and the amount of physical memory needed.

Valid-Invalid Bit

- We need some form of hardware support to distinguish between the pages that are in memory and the pages that are on the disk. When this bit is set to “valid,” the associated page is both legal and in memory.
- If the bit is set to “invalid,” the page is either not valid (that is, not in the logical address space of the process) or is valid but is currently on the disk. The page-table entry for a page that is brought into memory is set as usual, but the page-table entry for a page that is not currently in memory is either simply marked invalid or contains the address of the page on disk.

your roots to success...

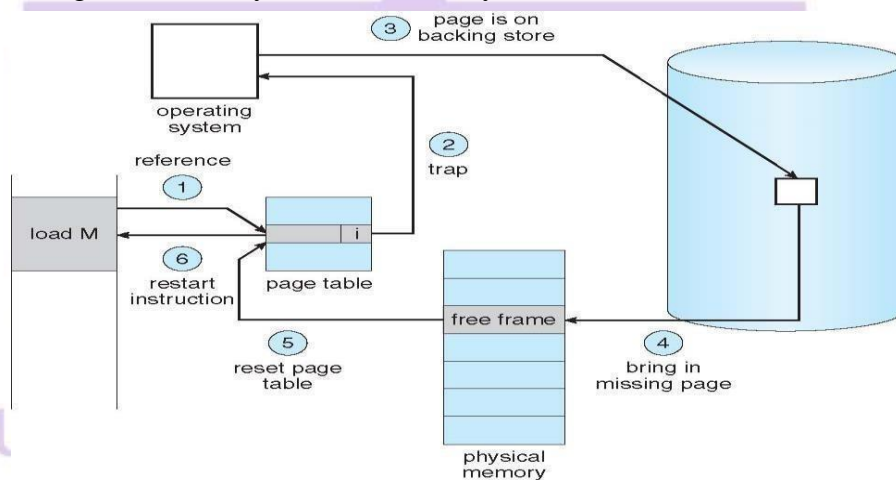


Page Fault

Access to a page marked invalid causes a **page fault**. The paging hardware, in translating the address through the page table, will notice that the invalid bit is set, causing a trap to the operating system. This trap is the result of the operating system's failure to bring the desired page into memory.

The procedure for handling this page fault is straightforward

1. We check an internal table (usually kept with the process control block) for this process to determine whether the reference was a valid or an invalid memory access.
2. If the reference was invalid, we terminate the process. If it was valid but we have not yet brought in that page, we now page it in.
3. We find a free frame (by taking one from the free-frame list, for example).
4. We schedule a disk operation to read the desired page into the newly allocated frame.
5. When the disk read is complete, we modify the internal table kept with the process and the page table to indicate that the page is now in memory.
6. We restart the instruction that was interrupted by the trap. The process can now access the page as though it had always been in memory.



Pure Demand Paging

In the extreme case, we can start executing a process with **no** pages in memory. When the operating system sets the instruction pointer to the first instruction of the process, which is on a non-memory-resident page, the process immediately faults for the page. After this page is brought into memory, the process continues to execute, faulting as necessary until

every page that it needs is in memory. At that point, it can execute with no more faults. This scheme is **pure demand paging**: never bring a page into memory until it is required.

Hardware to Support Demand Paging

- **Page table.** This table has the ability to mark an entry invalid through a valid–invalid bit or a special value of protection bits.
- **Secondary memory.** This memory holds those pages that are not present in main memory. The secondary memory is usually a high-speed disk. It is known as the swap device, and the section of disk used for this purpose is known as **swap space**.

A crucial requirement for demand paging is the ability to restart any instruction after a page fault. Because we save the state (registers, condition code, and instruction counter) of the interrupted process when the page fault occurs, we must be able to restart the process in *exactly* the same place and state, except that the desired page is now in memory and is accessible. In most cases, this requirement is easy to meet.

A page fault may occur at any memory reference. If the page fault occurs on the instruction fetch, we can restart by fetching the instruction again. If a page fault occurs while we are fetching an operand, we must fetch and decode the instruction again and then fetch the operand.

2. Performance of Demand Paging

Demand paging can significantly affect the performance of a computer system. Let's compute the **effective access time** for a demand-paged memory. For most computer systems, the memory-access time, denoted ma , ranges from 10 to 200 nanoseconds. As long as we have no page faults, the effective access time is equal to the memory access time. If, however, a page fault occurs, we must first read the relevant page from disk and then access the desired word.

Let p be the probability of a page fault ($0 \leq p \leq 1$). We would expect p to be close to zero—that is, we would expect to have only a few page faults.

The effective access time is then

$$\text{Effective Access Time} = (1-p) \times ma + p \times \text{page fault time}$$

To compute the effective access time, we must know how much time is needed to service a page fault.

Example:

With an average page-fault service time of 8 milliseconds and a memory access time of 200 nanoseconds, the effective access time in nanoseconds is

$$\begin{aligned} \text{Effective Access Time} &= (1-p) \times (200) + p(8 \text{ milliseconds}) \\ &= (1-p) \times 200 + p \times 8,000,000 \\ &= 200 + 7,999,800 \times p. \end{aligned}$$

We see, then, that the effective access time is directly proportional to the **page-fault rate**.

An additional aspect of demand paging is the handling and overall use of swap space. Disk I/O to swap space is generally faster than that to the file system. It is a faster file system because swap space is allocated in much larger blocks, and file lookups and indirect allocation methods are not used. However, swap space must still be used for pages not associated with a file (known as **anonymous memory**).

Mobile operating systems typically do not support swapping. Instead, these systems demand-page from the file system and reclaim read-only pages (such as code) from applications if memory becomes constrained. Such data can be demand-paged from the file system if it is later needed. Under iOS, anonymous memory pages are never reclaimed from an application unless the application is terminated or explicitly releases the memory.

Page Replacement

In Demand Paging, pages are only brought into memory only when needed. This has two benefits,

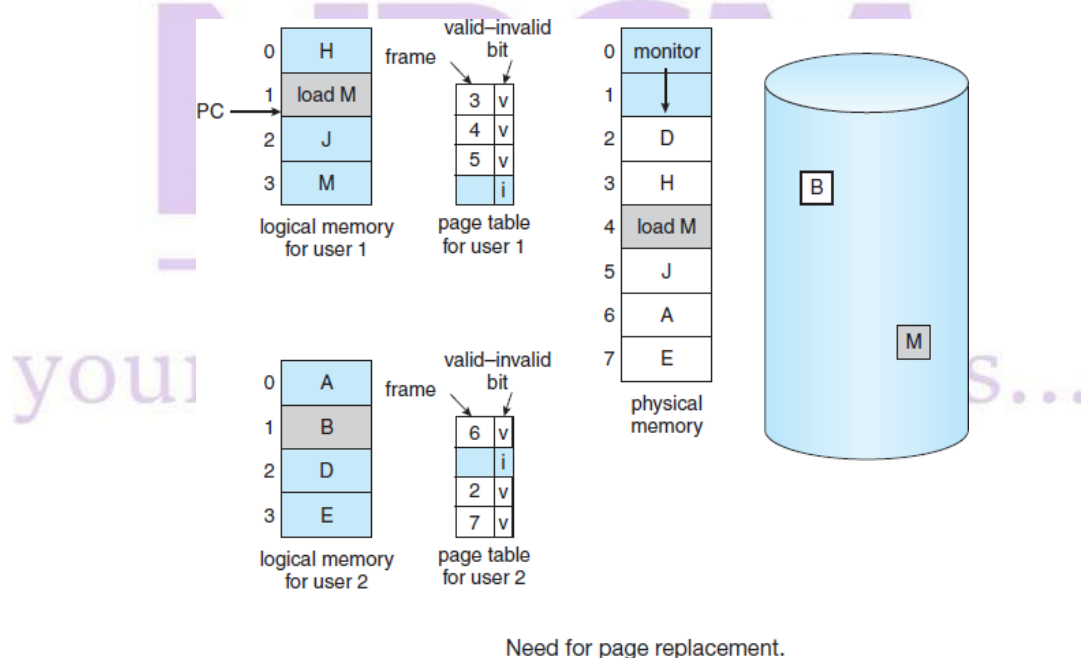
1. Saves I/O necessary to load unused pages.
2. Increases the degree of multiprogramming.

But increasing degree of multiprogramming may arise a new problem called “Over allocating of memory”.

Over-Allocating Memory

For example, there are 10 processes and each has 10 pages out of which only 5 may be used. If there are 50 frames then we can allocate only 5 processes if all the 10 pages are loaded. But by using demand paging (we load only used or demanded pages) we can accommodate 10 processes as only 5 pages are in demand. Problem arises when suddenly a process needs all 10 pages but no frames are free.

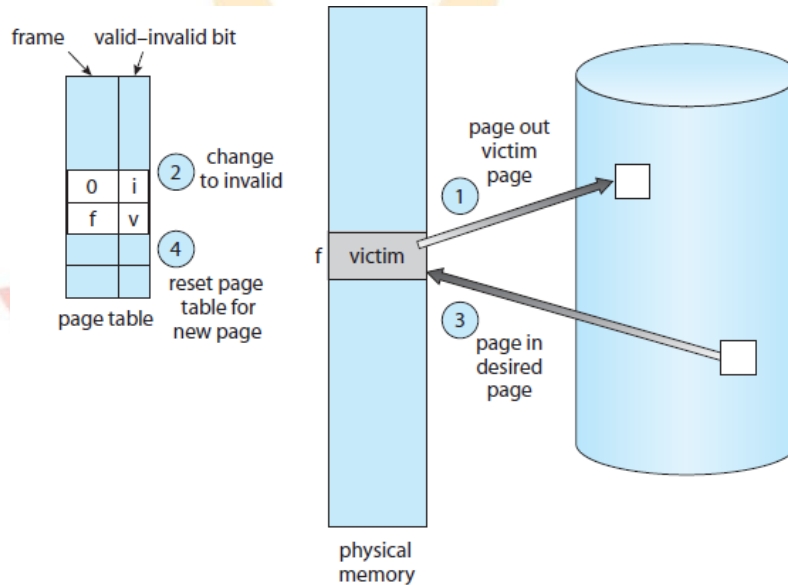
Over-allocation of memory manifests itself as follows. While a user process is executing, a page fault occurs. The operating system determines where the desired page is residing on the disk but then finds that there are *no* free frames on the free-frame list; all memory is in use. The operating system has several options at this point. It could terminate the user process. This option is not the best choice. The operating system could instead swap out a process, freeing all its frames and reducing the level of multiprogramming. This option is a good one but requires page replacement.



1. Basic Page Replacement

Page replacement takes the following approach,

1. Find the location of the desired page on the disk.
2. Find a free frame:
 - a. If there is a free frame, use it.
 - b. If there is no free frame, use a page-replacement algorithm to select a victim frame.
- c. Write the victim frame to the disk; change the page and frame tables accordingly.
3. Read the desired page into the newly freed frame; change the page and frame tables.
 - a. Write the victim frame to the disk; change the page and frame tables accordingly.
4. Continue the user process from where the page fault occurred.



Page replacement.

Modify Bit (or Dirty Bit).

- If no frames are free, *two* page transfers (one out and one in) are required. This situation effectively doubles the page-fault service time and increases the effective access time accordingly. We can reduce this overhead by using a **modify bit** (or **dirty bit**).
- When this scheme is used, each page or frame has a modify bit associated with it in the hardware. The modify bit for a page is set by the hardware whenever any byte in the page is written into, indicating that the page has been modified.
- When we select a page for replacement, we examine its modify bit. If the bit is set, we know that the page has been modified since it was read in from the disk. In this case, we must write the page to the disk. If the modify bit is not set, however, the page has *not* been modified since it was read into memory. In this case, we need not write the memory page to the disk: it is already there.

Major Problems to Implement Demand Paging

We must solve two major problems to implement demand paging: we must develop a **frame-allocation algorithm** and a **page-replacement algorithm**.

That is, if we have multiple processes in memory, we must decide how many frames to allocate to each process; and when page replacement is required, we must select the frames that are to be replaced.

Reference String

There are many different page-replacement algorithms. We evaluate an algorithm by running it on a particular string of memory references and computing the number of page faults. The string of memory references is called a **reference string**.

We can generate reference strings

- Artificially (by using a random-number generator, for example).
- We can trace a given system and record the address of each memory reference. But this produces large amount of data.

To reduce this, we use two facts

a. First, for a given page size (and the page size is generally fixed by the hardware or system), we need to consider only the page number, rather than the entire address.

b. Second, if we have a reference to a page *p*, then any reference to page *p* that **immediately** follow will never cause a page fault.

Example

If we trace a particular process, we might record the following address sequence: 0100, 0432, 0101, 0612, 0102, 0103, 0104, 0101, 0611, 0102, 0103, 0104, 0101, 0610, 0102, 0103, 0104, 0101, 0609, 0102, 0105

At 100 bytes per page, this sequence is reduced to the following reference string:

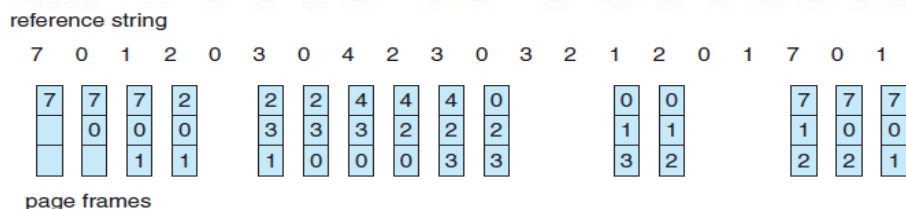
1, 4, 1, 6, 1, 6, 1, 6, 1, 6, 1

Page Replacement Algorithms

FIFO Page Replacement

- The simplest page-replacement algorithm is a first-in, first-out (FIFO) algorithm.
- A FIFO replacement algorithm associates with each page the time when that page was brought into memory. When a page must be replaced, the oldest page is chosen.
- We can create a FIFO queue to hold all pages in memory. We replace the page at the head of the queue. When a page is brought into memory, we insert it at the tail of the queue.
- The FIFO page-replacement algorithm is easy to understand and program.
- However, its performance is not always good. A bad replacement choice increases the page-fault rate and slows process execution. If we place an active page, some other page should be replaced to bring it back.

Example:



FIFO page-replacement algorithm.

Belady's anomaly: For some page-replacement algorithms, the page-fault rate may *increase* as the number of allocated frames increases.

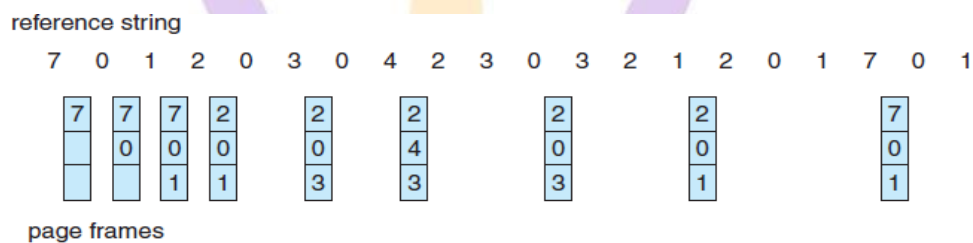
Consider the following reference string: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5

Number of faults for four frames (ten) is *greater* than the number of faults for three frames (nine)

Optimal Page Replacement

- It says that, Replace the page that will not be used for the longest period of time.
- It has the lowest page-fault rate of all algorithms and will never suffer from Belady's anomaly.
- Unfortunately, the optimal page-replacement algorithm is difficult to implement, because it requires future knowledge of the reference string.
- As a result, the optimal algorithm is used mainly for comparison studies.

Example:

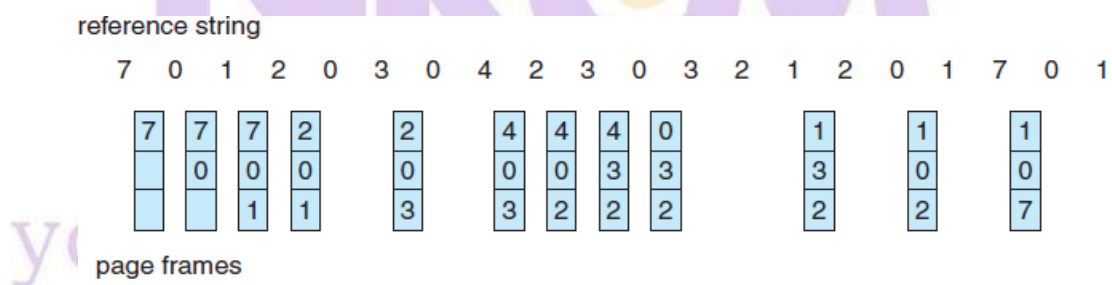


Optimal page-replacement algorithm.

LRU Page Replacement

- LRU replacement associates with each page the time of that page's last use.
- When a page must be replaced, LRU chooses the page that has not been used for the longest period of time.
- We can think of this strategy as the optimal page-replacement algorithm looking backward in time, rather than forward.
- Like optimal replacement, LRU replacement does not suffer from Belady's anomaly. Both belong to a class of page-replacement algorithms, called **stack algorithms**.

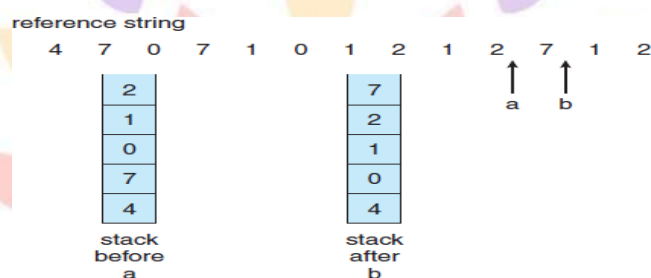
Example:



LRU page-replacement algorithm.

- The major problem is *how* to implement LRU replacement. An LRU page-replacement algorithm may require substantial hardware assistance. The problem is to determine an order for the frames defined by the time of last use.
- Two implementations are feasible:

- **Counters.** In the simplest case, we associate with each page-table entry a time-of-use field and add to the CPU a logical clock or counter. The clock is incremented for every memory reference. Whenever a reference to a page is made, the contents of the clock register are copied to the time-of-use field in the page-table entry for that page. In this way, we always have the “time” of the last reference to each page. We replace the page with the smallest time value. This scheme requires a search of the page table to find the LRU page and a write to memory (to the time-of-use field in the page table) for each memory access.
- **Stack.** Another approach to implementing LRU replacement is to keep a stack of page numbers. Whenever a page is referenced, it is removed from the stack and put on the top. In this way, the most recently used page is always at the top of the stack and the least recently used page is always at the bottom. Because entries must be removed from the middle of the stack, it is best to implement this approach by using a doubly linked list with a head pointer and a tail pointer.



Use of a stack to record the most recent page references

LRU-Approximation Page Replacement

- Few computer systems provide sufficient hardware support for true LRU page replacement. In fact, some systems provide no hardware support, and other page-replacement algorithms (such as a FIFO algorithm) must be used. Many systems provide some help, however, in the form of a **reference bit**.
- The reference bit for a page is set by the hardware whenever that page is referenced (either a read or a write to any byte in the page). Reference bits are associated with each entry in the page table.
- Initially, all bits are cleared (to 0) by the operating system. As a user process executes, the bit associated with each page referenced is set (to 1) by the hardware. After some time, we can determine which pages have been used and which have not been used by examining the reference bits, although we do not know the *order* of use. This information is the basis for many page-replacement algorithms that approximate LRU replacement.

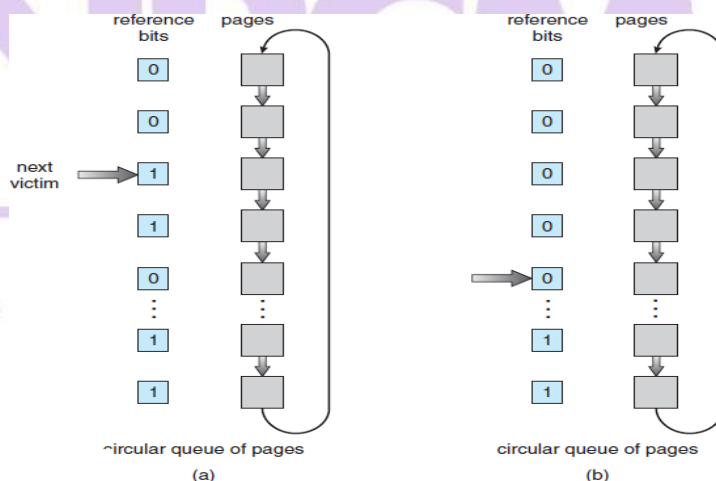
Additional-Reference-Bits Algorithm

- We can gain additional ordering information by recording the reference bits at regular intervals.
- We can keep an 8-bit byte for each page in a table in memory.
- At regular intervals (say, every 100 milliseconds), a timer interrupt transfers control to the operating system.

- The operating system shifts the reference bit for each page into the high-order bit of its 8-bit byte, shifting the other bits right by 1 bit and discarding the low-order bit. These 8-bit shift registers contain the history of page use for the last eight time periods.
- If the shift register contains 00000000, for example, then the page has not been used for eight time periods.
- A page that is used at least once in each period has a shift register value of 11111111. A page with a history register value of 11000100 has been used more recently than one with a value of 01110111.
- If we interpret these 8-bit bytes as unsigned integers, the page with the lowest number is the LRU page, and it can be replaced. Notice that the numbers are not guaranteed to be unique, however. We can either replace (swap out) all pages with the smallest value or use the FIFO method to choose among them.

Second-Chance Algorithm OR clock algorithm

- The basic algorithm of second-chance replacement is a FIFO replacement algorithm. When a page has been selected, we inspect its reference bit.
- If the value is 0, we proceed to replace this page; but if the reference bit is set to 1, we give the page a second chance and move on to select the next FIFO page.
- When a page gets a second chance, its reference bit is cleared, and its arrival time is reset to the current time. Thus, a page that is given a second chance will not be replaced until all other pages have been replaced (or given second chances).
- In addition, if a page is used often enough to keep its reference bit set, it will never be replaced.
- One way to implement the second-chance algorithm is as a circular queue. A pointer (that is, a hand on the clock) indicates which page is to be replaced next.
- When a frame is needed, the pointer advances until it finds a page with a 0 reference bit. As it advances, it clears the reference bits. Once a victim page is found, the page is replaced, and the new page is inserted in the circular queue in that position.



Second-chance (clock) page-replacement algorithm.

Enhanced Second-Chance Algorithm

- We can enhance the second-chance algorithm by considering the reference bit and the modify bit as an ordered pair. With these two bits, we have the following four possible classes:

- **(0,0)** neither recently used nor modified—best page to replace.
- **(0, 1)** not recently used but modified—not quite as good, because the page will need to be written out before replacement.
- **(1,0)** recently used but clean—probably will be used again soon.
- **(1, 1)** recently used and modified—probably will be used again soon, and the page will need to be written out to disk before it can be replaced.

Counting-Based Page Replacement

There are many other algorithms that can be used for page replacement. For example, we can keep a counter of the number of references that have been made to each page and develop the following two schemes,

Least Frequently Used (LFU)

- The **least frequently used (LFU)** page-replacement algorithm requires that the page with the smallest count be replaced. The reason for this selection is that an actively used page should have a large reference count.
- A problem arises, however, when a page is used heavily during the initial phase of a process but then is never used again.
- Since it was used heavily, it has a large count and remains in memory even though it is no longer needed.
- One solution is to shift the counts right by 1 bit at regular intervals, forming an exponentially decaying average usage count.

Most Frequently Used (MFU)

- The **most frequently used (MFU)** page-replacement algorithm is based on the argument that the page with the smallest count was probably just brought in and has yet to be used.



NRCM

your roots to success...

UNIT-5

File system interface and operation- Access methods, directory structure, Protection, File system structure, Allocation methods, Free space management, Usage of Open, Create, Read, Write, Close, lseek, Stat, ioctl System calls

File:-

A file is a named collection of related information that is recorded on secondary storage.

(or) A file is the smallest allotment of logical secondary storage.

(or) A file is a sequence of bits, bytes, lines, or records, the meaning of which is defined by the file's creator and user. Many different types of information may be stored in a file.

File Attributes

File attributes give the Operating System information about the file and how it is intended to use.

A file's attributes vary from one operating system to another but typically consist of these:

- **Name.** The symbolic filename is the only information kept in human-readable form.
- **Identifier.** This unique tag, usually a number, identifies the file within the file system; it is the non-human-readable name for the file.
- **Type.** This information is needed for systems that support different types of files.
- **Location.** This information is a pointer to a device and to the location of the file on that device.
- **Size.** The current size of the file (in bytes, words, or blocks) and possibly the maximum allowed size are included in this attribute.
- **Protection.** Access-control information determines who can do reading, writing, executing, and so on.
- **Time, date, and user identification.** This information may be kept for creation, last modification, and last use. These data can be useful for protection, security, and usage monitoring.

Some newer file systems also support **extended file attributes**, including character encoding of the file and security features such as a file checksum.

File Types

When we design a file system we always consider whether the operating system should recognize and support file types. If an operating system recognizes the type of a file, it can then operate on the file in reasonable ways. A common technique for implementing file types is to include the type as part of the file name. The name is split into two parts—a name and an extension, usually separated by a period. Examples include `resume.docx`, `server.c`, and `ReaderThread.cpp`.

file type	usual extension	function
executable	exe, com, bin or none	ready-to-run machine-language program
object	obj, o	compiled, machine language, not linked
source code	c, cc, java, perl, asm	source code in various languages
batch	bat, sh	commands to the command interpreter
markup	xml, html, tex	textual data, documents
word processor	xml, rtf, docx	various word-processor formats
library	lib, a, so, dll	libraries of routines for programmers
print or view	gif, pdf, jpg	ASCII or binary file in a format for printing or viewing
archive	rar, zip, tar	related files grouped into one file, sometimes compressed, for archiving or storage
multimedia	mpeg, mov, mp3, mp4, avi	binary file containing audio or A/V information

Common file types.

File Structure

File types also can be used to indicate the internal structure of the file. Some operating systems extend this idea by supporting their own file structures. But it has the following disadvantages

1. If operating system support multiple file structures: the resulting size of the operating system is large.
2. Some applications may require information structured in a way that is not supported by the OS some operating systems impose (and support) a minimal number of file structures. This approach has been adopted in UNIX, Windows, and others.

Internal File Structure

Block Structure

Disk systems typically have a well-defined block size determined by the size of a sector. All disk I/O is performed in units of one block (physical record), and all blocks are the same size.

Record Structure

Files contain a sequence of fixed length records. Physical records may or may not get exact match with the logical record. Logical records even vary in length.

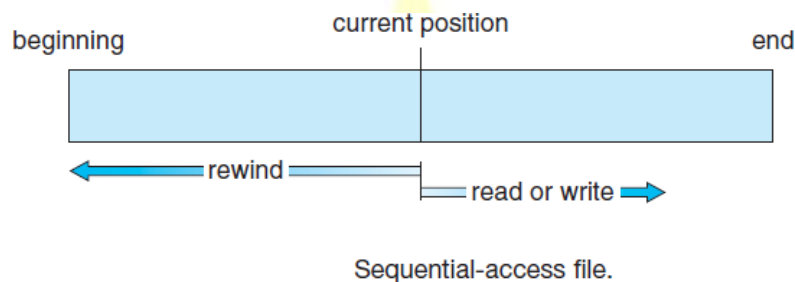
Access methods

Files store information. When it is used, this information must be accessed and read into computer memory. The information in the file can be accessed in the following ways,

1. Sequential Access

- The simplest access method is **sequential access**. Information in the file is processed in order, one record after the other. It is based on a tape model of a file and works as well on sequential-access devices.

- **Example:** Editors and Compilers usually access files in this fashion.
- **Operations**
 - A read operation—`read next ()`—reads the next portion of the file and automatically advances a file pointer, which tracks the I/O location. Similarly, the write operation—`write next ()`—appends to the end of the file and advances to the end of the newly written material (the new end of file). On some systems, a program may be able to skip forward or backward n records for some integer n —perhaps only for $n = 1$.



2. Direct Access (or Relative Access)

- Another method is **direct access** (or **relative access**). Here, a file is made up of fixed-length **logical records** that allow programs to read and write records rapidly in no particular order. The direct-access method is based on a disk model of a file, since disks allow random access to any file block.
- For direct access, the file is viewed as a numbered sequence of blocks or records. Thus, we may read block 14, then read block 53, and then write block 7. There are no restrictions on the order of reading or writing for a direct-access file.
- **Examples:**
 - Direct-access files are of great use for immediate access to large amounts of information. Databases are often of this type. When a query concerning a particular subject arrives, we compute which block contains the answer and then read that block directly to provide the desired information.
 - On an airline-reservation system, we might store all the information about a particular flight (for example, flight 713) in the block identified by the flight number. Thus, the number of available seats for flight 713 is stored in block 713 of the reservation file. To store information about a larger set, such as people, we might compute a hash function on the people's names or search a small in-memory index to determine a block to read and search.
- **Operations**
 - For the direct-access method, the file operations must be modified to include the block number as a parameter. Thus, we have `read (n)`, where n is the block number, rather than `read next ()`, and `write (n)` rather than `write next ()`.
 - An alternative approach is to retain `read next ()` and `write next ()`, as with sequential access, and to add an operation `position file (n)` where n is the block number. Then, to affect a `read (n)`, we would `position file (n)` and then `read next()`.

3. Indexed Access

- It involves the construction of an index for the file. The **index**, like an index in the back of a book, contains pointers to the various blocks. To find a record in the file, we first search the index and then use the pointer to access the file directly and to find the desired record.
- **Example:**
A retail-price file might list the universal product codes (UPCs) for items, with the associated prices. Each record consists of a 10-digit UPC and a 6-digit price, for a 16-byte record. If our disk has 1,024 bytes per block, we can store 64 records per block. A file of 120,000 records would occupy about 2,000 blocks (2 million bytes). By keeping the file sorted by UPC, we can define an index consisting of the first UPC in each block. This index would have 2,000 entries of 10 digits each, or 20,000 bytes, and thus could be kept in memory. To find the price of a particular item, we can make a binary search of the index. From this search, we learn exactly which block contains the desired record and access that block. This structure allows us to search a large file doing little I/O.

With large files, the index file itself may become too large to be kept in memory. One solution is to create an index for the index file. The primary index file contains pointers to secondary index files, which point to the actual data items. For example, IBM's indexed sequential-access method (ISAM) uses a small master index that points to disk blocks of a secondary index.

Directory Overview

A file system can be created on each of these parts of the disk. Any entity containing a file system is generally known as a **volume**. Each volume that contains a file system must also contain information about the files in the system. This information is kept in entries in a **device directory** or **volume table of contents**. The device directory (or **directory**) records information—such as name, location, size, and type—for all files on that volume.

The directory can be viewed as a symbol table that translates file names into their directory entries. The following are the operations that are to be performed on a directory:

- **Search for a file.**
- **Create a file.**
- **Delete a file.**
- **List a directory.**
- **Rename a file.**
- **Traverse the file system.**

Directory Structure

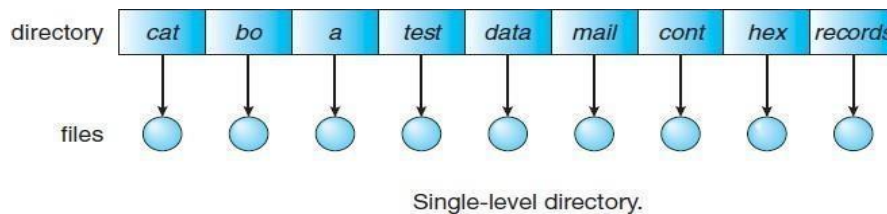
The most common schemes for defining the logical structure of a directory are the following,

1. Single-Level Directory

- The simplest directory structure is the single-level directory. All files are contained in the same directory, which is easy to support and understand.

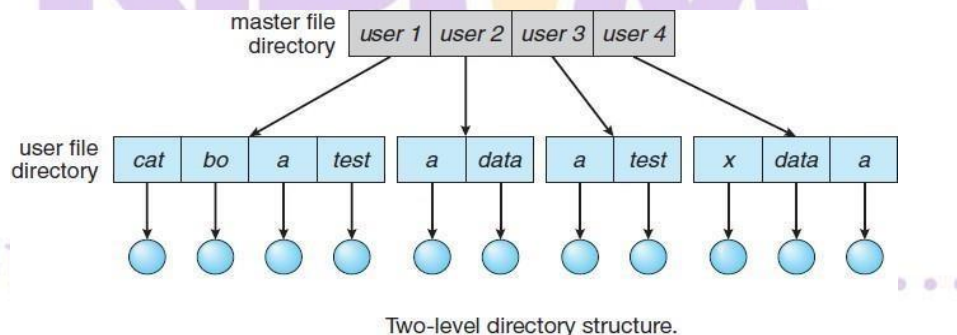
- **Limitations**

- All files are in the same directory, they must have unique names. If two users call their data file test.txt, then the unique-name rule is violated.
- Even a single user on a single-level directory may find it difficult to remember the names of all the files as the number of files increases. Keeping track of so many files is a problem.



2. Two-Level Directory

- The standard solution to eliminate confusion of file names among different users is to create a separate directory for each user.
- So the two-level directory structure contains 2 directories
 - Master File Directory (MFD) at the top level.
 - User File Directory (UFD) at the second level and
 - Actual files are at the third level.
- Each user has his own **user file directory (UFD)**. When a user job starts or a user logs in, the system's **master file directory (MFD)** is searched. The MFD is indexed by user name or account number, and each entry points to the UFD for that user
- When a user refers to a particular file, only his own UFD is searched.
- To create a file for a user, the operating system searches only that user's UFD to ascertain whether another file of that name exists.
- To delete a file, the operating system confines its search to the local UFD; thus, it cannot accidentally delete another user's file that has the same name.

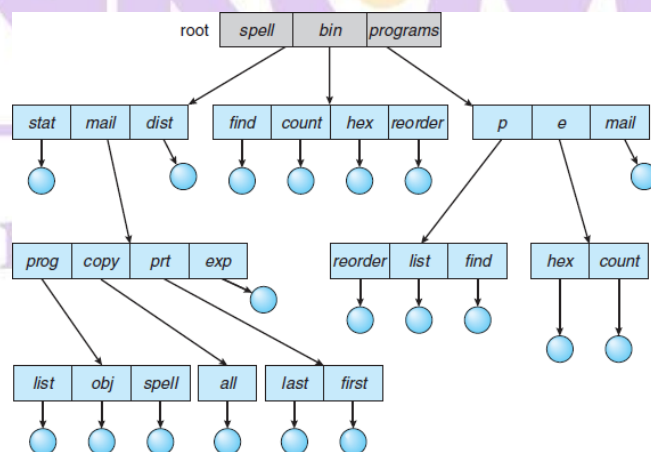


- Although the two-level directory structure solves the name-collision problem, it still has disadvantages.
- This structure effectively isolates one user from another. Isolation is an advantage when the users are completely independent but is a disadvantage when the users want to cooperate on some task and to access one another's files. Some systems simply do not allow local user files to be accessed by other users.

- If access is permitted one user must have the ability to name a file in another user's directory. To name a file uniquely, the user must give both user name and file name as Path Name.
- If user A wishes to access her own test file named test.txt, she can simply refer to test.txt. To access the file named test.txt of user B (with directory-entry name userb), however, she might have to refer to /userb/test.txt (windows) and /u/pbg/test (Unix, Linux).
- A special situation occurs with the system files. If a user wants them, they are searched in USD if found ok if not found we should copy the system files into each UFD but copying all the system files would waste an enormous amount of space.
- The standard solution is to use special user directory. Whenever a file name is given to be loaded, the operating system first searches the local UFD. If the file is found, it is used. If it is not found, the system automatically searches the special user directory that contains the system files.
- The sequence of directories searched when a file is named is called the **search path**.

3. Tree-Structured Directories

- A tree is the most common directory structure. The tree has a root directory, and every file in the system has a unique path name.
- A directory (or subdirectory) contains a set of files or subdirectories. A directory is simply another file, but it is treated in a special way. All directories have the same internal format. One bit in each directory entry defines the entry as a file (0) or as a subdirectory (1). Special system calls are used to create and delete directories.
- **Current Directory**
 - Each process has a current directory. The **current directory** should contain most of the files that are of current interest to the process.
 - When reference is made to a file, the current directory is searched. If a file is needed that is not in the current directory, then the user usually must either specify a path name or change the current directory (using change directory () system call) to be the directory holding that file.



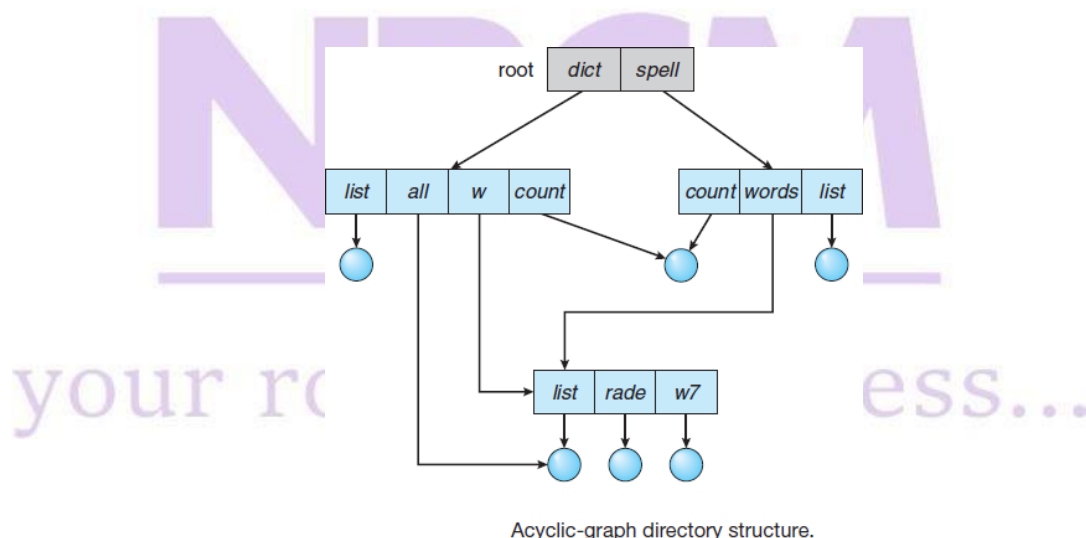
Tree-structured directory structure.

- **Path Names**

- It describes the path the OS must take to get to some point.
 - Path names can be of two types: absolute and relative.
 - An **absolute path name** begins at the root and follows a path down to the specified file, giving the directory names on the path.
 - A **relative path name** defines a path from the current directory.
 - If the current directory is root/spell/mail, then the relative path name prt/first refers to the same file as does the absolute path name root/spell/mail/prt/first.
- **Deletion of a directory**
 - If a directory is empty, its entry in the directory that contains it can simply be deleted.
 - However, suppose the directory to be deleted is not empty but contains several files or subdirectories.
 - One of two approaches can be taken. Some systems will not delete a directory unless it is empty. Thus, to delete a directory, the user must first delete all the files in that directory. If any subdirectories exist, this procedure must be applied recursively to them, so that they can be deleted also. This approach can result in a substantial amount of work.
 - An alternative approach, such as that taken by the UNIX rm command, is to provide an option: when a request is made to delete a directory, all that directory's files and subdirectories are also to be deleted.

4. Acyclic-Graph Directories

- A tree structure prohibits the sharing of files or directories. An **acyclic graph** i.e., a graph with no cycles which allows directories to share subdirectories and files.
- The same file or subdirectory may be in two different directories. An acyclic-graph directory structure is more flexible than a simple tree structure, but it is also more complex.



- **Implementation**
 - a. **Link:** A common way is to create a new directory entry called a link. A **link** is effectively a pointer to another file or subdirectory. A link may be implemented as an absolute or a relative path name. When a reference to a file is made, we search the directory. If the directory entry is marked as a link, then the name of

thereal file is included in the link information. We **resolve** the link by using that path name to locate the real file.

b. Duplication: Shared files duplicate all information about them in both sharing directories. Thus, both entries are identical and equal. A major problem with duplicate directory entries is maintaining consistency when a file is modified.

- **Problems**

- A file may now have multiple absolute path names creating problem in traversing.
- **Deletion:** When can the space allocated to a shared file be deallocated and reused?
 - One possibility is to remove the file whenever anyone deletes it, but this action may leave dangling pointers to the now-nonexistent file.
 - Another possibility occurs when symbolic links are used. The deletion of a link need not affect the original file; only the link is removed. If the file entry itself is deleted, the space for the file is deallocated, leaving the links dangling. We can search for these links and remove them as well, but unless a list of the associated links is kept with each file, this search can be expensive. Alternatively, we can leave the links until an attempt is made to use them. At that time, we can determine that the file of the name given by the link does not exist and can fail to resolve the link name; the access is treated just as with any other illegal file name.
 - Another approach to deletion is to preserve the file until all references to it are deleted. To implement this approach, we must have some mechanism for determining that the last reference to the file has been deleted. We could keep a list of all references to a file (directory entries or symbolic links). When a link or a copy of the directory entry is established, a new entry is added to the file-reference list. When a link or directory entry is deleted, we remove its entry on the list. The file is deleted when its file-reference list is empty.

Protection

When information is stored in a computer system, we want to keep it safe from physical damage (the issue of reliability) and improper access (the issue of protection).

Reliability is generally provided by duplicate copies of files. Many computers have systems programs that automatically (or through computer-operator intervention) copy disk files to tape at regular intervals (once per day or week or month) to maintain a copy should a file system be accidentally destroyed.

File systems can be damaged by hardware problems (such as errors in reading or writing), power surges or failures, head crashes, dirt, temperature extremes, and vandalism, accidental deletion, Bugs in the file-system software etc.,

Protection can be provided in many ways. For a single-user laptop system, we might provide protection by locking the computer in a desk drawer or file cabinet. In a larger multiuser system, however, other mechanisms are needed.

1. Types of Access

Protection mechanisms provide controlled access by limiting the types of file access that can be made. Access is permitted or denied depending on several factors, one of which is the type of access requested. Several different types of operations may be controlled:

- **Read.** Read from the file.
- **Write.** Write or rewrite the file.
- **Execute.** Load the file into memory and execute it.
- **Append.** Write new information at the end of the file.
- **Delete.** Delete the file and free its space for possible reuse.
- **List.** List the name and attributes of the file.

Other operations, such as renaming, copying, and editing the file, may also be controlled.

2. Access Control

The most common approach to the protection problem is to make access dependent on the identity of the user. Different users may need different types of access to a file or directory. The most general scheme to implement identity dependent access is to associate with each file and directory an **access-control list (ACL)** specifying user names and the types of access allowed for each user.

When a user requests access to a particular file, the operating system checks the access list associated with that file. If that user is listed for the requested access, the access is allowed.

Otherwise, a protection violation occurs, and the user's job is denied access to the file.

This approach has the advantage of enabling complex access methodologies. The main problem with access lists is their length. If we want to allow everyone to read a file, we must list all users with read access. This technique has two undesirable consequences:

- Constructing such a list may be a tedious and unrewarding task, especially if we do not know in advance the list of users in the system.
- The directory entry, previously of fixed size, now must be of variable size, resulting in more complicated space management.

These problems can be resolved by use of a condensed version of the access list. To condense the length of the access-control list, many systems recognize three classifications of users in connection with each file:

- **Owner.** The user who created the file is the owner.
- **Group.** A set of users who are sharing the file and need similar access is a group, or work group.
- **Universe.** All other users in the system constitute the universe.

The most common recent approach is to combine access-control lists with the more general (and easier to implement) owner, group, and universe access control scheme.

For this scheme to work properly, permissions and access lists must be controlled tightly. This control can be accomplished in several ways. For example, in the UNIX system, groups can be created and modified only by the manager of the facility (or by any superuser). Thus, control is achieved through human interaction.

With the more limited protection classification, only three fields are needed to define protection. Often, each field is a collection of bits, and each bit either allows or prevents the access associated with it. For example, the UNIX system defines three fields of 3 bits each—`rwx`, where `r` controls read access, `w` controls write access, and `x` controls execution. A separate field is kept for the file owner, for the file's group, and for all other users. In this scheme, 9 bits per file are needed to record protection information.

Example:

`19-rw-r--r--+1jimstaff130May2522:13 file1`

3. Other Protection Approaches

Another approach to the protection problem is to associate a password with each file. Just as access to the computer system is often controlled by a password, access to each file can be controlled in the same way. If the passwords are chosen randomly and changed often, this scheme may be effective in limiting access to a file. The use of passwords has a few disadvantages, however.

First, the number of passwords that a user needs to remember may become large, making the scheme impractical.

Second, if only one password is used for all the files, then once it is discovered, all files are accessible; protection is on an all-or-none basis. Some systems allow a user to associate a password with a subdirectory, rather than with an individual file, to address this problem.

File System Structure

Disks provide most of the secondary storage on which file systems are maintained. Two characteristics make them convenient for this purpose are,

1. A disk can be rewritten.
2. A disk can access directly any block of information it contains.

File systems provide efficient and convenient access to the disk by allowing data to be stored, located, and retrieved easily.

Design Problems in a File System

1. How the file system should look to the user, i.e. file, and the directory structure for organizing files.
2. Creating algorithms and data structures to map the logical file system onto the physical secondary-storage devices.

Layered design of file systems

Each level in the design uses the features of lower levels to create new features for use by higher levels.

Application Programs

It contains user code that is making a request.

Logical File System

The **logical file system** manages metadata information. Metadata includes all of the file-system structure except the actual data. The logical file system manages the directory

structure to provide the file-organization module with this information.

File-Organization Module

The **file-organization module** knows about files and their logical blocks and physical blocks. By knowing the type of file allocation used and the location of the file, the file organization module can translate logical block addresses to physical block addresses for the basic file system to transfer.

Basic File System

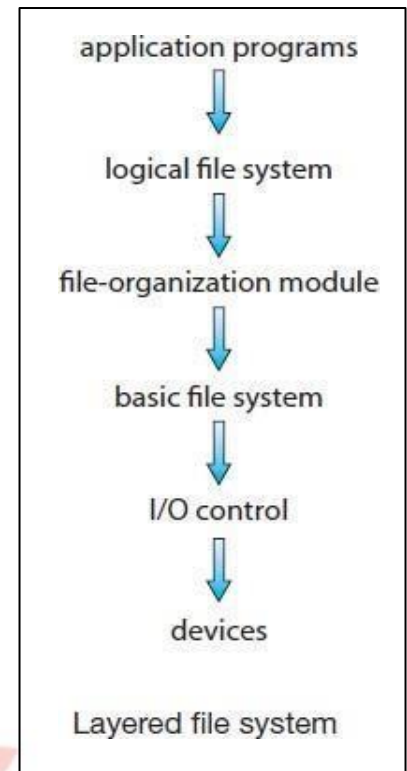
The **basic file system** needs only to issue generic commands to the appropriate device driver to read and write physical blocks on the disk. Each physical block is identified by its numeric disk address.

I/O Control

The **I/O control** level consists of device drivers and interrupt handlers to transfer information between the main memory and the disk system. It acts like a translator, inputting high-level commands such as “retrieve block 123.” And outputting low-level, hardware-specific instructions that are used by the hardware controller

Devices

These are the actual hardware devices like disk.

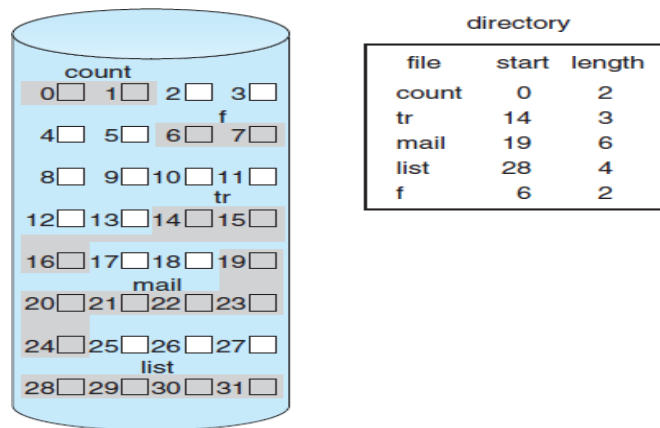


Allocation methods

Many files can be stored on the same disk. The main problem is how to allocate space to these files so that disk space is utilized effectively and files can be accessed quickly. The following are the three major methods of allocating disk space that are in wide use:

1. Contiguous Allocation

- **Contiguous allocation** requires that each file occupy a set of contiguous blocks on the disk. Disk addresses define a linear ordering on the disk.
- Contiguous allocation of a file is defined by the disk address and length (in block units) of the first block. If the file is n blocks long and starts at location b , then it occupies blocks $b, b + 1, b + 2, b + n - 1$. The directory entry for each file indicates the address of the starting block and the length of the area allocated for this file.
- **Accessing a file:**
Accessing a file that has been allocated contiguously is easy. It supports both sequential and random access. For sequential access, the file system remembers the disk address of the last block referenced and, when necessary, reads the next block. For direct access to block i of a file that starts at block b , we can immediately access block $b + i$.



Contiguous allocation of disk space.

- **Drawbacks**

- Finding space for a new file. The system chosen to manage free space determine show this task is accomplished. First fit and best fit are the most common strategies used to select a free hole from the set of available holes.

- **External Fragmentation**

As files are allocated and deleted, the free disk space is broken into little pieces. External fragmentation exists whenever free space is broken into chunks. It becomes a problem when the largest contiguous chunk is insufficient for a request; storage is fragmented into a number of holes, none of which is large enough to store the data.

- **Solution to external fragmentation**

Copy an entire file system onto another disk. The original disk is then freed completely, creating one large contiguous freespace. We then copy the files back onto the original disk by allocating contiguous space from this one large hole. This scheme effectively **compacts** all free space into one contiguous space, solving the fragmentation problem.

- Determining how much space is needed for a file. When the file is created, the total amount of space it will need must be found and allocated. If we allocate too little space to a file, we may find that the file cannot be extended.

Two possibilities then exist. First, the user program can be terminated, with an appropriate error message. The user must then allocate more space and run the program again. These repeated runs may be costly. To prevent them, the user will normally overestimate the amount of space needed, resulting in considerable wasted space. The other possibility is to find a larger hole, copy the contents of the file to the new space, and release the previous space. This series of actions can be repeated as long as space exists, although it can be time consuming. The user need never be informed explicitly about what is happening, however; the system continues despite the problem, although more and more slowly. Even if the total amount of space needed for a file is known in advance, preallocation may be inefficient. A file that will grow slowly over a long period.

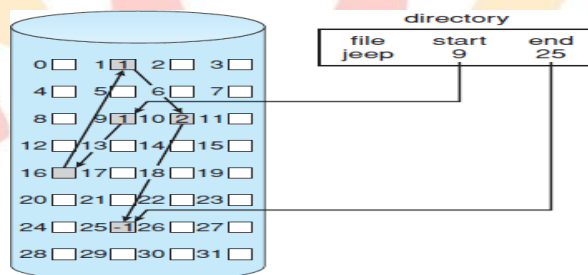
- **Modified Contiguous-Allocation**

- To minimize these drawbacks, some operating systems use a modified contiguous-allocation scheme. Here, a contiguous chunk of space is allocated

initially. Then, if that amount proves not to be large enough, another chunk of contiguous space, known as an **extent**, is added.

2. Linked Allocation

- **Linked allocation** solves all problems of contiguous allocation. With linked allocation, each file is a linked list of disk blocks; the disk blocks may be scattered anywhere on the disk. The directory contains a pointer to the first and last blocks of the file. For example, a file of five blocks might start at block 9 and continue at block 16, then block 1, then block 10, and finally block 25.
- To create a new file, we simply create a new entry in the directory. With linked allocation, each directory entry has a pointer to the first disk block of the file. This pointer is initialized to null (the end-of-list pointer value) to signify an empty file. The size field is also set to 0.
- A write to the file causes the free-space management system to find a free block, and this new block is written to and is linked to the end of the file.
- To read a file, we simply read blocks by following the pointers from block to block.



Linked allocation of disk space.

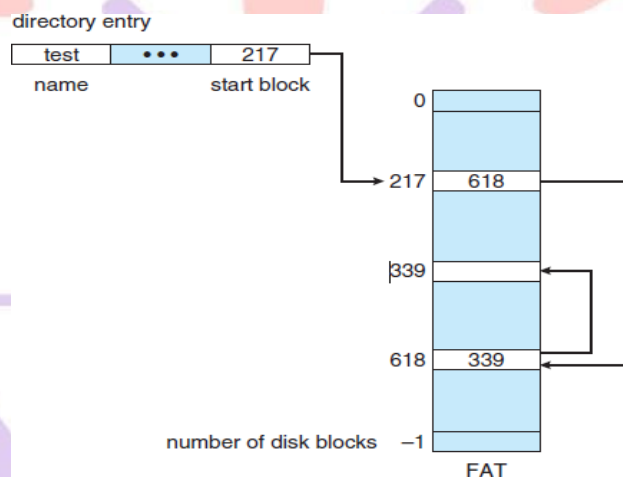
- **Advantages**
 - There is no external fragmentation with linked allocation, and any free block on the free-space list can be used to satisfy a request.
 - The size of a file need not be declared when the file is created. A file can continue to grow as long as free blocks are available. Consequently, it is never necessary to compact disk space.
- **Disadvantages**
 - The major problem is that it can be used effectively only for sequential-access files. To find the i^{th} block of a file, we must start at the beginning of that file and follow the pointers until we get to the i^{th} block. Each access to a pointer requires a disk read, and some require a disk seek. Consequently, it is inefficient to support a direct-access capability for linked-allocation files.
 - Another disadvantage is the space required for the pointers. If a pointer requires 4 bytes out of a 512-byte block, then 0.78 percent of the disk is being used for pointers, rather than for information. Each file requires slightly more space than it would otherwise.
 - The usual solution to this problem is to collect blocks into multiples, called **clusters**, and to allocate clusters rather than blocks. For instance, the file system may define a cluster as four blocks and operate on the disk only in cluster units.
 - Another problem of linked allocation is reliability: the files are linked together by pointers scattered all over the disk, and consider what would happen if a pointer

were lost or damaged.

- One partial solution is to use doubly linked lists, and another is to store the file name and relative block number in each block. However, these schemes require even more overhead for each file.

Variation on Linked Allocation

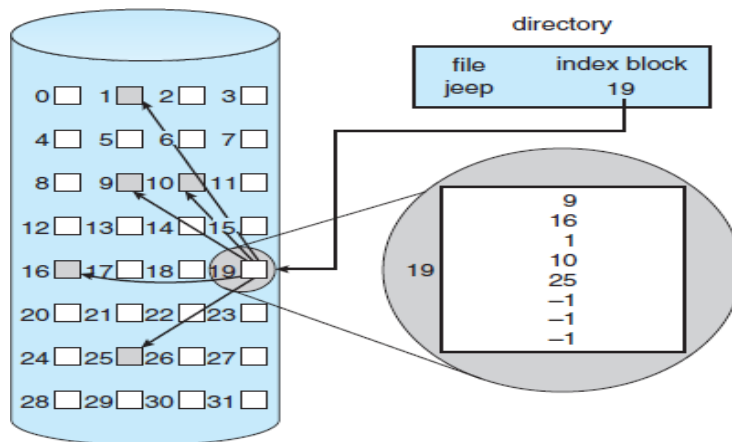
- An important variation on linked allocation is the use of a **file-allocation table (FAT)**. This simple but efficient method of disk-space allocation was used by the MS-DOS operating system.
- A section of disk at the beginning of each volume is set aside to contain the table. The table has one entry for each disk block and is indexed by block number.
- The FAT is used in much the same way as a linked list. The directory entry contains the block number of the first block of the file. The table entry indexed by that block number contains the block number of the next block in the file. This chain continues until it reaches the last block, which has a special end-of-file value as the table entry.
- An unused block is indicated by a table value of 0. Allocating a new block to a file is a simple matter of finding the first 0-valued table entry and replacing the previous end-of-file value with the address of the new block. The 0 is then replaced with the end-of-file value.



3. Indexed Allocation

- Linked allocation solves the external-fragmentation and size-declaration problems of contiguous allocation. However, in the absence of a FAT, linked allocation cannot support efficient direct access, since the pointers to the blocks are scattered with the blocks themselves all over the disk and must be retrieved in order.
- **Indexed allocation** solves this problem by bringing all the pointers together into one location: the **index block**.
- Each file has its own index block, which is an array of disk-block addresses. The i^{th} entry in the index block points to the i^{th} block of the file. The directory contains the address of the index block.
- To find and read the i^{th} block, we use the pointer in the i^{th} index-block entry.
- When the file is created, all pointers in the index block are set to null. When the i^{th} block is first written, a block is obtained from the free-space manager, and its address is put in the i^{th} index-block entry.

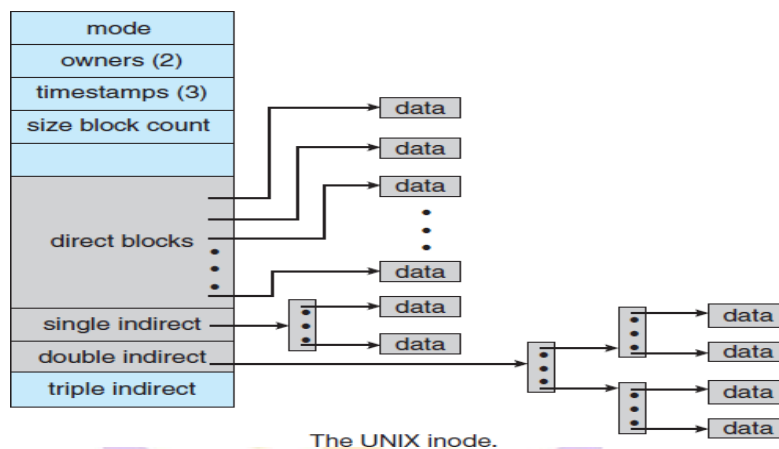
- Indexed allocation supports direct access, without suffering from external fragmentation, because any free block on the disk can satisfy a request for more space.



Indexed allocation of disk space.

- **Disadvantages**
 - Indexed allocation does suffer from wasted space.
 - The pointer overhead of the index block is generally greater than the pointer overhead of linked allocation.
- **Mechanisms for implementing Index Block**
 - **Linked scheme.** An index block is normally one disk block. Thus, it can be read and written directly by itself. To allow for large files, we can link together several index blocks. For example, an index block might contain a small header giving the name of the file and a set of the first 100 disk-block addresses. The next address (the last word in the index block) is null (for a small file) or is a pointer to another index block (for a large file).
 - **Multilevel index.** A variant of linked representation uses a first-level index block to point to a set of second-level index blocks, which in turn point to the file blocks. To access a block, the operating system uses the first-level index to find a second-level index block and then uses that block to find the desired data block. This approach could be continued to a third or fourth level, depending on the desired maximum file size. With 4,096-byte blocks, we could store 1,024 four-byte pointers in an index block. Two levels of indexes allow 1,048,576 data blocks and a file size of up to 4 GB.
 - **Combined scheme.** Another alternative, used in UNIX-based file systems, is to keep the first, say, 15 pointers of the index block in the file's inode. The first 12 of these pointers point to **direct blocks**; that is, they contain addresses of blocks that contain data of the file. Thus, the data for small files (of no more than 12 blocks) do not need a separate index block. If the block size is 4 KB, then up to 48 KB of data can be accessed directly. The next three pointers point to **indirect blocks**. The first points to a **single indirect block**, which is an index block containing not data but the addresses of blocks that do contain data. The second points to a **double indirect block**, which contains the address of a block that

contains the addresses of blocks that contain pointers to the actual data blocks. The last pointer contains the address of a **triple indirect block**.



Free-space Management

To keep track of free disk space, the system maintains a **free-space list**. The free-space list records all free disk blocks—those not allocated to some file or directory. The following are implementations of free space list.

1. BitVector

- Free-space list is frequently implemented as a **bit map** or **bit vector**. Each block is represented by 1 bit. If the block is free, the bit is 1; if the block is allocated, the bit is 0.
- For example, consider a disk where blocks 2, 3, 4, 5, 8, 9, 10, 11, 12, 13, 17, 18, 25, 26, and 27 are free and the rest of the blocks are allocated. The free-space bit map would be
001111001111110001100000011100000...

- **Advantage**

The main advantage of this approach is its relative simplicity and its efficiency in finding the first free block or n consecutive free blocks on the disk. Indeed, many computers supply bit-manipulation instructions that can be used effectively for that purpose. One technique for finding the first free block on a system that uses a bit-vector to allocate disk space is to sequentially check each word in the bit map to see whether that value is not 0, since a 0-valued word contains only 0 bits and represents a set of allocated blocks. The first non-0 word is scanned for the first 1 bit, which is the location of the first free block.

The calculation of the block number is

$$(\text{number of bits per word}) \times (\text{number of 0-valued words}) + \text{offset of first 1 bit.}$$

- **Disadvantage**

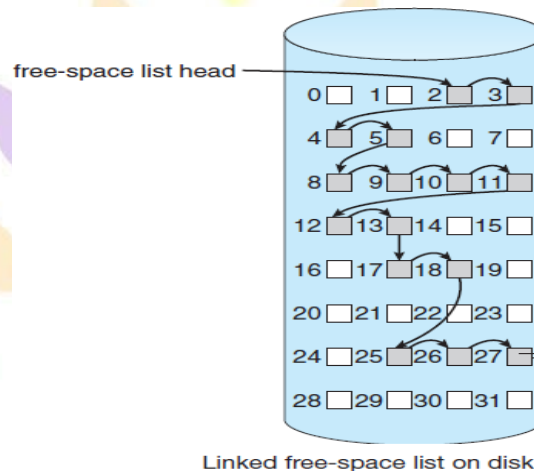
Bit vectors are inefficient unless the entire vector is kept in main memory (and is written to disk occasionally for recovery needs). Keeping it in main memory is possible for smaller disks but not necessarily for larger ones.

2. Linked List

- Another approach to free-space management is to link together all the free disk blocks, keeping a pointer to the first free block in a special location on the disk and

caching it in memory. This first block contains a pointer to the next free disk block, and so on.

- For example, blocks 2, 3, 4, 5, 8, 9, 10, 11, 12, 13, 17, 18, 25, 26, and 27 were free and the rest of the blocks were allocated. In this situation, we would keep a pointer to block 2 as the first free block. Block 2 would contain a pointer to block 3, which would point to block 4, which would point to block 5, which would point to block 8, and so on.
- **Disadvantages**
This scheme is not efficient; to traverse the list, we must read each block, which requires substantial I/O time



3. Grouping

- A modification of the free-list approach stores the addresses of n free blocks in the first free block. The first $n-1$ of these blocks are actually free. The last, block contains the addresses of another n free block, and so on. The addresses of a large number of free blocks can now be found quickly.

4. Counting

- Several contiguous blocks may be allocated or freed simultaneously, particularly when space is allocated with the contiguous-allocation algorithm or through clustering.
- Thus, rather than keeping a list of n free disk addresses, we can keep the address of the first free block and the number (n) of free contiguous blocks that follow the first block. Each entry in the free-space list then consists of a disk address and a count.

5. SpaceMaps

- Oracle's **ZFS** file system was designed to encompass huge numbers of files, directories, and even file systems.
- In its management of free space, ZFS creates **metaslabs** to divide the space on the device into chunks of manageable size. Each metaslab has an associated space map.
- The space map is a log of all block activity (allocating and freeing), in time order, in counting format. When ZFS decides to allocate or free space from a metaslab, it loads the associated space map into memory in a balanced-tree structure (for very efficient operation), indexed by offset, and replays the log to that structure.

- The in-memory space map is then an accurate representation of the allocated and free space in the meta slab.

System calls for file operations—`open()`, `read()`, `write()`, `close()`, `seek()`, `unlink()` (File Operations)

create()

This is used to create a file. Two steps are necessary to create a file. First, space in the file system must be found for the file. Second, an entry for the new file must be made in the directory.

open ()

Many systems require that an `open ()` system call be made before a file is first used. When a file has been opened its entry is added in the open file table. It also contains open count associated with each file to indicate how many processes have the file open.

read()

To read from a file, we use a system call that specifies the name of the file and **read pointer** to the location in the file where the next read is to take place. Once the read has taken place, the read pointer is updated.

write()

To write a file, we make a system call specifying both the name of the file and the information to be written to the file. Given the name of the file, the system searches the directory to find the file's location.

The system must keep a **write pointer** to the location in the file where the next write is to take place. The write pointer must be updated whenever a write occurs.

close()

This closes a file. Each `close ()` decrements the open count and when the count reaches zero, the file is no longer in use so it can be closed.

delete()

To delete a file, we search the directory for the named file. Having found the associated directory entry, we release all file space, so that it can be reused by other files, and erase the directory entry.

truncate()

The user may want to erase the contents of a file but keep its attributes. Rather than forcing the user to delete the file and then recreate it, this function allows all attributes to remain unchanged—except for file length—but lets the file be reset to length zero and its file space released.

seek()

It is also called as Reposition. The directory is searched for the appropriate entry, and the current-file-position pointer is repositioned to a given value. Repositioning within a file need not involve any actual I/O.

unlink()

Deletes a name from the file system. If that name was the last link to a file and no processes have the file open the file is deleted and the space it was using is made available for reuse.