# Operating Systems

*Narasimhulu M*<sub>M. Tech.</sub>
*Assistant Professor*
*Department of Computer Science & Engineering*

| S.No. | Course Outcomes | Cognitive Level |
|-------|-----------------|-----------------|
| 1 | Explain the fundamentals of operating systems like process, memory, storage, file system, security and protection. | Understand |
| 2 | Illustrate various operating System services, interfaces and system calls. | Apply |
| 3 | Demonstrate critics of process management and IPC. | Apply |
| 4 | Implement page replacement algorithms, memory management techniques and deadlock issues. | Apply |
| 5 | Illustrate architecture of file systems and I/O systems for mass storage structures. | Apply |
| 6 | Utilize the methods of operating system security and protection. | Apply |

12/12/2022    Prepared by M. Narasimhulu, CSE
Assistant Professor    2

# Unit 4: Storage Management & File System

- **Mass-Storage Structure:** Overview of Mass-Storage Structure, Disk Scheduling, Storage Attachment, RAID Structure.

- **I/O Systems:** I/O Hardware, Application I/O Interface, Kernel I/O Subsystem, Transforming I/O Requests to Hardware Operations.

- **File-System :** File Concept, Access Methods, Directory Structure, Protection, Memory-Mapped Files, File system structure and Implementation.

12/12/2022                        Prepared by: M. Narasimhulu, CSE,                       3
                                         Assistant Professor

# Unit 4 - Storage Management & File System

*Narasimhulu M*M. Tech.
*Assistant Professor*
*Department of Computer Science & Engineering*

# Chapter 1
## Mass-Storage Structure

*Narasimhulu M*<sub>M. Tech.</sub>
*Assistant Professor*
*Department of Computer Science & Engineering*

---

# *Overview of Mass-Storage Structure*

*Narasimhulu M*<sub>M. Tech.</sub>
*Assistant Professor*
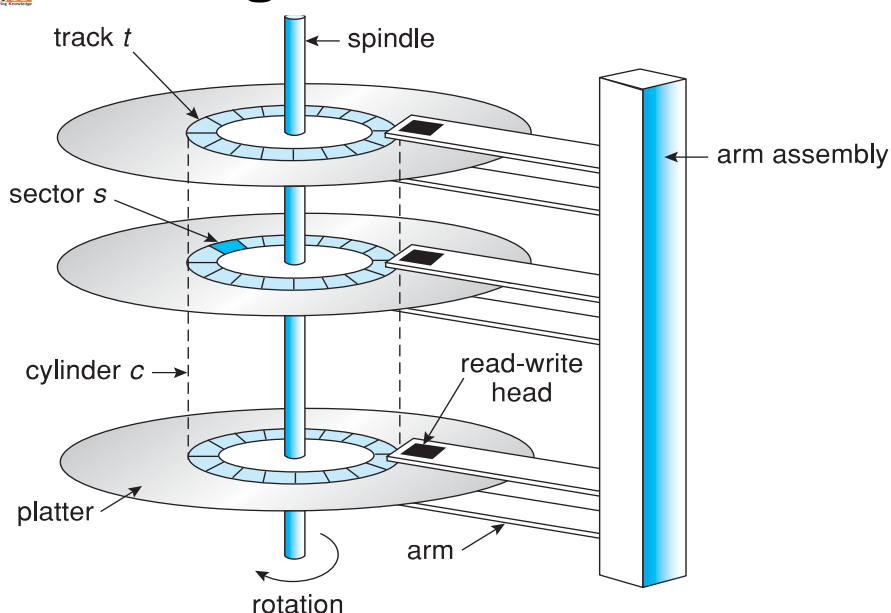*Department of Computer Science & Engineering*

# Overview of Mass Storage Structure

- **Magnetic disks** provide bulk of secondary storage of modern computers
    - Drives rotate at 60 to 250 times per second
    - **Transfer rate** is rate at which data flow between drive and computer
    - **Positioning time** (**random-access time**) is time to move disk arm to desired cylinder (**seek time**) and time for desired sector to rotate under the disk head (**rotational latency**)
    - **Head crash** results from disk head making contact with the disk surface -- That's bad
- Disks can be removable
- Drive attached to computer via **I/O bus**
    - Busses vary, including **EIDE**, **ATA**, **SATA**, **USB**, **Fibre Channel**, **SCSI, SAS, Firewire**
    - **Host controller** in computer uses bus to talk to **disk controller** built into drive or storage array

# Moving-head Disk Mechanism

# Hard Disks

- Platters range from .85" to 14" (historically)
  - Commonly 3.5", 2.5", and 1.8"
- Range from 30GB to 3TB per drive
- Performance
  - Transfer Rate – theoretical – 6 Gb/sec
  - Effective Transfer Rate – real – 1Gb/sec
  - Seek time from 3ms to 12ms – 9ms common for desktop drives
  - Average seek time measured or calculated based on 1/3 of tracks
  - Latency based on spindle speed
    - 1 / (RPM / 60) = 60 / RPM
  - Average latency = ½ latency

| Spindle [rpm] | Average latency [ms] |
|---------------|----------------------|
| 4200          | 7.14                 |
| 5400          | 5.56                 |
| 7200          | 4.17                 |
| 10000         | 3                    |
| 15000         | 2                    |

(From Wikipedia)

# Hard Disk Performance

- **Access Latency** = **Average access time** = average seek time + average latency
  - For fastest disk 3ms + 2ms = 5ms
  - For slow disk 9ms + 5.56ms = 14.56ms
- Average I/O time = average access time + (amount to transfer / transfer rate) + controller overhead
- For example to transfer a 4KB block on a 7200 RPM disk with a 5ms average seek time, 1Gb/sec transfer rate with a .1ms controller overhead =
  - 5ms + 4.17ms + 0.1ms + transfer time =
  - Transfer time = 4KB / 1Gb/s * 8Gb / GB * 1GB / $1024^2$KB = 32 / $(1024^2)$ = 0.031 ms
  - Average I/O time for 4KB block = 9.27ms + .031ms = 9.301ms

# The First Commercial Disk Drive

1956
IBM RAMDAC computer
included the IBM Model 350
disk storage system

5M (7 bit) characters
50 x 24" platters
Access time = < 1 second

# Solid-State Disks

- Nonvolatile memory used like a hard drive
  – Many technology variations
- Can be more reliable than HDDs
- More expensive per MB
- Maybe have shorter life span
- Less capacity
- But much faster
- Busses can be too slow -> connect directly to PCI for example
- No moving parts, so no seek time or rotational latency

## Magnetic Tape

- Was early secondary-storage medium
  - Evolved from open spools to cartridges
- Relatively permanent and holds large quantities of data
- Access time slow
- Random access ~1000 times slower than disk
- Mainly used for backup, storage of infrequently-used data, transfer medium between systems
- Kept in spool and wound or rewound past read-write head
- Once data under head, transfer rates comparable to disk
  - 140MB/sec and greater
- 200GB to 1.5TB typical storage
- Common technologies are LTO-{3,4,5} and T10000

# *Disk Scheduling*

*Narasimhulu M*$_{M. Tech.}$
*Assistant Professor*
*Department of Computer Science & Engineering*

# Disk Scheduling

- The operating system is responsible for using hardware efficiently — for the disk drives, this means having a fast access time and disk bandwidth

- Minimize seek time

- Seek time ≈ seek distance

- Disk **bandwidth** is the total number of bytes transferred, divided by the total time between the first request for service and the completion of the last transfer

# Disk Scheduling (Cont.)

- There are many sources of disk I/O request
  - OS
  - System processes
  - Users processes
- I/O request includes input or output mode, disk address, memory address, number of sectors to transfer
- OS maintains queue of requests, per disk or device
- Idle disk can immediately work on I/O request, busy disk means work must queue
  - Optimization algorithms only make sense when a queue exists

# Disk Scheduling (Cont.)

- Note that drive controllers have small buffers and can manage a queue of I/O requests (of varying "depth")

- Several algorithms exist to schedule the servicing of disk I/O requests

- The analysis is true for one or many platters

- We illustrate scheduling algorithms with a request queue (0-199)

        98, 183, 37, 122, 14, 124, 65, 67
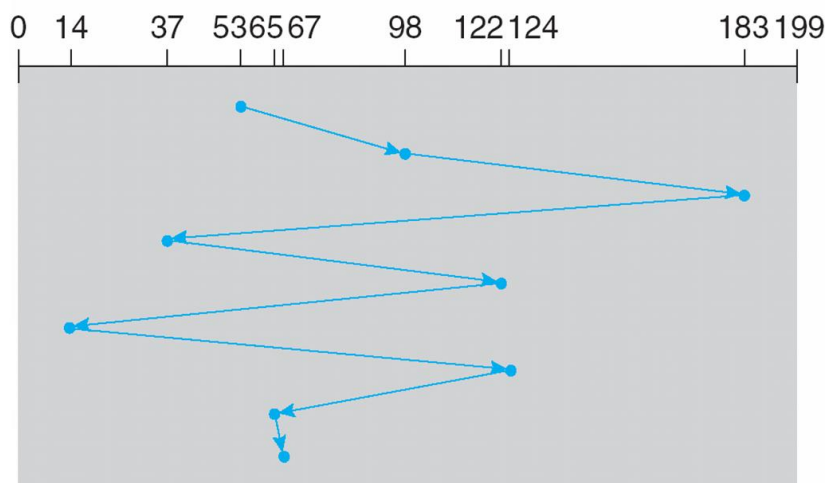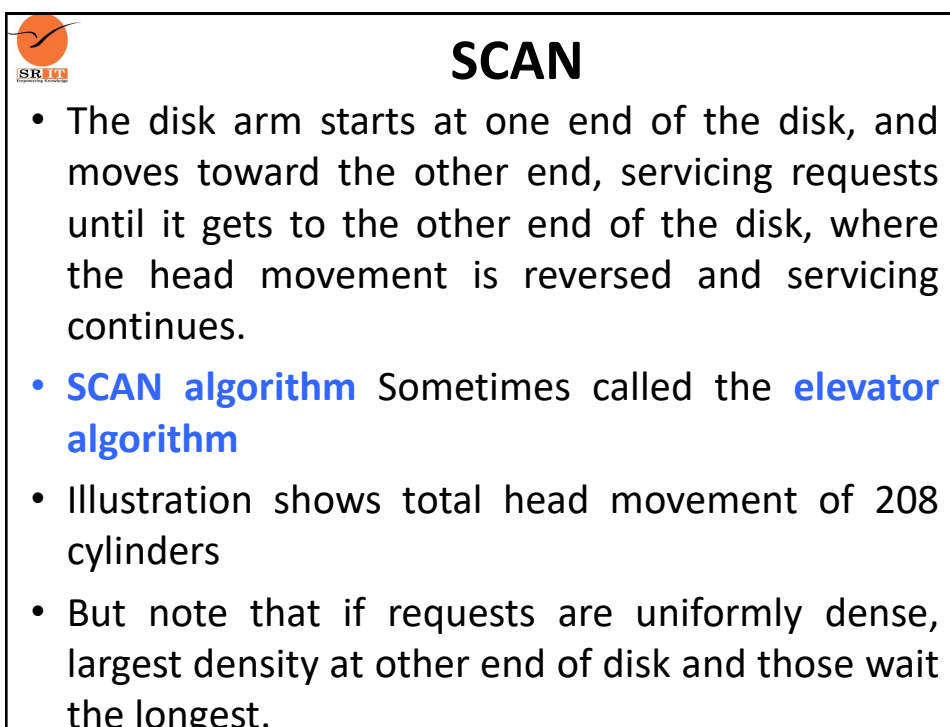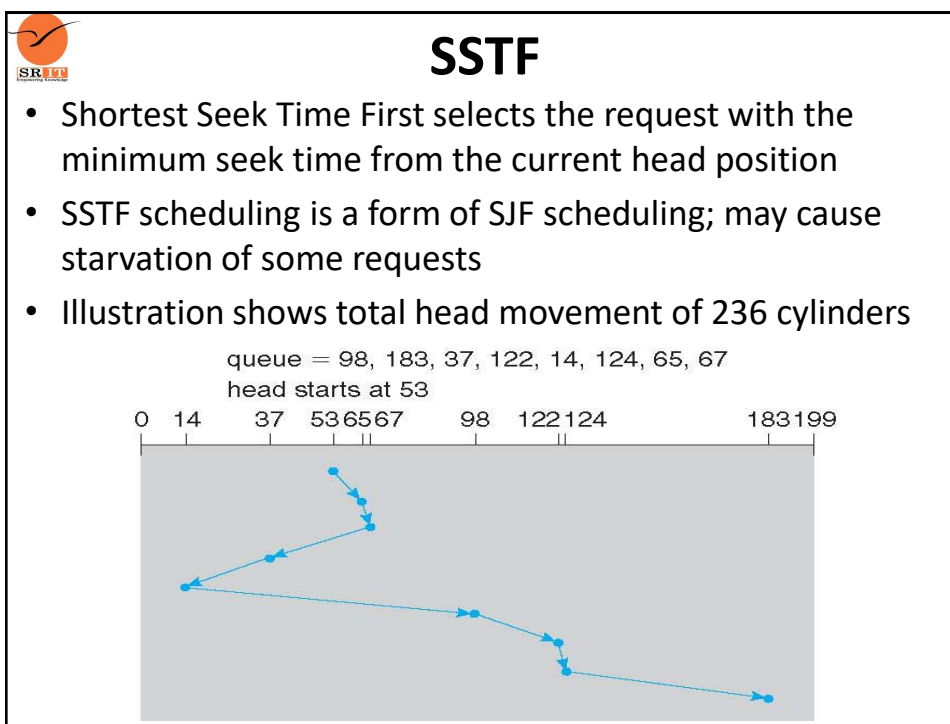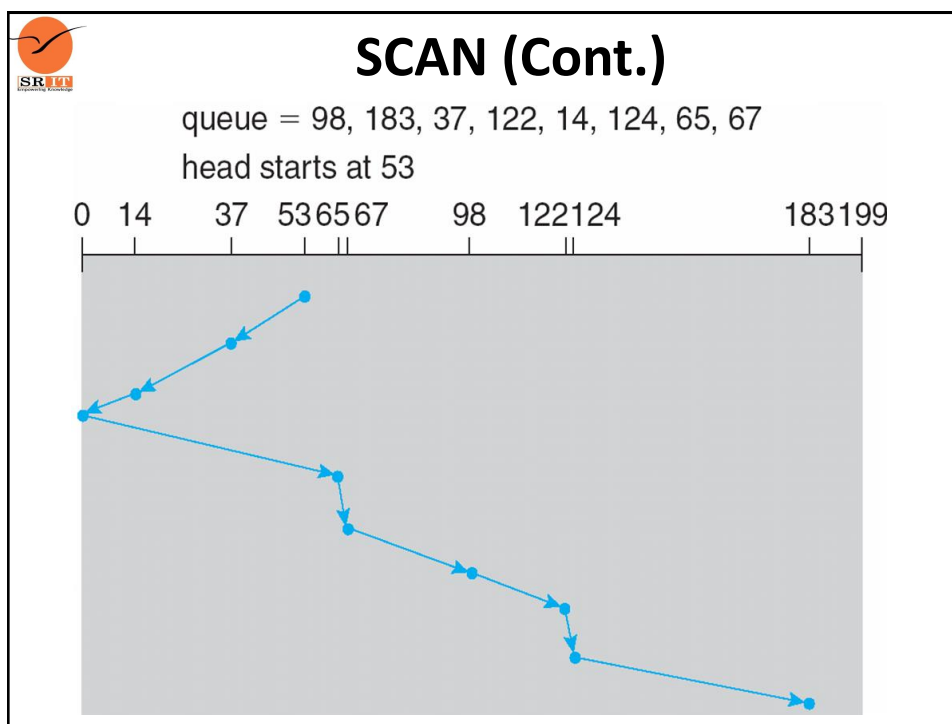
Head pointer 53

# FCFS

**Illustration shows total head movement of 640 cylinders**

queue = 98, 183, 37, 122, 14, 124, 65, 67
head starts at 53

# SSTF

- Shortest Seek Time First selects the request with the minimum seek time from the current head position
- SSTF scheduling is a form of SJF scheduling; may cause starvation of some requests
- Illustration shows total head movement of 236 cylinders

queue = 98, 183, 37, 122, 14, 124, 65, 67
head starts at 53



# SCAN

- The disk arm starts at one end of the disk, and moves toward the other end, servicing requests until it gets to the other end of the disk, where the head movement is reversed and servicing continues.

- **SCAN algorithm** Sometimes called the **elevator algorithm**

- Illustration shows total head movement of 208 cylinders

- But note that if requests are uniformly dense, largest density at other end of disk and those wait the longest.

## SCAN (Cont.)

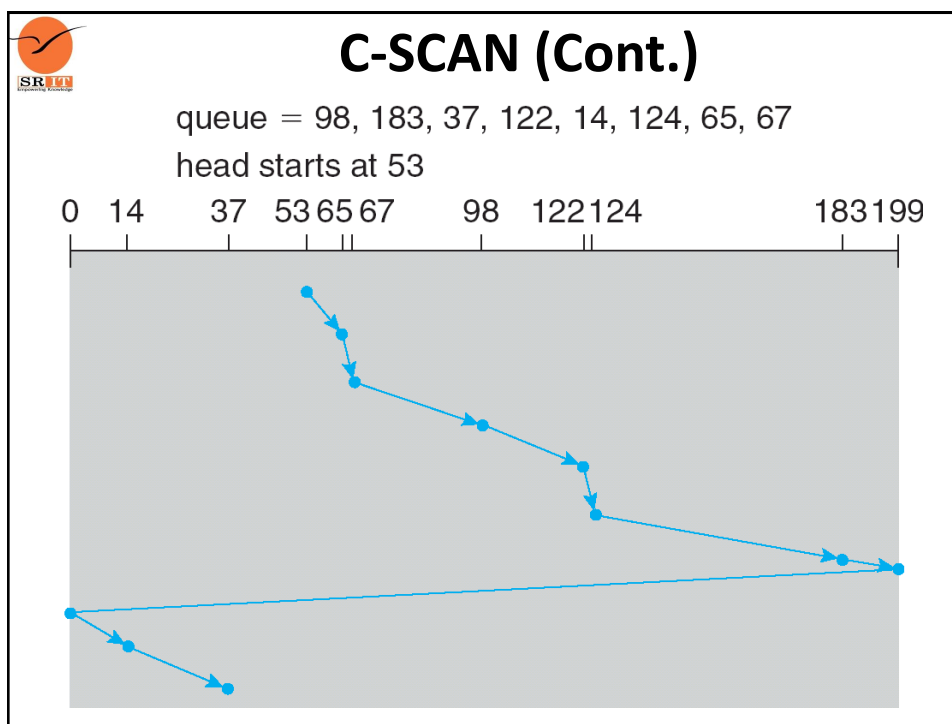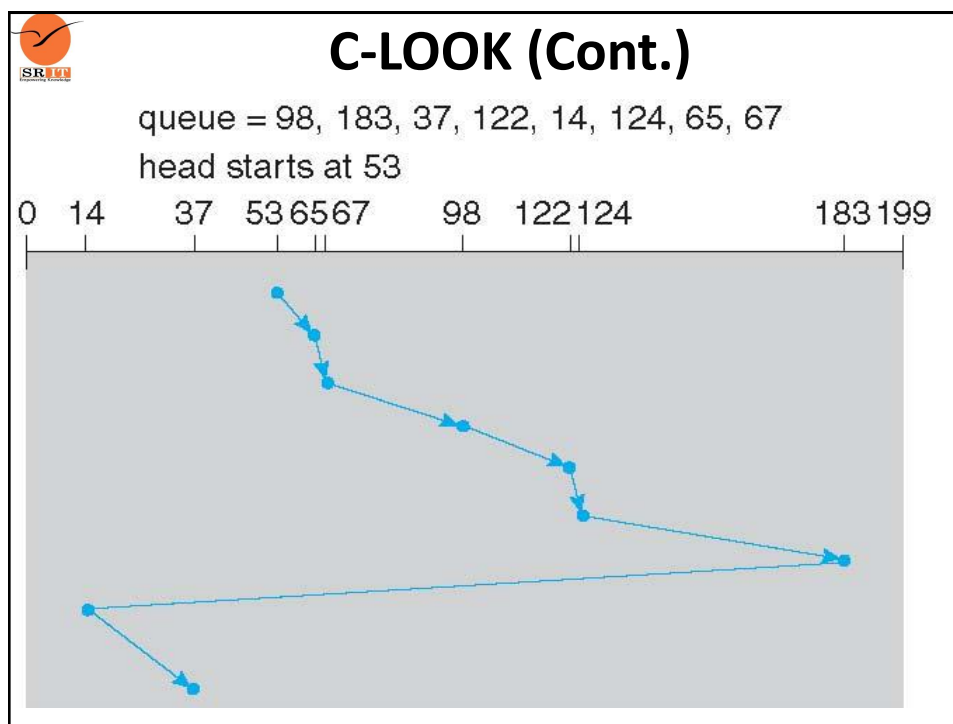queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53



## C-SCAN

- Provides a more uniform wait time than SCAN
- The head moves from one end of the disk to the other, servicing requests as it goes
  - When it reaches the other end, however, it immediately returns to the beginning of the disk, without servicing any requests on the return trip
- Treats the cylinders as a circular list that wraps around from the last cylinder to the first one
- Total number of cylinders?

# C-SCAN (Cont.)

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53



# C-LOOK

- LOOK a version of SCAN, C-LOOK a version of C-SCAN
- Arm only goes as far as the last request in each direction, then reverses direction immediately, without first going all the way to the end of the disk
- Total number of cylinders?

# C-LOOK (Cont.)

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53



---

# Selecting a Disk-Scheduling Algorithm

- SSTF is common and has a natural appeal
- SCAN and C-SCAN perform better for systems that place a heavy load on the disk
  - Less starvation
- Performance depends on the number and types of requests
- Requests for disk service can be influenced by the file-allocation method
  - And metadata layout
- The disk-scheduling algorithm should be written as a separate module of the operating system, allowing it to be replaced with a different algorithm if necessary
- Either SSTF or LOOK is a reasonable choice for the default algorithm
- What about rotational latency?
  - Difficult for OS to calculate
- How does disk-based queueing effect OS queue ordering efforts?

# *Storage Attachment*

*Narasimhulu M*<sub></sub>*M. Tech.*
*Assistant Professor*
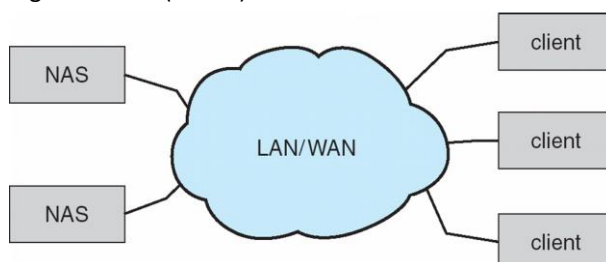*Department of Computer Science & Engineering*

---

# Storage Attachment

- Computers access storage in three ways
  - Host-attached
  - Network-attached
  - Cloud
- Host attached access through local I/O ports, using one of several technologies
  - To attach many devices, use storage busses such as USB, firewire, thunderbolt
  - High-end systems use **fibre channel** (**FC**)
    - High-speed serial architecture using fibre or copper cables
    - Multiple hosts and storage devices can connect to the FC fabric

# Network-Attached Storage

- Network-attached storage (**NAS**) is storage made available over a network rather than over a local connection (such as a bus)
  - Remotely attaching to file systems
- NFS and CIFS are common protocols
- Implemented via remote procedure calls (RPCs) between host and storage over typically TCP or UDP on IP network
- **iSCSI** protocol uses IP network to carry the SCSI protocol
  - Remotely attaching to devices (blocks)



# Cloud Storage

- Similar to NAS, provides access to storage across a network
  - Unlike NAS, accessed over the Internet or a WAN to remote data center
- NAS presented as just another file system, while cloud storage is API based, with programs using the APIs to provide access
  - Examples include Dropbox, Amazon S3, Microsoft OneDrive, Apple iCloud
  - Use APIs because of latency and failure scenarios (NAS protocols wouldn't work well)
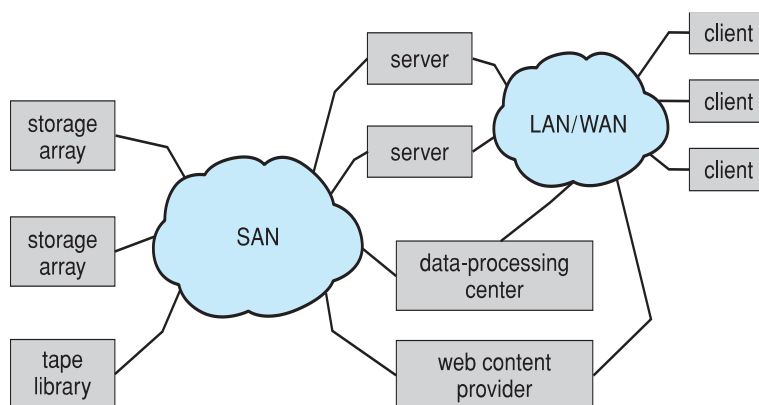
# Storage Array

- Can just attach disks, or arrays of disks
- Avoids the NAS drawback of using network bandwidth
- Storage Array has controller(s), provides features to attached host(s)
  - Ports to connect hosts to array
  - Memory, controlling software (sometimes NVRAM, etc.)
  - A few to thousands of disks
  - RAID, hot spares, hot swap (discussed later)
  - Shared storage -> more efficiency
  - Features found in some file systems
    - Snaphots, clones, thin provisioning, replication, deduplication, etc

# Storage Area Network

- Common in large storage environments
- Multiple hosts attached to multiple storage arrays
  - flexible

# Storage Area Network (Cont.)

- SAN is one or more storage arrays
  - Connected to one or more Fibre Channel switches or **InfiniBand** (**IB**) network
- Hosts also attach to the switches
- Storage made available via **LUN Masking** from specific arrays to specific servers
- Easy to add or remove storage, add new host and allocate it storage
- Why have separate storage networks and communications networks?
  - Consider iSCSI, FCOE



A Storage Array

---

# *RAID Structure*

*Narasimhulu M*~M. Tech.~
*Assistant Professor*
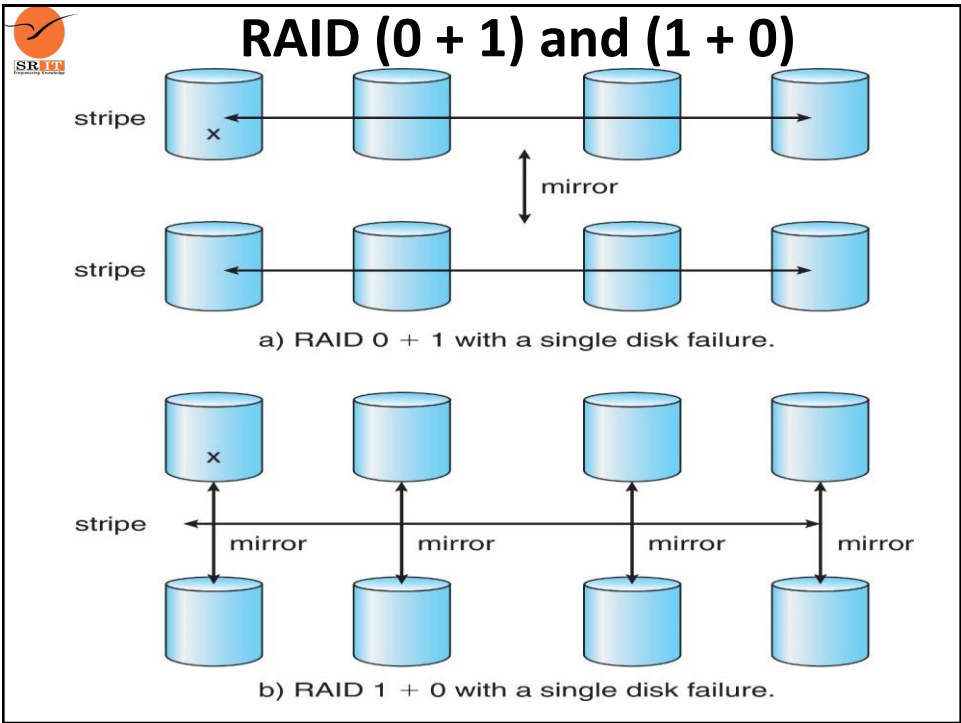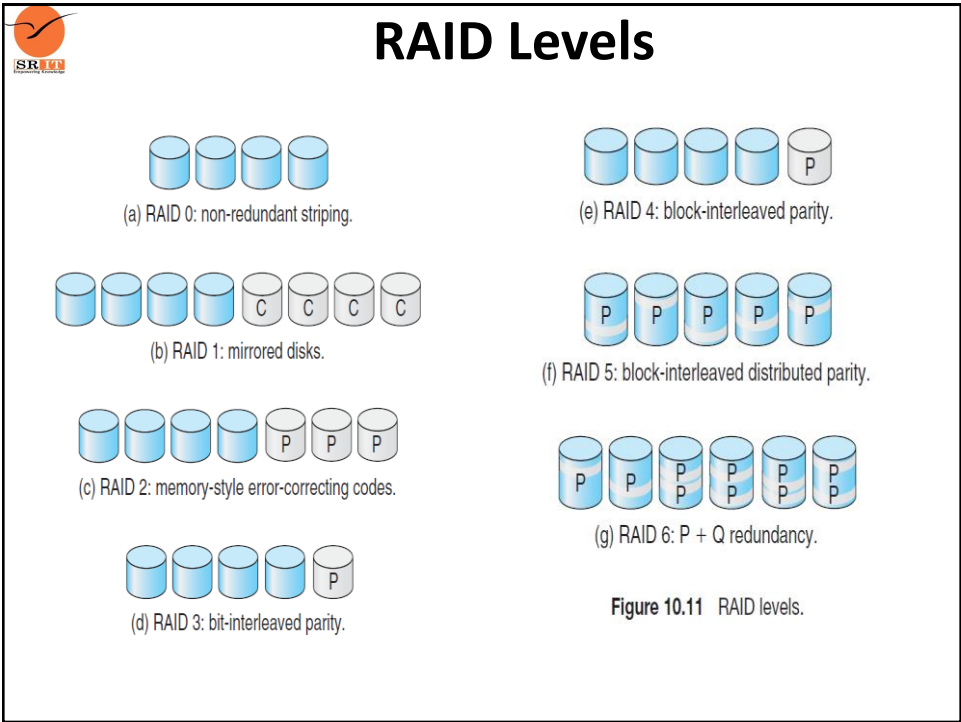*Department of Computer Science & Engineering*

# RAID Structure

- **RAID** – **redundant array of inexpensive disks**
  - multiple disk drives provides reliability via **redundancy**
- Increases the **mean time to failure**
- **Mean time to repair –** exposure time when another failure could cause data loss
- **Mean time to data loss** based on above factors
- If mirrored disks fail independently, consider disk with 100,000 **mean time to failure** and 10 hour mean time to repair
  - Mean time to data loss is $100,000^2 / (2 * 10) = 500 * 10^6$ hours, or 57,000 years!
- Frequently combined with **NVRAM** to improve write performance
- Several improvements in disk-use techniques involve the use of multiple disks working cooperatively

# RAID (Cont.)

- Disk **striping** uses a group of disks as one storage unit
- RAID is arranged into six different levels
- RAID schemes improve performance and improve the reliability of the storage system by storing redundant data
  - **Mirroring** or **shadowing** (**RAID 1**) keeps duplicate of each disk
  - Striped mirrors (**RAID 1+0**) or mirrored stripes (**RAID 0+1**) provides high performance and high reliability
  - **Block interleaved parity** (**RAID 4, 5, 6**) uses much less redundancy
- RAID within a storage array can still fail if the array fails, so automatic **replication** of the data between arrays is common
- Frequently, a small number of **hot-spare** disks are left unallocated, automatically replacing a failed disk and having data rebuilt onto them

# RAID Levels



(a) RAID 0: non-redundant striping.

(b) RAID 1: mirrored disks.

(c) RAID 2: memory-style error-correcting codes.

(d) RAID 3: bit-interleaved parity.

(e) RAID 4: block-interleaved parity.

(f) RAID 5: block-interleaved distributed parity.

(g) RAID 6: P + Q redundancy.

**Figure 10.11**   RAID levels.

# RAID (0 + 1) and (1 + 0)



a) RAID 0 + 1 with a single disk failure.

b) RAID 1 + 0 with a single disk failure.

# END of Chapter - 1

Prepared by: M. Narasimhulu, CSE,            49
                                  Assistant Professor

# Chapter 2
# I/O Systems

*Narasimhulu M*<sub></sub>*M. Tech.*
*Assistant Professor*
*Department of Computer Science & Engineering*

# *I/O Hardware*

*Narasimhulu M*<sub></sub> *M. Tech.*
**Assistant Professor**
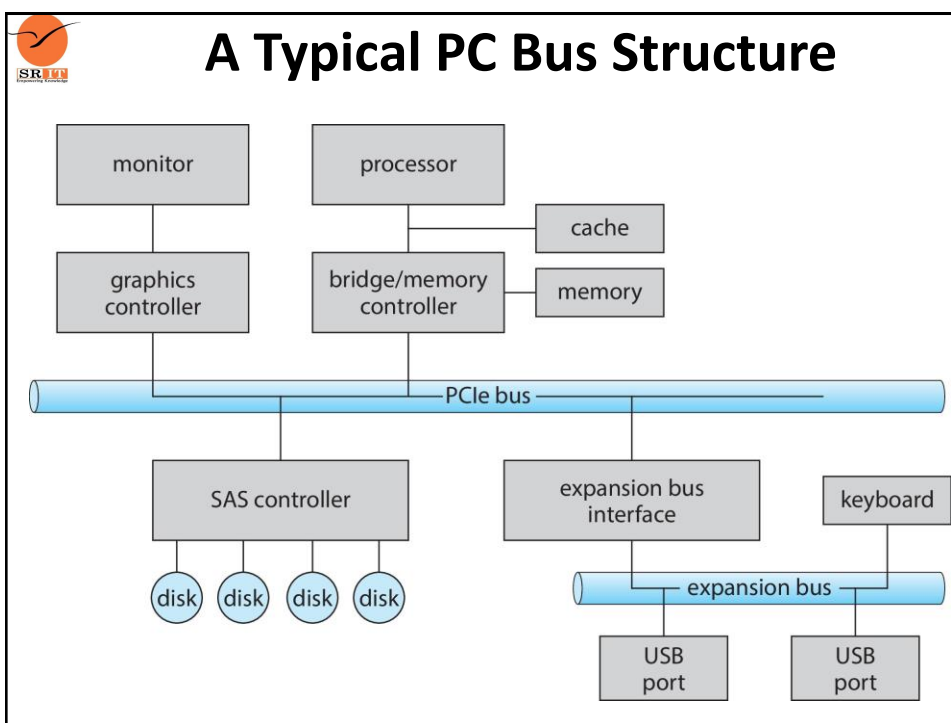**Department of Computer Science & Engineering**

---

# Overview

- I/O management is a major component of operating system design and operation
  - Important aspect of computer operation
  - I/O devices vary greatly
  - Various methods to control them
  - Performance management
  - New types of devices frequent
- Ports, busses, device controllers connect to various devices
- **Device drivers** encapsulate device details
  - Present uniform device-access interface to I/O subsystem

# I/O Hardware

- Incredible variety of I/O devices
  - Storage
  - Transmission
  - Human-interface
- Common concepts – signals from I/O devices interface with computer
  - **Port** – connection point for device
  - **Bus** - **daisy chain** or shared direct access
    - **PCI** bus common in PCs and servers, PCI Express (**PCIe**)
    - **expansion bus** connects relatively slow devices
    - **Serial-attached SCSI** (**SAS**) common disk interface
  - **Controller** (**host adapter**) – electronics that operate port, bus, device
    - Sometimes integrated
    - Sometimes separate circuit board (host adapter)
    - Contains processor, microcode, private memory, bus controller, etc.
      - Some talk to per-device controller with bus controller, microcode, memory, etc.

# A Typical PC Bus Structure

# I/O Hardware (Cont.)

- **Fibre channel** (**FC**) is complex controller, usually separate circuit board (**host-bus adapter**, **HBA**) plugging into bus
- I/O instructions control devices
- Devices usually have registers where device driver places commands, addresses, and data to write, or read data from registers after command execution
  - Data-in register, data-out register, status register, control register
  - Typically 1-4 bytes, or FIFO buffer
- Devices have addresses, used by
  - Direct I/O instructions
  - **Memory-mapped I/O**
    - Device data and command registers mapped to processor address space
    - Especially for large address spaces (graphics)

# Device I/O Port Locations on PCs (partial)

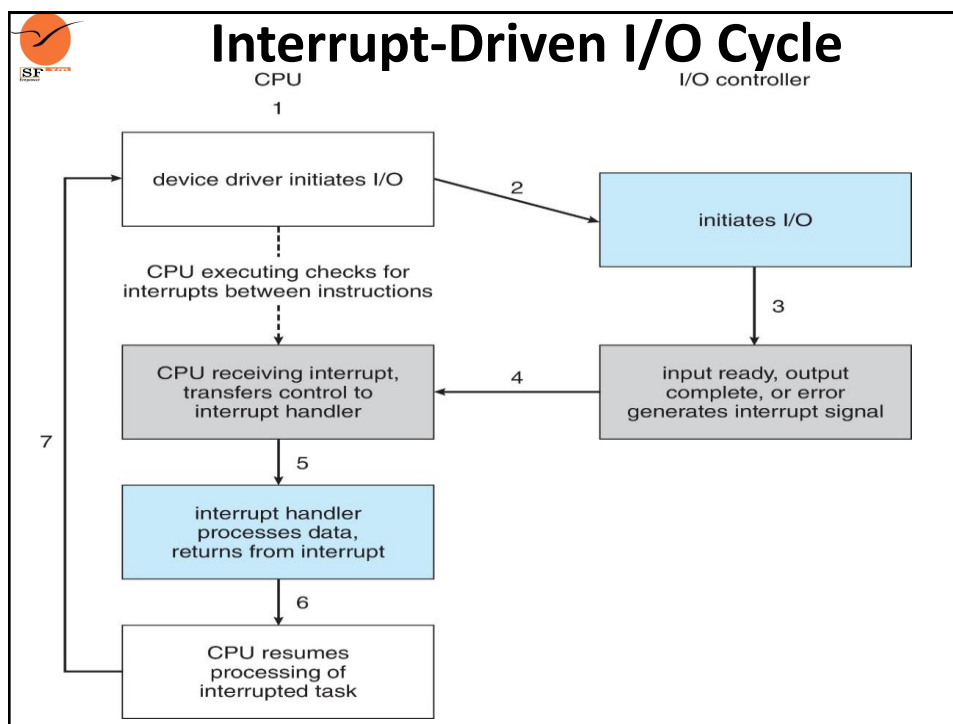| I/O address range (hexadecimal) | device |
|---|---|
| 000–00F | DMA controller |
| 020–021 | interrupt controller |
| 040–043 | timer |
| 200–20F | game controller |
| 2F8–2FF | serial port (secondary) |
| 320–32F | hard-disk controller |
| 378–37F | parallel port |
| 3D0–3DF | graphics controller |
| 3F0–3F7 | diskette-drive controller |
| 3F8–3FF | serial port (primary) |

# Polling

- For each byte of I/O
    1. Read busy bit from status register until 0
    2. Host sets read or write bit and if write copies data into data-out register
    3. Host sets command-ready bit
    4. Controller sets busy bit, executes transfer
    5. Controller clears busy bit, error bit, command-ready bit when transfer done
- Step 1 is **busy-wait** cycle to wait for I/O from device
    - Reasonable if device is fast
    - But inefficient if device slow
    - CPU switches to other tasks?
        - But if miss a cycle data overwritten / lost

# Interrupts

- Polling can happen in 3 instruction cycles
    - Read status, logical-and to extract status bit, branch if not zero
    - How to be more efficient if non-zero infrequently?
- CPU **Interrupt-request line** triggered by I/O device
    - Checked by processor after each instruction
- **Interrupt handler** receives interrupts
    - **Maskable** to ignore or delay some interrupts
- **Interrupt vector** to dispatch interrupt to correct handler
    - Context switch at start and end
    - Based on priority
    - Some **nonmaskable**
    - Interrupt chaining if more than one device at same interrupt number

## Interrupt-Driven I/O Cycle

CPU                                                    I/O controller

1

device driver initiates I/O ———2———→ initiates I/O

CPU executing checks for
interrupts between instructions

3

CPU receiving interrupt,  ←———4——— input ready, output
transfers control to                            complete, or error
interrupt handler                               generates interrupt signal

7

5

interrupt handler
processes data,
returns from interrupt

6

CPU resumes
processing of
interrupted task

## Interrupts (Cont.)

- Interrupt mechanism also used for **exceptions**
  - Terminate process, crash system due to hardware error
- Page fault executes when memory access error
- System call executes via **trap** to trigger kernel to execute request
- Multi-CPU systems can process interrupts concurrently
  - If operating system designed to handle it
- Used for time-sensitive processing, frequent, must be fast

# Latency

- Stressing interrupt management because even single-user systems manage hundreds or interrupts per second and servers hundreds of thousands

- For example, a quiet macOS desktop generated 23,000 interrupts over 10 seconds

```
Fri Nov 25 13:55:59                              0:00:10
                        SCHEDULER      INTERRUPTS
-------------------------------------------------
total_samples                 13           22998

delays <  10 usecs            12           16243
delays <  20 usecs             1            5312
delays <  30 usecs             0             473
delays <  40 usecs             0             590
delays <  50 usecs             0              61
delays <  60 usecs             0             317
delays <  70 usecs             0               2
delays <  80 usecs             0               0
delays <  90 usecs             0               0
delays < 100 usecs             0               0
total   < 100 usecs           13           22998
```

# Intel Pentium Processor Event-Vector Table

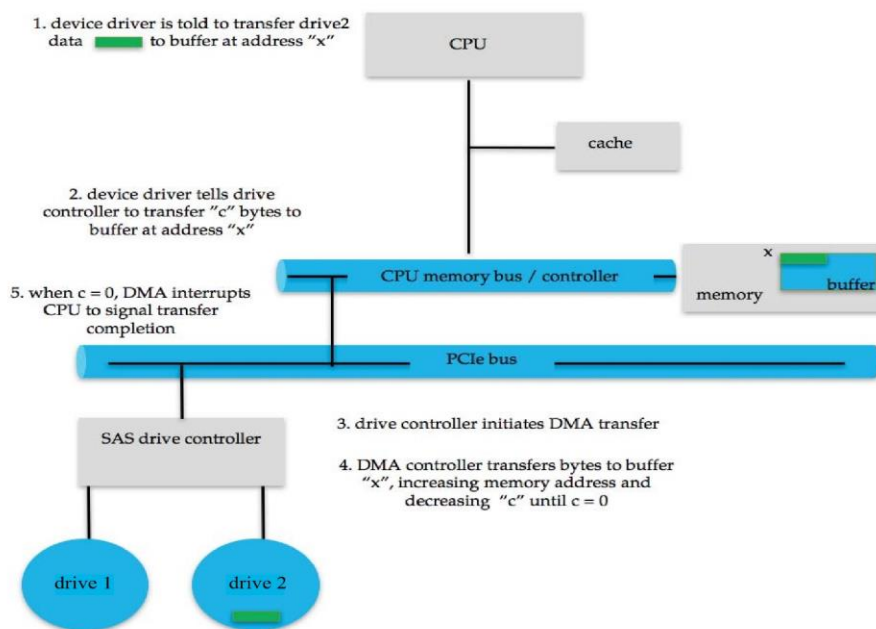| vector number | description |
| --- | --- |
| 0 | divide error |
| 1 | debug exception |
| 2 | null interrupt |
| 3 | breakpoint |
| 4 | INTO-detected overflow |
| 5 | bound range exception |
| 6 | invalid opcode |
| 7 | device not available |
| 8 | double fault |
| 9 | coprocessor segment overrun (reserved) |
| 10 | invalid task state segment |
| 11 | segment not present |
| 12 | stack fault |
| 13 | general protection |
| 14 | page fault |
| 15 | (Intel reserved, do not use) |
| 16 | floating-point error |
| 17 | alignment check |
| 18 | machine check |
| 19–31 | (Intel reserved, do not use) |
| 32–255 | maskable interrupts |

# Direct Memory Access

- Used to avoid **programmed I/O** (one byte at a time) for large data movement
- Requires **DMA** controller
- Bypasses CPU to transfer data directly between I/O device and memory
- OS writes DMA command block into memory
  - Source and destination addresses
  - Read or write mode
  - Count of bytes
  - Writes location of command block to DMA controller
  - Bus mastering of DMA controller – grabs bus from CPU
    - **Cycle stealing** from CPU but still much more efficient
  - When done, interrupts to signal completion
- Version that is aware of virtual addresses can be even more efficient - **DVMA**

# Six Step Process to Perform DMA Transfer



1. device driver is told to transfer drive2 data to buffer at address "x"

CPU

cache

2. device driver tells drive controller to transfer "c" bytes to buffer at address "x"

5. when c = 0, DMA interrupts CPU to signal transfer completion

CPU memory bus / controller

x

memory            buffer

PCIe bus

SAS drive controller

3. drive controller initiates DMA transfer

4. DMA controller transfers bytes to buffer "x", increasing memory address and decreasing "c" until c = 0

drive 1          drive 2

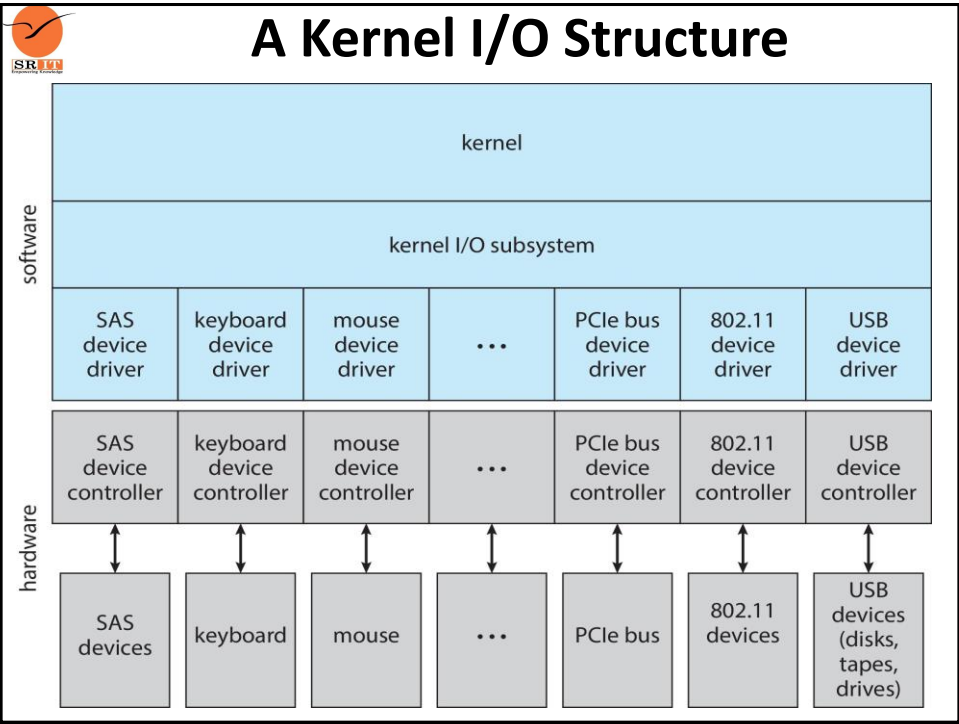# *Application I/O Interface*

*Narasimhulu M*$_{M. Tech.}$
*Assistant Professor*
*Department of Computer Science & Engineering*

## Application I/O Interface

- I/O system calls encapsulate device behaviors in generic classes
- Device-driver layer hides differences among I/O controllers from kernel
- New devices talking already-implemented protocols need no extra work
- Each OS has its own I/O subsystem structures and device driver frameworks
- Devices vary in many dimensions
  - **Character-stream** or **block**
  - **Sequential** or **random-access**
  - **Synchronous** or **asynchronous** (or both)
  - **Sharable** or **dedicated**
  - **Speed of operation**
  - **read-write**, **read only**, or **write only**

# A Kernel I/O Structure

| software | kernel | | | | | | |
|---|---|---|---|---|---|---|---|
| | kernel I/O subsystem | | | | | | |
| | SAS device driver | keyboard device driver | mouse device driver | ... | PCIe bus device driver | 802.11 device driver | USB device driver |
| hardware | SAS device controller | keyboard device controller | mouse device controller | ... | PCIe bus device controller | 802.11 device controller | USB device controller |
| | SAS devices | keyboard | mouse | ... | PCIe bus | 802.11 devices | USB devices (disks, tapes, drives) |

# Characteristics of I/O Devices

| aspect | variation | example |
|---|---|---|
| data-transfer mode | character block | terminal disk |
| access method | sequential random | modem CD-ROM |
| transfer schedule | synchronous asynchronous | tape keyboard |
| sharing | dedicated sharable | tape keyboard |
| device speed | latency seek time transfer rate delay between operations | |
| I/O direction | read only write only read–write | CD-ROM graphics controller disk |

# Characteristics of I/O Devices (Cont.)

- Subtleties of devices handled by device drivers
- Broadly I/O devices can be grouped by the OS into
  - Block I/O
  - Character I/O (Stream)
  - Memory-mapped file access
  - Network sockets
- For direct manipulation of I/O device specific characteristics, usually an escape / back door
  - Unix `ioctl()` call to send arbitrary bits to a device control register and data to device data register
- UNIX and Linux use tuple of "major" and "minor" device numbers to identify type and instance of devices (here major 8 and minors 0-4)

```
% ls -l /dev/sda*

brw-rw---- 1 root disk 8, 0 Mar 16 09:18 /dev/sda
brw-rw---- 1 root disk 8, 1 Mar 16 09:18 /dev/sda1
brw-rw---- 1 root disk 8, 2 Mar 16 09:18 /dev/sda2
brw-rw---- 1 root disk 8, 3 Mar 16 09:18 /dev/sda3
```

# Block and Character Devices

- Block devices include disk drives
  - Commands include read, write, seek
  - **Raw I/O**, **direct I/O**, or file-system access
  - Memory-mapped file access possible
    - File mapped to virtual memory and clusters brought via demand paging
  - DMA
- Character devices include keyboards, mice, serial ports
  - Commands include `get(), put()`
  - Libraries layered on top allow line editing

# Network Devices

- Varying enough from block and character to have own interface
- Linux, Unix, Windows and many others include **socket** interface
  - Separates network protocol from network operation
  - Includes **select()** functionality
- Approaches vary widely (pipes, FIFOs, streams, queues, mailboxes)

# Clocks and Timers

- Provide current time, elapsed time, triggering timer.
- Normal resolution about 1/60 second
- Some systems provide higher-resolution timers
- **Programmable interval timer** used for timings, periodic interrupts
- **ioctl()** (on UNIX) covers odd aspects of I/O such as clocks and timers

# Nonblocking and Asynchronous I/O

- **Blocking** - process suspended until I/O completed
  - Easy to use and understand
  - Insufficient for some needs
- **Non blocking** - I/O call returns as much as available
  - User interface, data copy (buffered I/O)
  - Implemented via multi-threading
  - Returns quickly with count of bytes read or written
  - `select()` to find if data ready then `read()` or `write()` to transfer
- **Asynchronous** - process runs while I/O executes
  - Difficult to use
  - I/O subsystem signals process when I/O completed

# Two I/O Methods

# Vectored I/O

- **Vectored I/O** allows one system call to perform multiple I/O operations
- For example, Unix `readve()` accepts a vector of multiple buffers to read into or write from.
- This scatter-gather method better than multiple individual I/O calls
  - Decreases **context switching and system call overhead**
  - Some versions provide **atomicity**
    - Avoid for example worry about multiple threads changing data as reads / writes occurring

# *Kernel I/O Subsystem*

*Narasimhulu M*<sub>M. Tech.</sub>
*Assistant Professor*
*Department of Computer Science & Engineering*

# Kernel I/O Subsystem

- **Scheduling**
  - Some I/O request ordering via per-device queue
  - Some OSs try fairness
  - Some implement Quality Of Service (i.e. IPQOS)
- **Buffering** - store data in memory while transferring between devices
  - To manage with **device speed mismatch**
  - To manage with **device transfer size mismatch**
  - To Support "**copy semantics**" for an Application I/O
  - **Double buffering** – **This double buffering** decouples the producer of data from the consumer, thus relaxing timing requirements between them.

# Device-status Table

| device: keyboard<br>status: idle | |
| device: laser printer<br>status: busy | → | request for<br>laser printer<br>address: 38546<br>length: 1372 |
| device: mouse<br>status: idle | |
| device: disk unit 1<br>status: idle | |
| device: disk unit 2<br>status: busy | → | request for<br>disk unit 2<br><br>file: xxx<br>operation: read<br>address: 43046<br>length: 20000 | → | request for<br>disk unit 2<br><br>file: yyy<br>operation: write<br>address: 03458<br>length: 500 |
| ⋮ | |

Common PC and Data-center I/O devices and Interface Speeds

# Kernel I/O Subsystem

- **Caching** - faster device holding copy of data
  – Always just a copy
  – Key to performance
  – Sometimes combined with buffering
- **Spooling** - **hold output for a device**
  – If device can serve only one request at a time
  – i.e., Printing
- **Device reservation** - **provides exclusive access to a device**
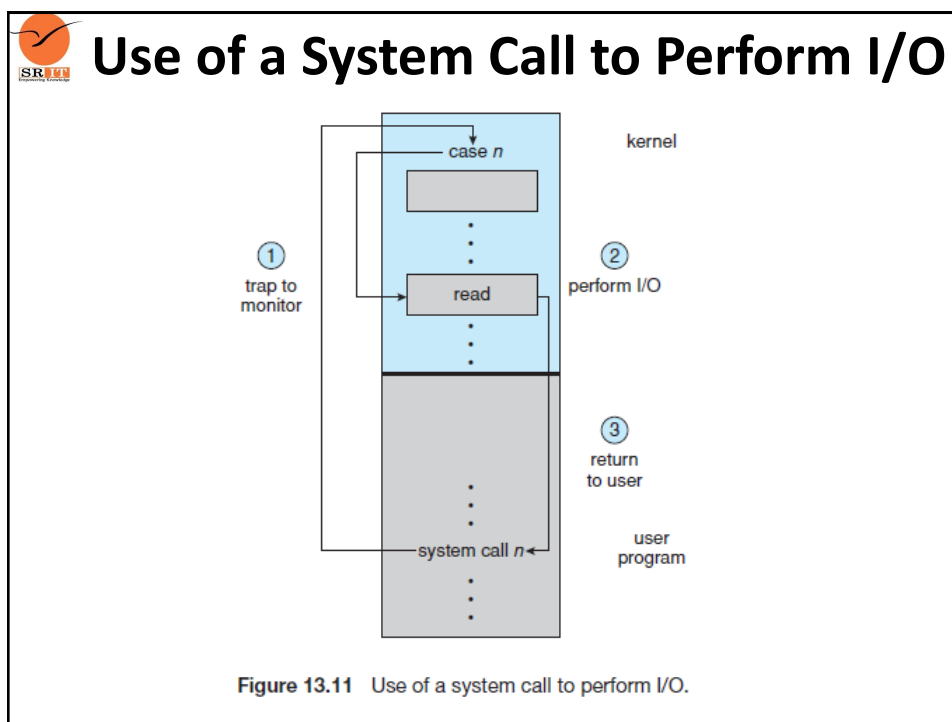  – System calls for allocation and de-allocation
  – Watch out for deadlock

# Error Handling

- OS can recover from **disk read**, **device unavailable**, **transient write failures.**
  - Retry a read or write, for example
  - Some systems more advanced – Solaris FMA, AIX
    - Track error frequencies, stop using device with increasing frequency of retry-able errors
- Most return an error number or code when I/O request fails
- System error logs hold problem reports

# I/O Protection

- User process may accidentally or purposefully attempt to disrupt normal operation via illegal I/O instructions
  - All I/O instructions defined to be privileged
  - I/O must be performed via system calls
    - Memory-mapped and I/O port memory locations must be protected too

## Use of a System Call to Perform I/O



Figure 13.11   Use of a system call to perform I/O.

## Kernel Data Structures

- Kernel keeps state information for I/O components, including open file tables, network connections, character device state
- Many, many complex data structures to track buffers, memory allocation, "dirty" blocks
- Some use object-oriented methods and message passing to implement I/O
  - Windows uses message passing
    - Message with I/O information passed from user mode into kernel
    - Message modified as it flows through to device driver and back to process
    - Pros / cons?

# UNIX I/O Kernel Structure



# Power Management

- Not strictly domain of I/O, but much is I/O related
- Computers and devices use electricity, generate heat, frequently require cooling
- OSes can help manage and improve use
  - Cloud computing environments move virtual machines between servers
    - Can end up evacuating whole systems and shutting them down
- Mobile computing has power management as first class OS aspect

# Power Management (Cont.)

- For example, Android implements
  - Component-level power management
    - Understands relationship between components
    - Build device tree representing physical device topology
    - System bus -> I/O subsystem -> {flash, USB storage}
    - Device driver tracks state of device, whether in use
    - Unused component – turn it off
    - All devices in tree branch unused – turn off branch
  - Wake locks – like other locks but prevent sleep of device when lock is held
  - Power collapse – put a device into very deep sleep
    - Marginal power use
    - Only awake enough to respond to external stimuli (button press, incoming call)
- Modern systems use **advanced configuration and power interface** (**ACPI**) firmware providing code that runs as routines called by kernel for device discovery, management, error and power management

# Kernel I/O Subsystem Summary

- In summary, the I/O subsystem coordinates an extensive collection of services that are available to applications and to other parts of the kernel
  - Management of the name space for files and devices
  - Access control to files and devices
  - Operation control (for example, a modem cannot seek())
  - File-system space allocation
  - Device allocation
  - Buffering, caching, and spooling
  - I/O scheduling
  - Device-status monitoring, error handling, and failure recovery
  - Device-driver configuration and initialization
  - Power management of I/O devices
- The upper levels of the I/O subsystem access devices via the uniform interface provided by the device drivers

# *Transforming I/O Requests to Hardware Operations*
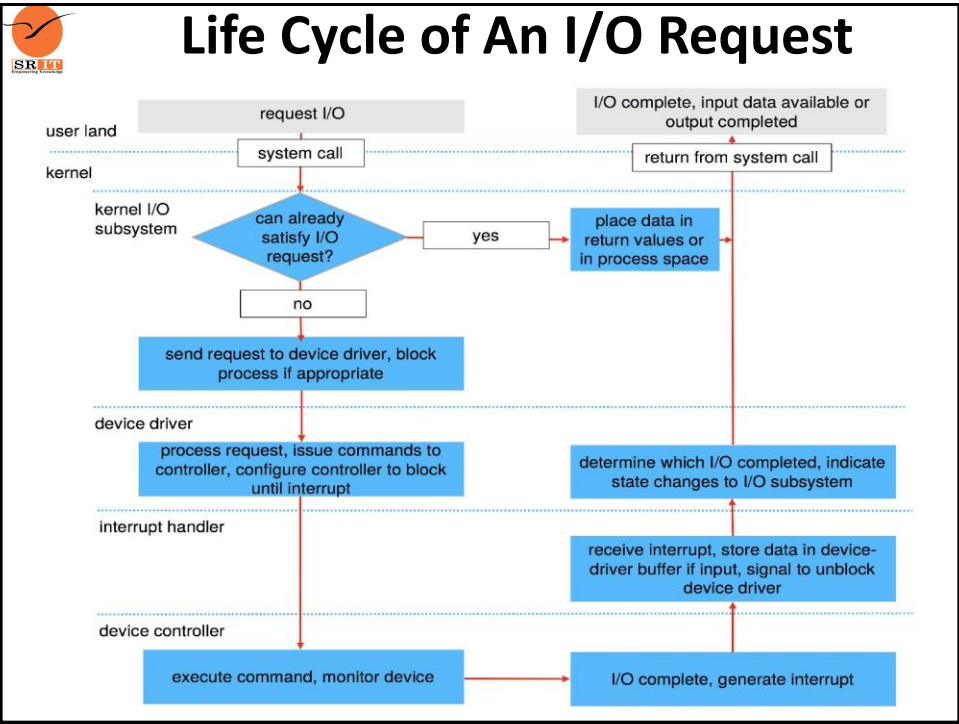
*Narasimhulu M$_{M. Tech.}$*
*Assistant Professor*
*Department of Computer Science & Engineering*

---

**Transforming I/O Requests to Hardware Operations**

- Consider reading a file from disk for a process:
  - Determine device holding file
  - Translate name to device representation
  - Physically read data from disk into buffer
  - Make data available to requesting process
  - Return control to process

# Life Cycle of An I/O Request





# END of Chapter - 2

# Chapter 3
# File-System

*Narasimhulu M*$_{M. Tech.}$
*Assistant Professor*
*Department of Computer Science & Engineering*

# *File Concept*

*Narasimhulu M*$_{M. Tech.}$
*Assistant Professor*
*Department of Computer Science & Engineering*

# File Concept

- The operating system abstracts from the physical properties of its storage devices to define a logical storage unit, **the file.**
- A file is a named collection of related information that is recorded on secondary storage
- Contiguous logical address space
- Types:
  - Data : Numeric, Character, Binary
  - Program
- Contents defined by file's creator
  - Many types like **text file, source file, executable file**

# File Attributes

- **Name** – only information kept in human-readable form
- **Identifier** – unique tag (number) identifies file within file system
- **Type** – needed for systems that support different types
- **Location** – pointer to a device and location of file on that device
- **Size** – current file size
- **Protection** – controls who can do reading, writing, executing
- **Time, date, and user identification** – data for protection, security, and usage monitoring.
- Many variations, including extended file attributes such as file checksum, type of Encoding etc.,
- Information about files are kept in the directory structure, which is maintained on the disk

# File info Window on Mac OS X



# Directory Structure

- A collection of nodes containing information about all files



- Both the directory structure and the files reside on disk

# File Operations

- **Create**
- **Write –** at **write pointer** location
- **Read –** at **read pointer** location
- **Reposition within file - seek**
- **Delete**
- **Truncate**
- *Open* (*F<sub>i</sub>*) – search the directory structure on disk for entry $F_i$, and move the content of entry to memory
- *Close* (*F<sub>i</sub>*) – move the content of entry $F_i$ in memory to directory structure on disk

# Open Files

- Several pieces of information  are associated with an  open files:
  - **Open-file table**: tracks open files
  - **File pointer**:  pointer to last read/write location, per process that has the file open
  - **File-open count**: counter of number of times a file is open – to allow removal of data from open-file table when last processes closes it.
  - **Disk location of the file:** cache of data access information
  - **Access rights**: per-process access mode information

# File Locking

- Provided by some operating systems and file systems
  - Similar to reader-writer locks
  - **Shared lock** similar to reader lock – several processes can acquire concurrently
  - **Exclusive lock** similar to writer lock
- Mediates access to a file
- Mandatory or advisory:(File-locking Mechanisms)
  - **Mandatory** – access is denied depending on locks held and requested
  - **Advisory** – processes can find status of locks and decide what to do

# File Locking Example – Java API

```
import java.io.*;
import java.nio.channels.*;
public class LockingExample {
    public static final boolean EXCLUSIVE = false;
    public static final boolean SHARED = true;
    public static void main(String arsg[]) throws IOException {
        FileLock sharedLock = null;
        FileLock exclusiveLock = null;
        try {
                RandomAccessFile raf = new RandomAccessFile("file.txt",
"rw");

                // get the channel for the file
                FileChannel ch = raf.getChannel();
                // this locks the first half of the file - exclusive
                exclusiveLock = ch.lock(0, raf.length()/2, EXCLUSIVE);
                /** Now modify the data . . . */
                // release the lock
                exclusiveLock.release();
```

### File Locking Example – Java API (Cont.)

```
        // this locks the second half of the file - shared
        sharedLock = ch.lock(raf.length()/2+1, raf.length(),
                SHARED);
        /** Now read the data . . . */
        // release the lock
        sharedLock.release();
} catch (java.io.IOException ioe) {
        System.err.println(ioe);
}finally {
        if (exclusiveLock != null)
        exclusiveLock.release();
        if (sharedLock != null)
        sharedLock.release();
    }
  }
}
```

# File Types – Name, Extension

| 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, pas, asm, a | source code in various languages |
| batch | bat, sh | commands to the command interpreter |
| text | txt, doc | textual data, documents |
| word processor | wp, tex, rtf, doc | various word-processor formats |
| library | lib, a, so, dll | libraries of routines for programmers |
| print or view | ps, pdf, jpg | ASCII or binary file in a format for printing or viewing |
| archive | arc, zip, tar | related files grouped into one file, sometimes compressed, for archiving or storage |
| multimedia | mpeg, mov, rm, mp3, avi | binary file containing audio or A/V information |

# File Structure

- File types also can be used to indicate the internal structure of the file.
- None - sequence of words, bytes
- Simple record structure
  - Lines
  - Fixed length
  - Variable length
- Complex Structures
  - Formatted document
  - Relocatable load file
- Can simulate last two with first method by inserting appropriate control characters
- Who decides:
  - Operating system
  - Program

# *Access Methods*

*Narasimhulu M*<sub>M. Tech.</sub>
*Assistant Professor*
*Department of Computer Science & Engineering*

# Access Methods

- A file is fixed length **logical records**
- **Sequential Access**
- **Direct Access**
- **Other Access Methods**

# Sequential Access

- Operations
  - **Read_next**
  - **Write_next**
  - **Reset**
  - no read after last write  (rewrite)

  Editors and Comipilers need sequential access
- Figure

# Direct Access

- Operations
  - **read *n***
  - **write *n***
  - **position to *n***
    - **Read_next**
    - **Write_next**
    - **rewrite *n***

  *n* = **relative block number**

- Relative block numbers allow OS to decide where file should be placed
- Databases need direct acess

## Simulation of Sequential Access on Direct-access File

| sequential access | implementation for direct access |
|---|---|
| reset | cp = 0; |
| read_next | read cp ;<br>cp = cp + 1; |
| write_next | write cp;<br>cp = cp + 1; |

**Figure 11.5**   Simulation of sequential access on a direct-access file.

# Other Access Methods

- Can be other access methods built on top of direct-acesss methods.
- General involve creation of an **index** for the file
- Keep index in memory for fast determination of location of data to be operated on (consider Universal Produce Code (UPC code) plus record of data about that item)
- If the index is too large, create an in-memory index, which an index of a disk index
- IBM indexed sequential-access method (ISAM)
  - Small master index, points to disk blocks of secondary index
  - File kept sorted on a defined key
  - All done by the OS
- VMS operating system provides index and relative files as another example (see next slide)
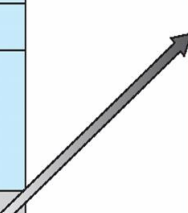
# Example of Index and Relative Files

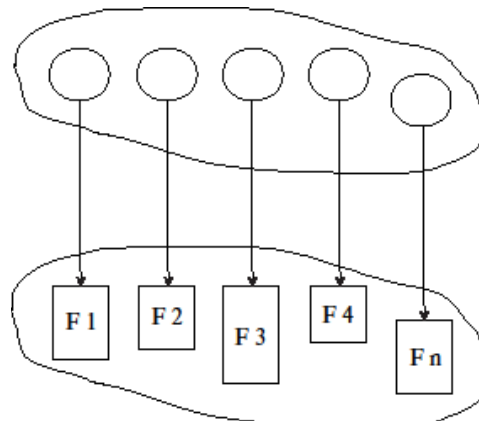# *Directory Structure*

*Narasimhulu M*<sub>M. Tech.</sub>
*Assistant Professor*
*Department of Computer Science & Engineering*

---

# Directory Structure

- A collection of nodes containing information about all files

# Operations Performed on Directory

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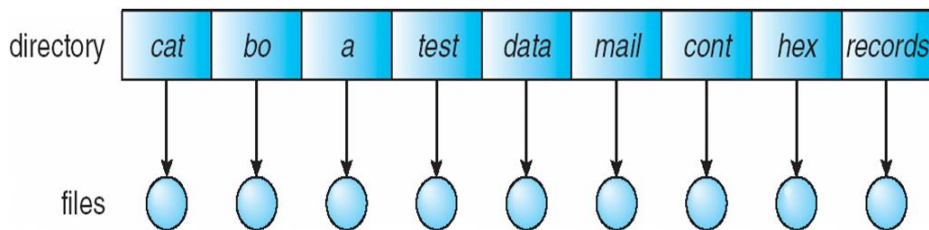- Traverse the file system

# Directory Organization

**The directory is organized logically to obtain**

- Efficiency – locating a file quickly
- Naming – convenient to users
  - Two users can have same name for different files
  - The same file can have several different names
- Grouping – logical grouping of files by properties, (e.g., all Java programs, all games, ...)

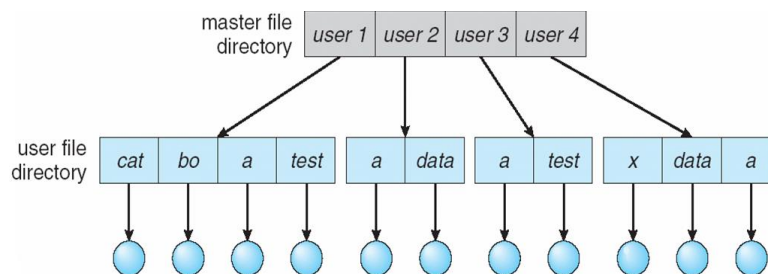# Single-Level Directory
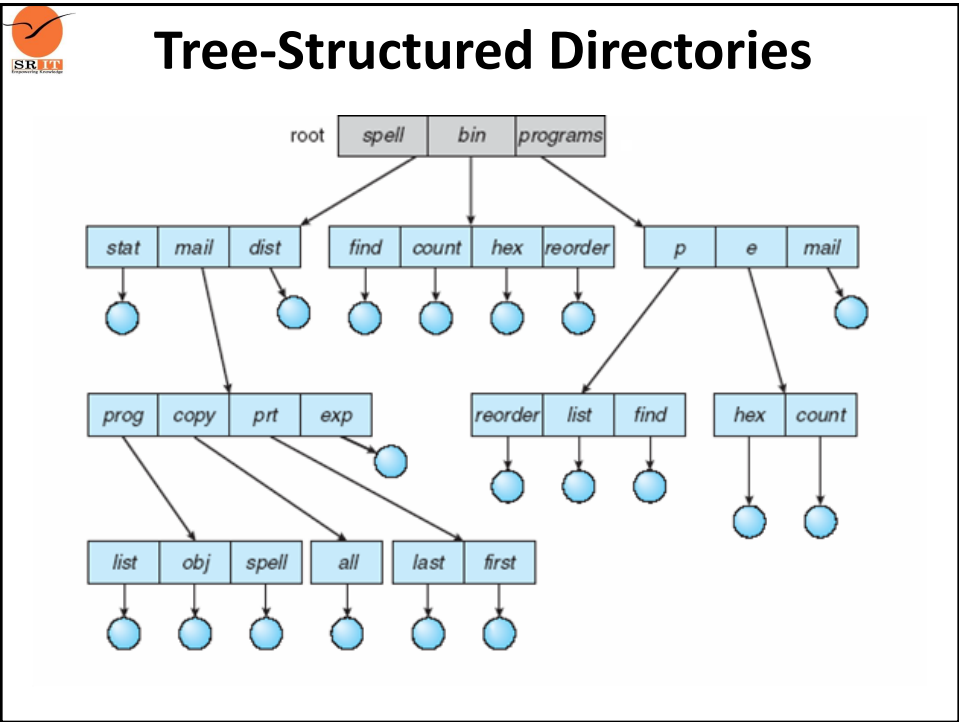
- A single directory for all users

| directory | cat | bo | a | test | data | mail | cont | hex | records |

files

- Naming problem
- Grouping problem

# Two-Level Directory

- **Separate directory for each user**

master file directory: | user 1 | user 2 | user 3 | user 4 |

user file directory: | cat | bo | a | test | a | data | a | test | x | data | a |

- Path name
- Can have the same file name for different user
- Efficient searching
- No grouping capability

# Tree-Structured Directories



# Acyclic-Graph Directories

- **Have shared subdirectories and files**
- **Example**

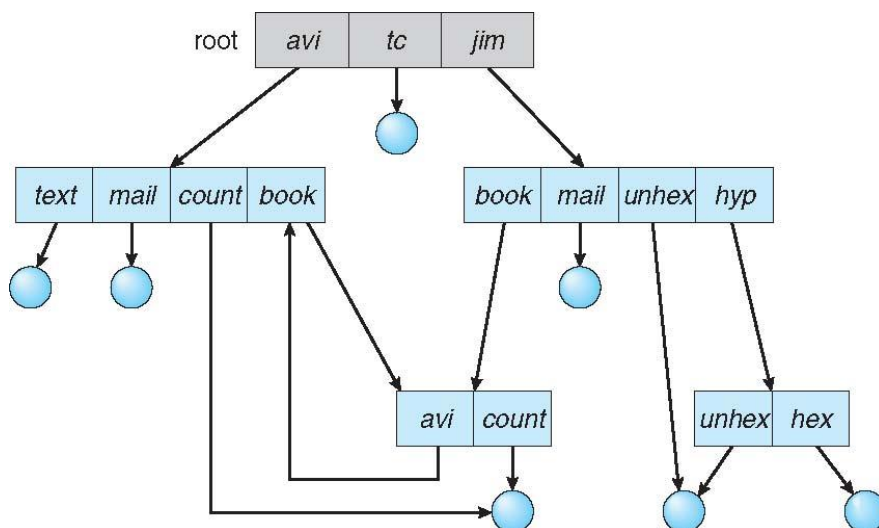# Acyclic-Graph Directories (Cont.)

- Two different names (aliasing)
- If **dict** deletes **w/list** $\Rightarrow$ dangling pointer
  Solutions:
  - Backpointers, so we can delete all pointers.
    - Variable size records a problem
  - Backpointers using a daisy chain organization
  - Entry-hold-count solution
- New directory entry type
  - **Link** – another name (pointer) to an existing file
  - **Resolve the link** – follow pointer to locate the file

# General Graph Directory

# General Graph Directory (Cont.)

- How do we guarantee no cycles?
  - Allow only links to files not subdirectories
  - **Garbage collection**
  - Every time a new link is added use a cycle detection algorithm to determine whether it is OK
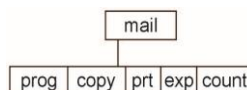
# Current Directory

- Can designate one of the directories as the current (working) directory
  - `cd /spell/mail/prog`
  - `type list`
- Creating and deleting a file is done in current directory
- Example of creating a new file
  - If in current directory is `/mail`
  - The command

    **mkdir <dir-name>**

  - Results in:

    | mail | | | | |
    |------|------|-----|-----|-------|
    | prog | copy | prt | exp | count |

  - Deleting "mail" $\Rightarrow$ deleting the entire subtree rooted by "mail"

# *Protection*

*Narasimhulu M*<sub>M. Tech.</sub>
*Assistant Professor*
*Department of Computer Science & Engineering*

---

# 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)

- 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.

# Protection

- File owner/creator should be able to control:
  - What can be done
  - By whom
- Types of access
  - **Read**
  - **Write**
  - **Execute**
  - **Append**
  - **Delete**
  - **List**

# Access Lists and Groups in Unix

- Mode of access:  read, write, execute
- Three classes of users on Unix / Linux

|                  |     |               | RWX   |
|------------------|-----|---------------|-------|
| a) **owner access** | 7 | $\Rightarrow$ | 1 1 1 |
| b) **group access** | 6 | $\Rightarrow$ | 1 1 0 |
| c) **public access** | 1 | $\Rightarrow$ | 0 0 1 |

- Ask manager to create a group (unique name), say G, and add some users to the group.

```
chgrp    G    game
```

- For a file (say *game*) or subdirectory, define an appropriate access.

```
owner  group  public

chmod  761  game
```

- Attach a group to a file
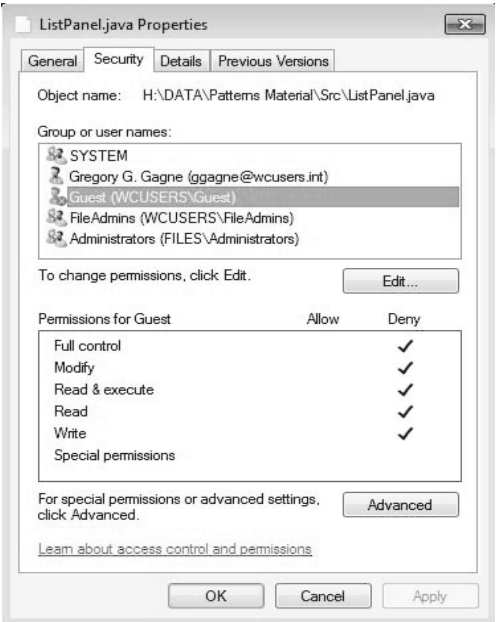
# A Sample UNIX Directory Listing

```
-rw-rw-r--      1 pbg    staff      31200   Sep 3 08:30    intro.ps
drwx------      5 pbg    staff        512   Jul 8 09.33    private/
drwxrwxr-x      2 pbg    staff        512   Jul 8 09:35    doc/
drwxrwx---      2 pbg    student      512   Aug 3 14:13    student-proj/
-rw-r--r--      1 pbg    staff       9423   Feb 24 2003    program.c
-rwxr-xr-x      1 pbg    staff      20471   Feb 24 2003    program
drwx--x--x      4 pbg    faculty      512   Jul 31 10:31   lib/
drwx------      3 pbg    staff       1024   Aug 29 06:52   mail/
drwxrwxrwx      3 pbg    staff        512   Jul 8 09:35    test/
```

# Windows 7 Access-Control List Management

# *Memory-Mapped Files*

*Narasimhulu M*$_{M. Tech.}$
*Assistant Professor*
*Department of Computer Science & Engineering*

---

## *Memory-Mapped Files*

- A part of the virtual Address space to be allocated logically associated with the file is referred as Memory-mapped File.

- Memory mapping a file is accomplished by mapping a disk block to a page (or pages) in memory.

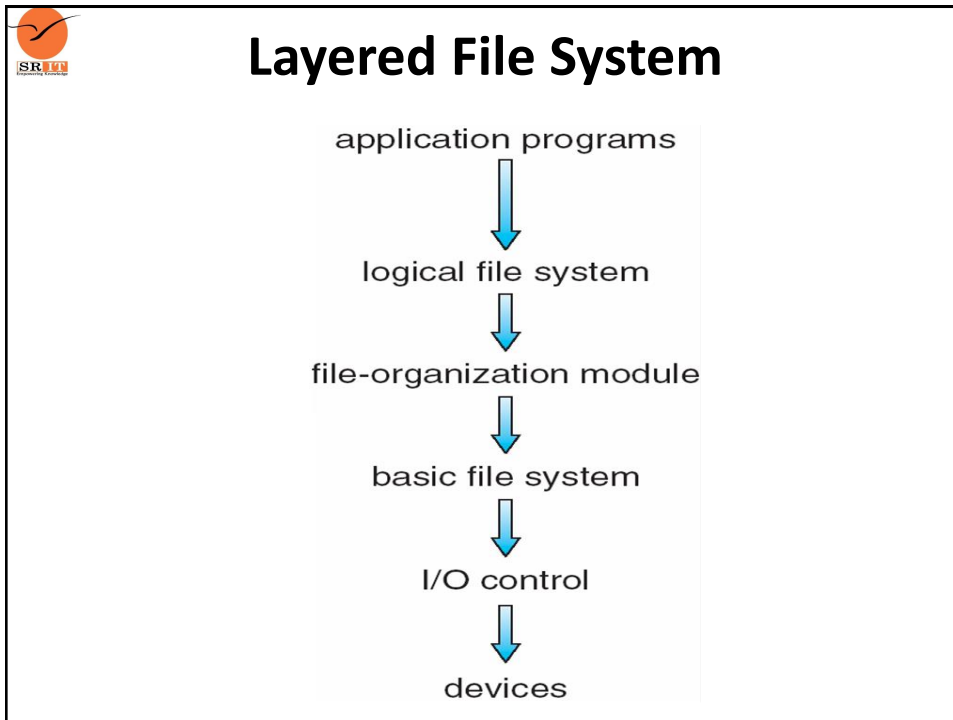# *File system structure and Implementation*

*Narasimhulu M*$_{M. Tech.}$
*Assistant Professor*
*Department of Computer Science & Engineering*

---

# File-System Structure

- File structure
  - Logical storage unit
  - Collection of related information
- **File system** resides on secondary storage (disks)
  - Provided user interface to storage, mapping logical to physical
  - Provides efficient and convenient access to disk by allowing data to be stored, located retrieved easily
- Disk provides in-place rewrite and random access
  - I/O transfers performed in **blocks** of **sectors** (usually 512 bytes)
- **File control block** – storage structure consisting of information about a file
- **Device driver** controls the physical device
- File system organized into layers

# Layered File System



application programs

↓

logical file system

↓

file-organization module

↓

basic file system

↓

I/O control

↓

devices

# File System Layers

- **Device drivers** manage I/O devices at the I/O control layer
  - Given commands like "read drive1, cylinder 72, track 2, sector 10, into memory location 1060" outputs low-level hardware specific commands to hardware controller
- **Basic file system** given command like "retrieve block 123" translates to device driver
- Also manages memory buffers and caches (allocation, freeing, replacement)
  - Buffers hold data in transit
  - Caches hold frequently used data
- **File organization module** understands files, logical address, and physical blocks
- ■ Translates logical block # to physical block #
- ■ Manages free space, disk allocation

# File System Layers (Cont.)

- **Logical file system** manages metadata information
  - Translates file name into file number, file handle, location by maintaining file control blocks (**inodes** in UNIX)
  - Directory management
  - Protection
- Layering useful for reducing complexity and redundancy, but adds overhead and can decrease Performance Translates file name into file number, file handle, location by maintaining file control blocks (**inodes** in UNIX)
  - Logical layers can be implemented by any coding method according to OS designer

# File-System Implementation

- We have system calls at the API level, but how do we implement their functions?
  - On-disk and in-memory structures
- **Boot control block** contains info needed by system to boot OS from that volume
  - Needed if volume contains OS, usually first block of volume
- **Volume control block (superblock(unix), master file table(windows))** contains volume details
  - Total # of blocks, # of free blocks, block size, free block pointers or array
- Directory structure organizes the files
  - (unix)Names and inode numbers, (windows)master file table

# File-System Implementation (Cont.)

- Per-file **File Control Block** (**FCB**) contains many details about the file
  - inode number, permissions, size, dates
  - NFTS stores into in master file table using relational DB structures

| |
|---|
| file permissions |
| file dates (create, access, write) |
| file owner, group, ACL |
| file size |
| file data blocks or pointers to file data blocks |

# In-Memory File System Structures

- The in-memory information is used for both file-system management and performance improvement via caching.

- Several types of structures may be included.

- An in-memory mount table contains information about each mounted volume.

- An in-memory directory-structure cache holds the directory information of recently accessed directories.

- The **system-wide open-file table** contains a copy of the FCB of each open file.
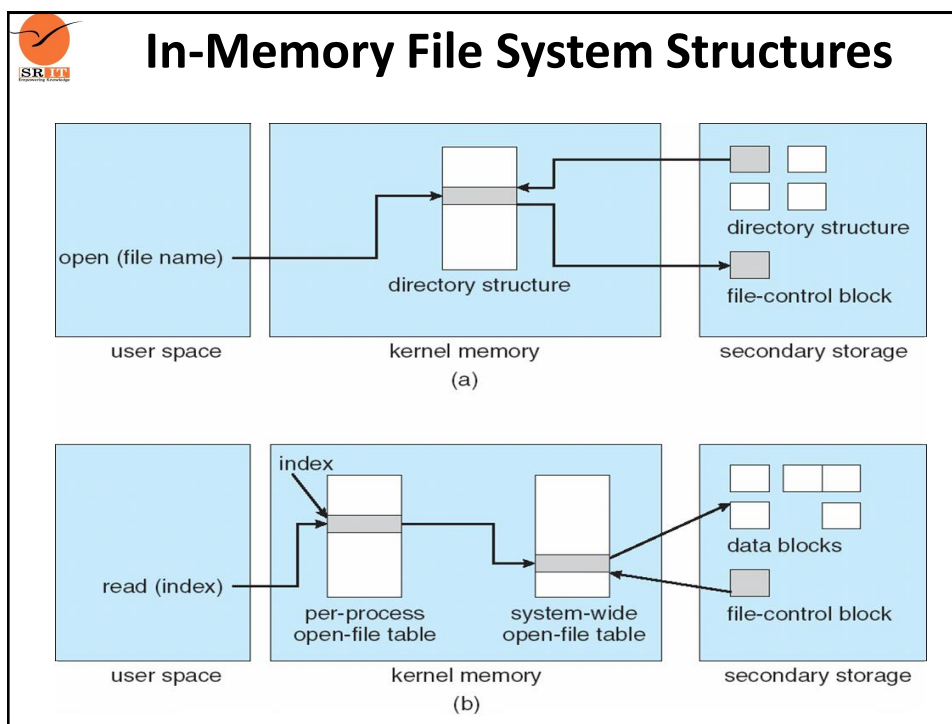
# In-Memory File System Structures

- The **per-process open-file table** contains a pointer to the appropriate entry in the system-wide open-file table
- Buffers hold file-system blocks when they are being read from disk or written to disk.

# In-Memory File System Structures

- Mount table storing file system mounts, mount points, file system types
- The following figure illustrates the necessary file system structures provided by the operating systems
- Figure 12-3(a) refers to opening a file
- Figure 12-3(b) refers to reading a file
- Plus buffers hold data blocks from secondary storage
- Open returns a file handle for subsequent use
- Data from read eventually copied to specified user process memory address

# In-Memory File System Structures



# Partitions and Mounting

- Partition can be a volume containing a file system ("cooked") or **raw** – just a sequence of blocks with no file system
- Boot block can point to boot volume or boot loader set of blocks that contain enough code to know how to load the kernel from the file system
  - Or a boot management program for multi-os booting
- **Root partition** contains the OS, other partitions can hold other Oses, other file systems, or be raw
  - Mounted at boot time
  - Other partitions can mount automatically or manually
- At mount time, file system consistency checked
  - Is all metadata correct?
    - If not, fix it, try again
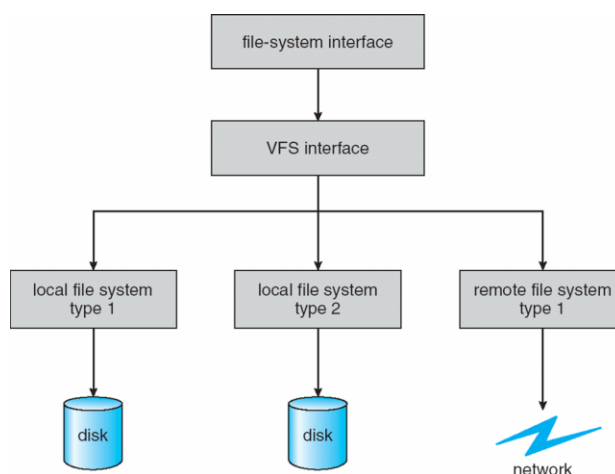    - If yes, add to mount table, allow access

# Virtual File Systems

- **Virtual File Systems** (**VFS**) on Unix provide an object-oriented way of implementing file systems
- VFS allows the same system call interface (the API) to be used for different types of file systems
  - Separates file-system generic operations from implementation details
  - Implementation can be one of many file systems types, or network file system
    - Implements **vnodes** which hold inodes or network file details
  - Then dispatches operation to appropriate file system implementation routines

# Virtual File Systems (Cont.)

- The API is to the VFS interface, rather than any specific type of file system

# Virtual File System Implementation

- For example, Linux has four object types:
  - inode, file, superblock, dentry
- VFS defines set of operations on the objects that must be implemented
  - Every object has a pointer to a function table
    - Function table has addresses of routines to implement that function on that object
    - For example:
    - • `int open(. . .)` —Open a file
    - • `int close(. . .)` —Close an already-open file
    - • `ssize t read(. . .)` —Read from a file
    - • `ssize t write(. . .)` —Write to a file
    - • `int mmap(. . .)` —Memory-map a file

# END of Unit-4