

# Modulo VI

## Sistemas de Entrada e Saída

*Prof. Ismael H F Santos*

## Ementa

- Gerência de Entrada e Saída – Armazenamento
  - Disk Structure
  - Disk Scheduling
  - Disk Management
  - RAID Structure
  - Tertiary Storage Devices
  - Operating System Issues
- Gerência de Entrada e Saída – Implementação
  - I/O Hardware - DMA
  - Application I/O Interface
  - Kernel I/O Subsystem

## SOP – CO023

Gerência de  
E/S



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

3

## Objectives

- Describe the physical structure of secondary and tertiary storage devices and the resulting effects on the uses of the devices
- Explain the performance characteristics of mass-storage devices
- Discuss operating-system services provided for mass storage, including **RAID** and **HSM**

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

4

## Overview of Mass Storage Structure

- Magnetic disks provide bulk of secondary storage of modern computers
  - Drives rotate at 60 to 200 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**

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

5

## Overview of Mass Storage Structure

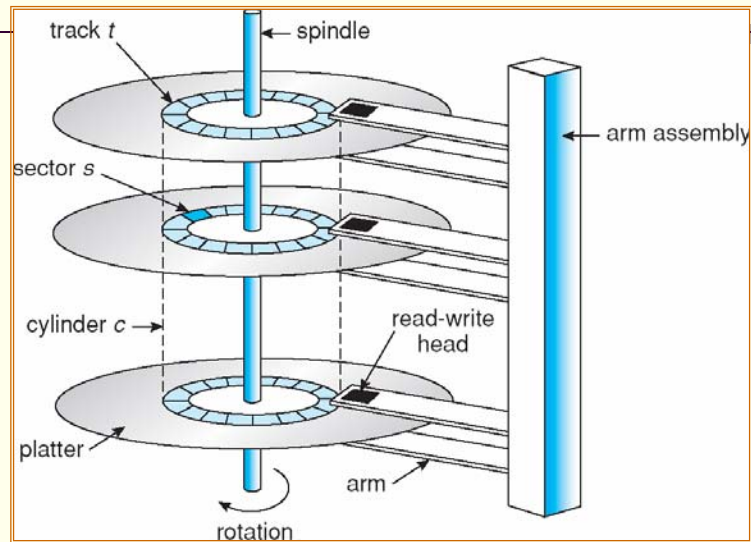
- Disks can be removable
- Drive attached to computer via **I/O bus**
  - Busses vary, including **EIDE, ATA, SATA, USB, Fibre Channel, SCSI**
  - **Host controller** in computer uses bus to talk to **disk controller** built into drive or storage array

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

6

## Moving-head Disk Mechanism



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

7

## Overview of Mass Storage Structure (Cont.)

### ■ Magnetic tape

- Was early secondary-storage medium
- 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
- 20-200GB typical storage
- Common technologies are 4mm, 8mm, 19mm, LTO-2 and SDLT

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

8

## SOP – CO023

*Estrutura  
Do  
Disco*



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

9

## Disk Structure

- Disk drives are addressed as large 1-dimensional *arrays of logical blocks*, where the logical block is the smallest unit of transfer.
- The 1-dimensional array of logical blocks is mapped into the sectors of the disk sequentially.
  - Sector 0 is the first sector of the first track on the outermost cylinder.
  - Mapping proceeds in order through that track, then the rest of the tracks in that cylinder, and then through the rest of the cylinders from outermost to innermost.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

10

## Disk Attachment

- Host-attached storage accessed through I/O ports talking to I/O busses
- SCSI itself is a bus, up to 16 devices on one cable, **SCSI initiator** requests operation and **SCSI targets** perform tasks
  - Each target can have up to 8 **logical units** (disks attached to device controller)
- FC is high-speed serial architecture
  - Can be switched fabric with 24-bit address space – the basis of **storage area networks (SANs)** in which many hosts attach to many storage units
  - Can be **arbitrated loop (FC-AL)** of 126 devices

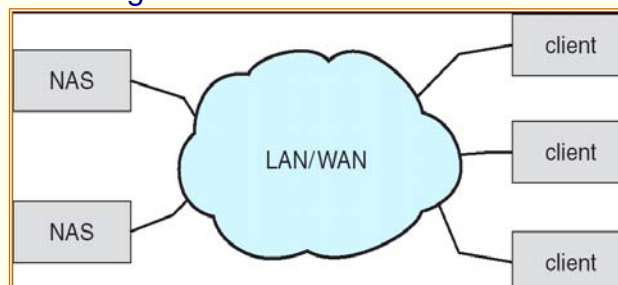
April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

11

## Network-Attached Storage

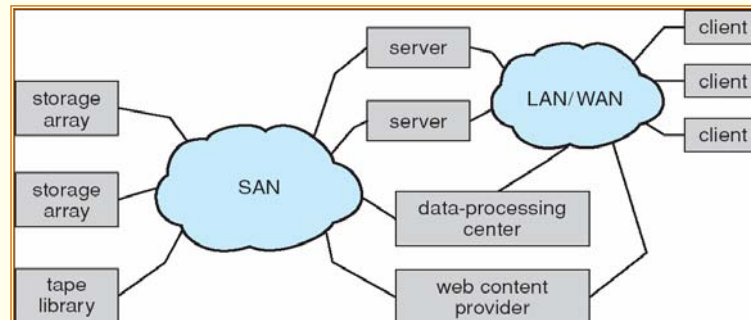
- Network-attached storage (**NAS**) is storage made available over a network rather than over a local connection (such as a bus)
- **NFS** and **CIFS** are common protocols
- Implemented via remote procedure calls (RPCs) between host and storage
- New **iSCSI** protocol uses IP network to carry the SCSI protocol



April 05

# Storage Area Network

- Common in large storage environments (and becoming more common)
- Multiple hosts attached to multiple storage arrays - flexible



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

13

# SOP – CO023

*Escalonamento  
Do  
Disco*



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

14

## 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.
- **Access time** has two major components
  - **Seek time** is the time for the disk are to move the heads to the cylinder containing the desired sector.
  - **Rotational latency** is the additional time waiting for the disk to rotate the desired sector to the disk head.
- Minimize seek time
- Seek time  $\approx$  seek distance

## Disk Scheduling (Cont.)

- 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.
- Several algorithms exist to schedule the servicing of disk I/O requests. We illustrate them with a request queue (0-199).

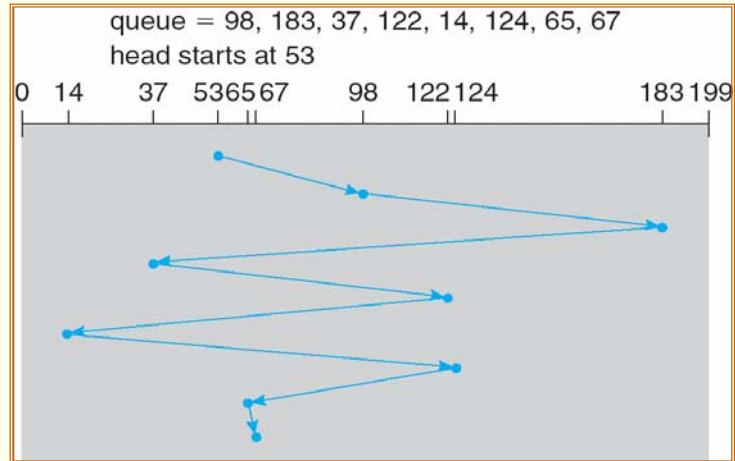
98, 183, 37, 122, 14, 124, 65, 67

Head pointer 53



# FCFS

Total head movement of 640 cylinders.



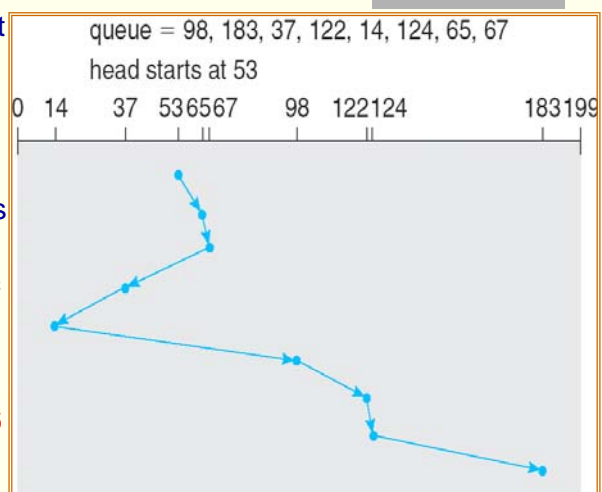
April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

17

# SSTF

- 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.
- **Total head movement of 236 cylinders !**



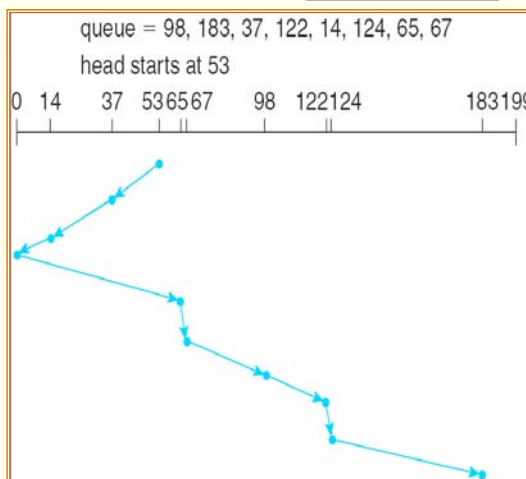
April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

18

## 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.
- Sometimes called the *elevator algorithm*.
- **Total head movement of 208 cylinders !**



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

19

## C-SCAN

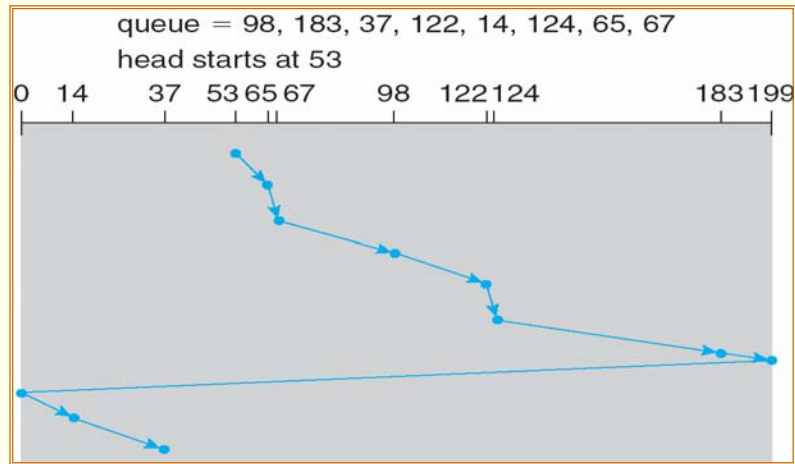
- Provides a **more uniform wait time** than SCAN.
- The head moves from one end of the disk to the other servicing requests as it goes. When it reaches the other end, however, it immediately returns to the beginning of the disk, without servicing any requests on the return trip.
- **Treats the cylinders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

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

20

## C-SCAN (Cont.)



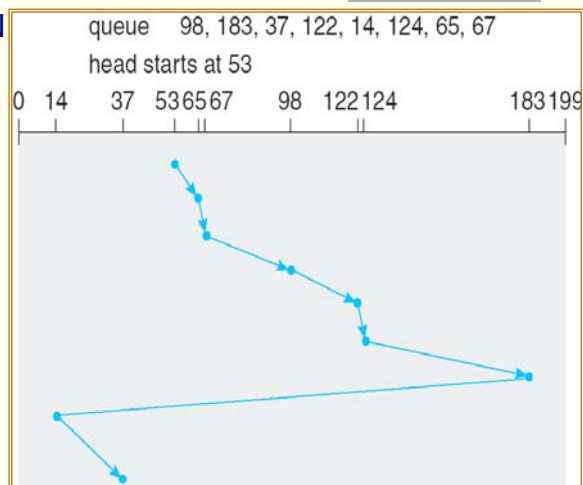
April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

21

## C-LOOK

- 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.



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

22

## 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.
- Performance depends on the number and types of requests.
- Requests for disk service can be influenced by the file-allocation method.
- 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 **C-LOOK** is a reasonable choice for the default algorithm.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

23

## SOP – CO023

Gerenciamento  
Do  
Disco



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

24

## Disk Management

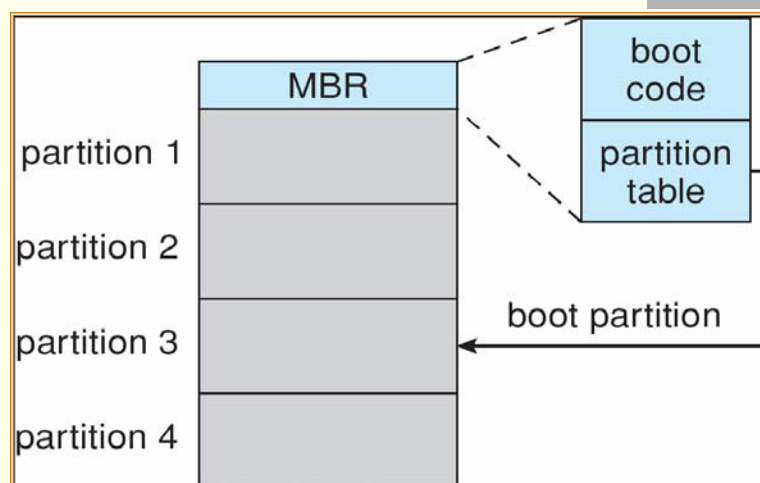
- *Low-level formatting, or physical formatting* — Dividing a disk into sectors that the disk controller can read and write.
- To use a disk to hold files, the operating system still needs to record its own data structures on the disk.
  - *Partition* the disk into one or more groups of cylinders.
  - *Logical formatting* or “making a file system”.
- Boot block initializes system.
  - The bootstrap is stored in ROM.
  - *Bootstrap loader* program.
- Methods such as *sector sparing* used to handle bad blocks.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

25

## Booting from a Disk in Windows 2000



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

26

## Swap-Space Management

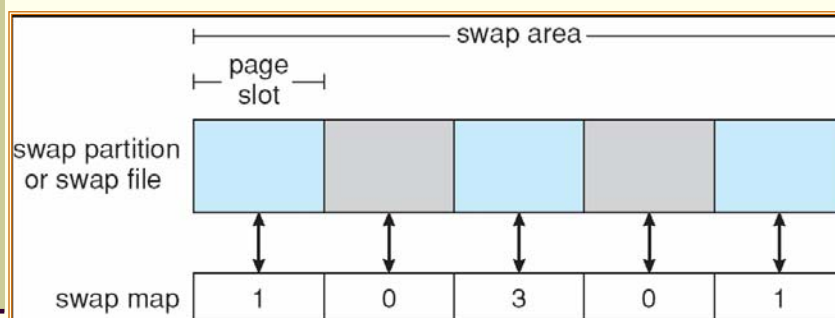
- **Swap-space** — Virtual memory uses disk space as an extension of main memory. Can be carved out of the normal file system, or, more commonly, it can be in a separate disk partition.
- **Swap-space management**
  - 4.3BSD allocates swap space when process starts; holds *text segment* (the program) and *data segment*.
  - Kernel uses *swap maps* to track swap-space use.
  - Solaris 2 allocates swap space only when a page is forced out of physical memory, not when the virtual memory page is first created.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

27

## Data Structures for Swapping on Linux Systems

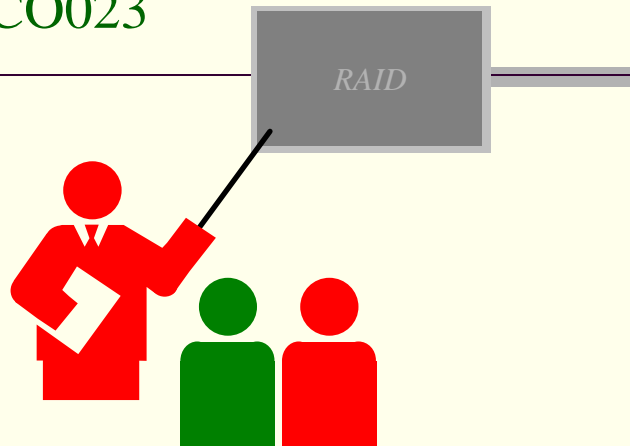


April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

28

## SOP – CO023



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

29

## RAID Structure

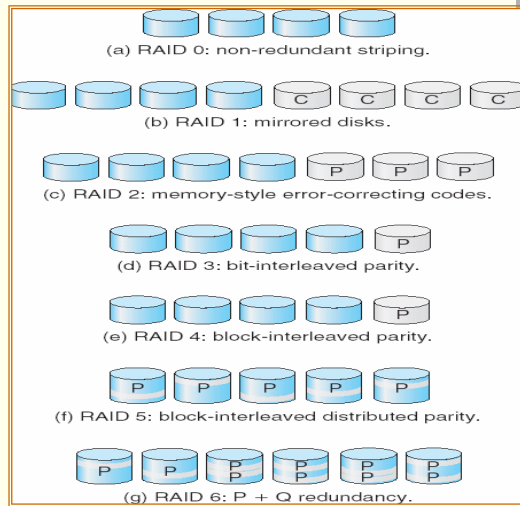
- **RAID** – multiple disk drives provides **reliability** via **redundancy**. RAID is arranged into six different levels.
- Several improvements in disk-use techniques involve the use of multiple disks working cooperatively.
- **Disk striping** uses a group of disks as one storage unit.
- RAID schemes improve **performance** and improve the **reliability** of the storage system by storing redundant data.
  - **Mirroring** or **shadowing** keeps duplicate of each disk.
  - **Block interleaved parity** uses much less redundancy.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

30

# RAID Levels

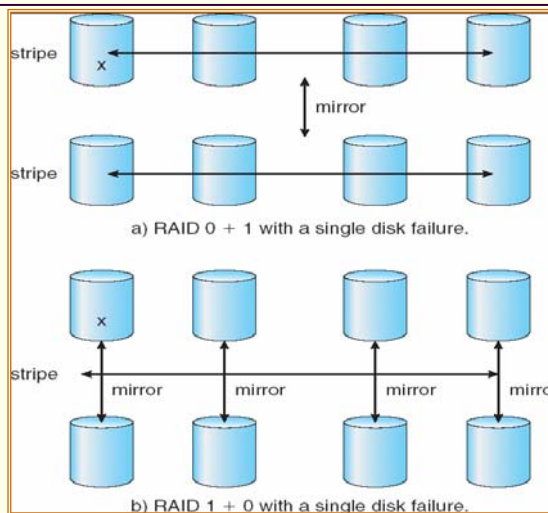


April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

31

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



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

32



## Stable-Storage Implementation

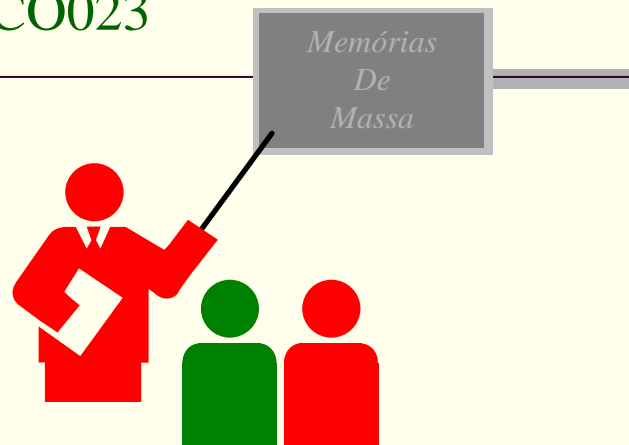
- Write-ahead log scheme requires stable storage.
- To implement stable storage:
  - Replicate information on more than one nonvolatile storage media with independent failure modes.
  - Update information in a controlled manner to ensure that we can recover the stable data after any failure during data transfer or recovery.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

33

## SOP – CO023



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

34

## Tertiary Storage Devices

---

- Low cost is the defining characteristic of tertiary storage.
- Generally, tertiary storage is built using *removable media*
- Common examples of removable media are floppy disks and CD-ROMs; other types are available.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

35

## Removable Disks

---

- Floppy disk — thin flexible disk coated with magnetic material, enclosed in a protective plastic case.
  - Most floppies hold about 1 MB; similar technology is used for removable disks that hold more than 1 GB.
  - Removable magnetic disks can be nearly as fast as hard disks, but they are at a greater risk of damage from exposure.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

36

## Removable Disks (Cont.)

- A magneto-optic disk records data on a rigid platter coated with magnetic material.
  - Laser heat is used to amplify a large, weak magnetic field to record a bit.
  - Laser light is also used to read data (Kerr effect).
  - The magneto-optic head flies much farther from the disk surface than a magnetic disk head, and the magnetic material is covered with a protective layer of plastic or glass; resistant to head crashes.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

37

## WORM Disks

- The data on read-write disks can be modified over and over.
- **WORM** ("Write Once, Read Many Times") disks can be written only once.
- Thin aluminum film sandwiched between two glass or plastic platters.
- To write a bit, the drive uses a laser light to burn a small hole through the aluminum; information can be destroyed by not altered.
- Very durable and reliable.
- **Read Only disks**, such as CD-ROM and DVD, come from the factory with the data pre-recorded.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

38

# Tapes

- Compared to a disk, a tape is less expensive and holds more data, but random access is much slower.
- Tape is an economical medium for purposes that do not require fast random access, e.g., **backup** copies of disk data, holding huge volumes of data.

# Tapes (cont.)

- Large tape installations typically use robotic tape changers that move tapes between tape drives and storage slots in a tape library.
  - stacker – library that holds a few tapes
  - silo – library that holds thousands of tapes
- A disk-resident file can be *archived* to tape for low cost storage; the computer can *stage* it back into disk storage for active use.

## Operating System Issues

---

- Major OS jobs are to manage physical devices and to present a virtual machine abstraction to applications
- For hard disks, the OS provides two abstraction:
  - **Raw device** – an array of data blocks.
  - **File system** – the OS queues and schedules the interleaved requests from several applications.

## Application Interface

---

- Most OSs handle removable disks almost exactly like fixed disks — a new cartridge is formatted and an empty file system is generated on the disk.
- Tapes are presented as a raw storage medium, i.e., and application does not open a file on the tape, it opens the whole tape drive as a raw device.
- Usually the tape drive is reserved for the exclusive use of that application.

## Application Interface (cont.)

---

- Since the OS does not provide file system services, the application must decide how to use the array of blocks.
- Since every application makes up its own rules for how to organize a tape, a tape full of data can generally only be used by the program that created it.

## Tape Drives

---

- The basic operations for a tape drive differ from those of a disk drive.
- **locate** positions the tape to a specific logical block, not an entire track (corresponds to **seek**).
- The **read position** operation returns the logical block number where the tape head is.
- The **space** operation enables relative motion.
- Tape drives are “append-only” devices; updating a block in the middle of the tape also effectively erases everything beyond that block.
- An **EOT** mark is placed after a block that is written.

## File Naming

- The issue of **naming files** on removable media is especially difficult when we want to write data on a removable cartridge on one computer, and then use the cartridge in another computer.
- Contemporary OSs generally leave the name space problem unsolved for removable media, and depend on applications and users to figure out how to access and interpret the data.
- Some kinds of removable media (e.g., CDs) are so well standardized that all computers use them the same way.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

45

## Hierarchical Storage Management (HSM)

- A **hierarchical storage system** extends the storage hierarchy beyond primary memory and secondary storage to incorporate tertiary storage — usually implemented as a jukebox of tapes or removable disks.
- Usually incorporate tertiary storage by extending the file system.
  - Small and frequently used files remain on disk.
  - Large, old, inactive files are archived to the jukebox.
- HSM is usually found in supercomputing centers and other large installations that have enormous volumes of data.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

46

## Speed

- Two aspects of speed in tertiary storage are bandwidth and latency.
- Bandwidth is measured in bytes per second.
  - **Sustained bandwidth** – average data rate during a large transfer; # of bytes/transfer time.  
Data rate when the data stream is actually flowing.
  - **Effective bandwidth** – average over the entire I/O time, including **seek** or **locate**, and cartridge switching.  
Drive's overall data rate.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

47

## Speed (Cont.)

- Access latency – amount of time needed to locate data.
  - Access time for a disk – move the arm to the selected cylinder and wait for the rotational latency; < 35 milliseconds.
  - Access on tape requires winding the tape reels until the selected block reaches the tape head; tens or hundreds of seconds.
  - Generally say that random access within a tape cartridge is about a thousand times slower than random access on disk.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

48



## Speed (Cont.)

---

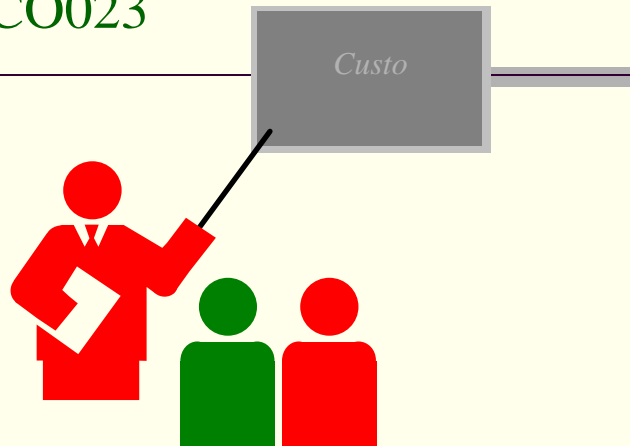
- The low cost of tertiary storage is a result of having many cheap cartridges share a few expensive drives.
- A removable library is best devoted to the storage of infrequently used data, because the library can only satisfy a relatively small number of I/O requests per hour.

## Reliability

---

- A fixed disk drive is likely to be more reliable than a removable disk or tape drive.
- An optical cartridge is likely to be more reliable than a magnetic disk or tape.
- A head crash in a fixed hard disk generally destroys the data, whereas the failure of a tape drive or optical disk drive often leaves the data cartridge unharmed.

## SOP – CO023



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

51

## Cost

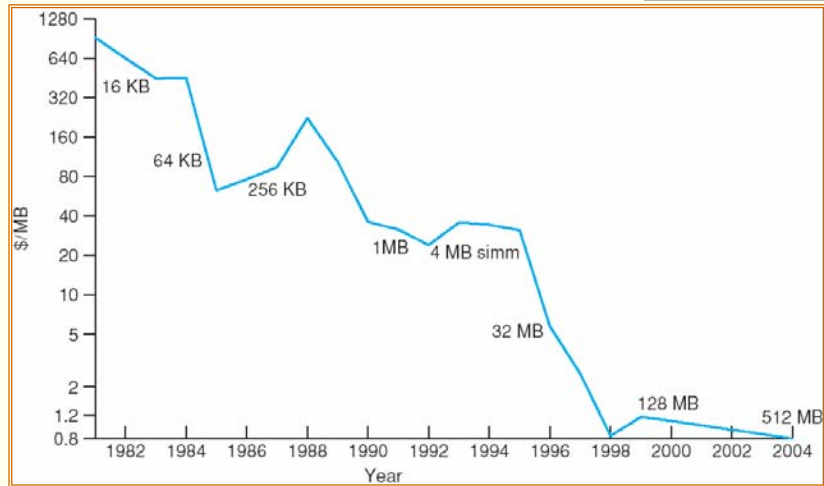
- Main memory is much more expensive than disk storage
- The cost per megabyte of hard disk storage is competitive with magnetic tape if only one tape is used per drive.
- The cheapest tape drives and the cheapest disk drives have had about the same storage capacity over the years.
- Tertiary storage gives a cost savings only when the number of cartridges is considerably larger than the number of drives.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

52

## Price per Megabyte of DRAM, From 1981 to 2004

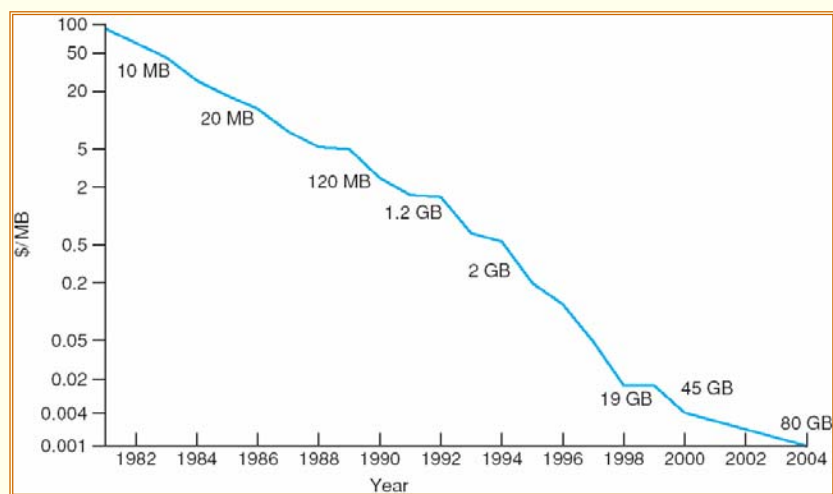


April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

53

## Price per Megabyte of Magnetic Hard Disk, From 1981 to 2004

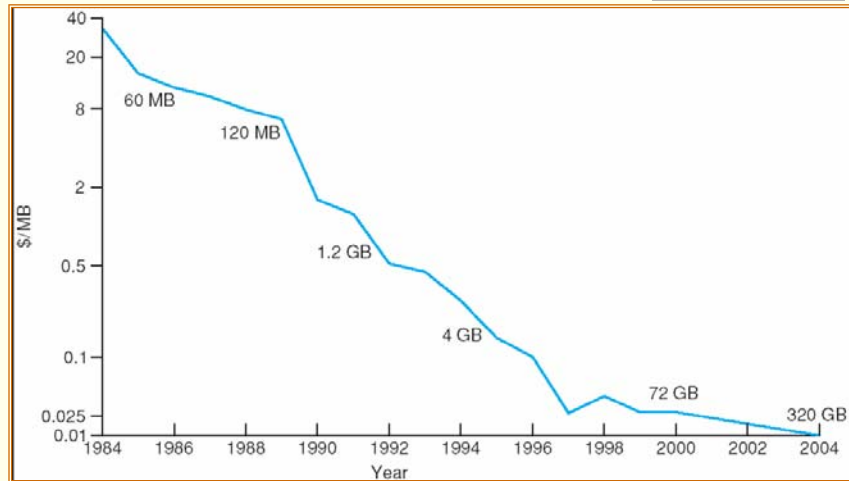


April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

54

## Price per Megabyte of a Tape Drive, From 1984-2000



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

55

## SOP – CO023

Gerência de  
E/S



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

56

## Objectives

---

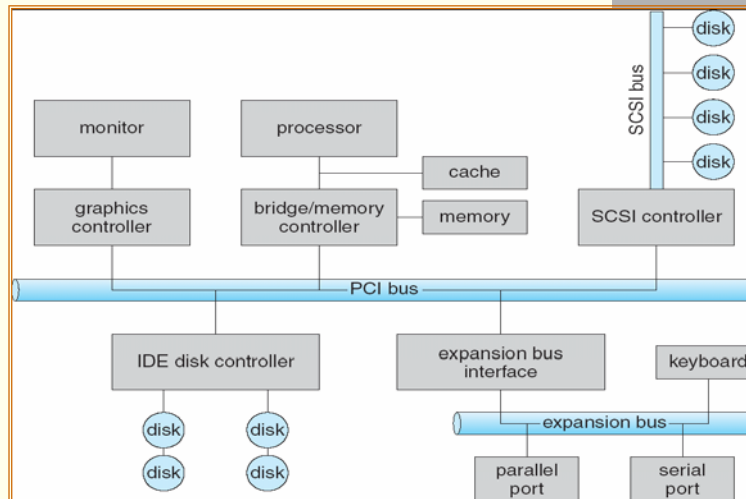
- Explore the structure of an operating system's I/O subsystem
- Discuss the principles of I/O hardware and its complexity
- Provide details of the performance aspects of I/O hardware and software

## I/O Hardware

---

- Incredible variety of I/O devices
- Common concepts
  - **Port**
  - **Bus (daisy chain** or shared direct access)
  - **Controller (host adapter)**
- I/O instructions control devices
- Devices have addresses, used by
  - Direct I/O instructions
  - **Memory-mapped I/O**

## A Typical PC Bus Structure



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

59

## 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)

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

60

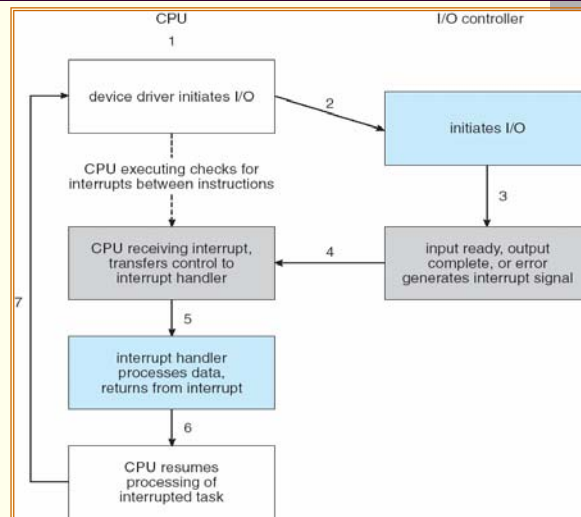
## Polling

- Determines state of device
  - command-ready
  - busy
  - Error
- **Busy-wait** cycle to wait for I/O from device

## Interrupts

- CPU **Interrupt-request line** triggered by I/O device  
**Interrupt handler** receives interrupts
- **Maskable** to ignore or delay some interrupts
- Interrupt vector to dispatch interrupt to correct handler
  - Based on priority
  - Some **nonmaskable**
- Interrupt mechanism also used for exceptions

## Interrupt-Driven I/O Cycle



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

63

## 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

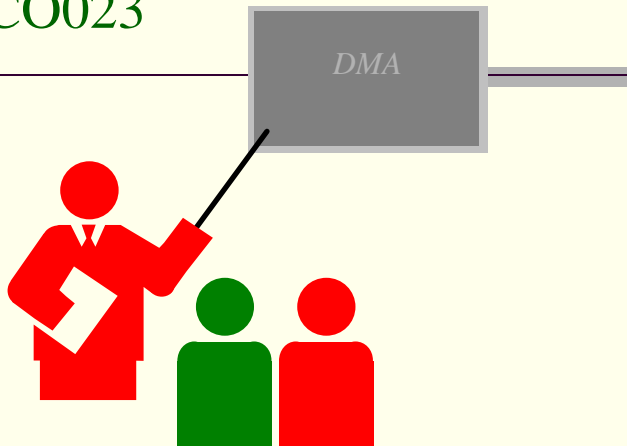
April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

64



## SOP – CO023



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

65

## Direct Memory Access

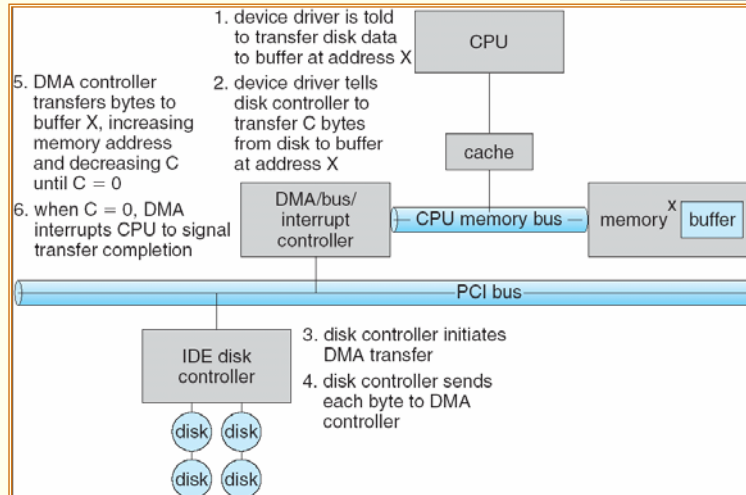
- Used to avoid **programmed I/O** for large data movement
- Requires **DMA** controller
- Bypasses CPU to transfer data directly between I/O device and memory

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

66

## Six Step Process to Perform DMA Transfer

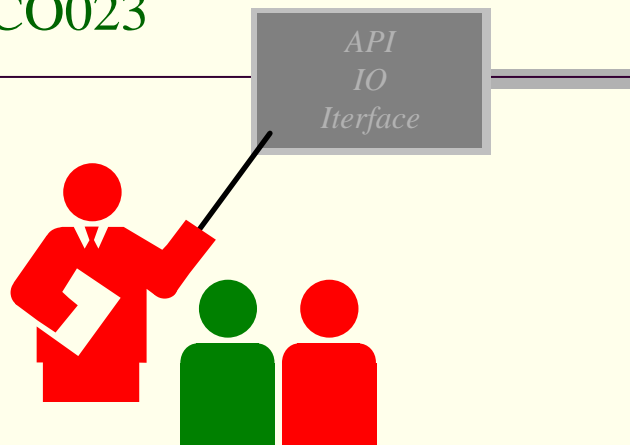


April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

67

## SOP – CO023



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

68

## Application I/O Interface

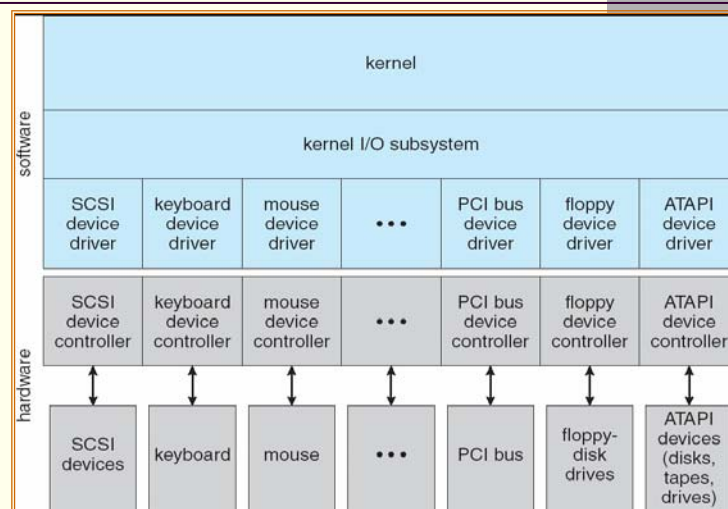
- I/O system calls encapsulate device behaviors in generic classes
- Device-driver layer hides differences among I/O controllers from kernel
- Devices vary in many dimensions
  - **Character-stream or block**
  - **Sequential or random-access**
  - **Sharable or dedicated**
  - **Speed of operation**
  - **read-write, read only, or write only**

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

69

## A Kernel I/O Structure



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

70

## 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

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

71

## Block and Character Devices

- **Block devices include disk drives**
  - Commands include read, write, seek
  - Raw I/O or file-system access
  - Memory-mapped file access possible
- **Character devices include keyboards, mice, serial ports**
  - Commands include `get`, `put`
  - Libraries layered on top allow line editing

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

72

## Network Devices

---

- Varying enough from block and character to have own interface
- Unix and Windows NT/9x/2000 include socket interface
  - Separates network protocol from network operation
  - Includes `select` functionality
- Approaches vary widely (pipes, FIFOs, streams, queues, mailboxes)

April 05

Prof. Ismael H. F. Santos - [ismael@tecgraf.puc-rio.br](mailto:ismael@tecgraf.puc-rio.br)

73

## Clocks and Timers

---

- Provide current time, elapsed time, timer
- **Programmable interval timer** used for timings, periodic interrupts
- `ioctl` (on UNIX) covers odd aspects of I/O such as clocks and timers

April 05

Prof. Ismael H. F. Santos - [ismael@tecgraf.puc-rio.br](mailto:ismael@tecgraf.puc-rio.br)

74

## Blocking and Nonblocking I/O

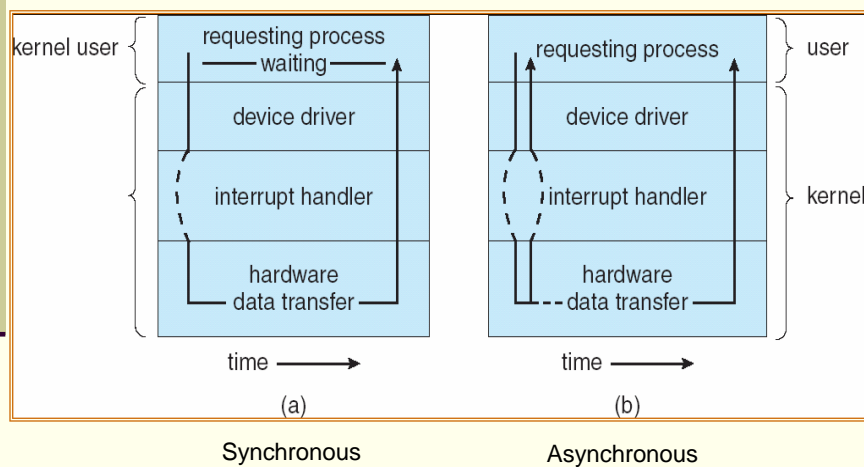
- **Blocking** - process suspended until I/O completed
  - Easy to use and understand
  - Insufficient for some needs
- **Nonblocking** - I/O call returns as much as available
  - User interface, data copy (buffered I/O)
  - Implemented via multi-threading
  - Returns quickly with count of bytes read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

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

75

## Two I/O Methods

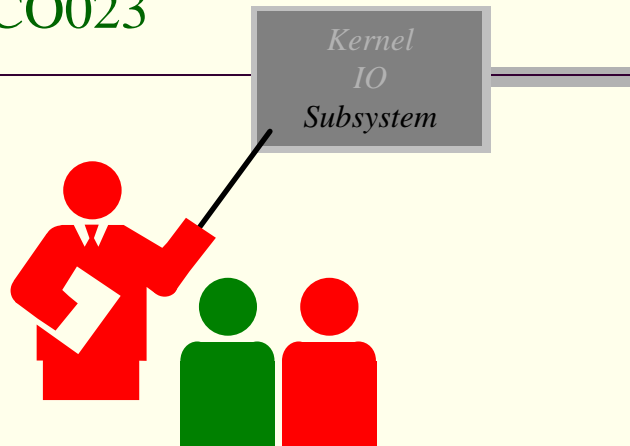


April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

76

## SOP – CO023



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

77

## Kernel I/O Subsystem

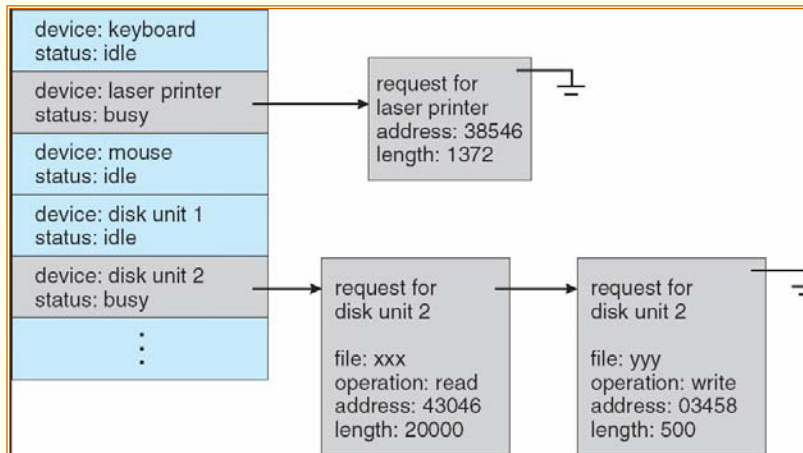
- **Scheduling**
  - Some I/O request ordering via per-device queue
  - Some OSs try fairness
- **Buffering** - store data in memory while transferring between devices
  - To cope with device speed mismatch
  - To cope with device transfer size mismatch
  - To maintain “copy semantics”

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

78

## Device-status Table

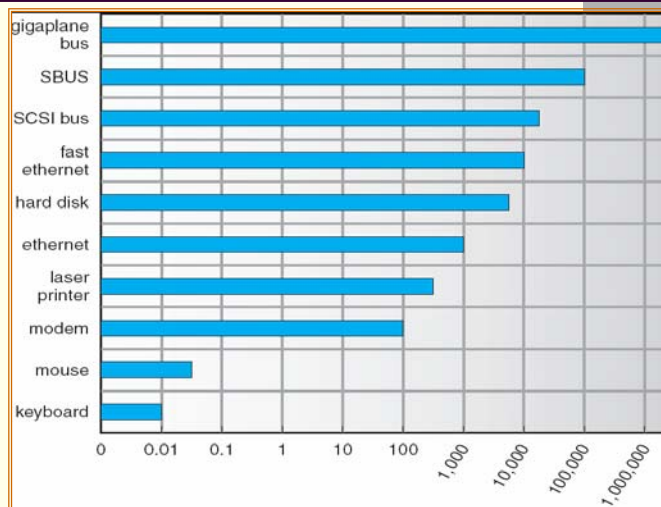


April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

79

## Sun Enterprise 6000 Device-Transfer Rates



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

80



## Kernel I/O Subsystem

---

- **Caching** - fast memory holding copy of data
  - Always just a copy
  - Key to performance
- **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 deallocation
  - Watch out for deadlock

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

81

## Error Handling

---

- OS can recover from disk read, device unavailable, transient write failures
- Most return an error number or code when I/O request fails
- System error logs hold problem reports

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

82

## I/O Protection

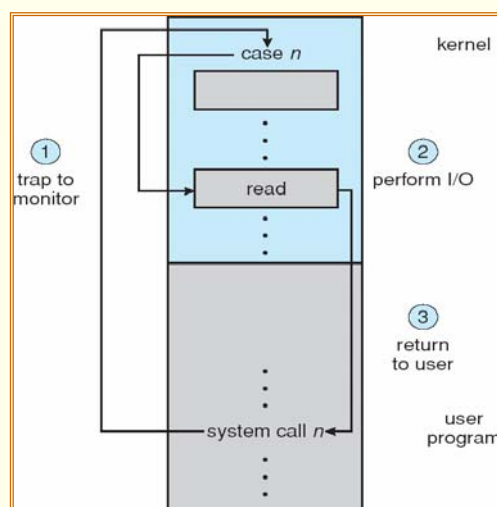
- User process may accidentally or purposefully attempt to disrupt normal operation via illegal I/O instructions
  - All I/O instructions defined to be privileged
  - I/O must be performed via system calls
    - Memory-mapped and I/O port memory locations must be protected too

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

83

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



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

84

## Kernel Data Structures

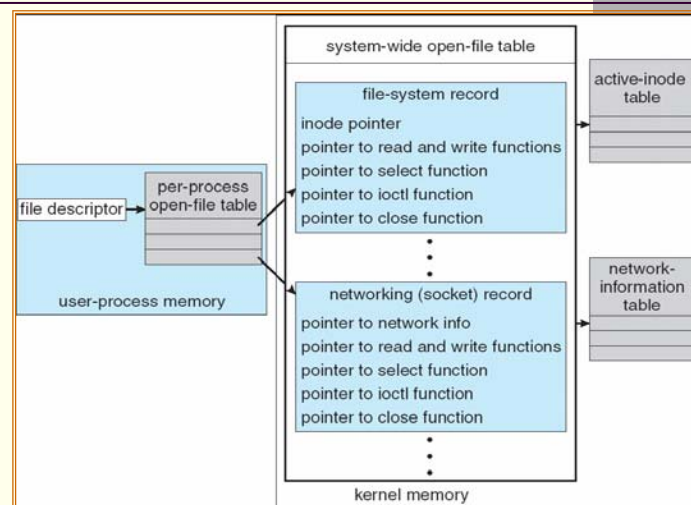
- Kernel keeps state info for I/O components, including open file tables, network connections, character device state
- Many, many complex data structures to track buffers, memory allocation, “dirty” blocks
- Some use object-oriented methods and message passing to implement I/O

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

85

## UNIX I/O Kernel Structure



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

86

# I/O Requests to Hardware Operations

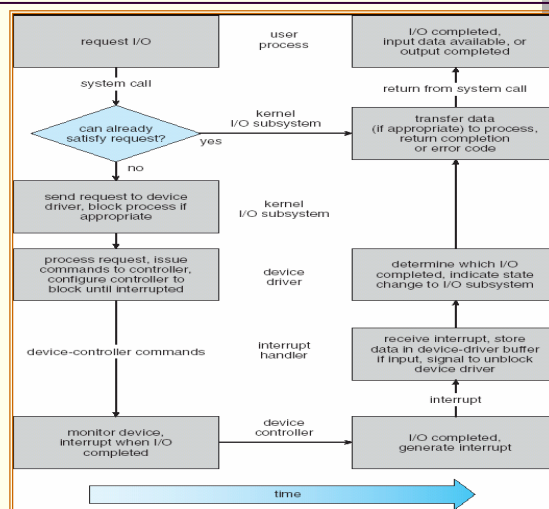
- Consider reading a file from disk for a process:
  - Determine device holding file
  - Translate name to device representation
  - Physically read data from disk into buffer
  - Make data available to requesting process
  - Return control to process

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

87

# Life Cycle of An I/O Request



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

88

# STREAMS

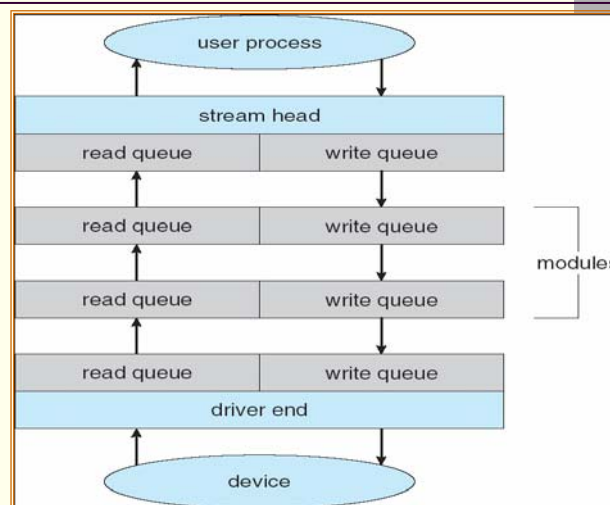
- **STREAM** – a full-duplex communication channel between a user-level process and a device in Unix System V and beyond
- A STREAM consists of:
  - STREAM head interfaces with the user process
  - driver end interfaces with the device
  - zero or more STREAM modules between them.
- Each module contains a **read queue** and a **write queue**. Message passing is used to communicate between queues

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

89

## The STREAMS Structure



April 05

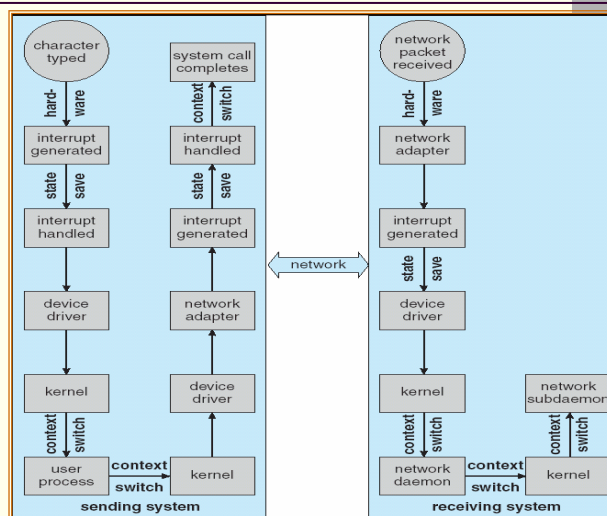
Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

90

# Performance

- I/O a major factor in system performance:
  - Demands CPU to execute device driver, kernel I/O code
  - Context switches due to interrupts
  - Data copying
  - Network traffic especially stressful

# Intercomputer Communications



## Improving Performance

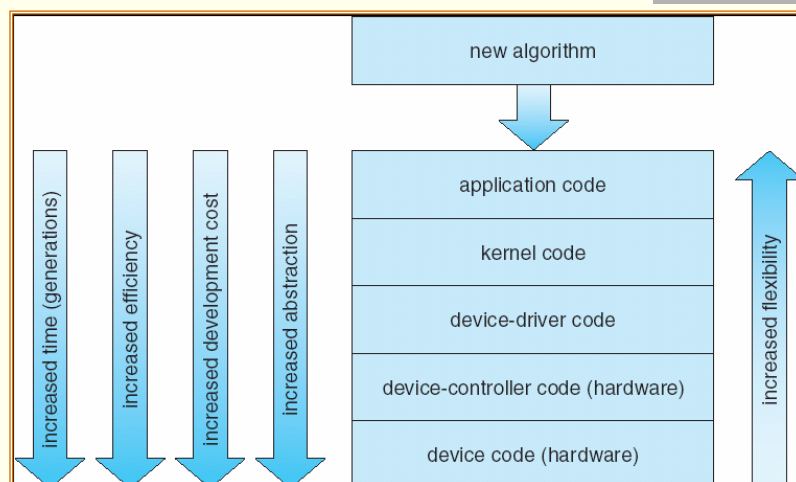
- Reduce number of context switches
- Reduce data copying
- Reduce interrupts by using large transfers, smart controllers, polling
- Use DMA
- Balance CPU, memory, bus, and I/O performance for highest throughput

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

93

## Device-Functionality Progression



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

94