# File Management

Ref:-[Galvin][Stallings][Tanenbaum]

# Objectives

- To describe the details of implementing local file systems and directory structures

- To describe the implementation of remote file systems

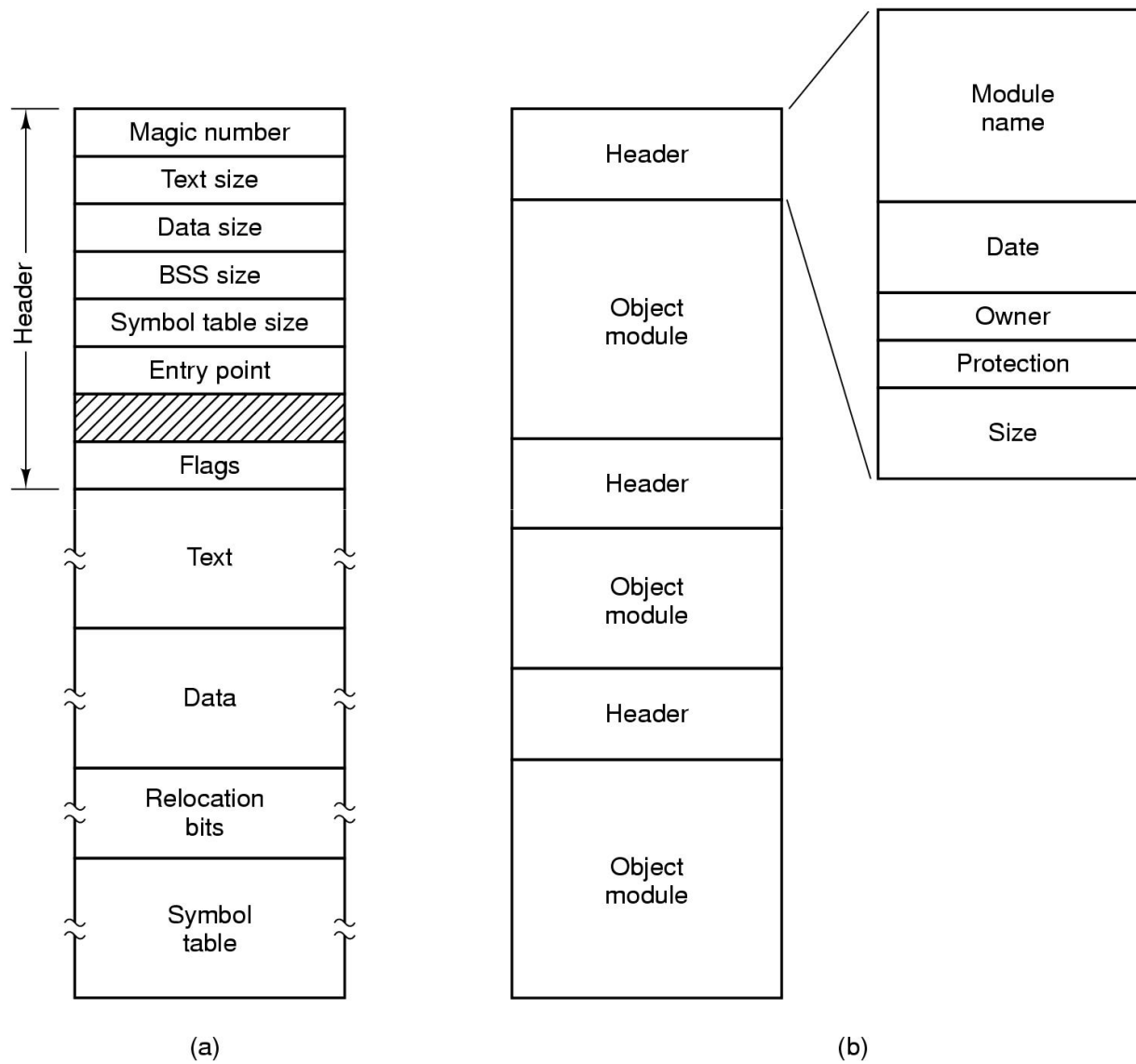- To discuss block allocation and free-block algorithms and trade-offs

# Long-term Information Storage

- Must store large amounts of data

- Must be persistent

- Multiple processes must be able to access the information concurrently

# File Naming

| Extension | Meaning |
|---|---|
| file.bak | Backup file |
| file.c | C source program |
| file.gif | Compuserve Graphical Interchange Format image |
| file.hlp | Help file |
| file.html | World Wide Web HyperText Markup Language document |
| file.jpg | Still picture encoded with the JPEG standard |
| file.mp3 | Music encoded in MPEG layer 3 audio format |
| file.mpg | Movie encoded with the MPEG standard |
| file.o | Object file (compiler output, not yet linked) |
| file.pdf | Portable Document Format file |
| file.ps | PostScript file |
| file.tex | Input for the TEX formatting program |
| file.txt | General text file |
| file.zip | Compressed archive |

# File Types



(a) An executable file   (b) An archive

# File Access

- Sequential access
  - read all bytes/records from the beginning
  - cannot jump around, could rewind or back up
  - convenient when medium was magnetic tape

- Random access
  - bytes/records read in any order
  - essential for data base systems

# File Attributes

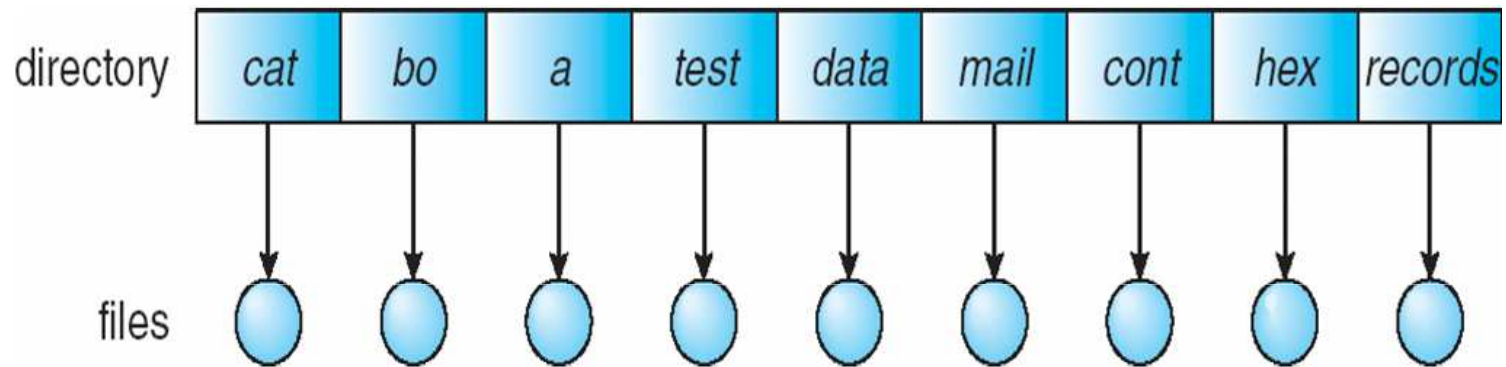| Attribute | Meaning |
|---|---|
| Protection | Who can access the file and in what way |
| Password | Password needed to access the file |
| Creator | ID of the person who created the file |
| Owner | Current owner |
| Read-only flag | 0 for read/write; 1 for read only |
| Hidden flag | 0 for normal; 1 for do not display in listings |
| System flag | 0 for normal files; 1 for system file |
| Archive flag | 0 for has been backed up; 1 for needs to be backed up |
| ASCII/binary flag | 0 for ASCII file; 1 for binary file |
| Random access flag | 0 for sequential access only; 1 for random access |
| Temporary flag | 0 for normal; 1 for delete file on process exit |
| Lock flags | 0 for unlocked; nonzero for locked |
| Record length | Number of bytes in a record |
| Key position | Offset of the key within each record |
| Key length | Number of bytes in the key field |
| Creation time | Date and time the file was created |
| Time of last access | Date and time the file was last accessed |
| Time of last change | Date and time the file has last changed |
| Current size | Number of bytes in the file |
| Maximum size | Number of bytes the file may grow to |

Possible file attributes

# File Operations

1. Create
2. Delete
3. Open
4. Close
5. Read
6. Write

7. Append
8. Seek
9. Get attributes
10. Set Attributes
11. Rename

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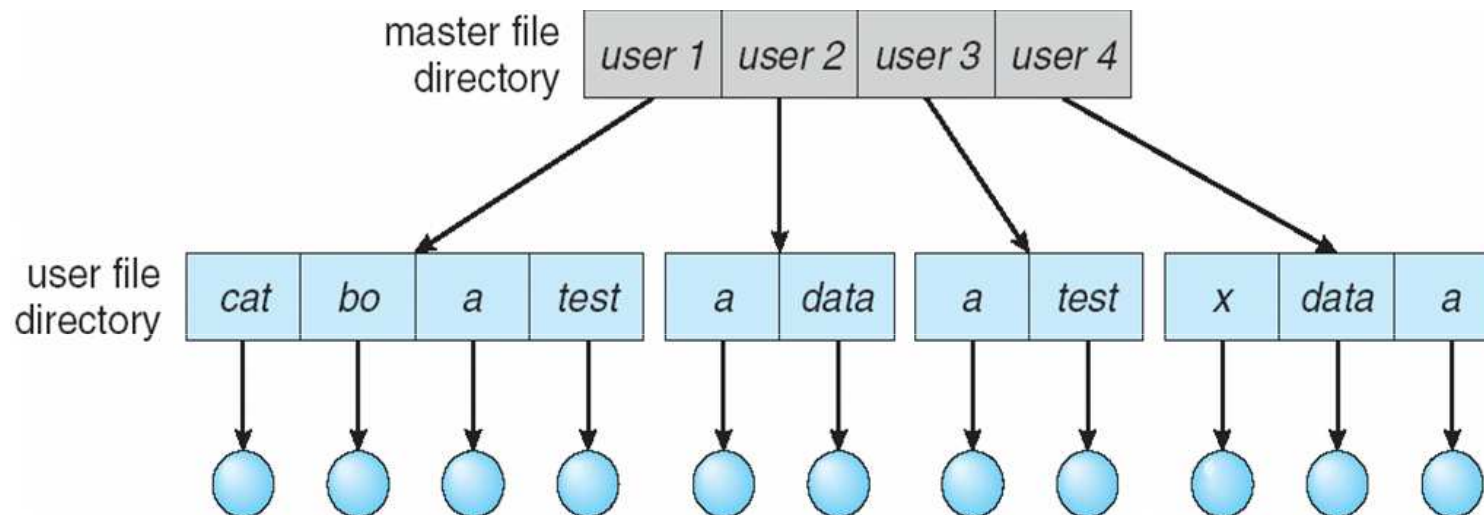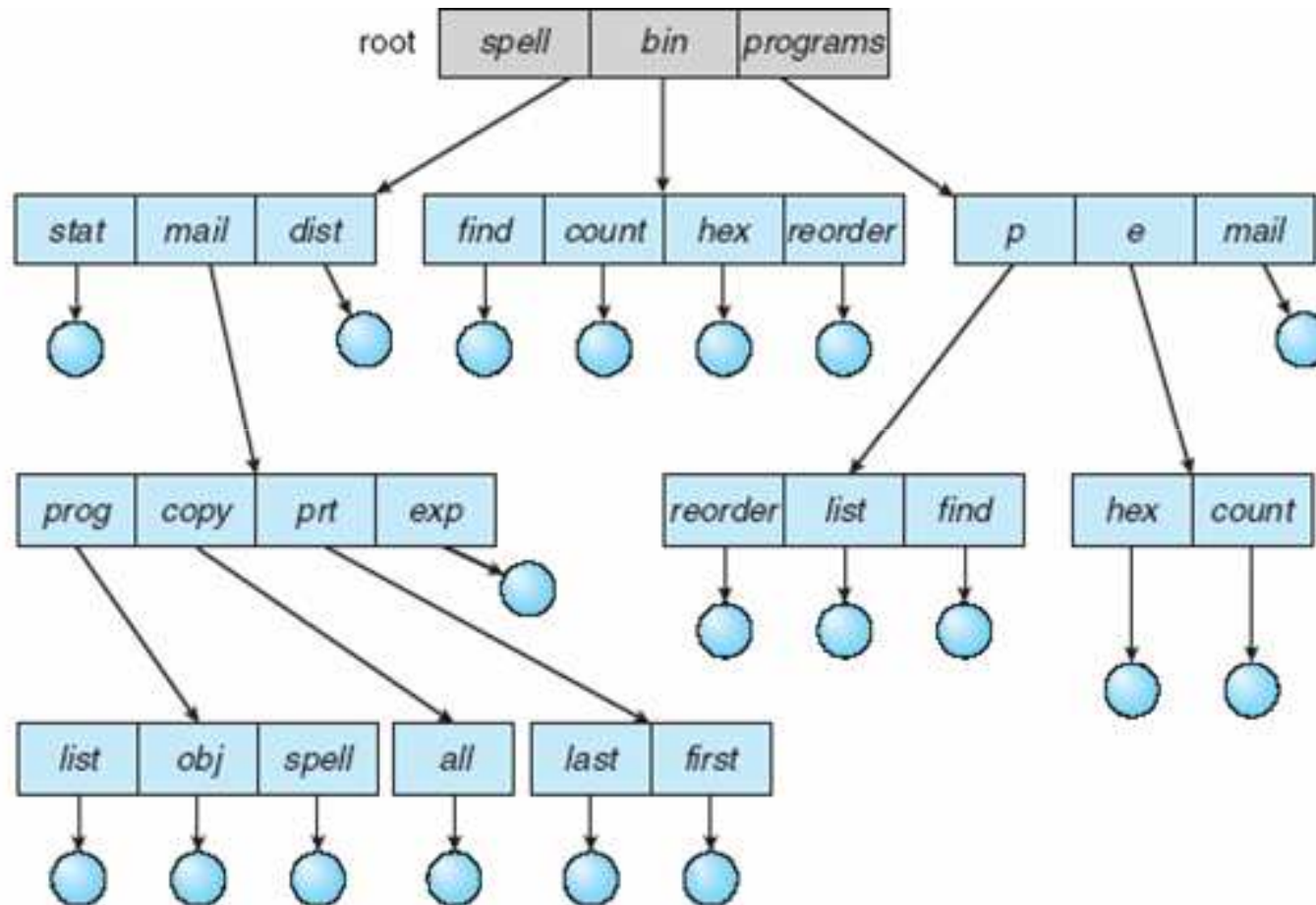
- A single directory for all users

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

files ○ ○ ○ ○ ○ ○ ○ ○ ○

Naming problem

Grouping problem

# Two-Level Directory

- Separate directory for each user



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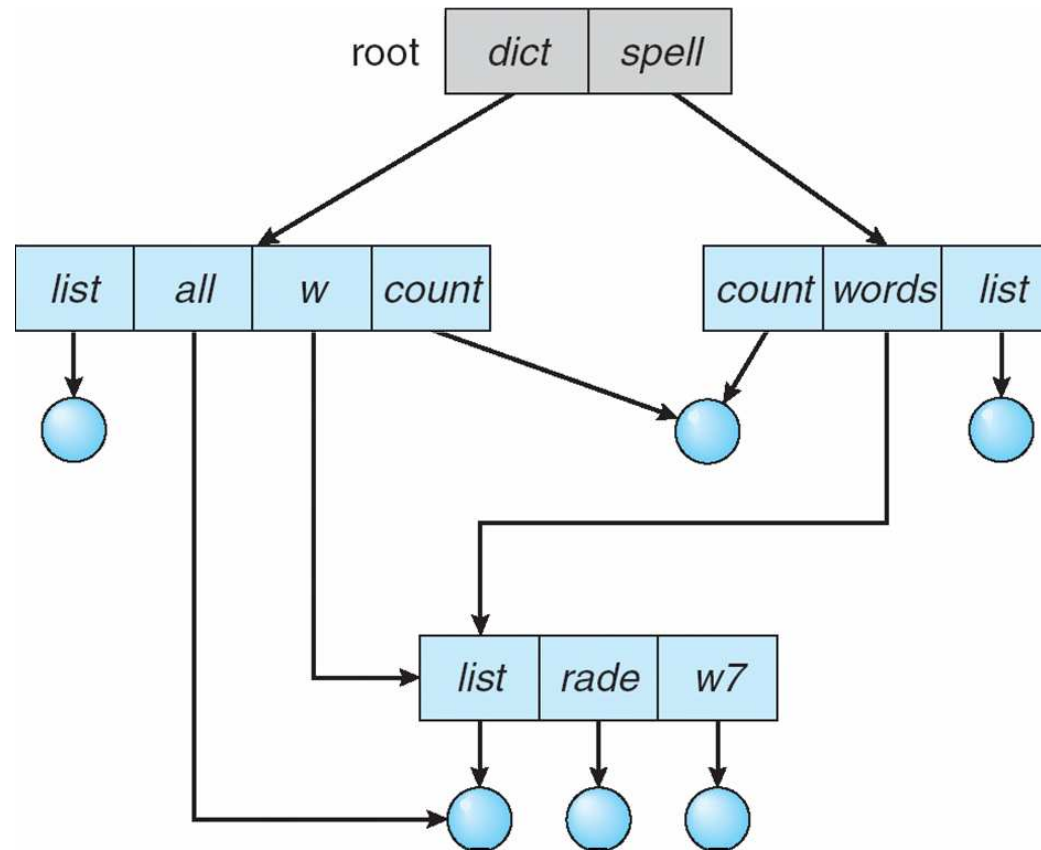- Efficient searching
- No grouping capability

# Tree-Structured Directories



- Efficient searching
- Grouping Capability
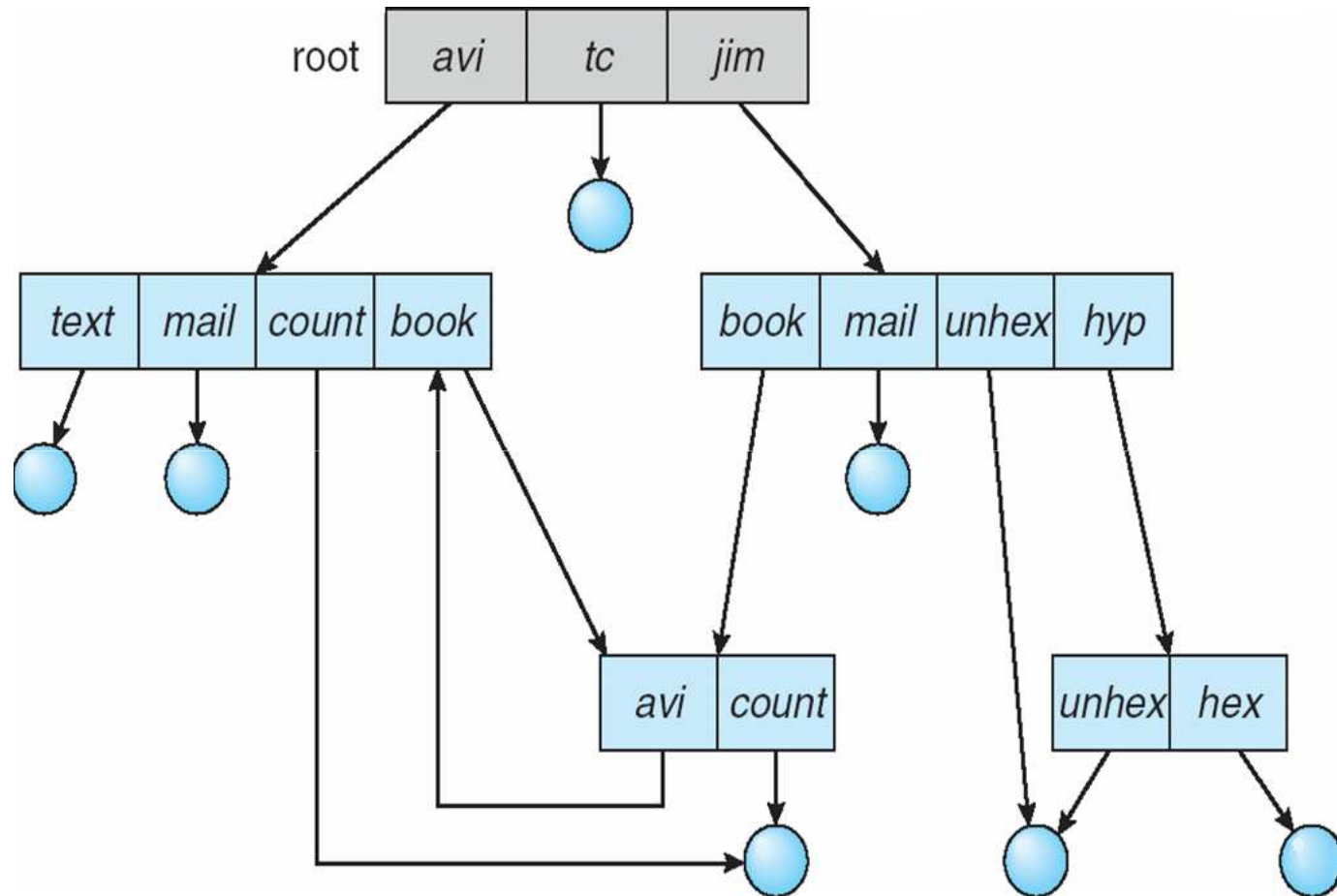- Absolute or relative path name

# Acyclic-Graph Directories

- Have shared subdirectories and files

# Acyclic-Graph Directories (Cont.)

- Two different names (aliasing)

- If *dict* deletes *list* $\Rightarrow$ dangling pointer
  Solutions:
  - Backpointers, so we can delete all pointers

- New directory entry type
  - **Link** – another name (pointer) to an existing file
  - **Resolve the link** – follow pointer to locate the file
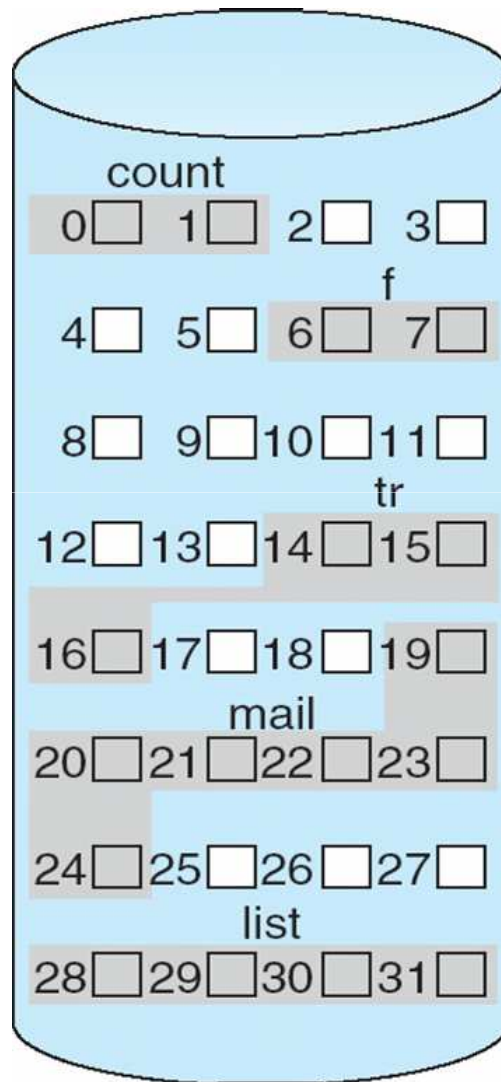
# General Graph Directory

# Allocation Methods

- An allocation method refers to how disk blocks are allocated for files:

- Contiguous allocation

- Linked allocation

- Indexed allocation

# Contiguous Allocation

- Each file occupies a set of contiguous blocks on the disk

- Simple – only starting location (block #) and length (number of blocks) are required

- Random access

- Wasteful of space (dynamic storage-allocation problem)

- Files cannot grow
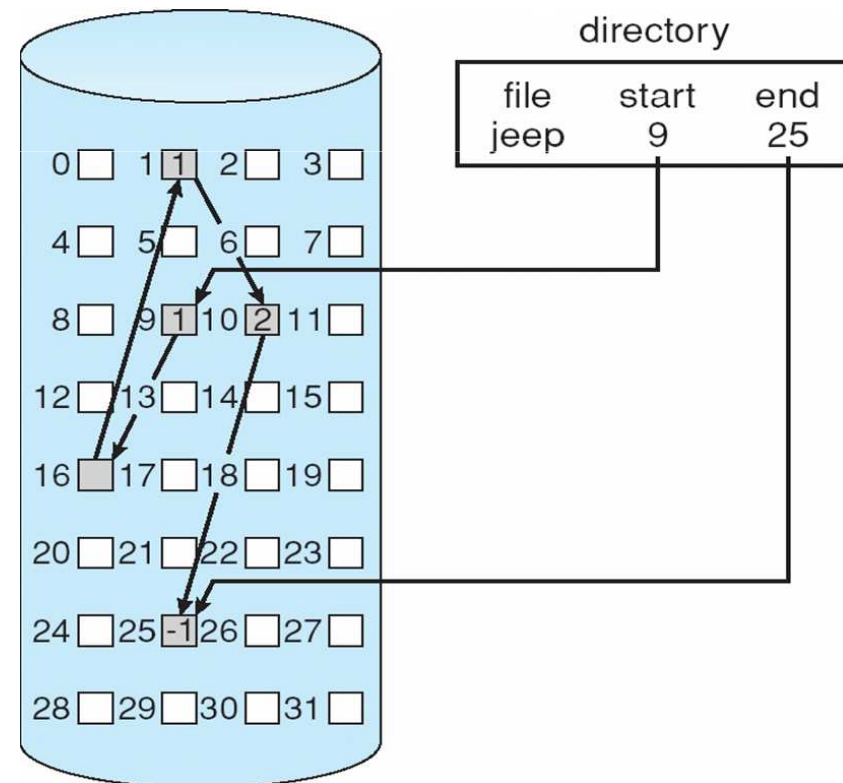
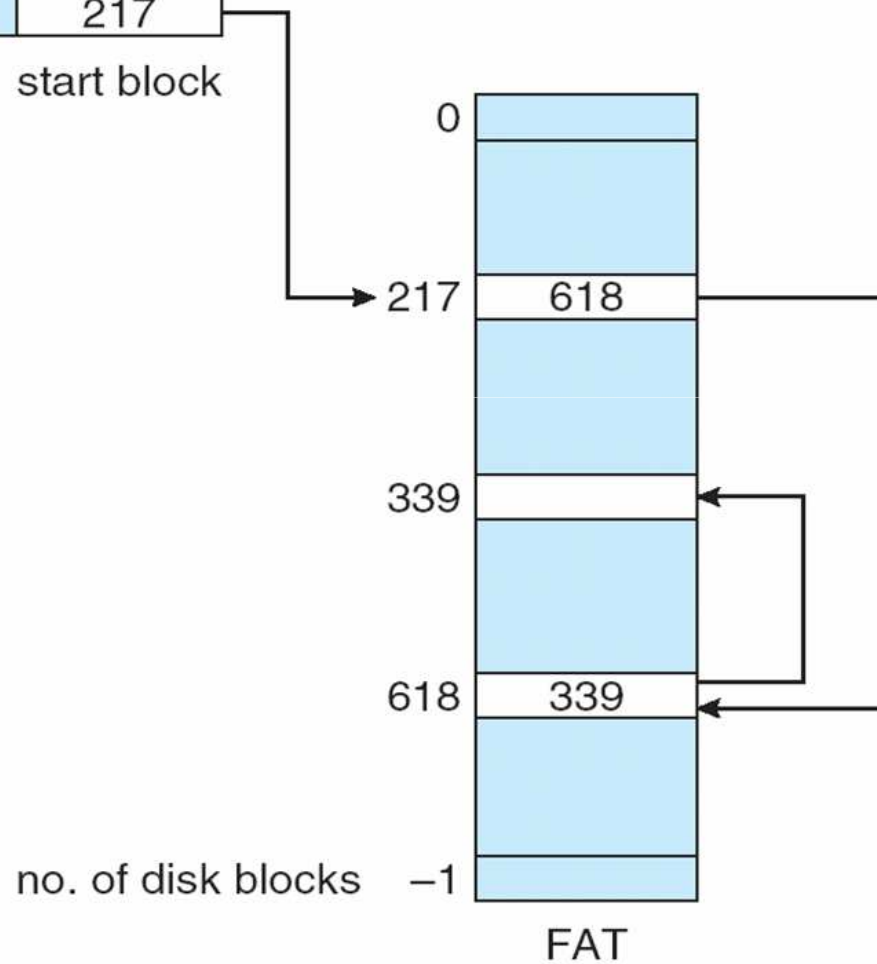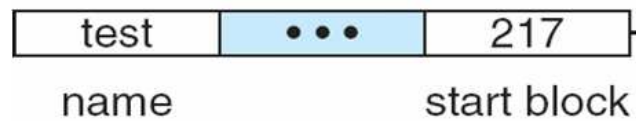# Contiguous Allocation of Disk Space

# Linked Allocation

- Each file is a linked list of disk blocks: blocks may be scattered anywhere on the disk.

- Simple – need only starting address

- Free-space management system – no waste of space

- No random access

# File-Allocation Table
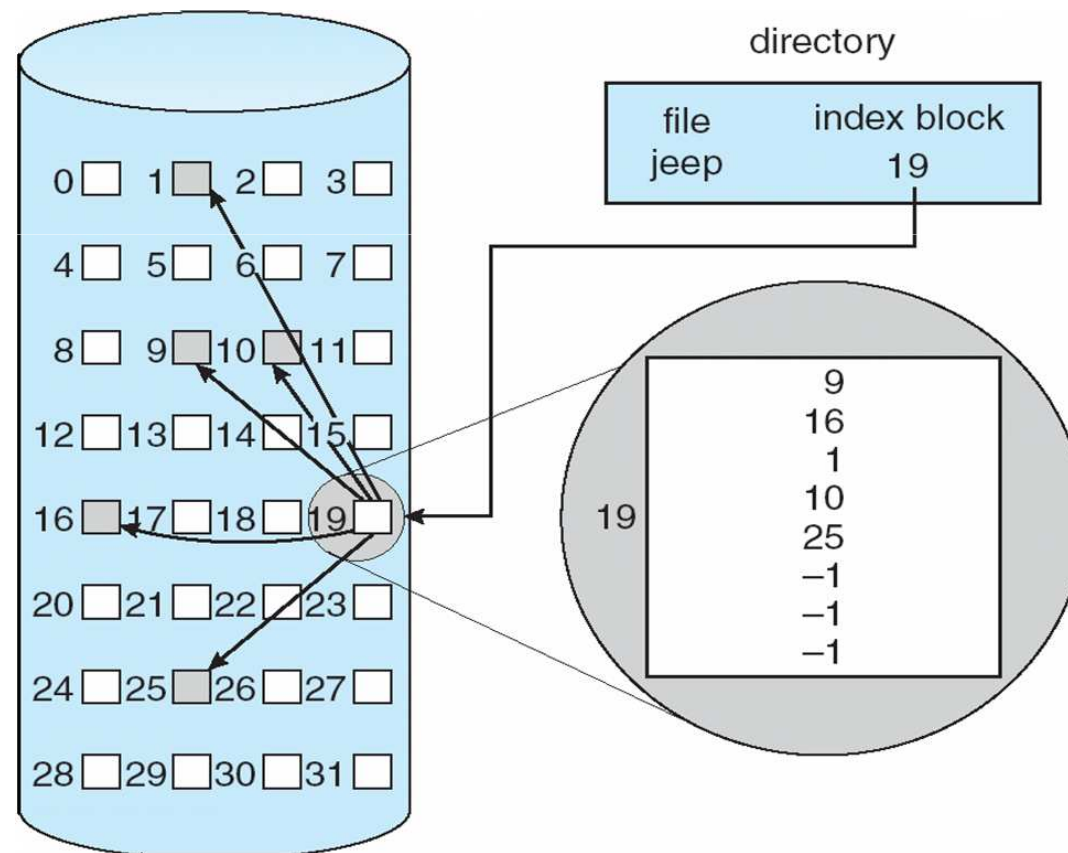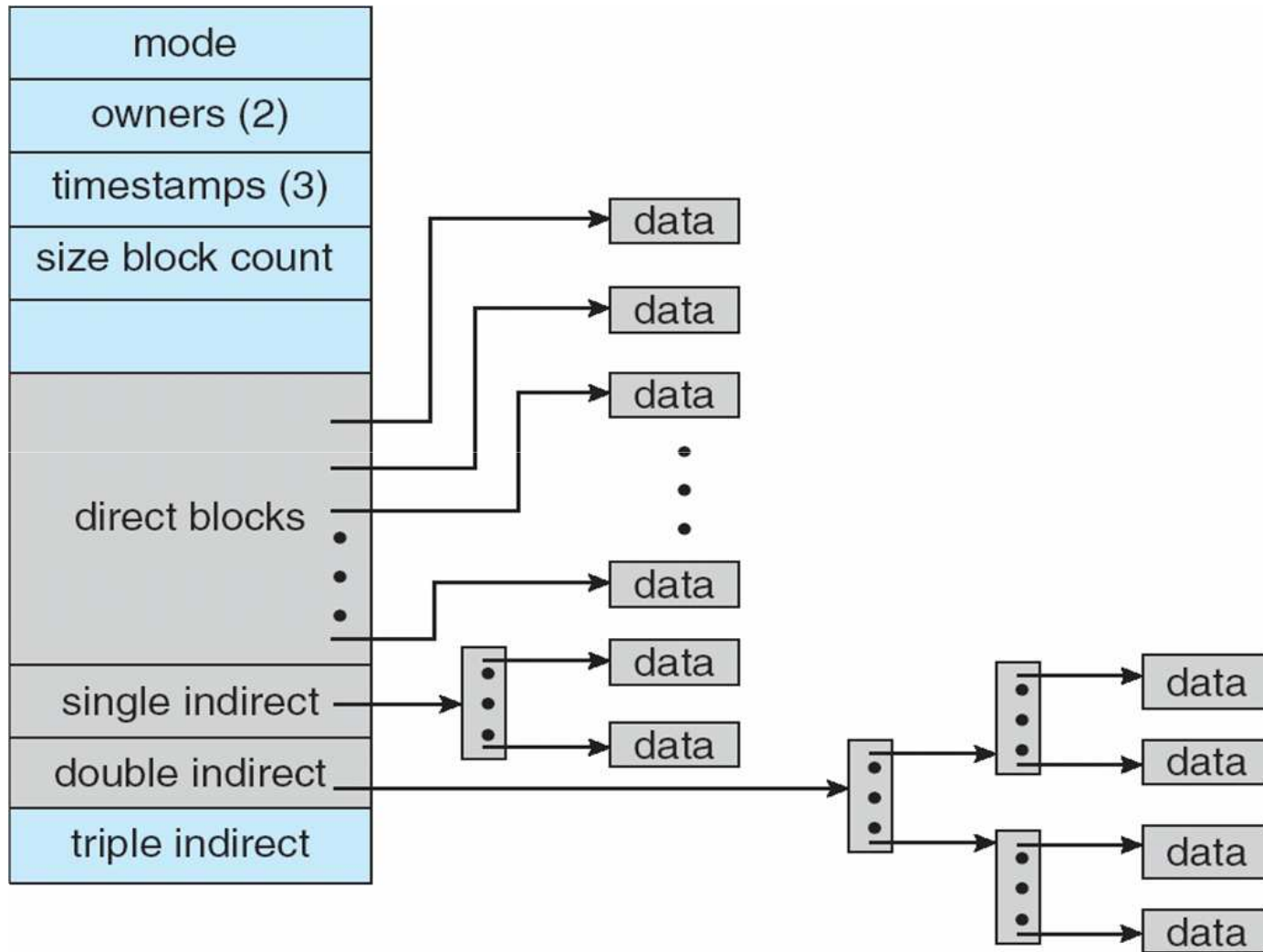
# Indexed Allocation

- Brings all pointers together into the index block
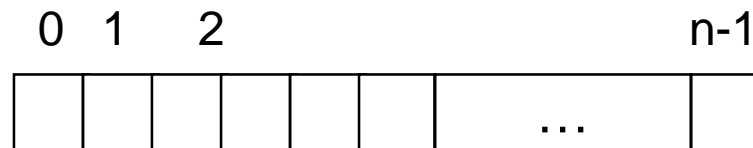- Logical view

# Indexed Allocation (Cont.)

- Need index table
- Random access
- Dynamic access without external fragmentation, but have overhead of index block

# Combined Scheme: UNIX UFS (4K bytes per block)

# Free-Space Management

- Bit vector   (*n* blocks)

0   1     2                                                    n-1

|   |   |   |   |   |   | … |   |

$$bit[i] = \begin{cases} 1 \Rightarrow block[i] \text{ free} \\ 0 \Rightarrow block[i] \text{ occupied} \end{cases}$$

Block number calculation

(number of bits per word) *(number of 0-value words) + offset of first 1 bit

# Free-Space Management (Cont.)

- Bit map requires extra space
  - Example:
- $\qquad$ block size $= 2^{12}$ bytes
- $\qquad$ disk size $= 2^{30}$ bytes (1 gigabyte)
- $\qquad$ $n = 2^{30}/2^{12} = 2^{18}$ bits (or 32K bytes)
- Easy to get contiguous files
- Linked list (free list)
  - Cannot get contiguous space easily
  - No waste of space
- *Grouping*
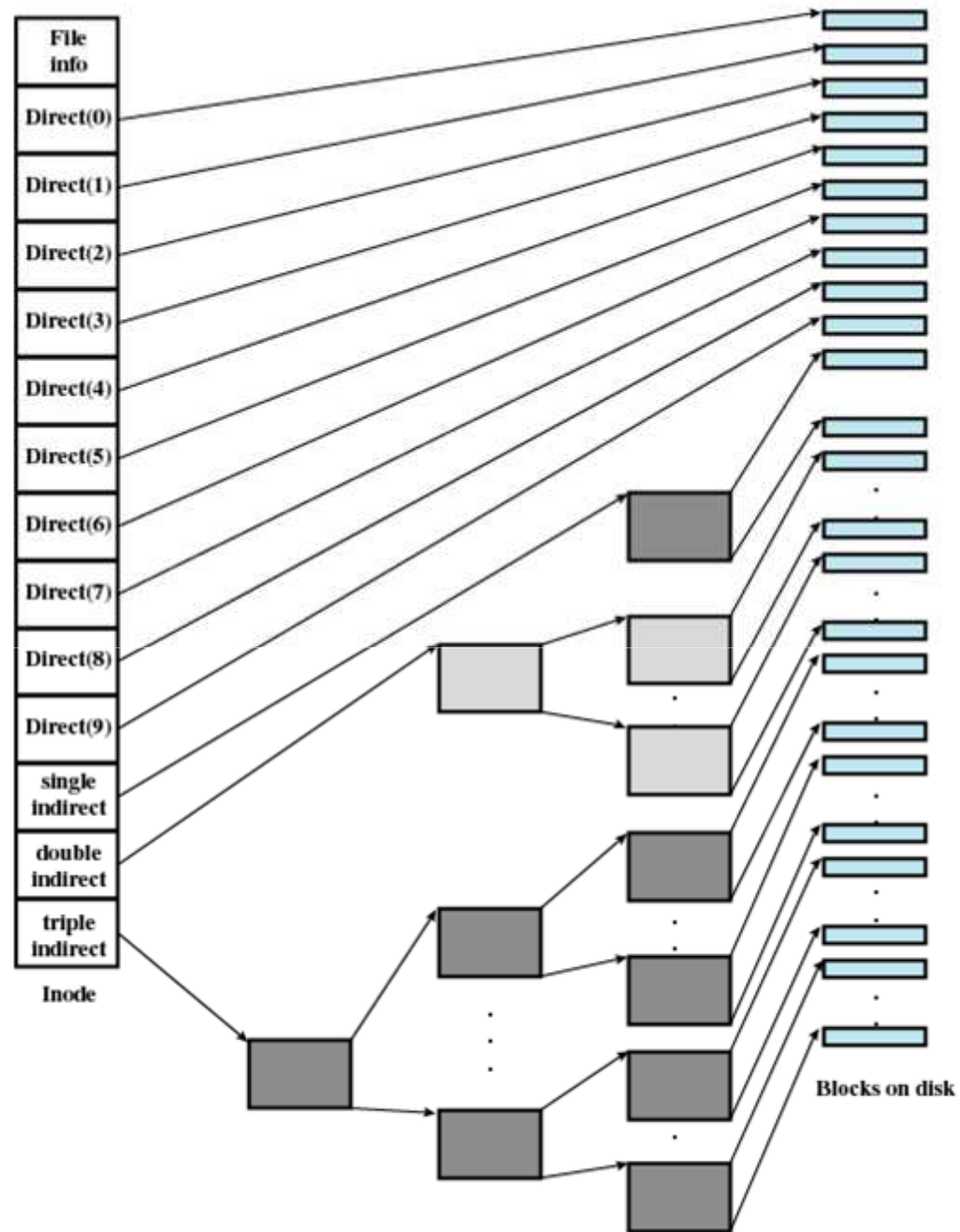- *Counting*

# UNIX File Management

- Types of files
  - Regular, or ordinary
  - Directory
  - Special
  - Named pipes
  - Links
  - Symbolic links

# Inodes

- Index node

- Control structure that contains key information for a particular file

# Unix Inodes

| | |
|---|---|
| **File Mode** | 16-bit flag that stores access and execution permissions associated with the file. |
| | 12-14 File type (regular, directory, character or block special, FIFO pipe |
| | 9-11 Execution flags |
| | 8 Owner read permission |
| | 7 Owner write permission |
| | 6 Owner execute permission |
| | 5 Group read permission |
| | 4 Group write permission |
| | 3 Group execute permission |
| | 2 Other read permission |
| | 1 Other write permission |
| | 0 Other execute permission |
| **Link Count** | Number of directory references to this inode |
| **Owner ID** | Individual owner of file |
| **Group ID** | Group owner associated with this file |
| **File Size** | Number of bytes in file |
| **File Addresses** | 39 bytes of address information |
| **Last Accessed** | Time of last file access |
| **Last Modified** | Time of last file modification |
| **Inode Modified** | Time of last inode modification |

**Layout of a UNIX File on Disk**

# Linux Virtual File System

- Uniform file system interface to user processes

- Represents any conceivable file system's general feature and behavior

- Assumes files are objects that share basic properties regardless of the target file system
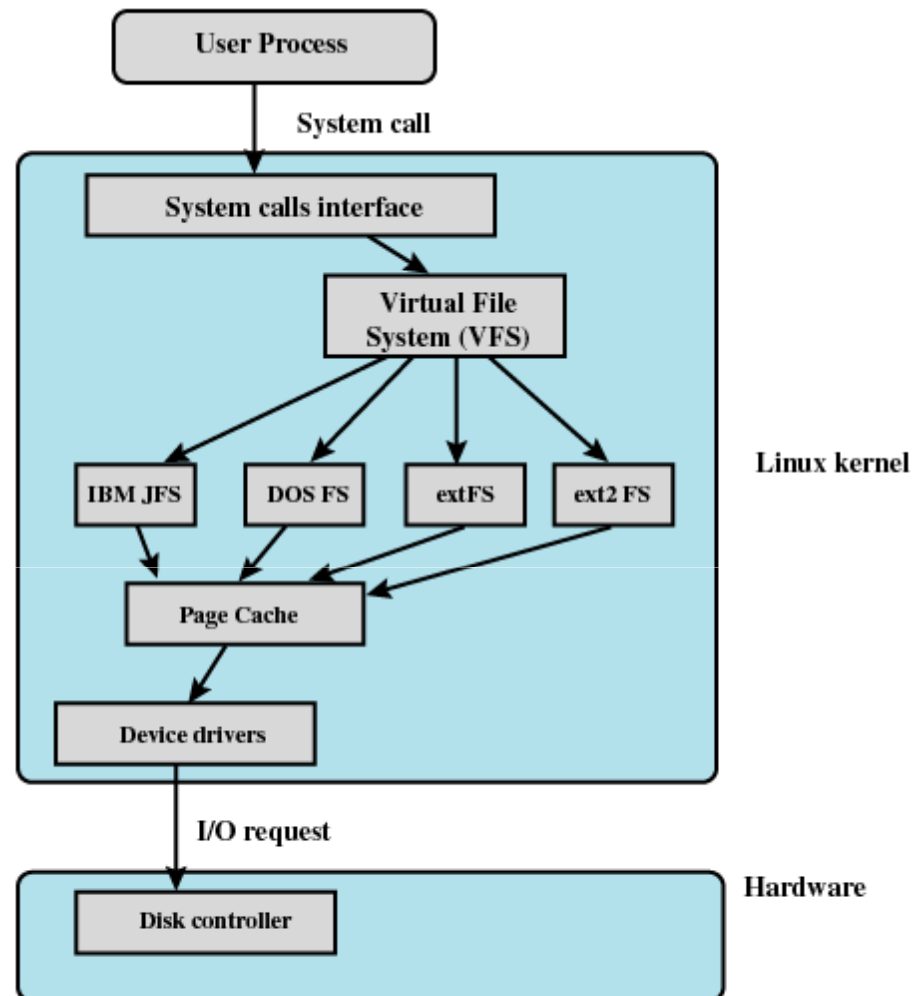
# VFS

What's VFS?

The Virtual File system is a kernel software layer that handles all system calls related to file systems. Its main strength is providing a common interface to several kinds of file systems.

What's Linux VFS's key idea?

For each read, write, or other function called, the kernel substitutes the actual function that supports a native Linux file system, the NT file system, or whatever other file system the file is on.
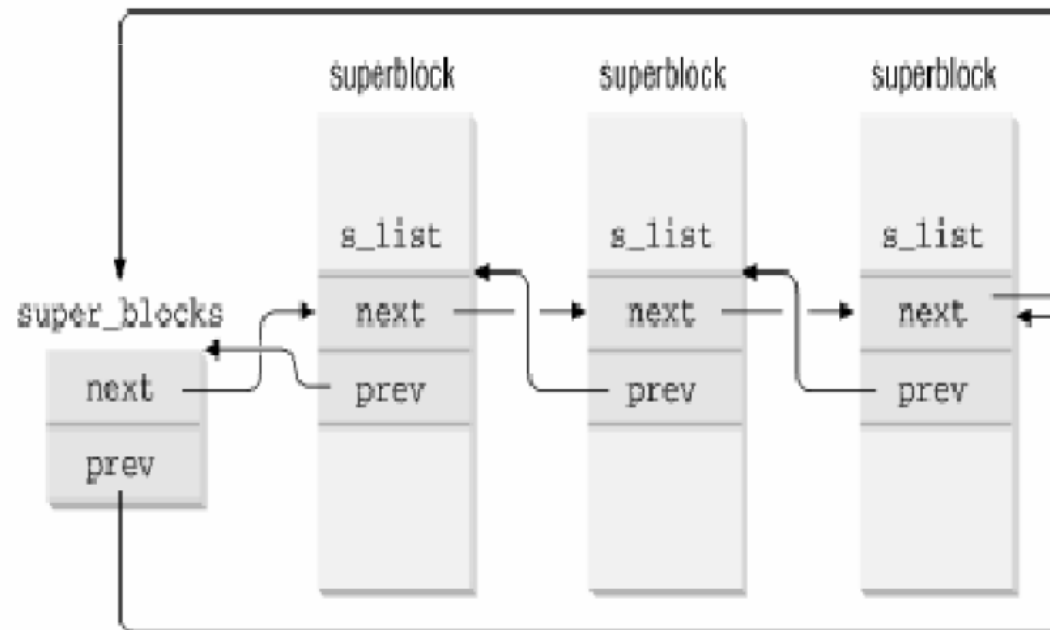
**User Process**

System call

**System calls interface**

**Virtual File System (VFS)**

**IBM JFS**   **DOS FS**   **extFS**   **ext2 FS**

Linux kernel

**Page Cache**

**Device drivers**

I/O request

**Disk controller**

Hardware

**Linux Virtual File System Context**

# Primary Objects in VFS

- Superblock object
  - Represents a specific mounted file system
- Inode object
  - Represents a specific file
- Dentry object
  - Represents a specific directory entry
- File object
  - Represents an open file associated with a process

# Superblock Object

All superblock objects (one per mounted file system) are linked together in a circular doubly linked list.

# Superblock Operations

```
struct super_operations {
    void (*read_inode) (struct inode *);   /* fill the structure */
    int (*notify_change) (struct inode *, struct iattr *);
    void (*write_inode) (struct inode *);
    void (*put_inode) (struct inode *);
    void (*put_super) (struct super_block *);
    void (*write_super) (struct super_block *);
    void (*statfs) (struct super_block *, struct statfs *, int);
    int (*remount_fs) (struct super_block *, int *, char *);
};
```

# Inode operations

```
struct inode_operations {
    struct file_operations * default_file_ops;
    int (*create) (struct inode *,const char *,int,int,struct inode **);
    int (*lookup) (struct inode *,const char *,int,struct inode **);
    int (*link) (struct inode *,struct inode *,const char *,int);
    int (*unlink) (struct inode *,const char *,int);
    int (*symlink) (struct inode *,const char *,int,const char *);
    int (*mkdir) (struct inode *,const char *,int,int);
    int (*rmdir) (struct inode *,const char *,int);
    int (*mknod) (struct inode *,const char *,int,int,int);
    int (*rename) (struct inode *,const char *,int, struct inode *,
                    const char *,int, int); /* this from 2.0.1 onwards */
    int (*readlink) (struct inode *,char *,int);
    int (*follow_link) (struct inode *,struct inode *,int,int,struct inode **);
    int (*readpage) (struct inode *, struct page *);
    int (*writepage) (struct inode *, struct page *);
    int (*bmap) (struct inode *,int);
    void (*truncate) (struct inode *);
    int (*permission) (struct inode *, int);
    int (*smap) (struct inode *,int);
};
```
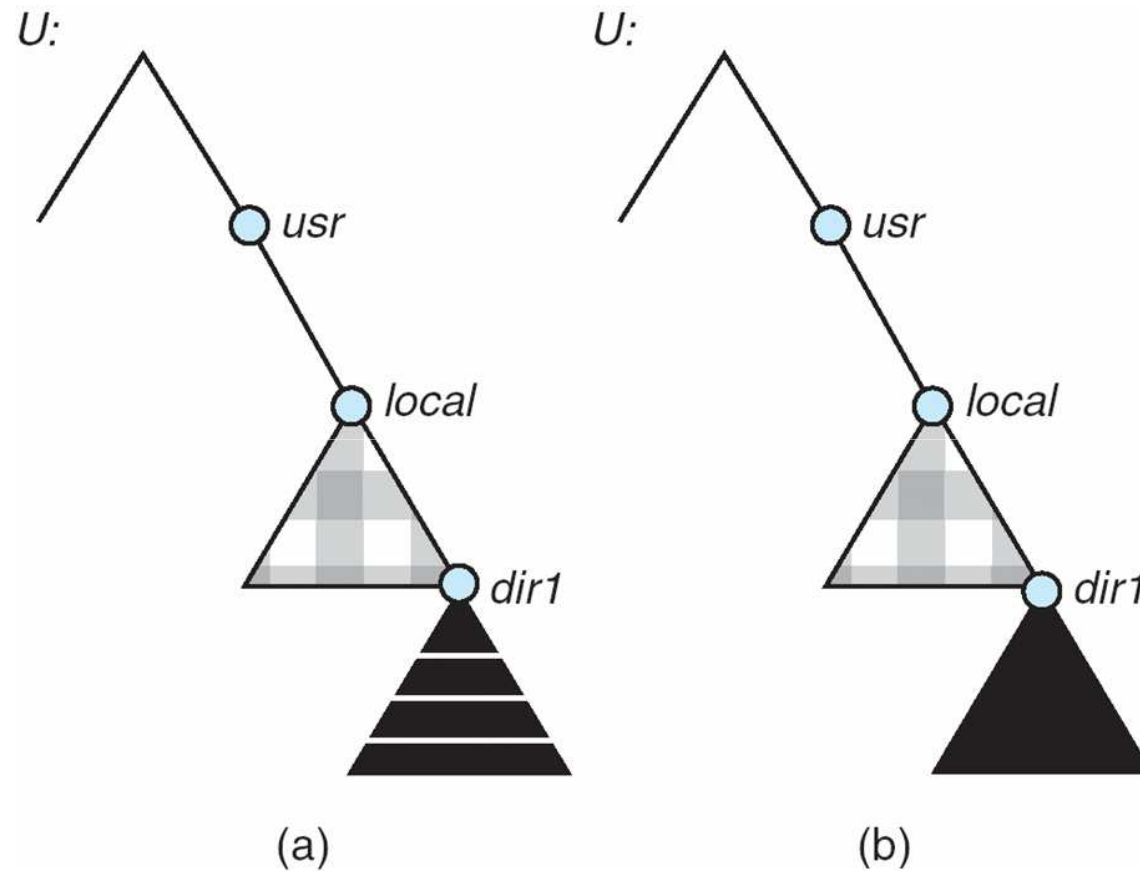
# The Sun Network File System (NFS)

- An implementation and a specification of a software system for accessing remote files across LANs (or WANs)

- The implementation is part of the Solaris and SunOS operating systems running on Sun workstations using an unreliable datagram protocol (UDP/IP protocol and Ethernet)

# NFS Mount Protocol

- Establishes initial logical connection between server and client
- Mount operation includes name of remote directory to be mounted and name of server machine storing it
  - Mount request is mapped to corresponding RPC and forwarded to mount server running on server machine
  - Export list – specifies local file systems that server exports for mounting, along with names of machines that are permitted to mount them
- Following a mount request that conforms to its export list, the server returns a file handle—a key for further accesses
- File handle – a file-system identifier, and an inode number to identify the mounted directory within the exported file system
- The mount operation changes only the user's view and does not affect the server side
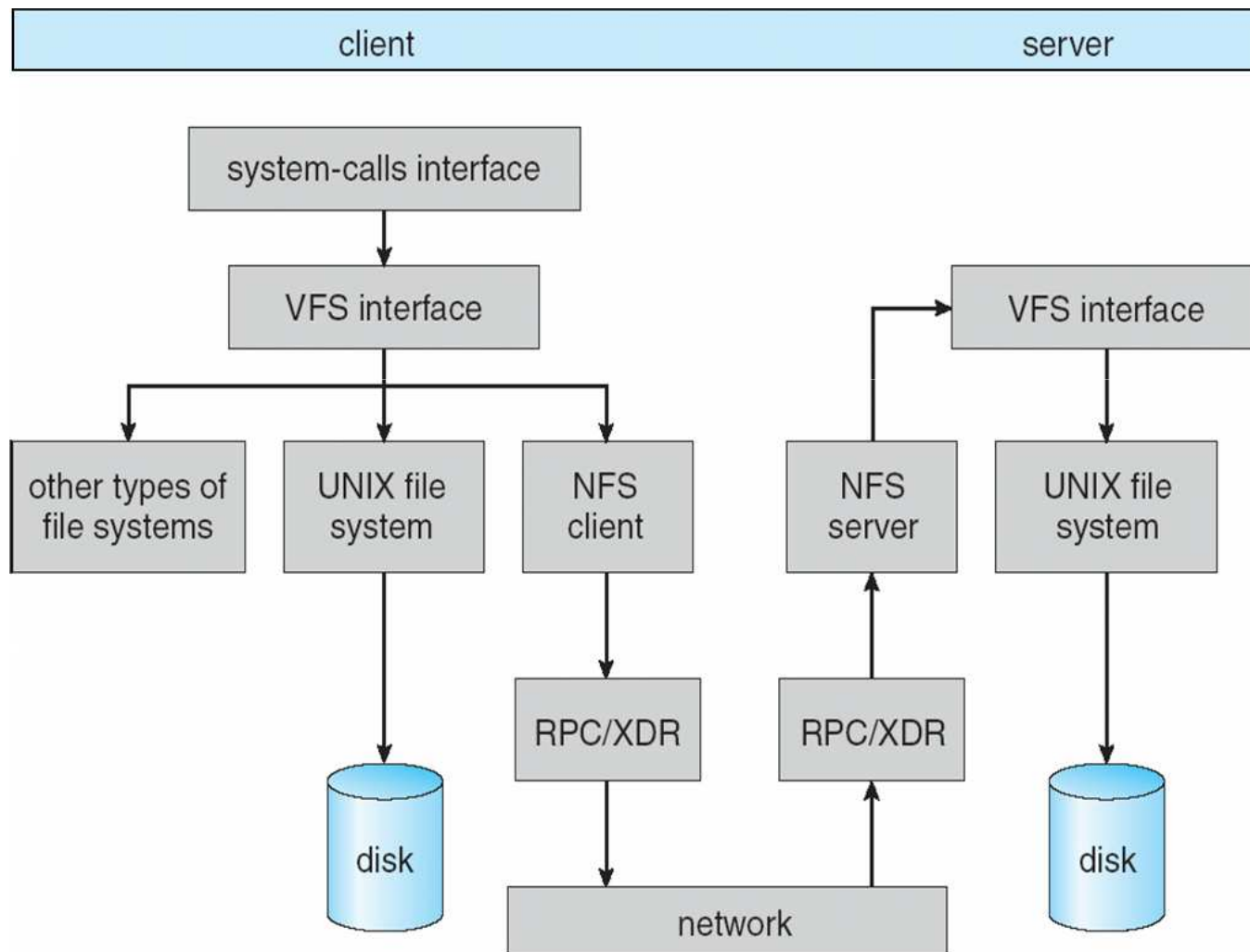
# Mounting in NFS



(a) Mounts

(b) Cascading mounts

# Schematic View of NFS Architecture

# NFS Path-Name Translation

- Performed by breaking the path into component names and performing a separate NFS lookup call for every pair of component name and directory vnode

- To make lookup faster, a directory name lookup cache on the client's side holds the vnodes for remote directory names
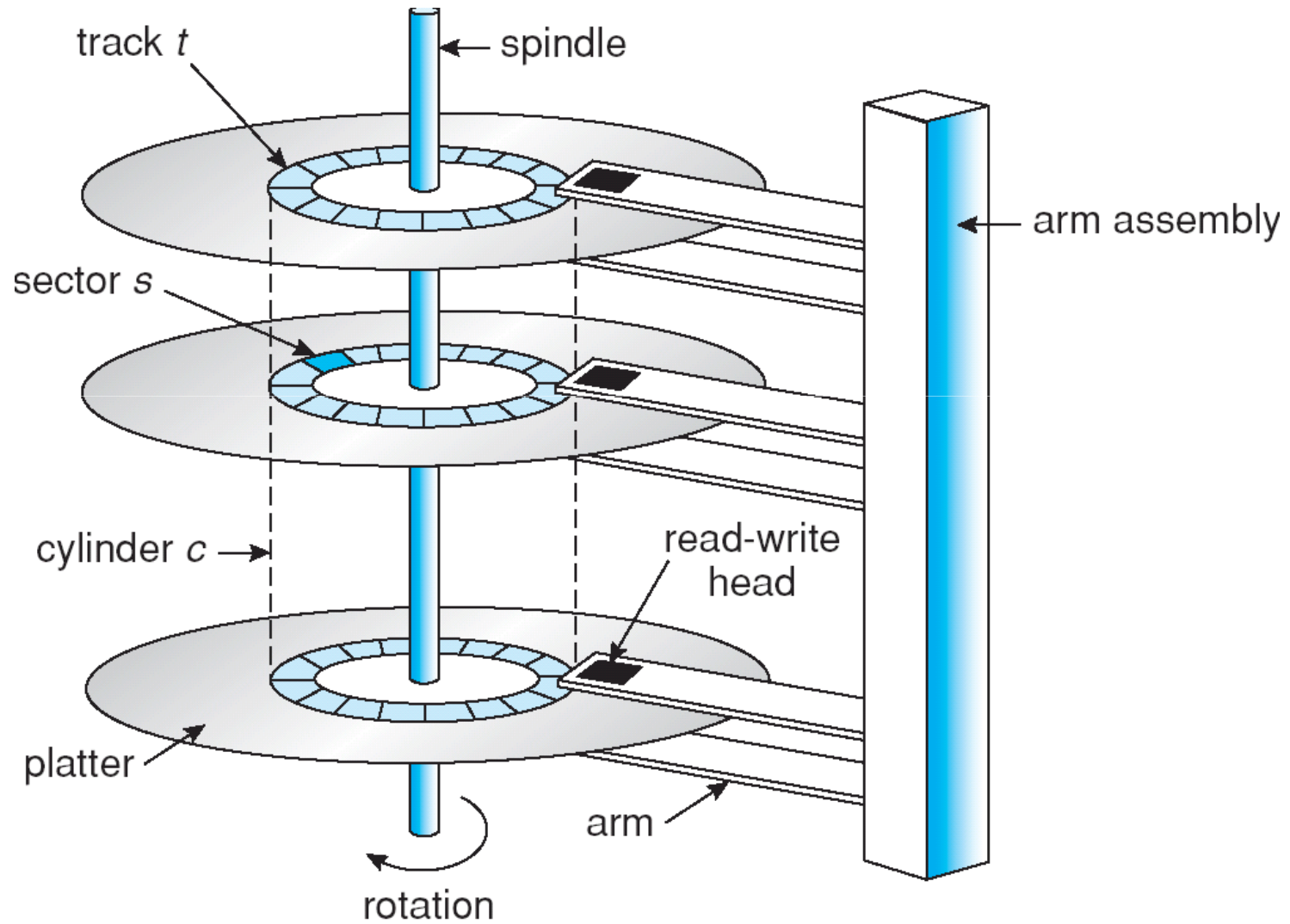
# Mass Storage

# Mass-Storage Systems

- Overview of Mass Storage Structure
- Disk Structure
- Disk Attachment
- Disk Scheduling
- Disk Management
- Swap-Space Management
- RAID Structure
- Disk Attachment
- Stable-Storage Implementation
- Tertiary Storage Devices
- Operating System Support
- Performance Issues

# 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

# Moving-head Disk Mechanism

# 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

- 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

# 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
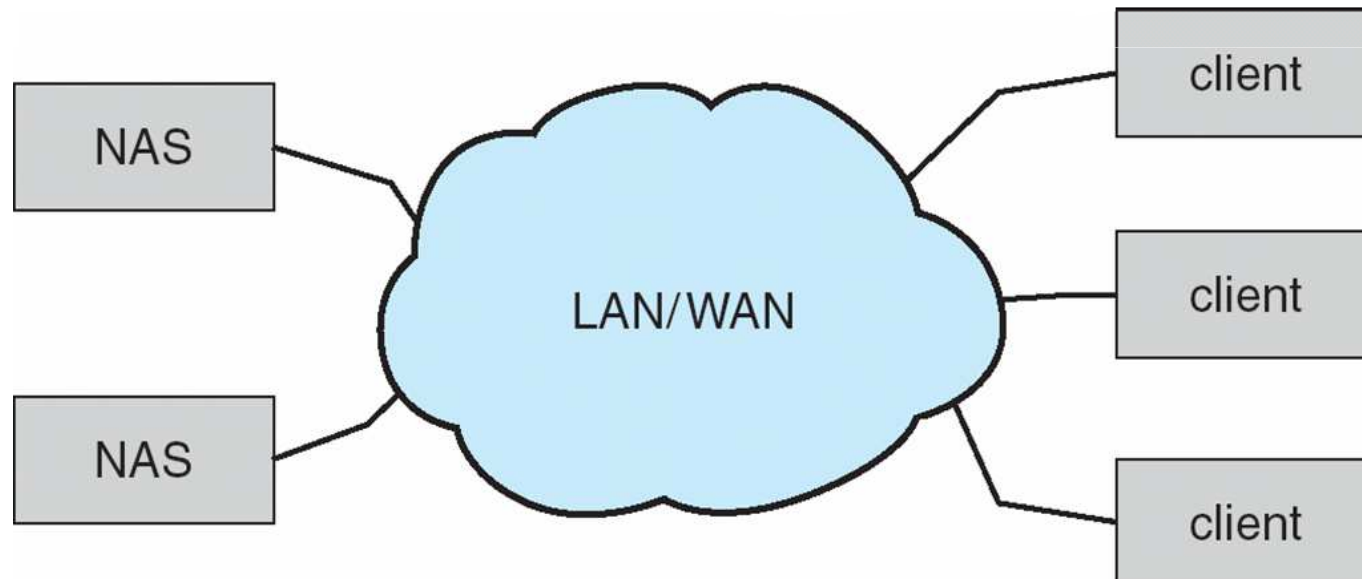  - 20-200GB typical storage

# Disk Structure

- Disk drives are addressed as large 1-dimensional arrays of <span style="color:blue">logical blocks</span>, 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

# 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

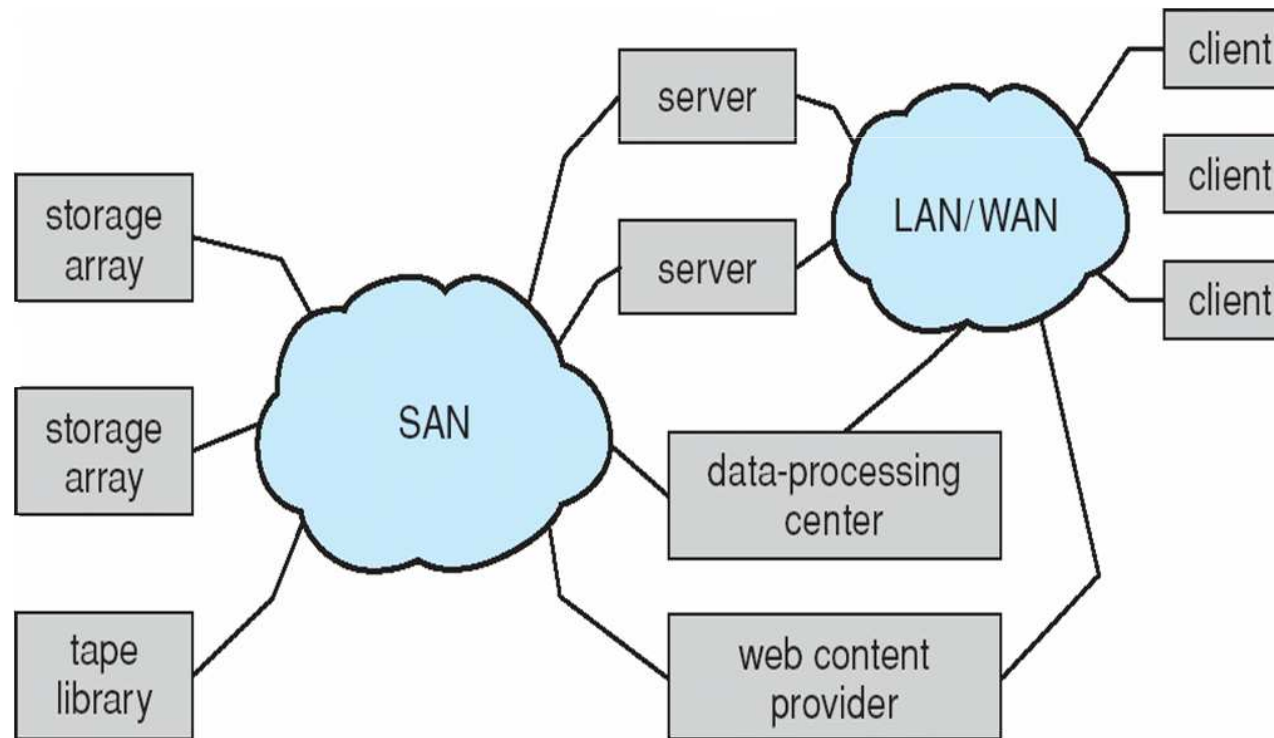  - Each target can have up to 8 logical units (disks attached to device controller

# 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

# Storage Area Network

- Common in large storage environments (and becoming more common)
- Multiple hosts attached to multiple storage arrays – flexible
- Nothing about file systems, they are just blocks.

# 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 ar 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 ≈ 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
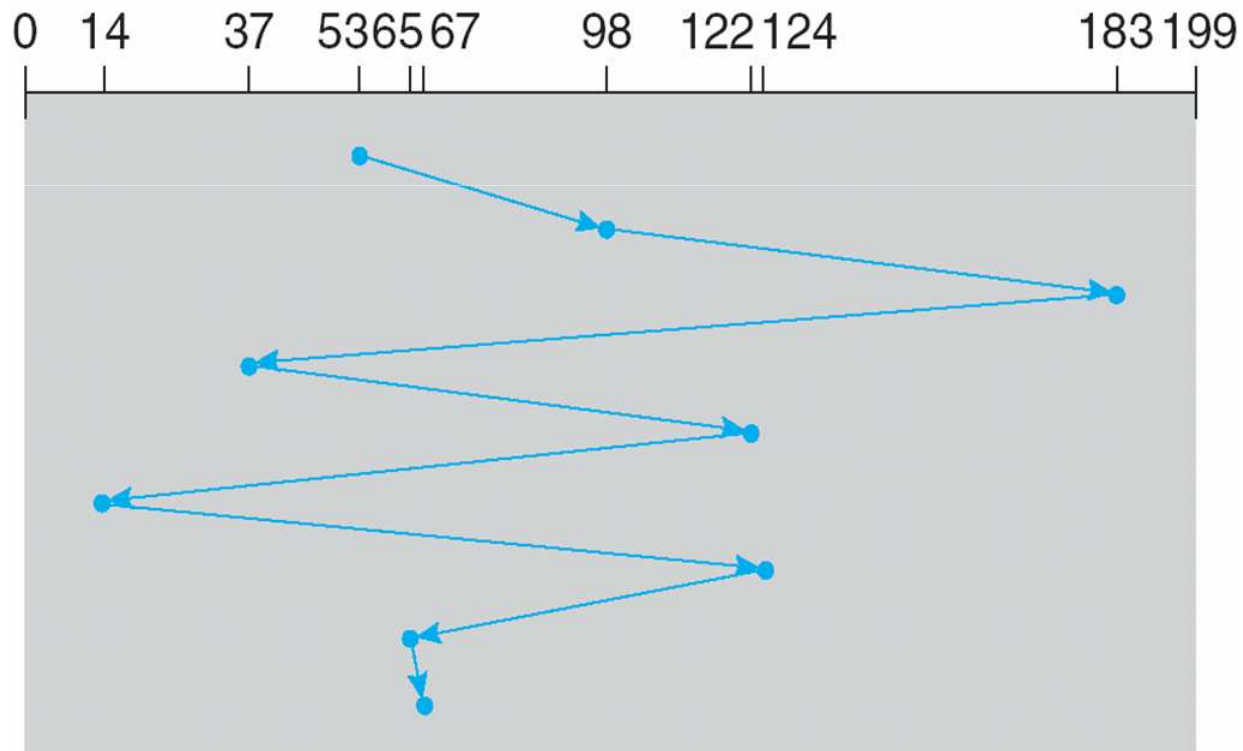
# Disk Scheduling (Cont)

- 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

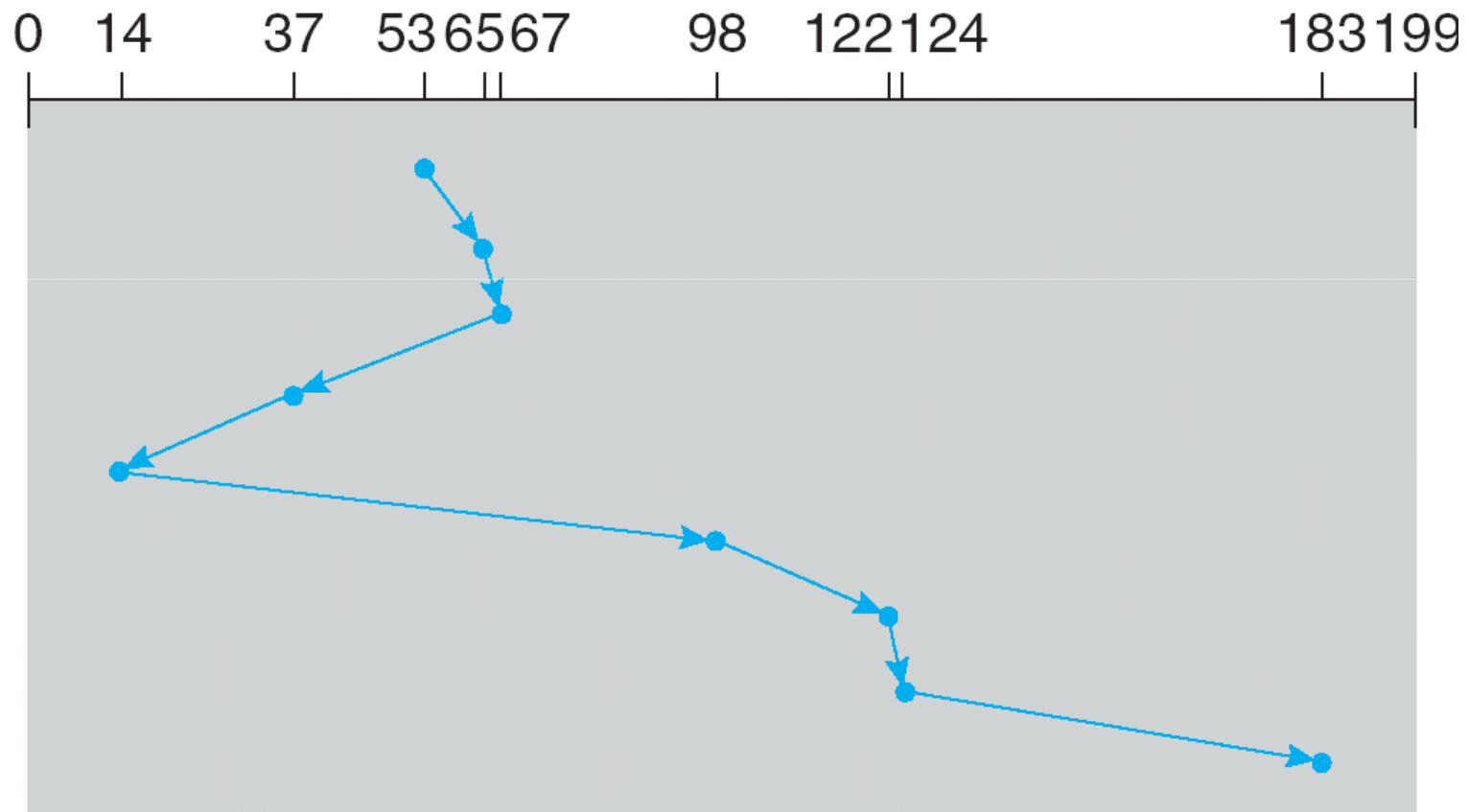Illustration shows total head movement of 640 cylinders

# 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

- Illustration shows total head movement of 236 cylinders

# SSTF (Cont)

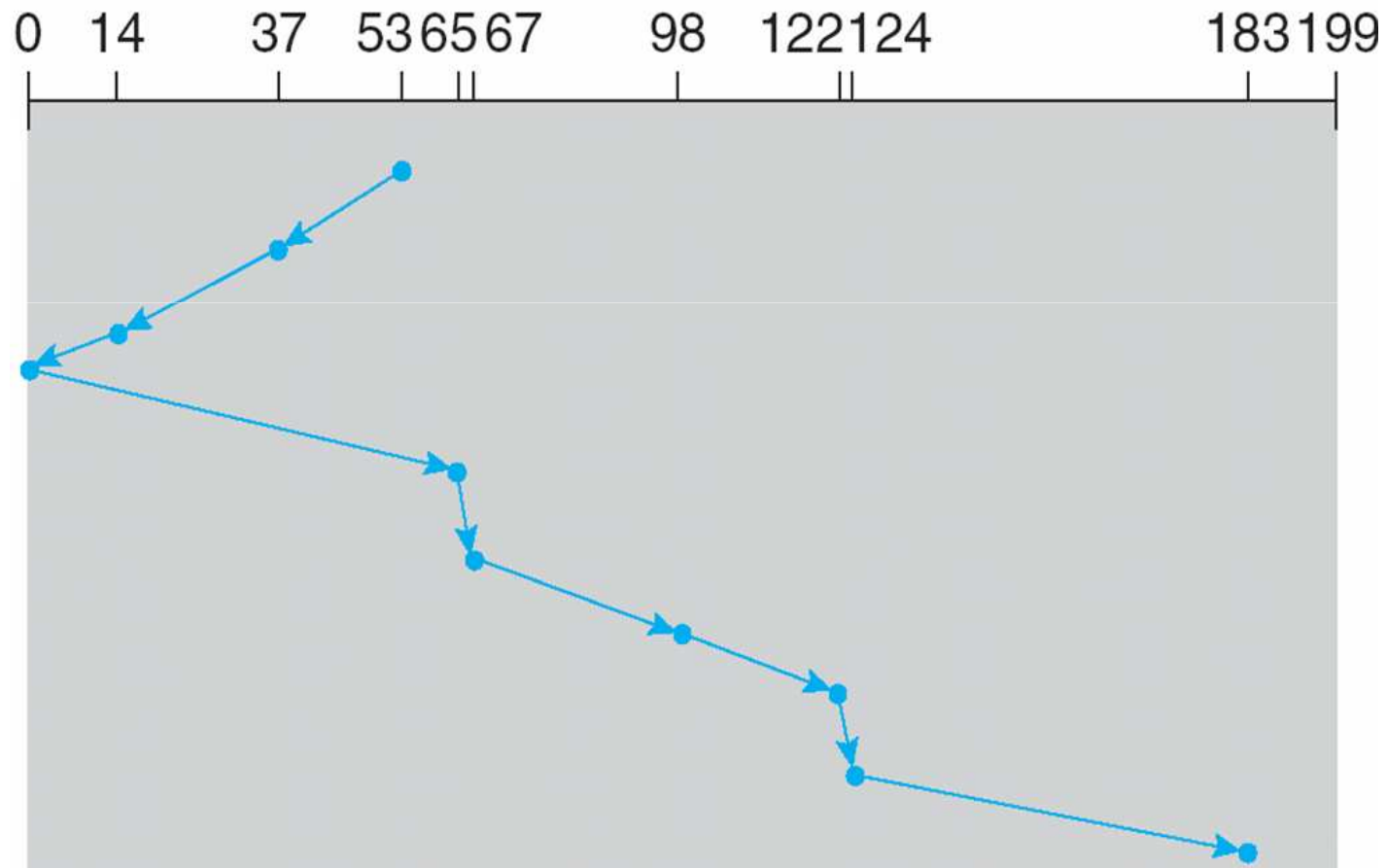queue = 98, 183, 37, 122, 14, 124, 65, 67

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

# SCAN (Cont.)

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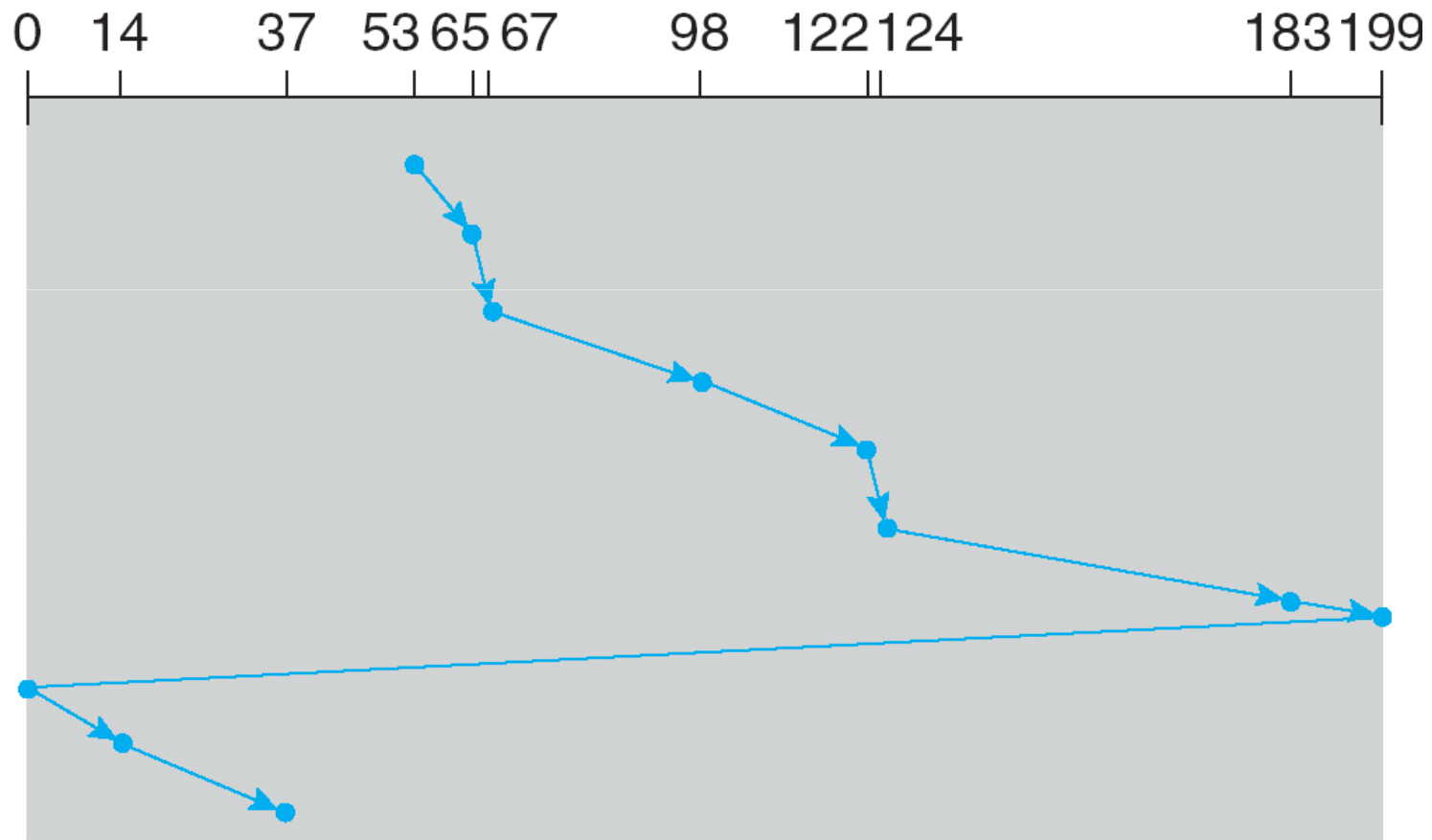
head starts at 53

# 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

# C-SCAN (Cont)

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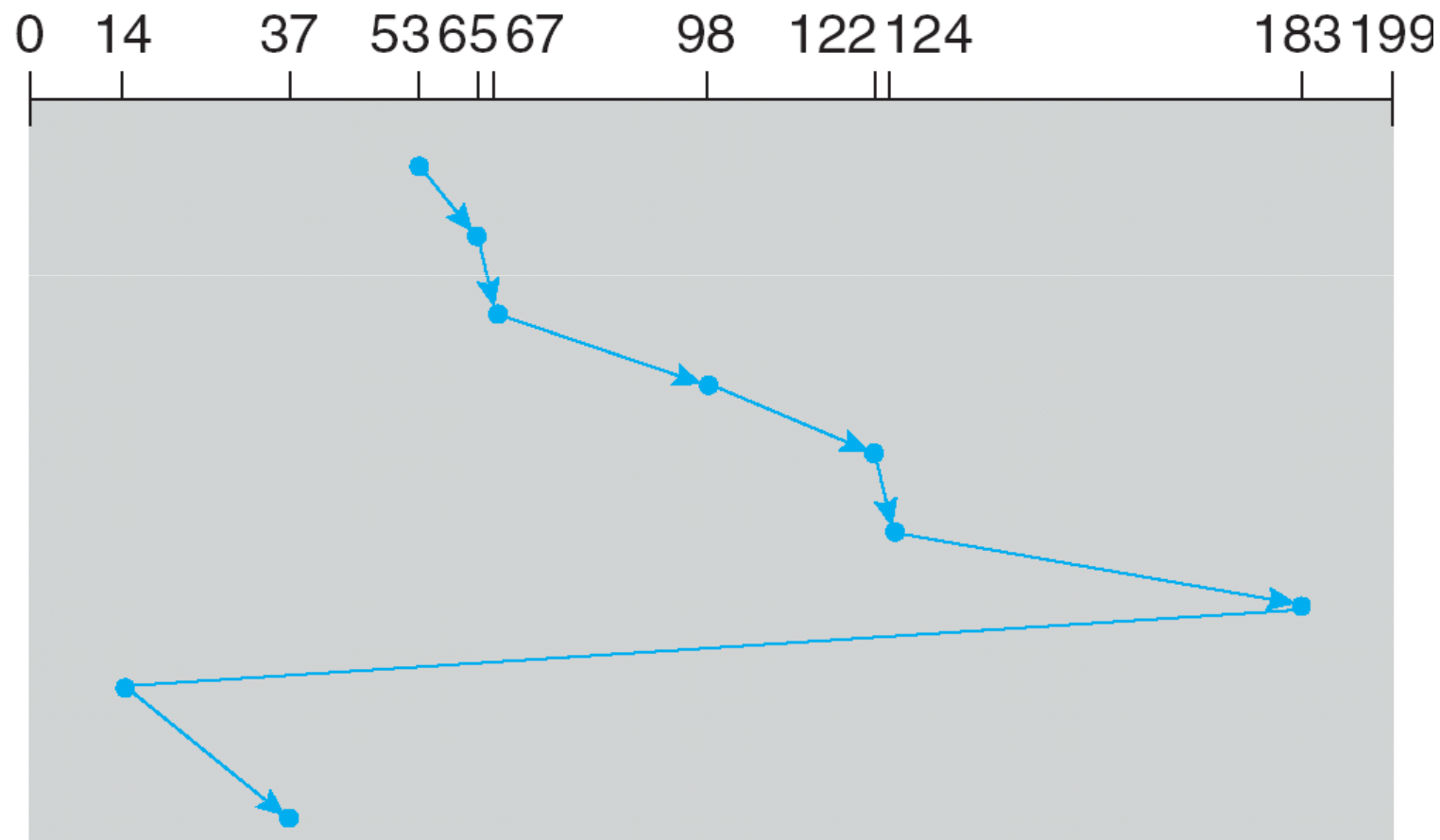
head starts at 53

# 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

# C-LOOK (Cont)

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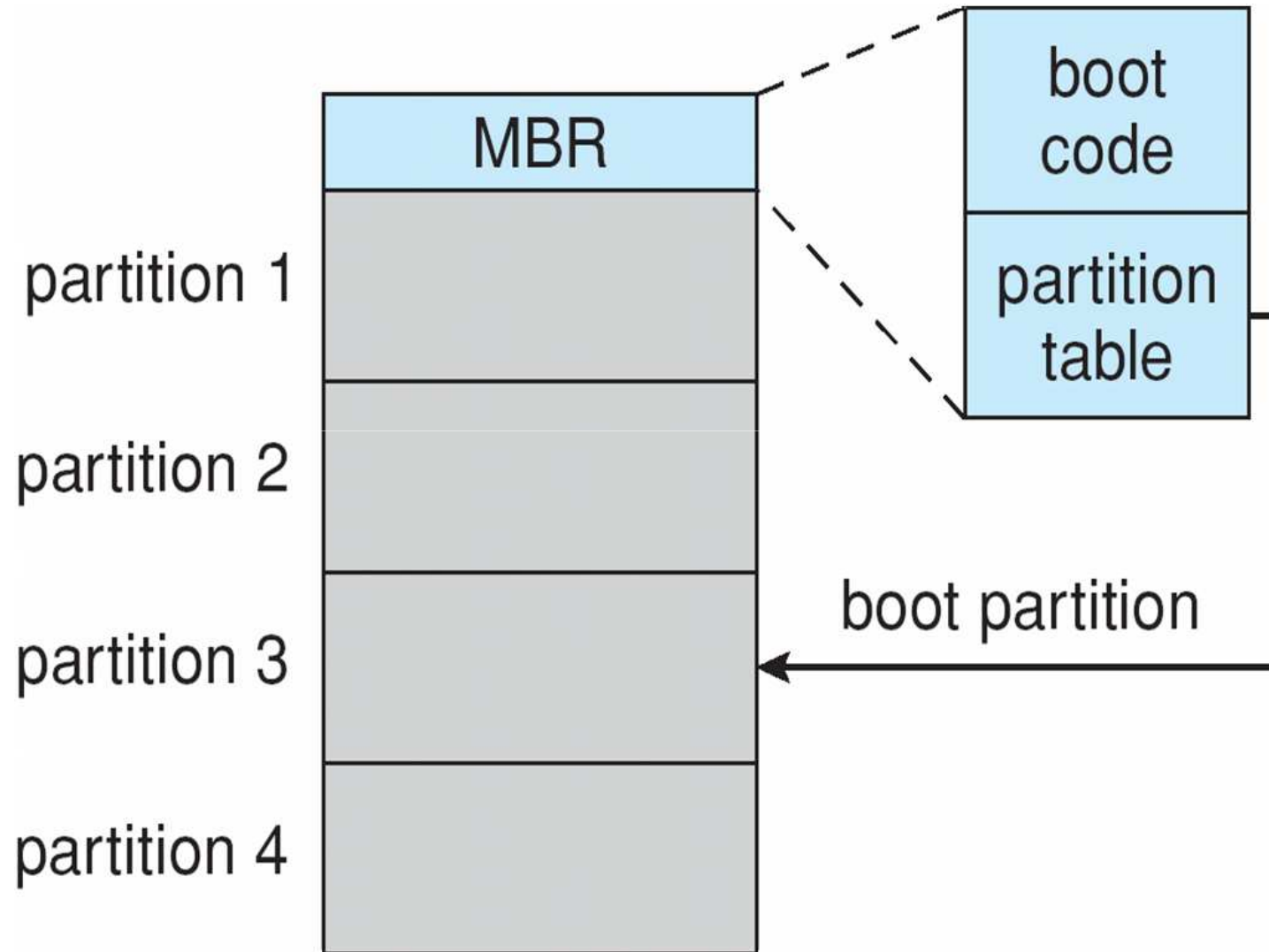
head starts at 53

# 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
- 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 LOOK is a reasonable choice for the default algorithm

# 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"
  - To increase efficiency most file systems group blocks into clusters
    - Disk I/O done in blocks
    - File I/O done in clusters
- Boot block initializes system
  - The bootstrap is stored in ROM
  - Bootstrap loader program
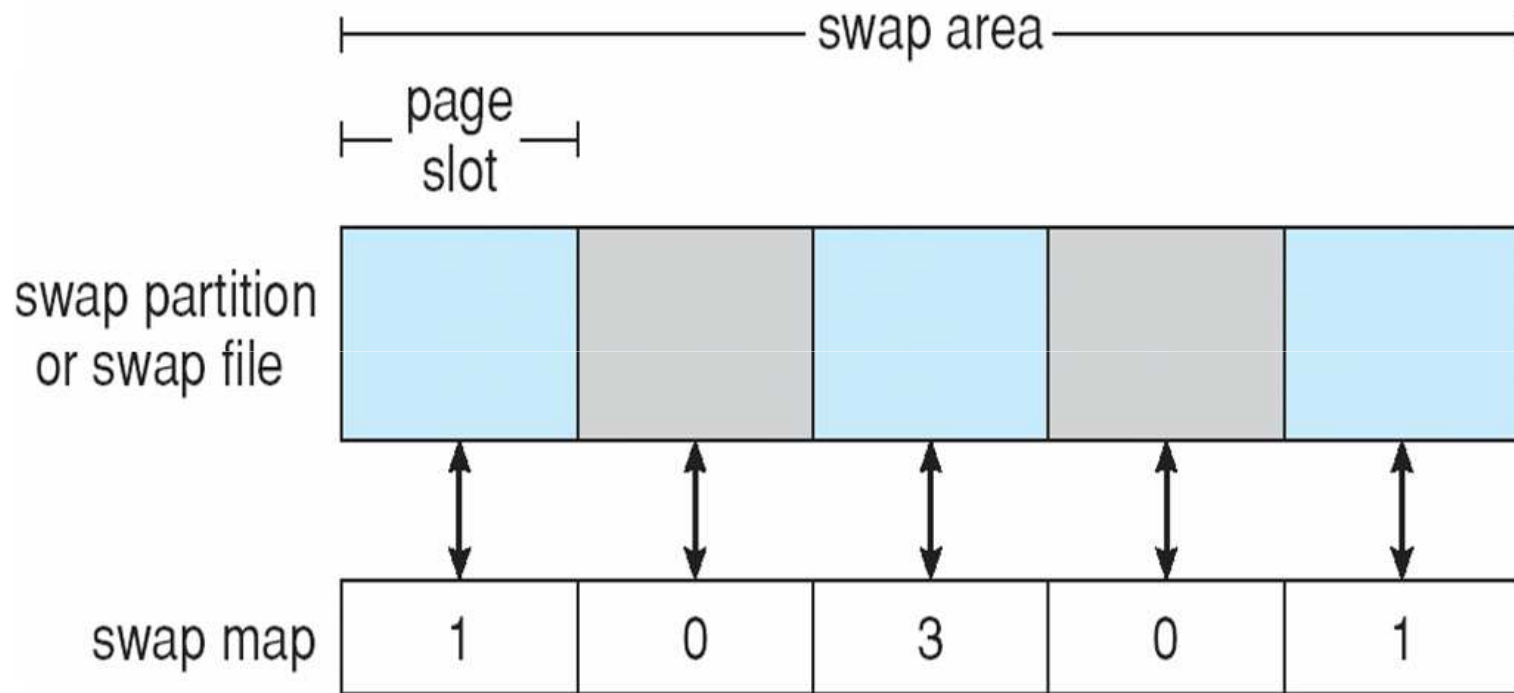- Methods such as sector sparing used to handle bad blocks

# Booting from a Disk in Windows 2000

# Swap-Space Management

- Swap-space — Virtual memory uses disk space as an extension of main memory

- Swap-space 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
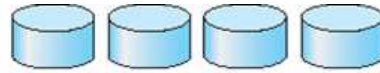
# Data Structures for Swapping on Linux Systems

# RAID Structure

- Redundant Array of Inexpensive (Independent)Disks

- RAID – multiple disk drives provides reliability via redundancy

- Increases the mean time to failure

- Frequently combined with NVRAM to improve write performance

- RAID is arranged into six different levels

# RAID (Cont)

- Mirroring and Striping

- 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  (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.
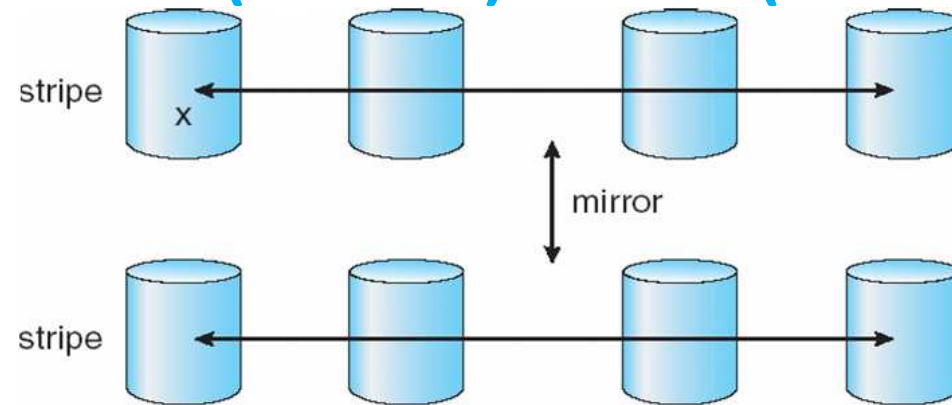
(e) RAID 4: block-interleaved parity.
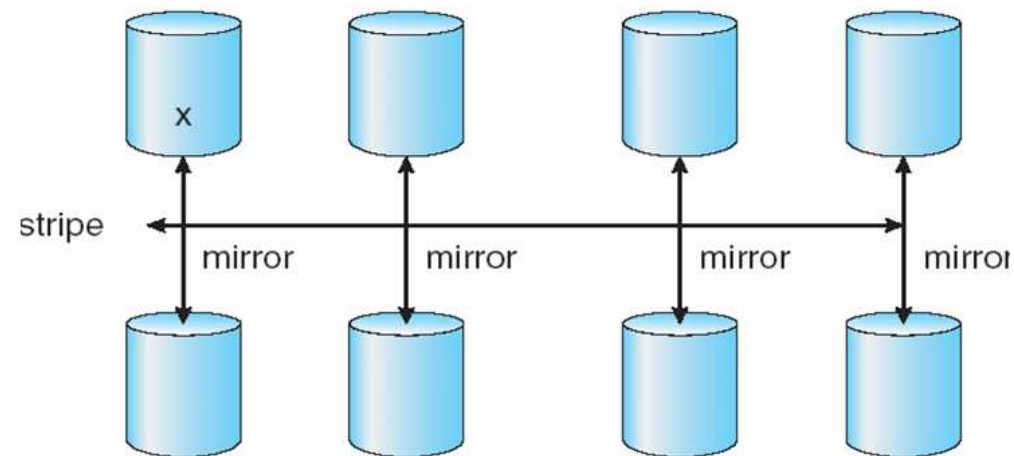
(f) RAID 5: block-interleaved distributed parity.

(g) RAID 6: P + Q redundancy.

69

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



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

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

# Thanks