# Module 5 :

# File-System, Implementation of File System, Secondary Storage Structure, Protection

By: Girish Kumar B C

# File-System

# File-System

- File Concept
- Access Methods
- Directory Structure
- File-System Mounting
- File Sharing
- Protection

# Objectives

- To explain the function of file systems

- To describe the interfaces to file systems

- To discuss file-system design tradeoffs, including access methods, file sharing, file locking, and directory structures

- To explore file-system protection

# File Concept

- Contiguous logical address space

- Types:
  - Data
    - numeric
    - character
    - binary
  - Program

# File Structure

- 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

# 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 file location on 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
- Information about files are kept in the directory structure, which is maintained on the disk

# File Operations

- File is an **abstract data type**
- **Create**
- **Write**
- **Read**
- **Reposition within file**
- **Delete**
- **Truncate**
- *Open($F_i$)* – search the directory structure on disk for entry $F_i$, and move the content of entry to memory
- *Close ($F_i$)* – move the content of entry $F_i$ in memory to directory structure on disk

# Open Files

- Several pieces of data are needed to manage 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

# Open File Locking

- Provided by some operating systems and file systems

- Mediates access to a file

- Mandatory or advisory:

  - **Mandatory** – access is denied depending on locks held and requested

  - **Advisory** – processes can find status of locks and decide what to do

| 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 |

# Access Methods

- **Sequential Access**

  read next
  write next
  reset
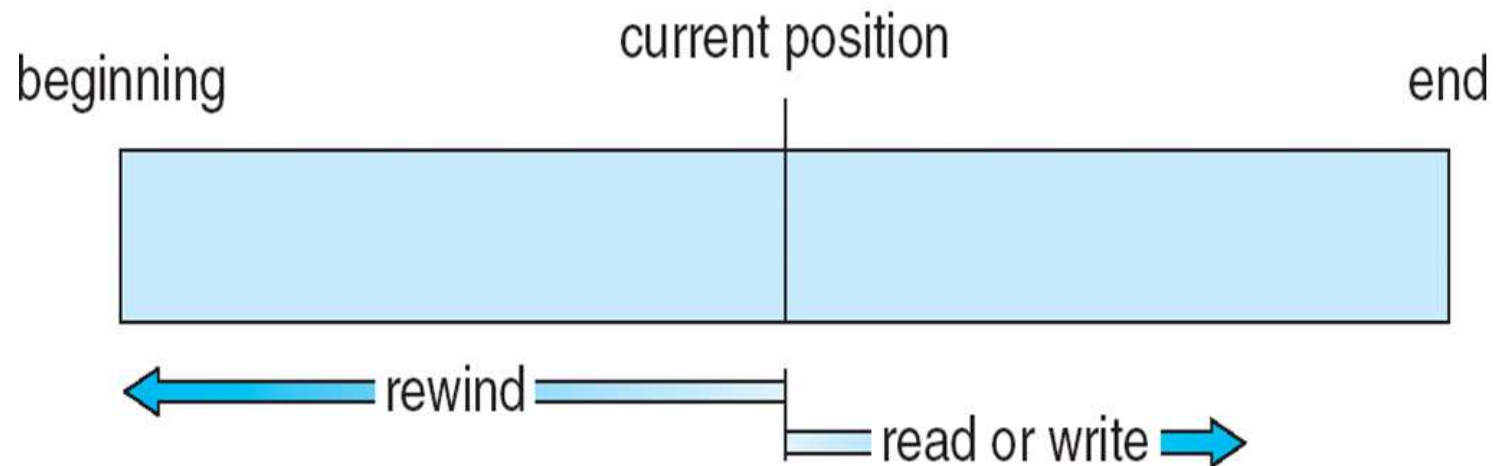  no read after last write
  (rewrite)

- **Direct Access**

  read $n$
  write $n$
  position to $n$
  read next
  write next
  rewrite $n$

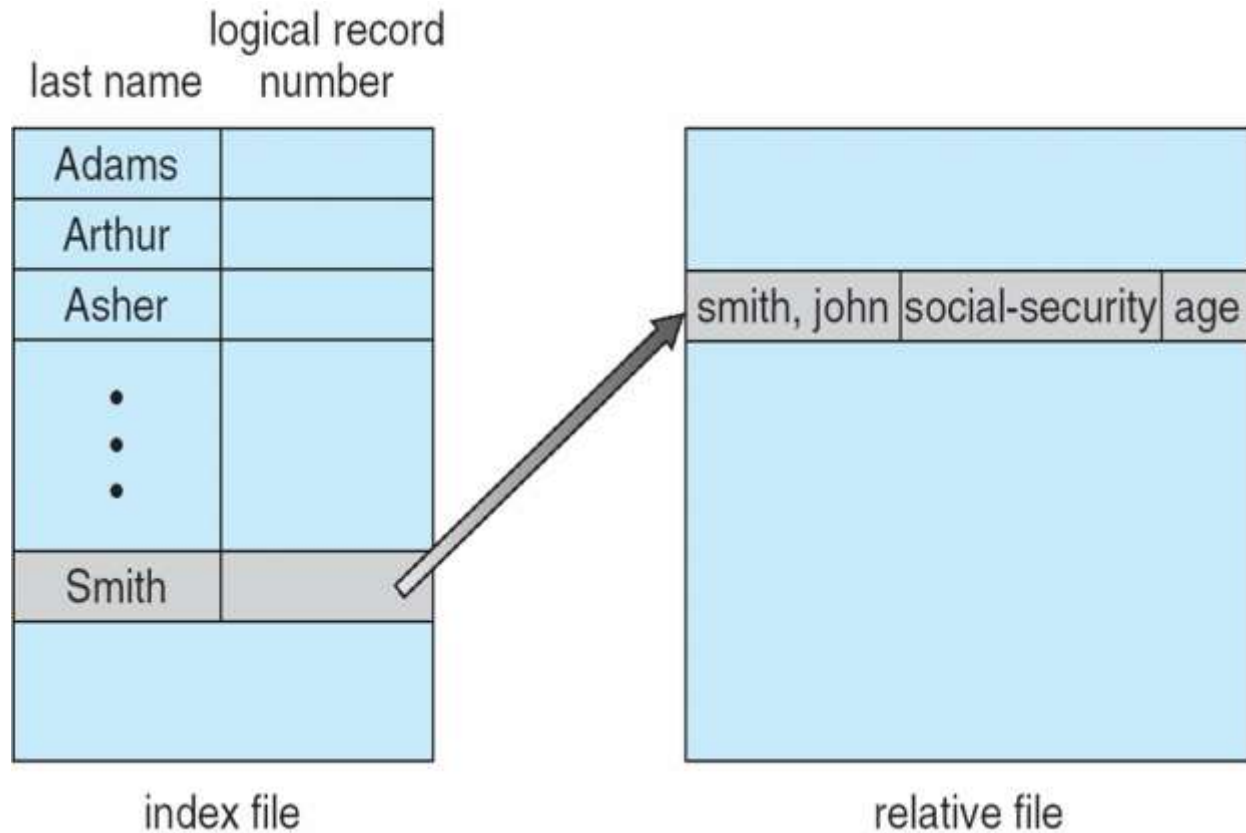$n$ = relative block number

# Sequential-access File

# Simulation of Sequential Access on Direct-access File

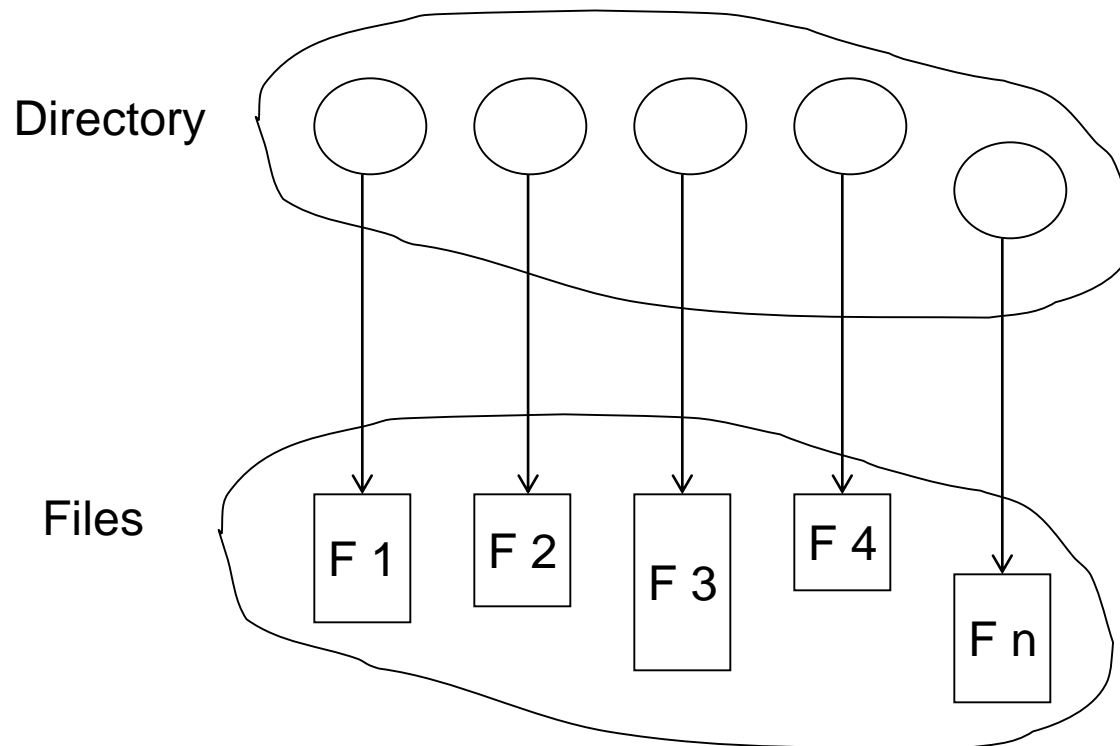| sequential access | implementation for direct access |
|---|---|
| reset | $cp = 0;$ |
| read next | read $cp$; <br> $cp = cp + 1;$ |
| write next | write $cp$; <br> $cp = cp + 1;$ |

# Example of Index and Relative Files

# Directory Structure

■ A collection of nodes containing information about all files

Directory

Files

F 1   F 2   F 3   F 4   F n

Both the directory structure and the files reside on disk
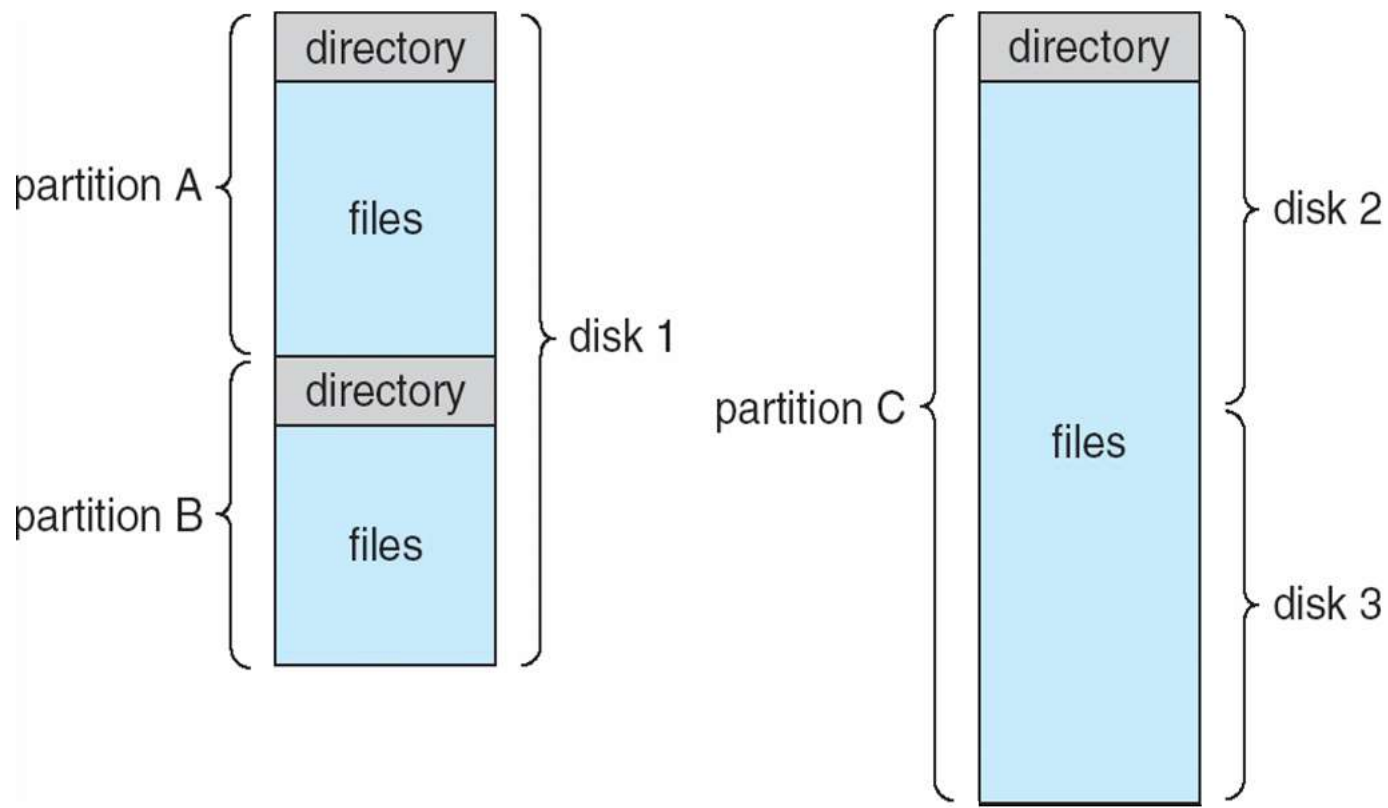Backups of these two structures are kept on tapes

# Disk Structure

Go, change the world

- Disk can be subdivided into partitions

- Disks or partitions can be RAID protected against failure

- Disk or partition can be used raw – without a file system, or formatted with a file system

- Partitions also known as minidisks, slices

- Entity containing file system known as a volume

- Each volume containing file system also tracks that file system's info in device directory or volume table of contents

- As well as general-purpose file systems there are many special-purpose file systems, frequently all within the same operating system or computer

# A Typical File-system Organization

# Operations Performed on Directory

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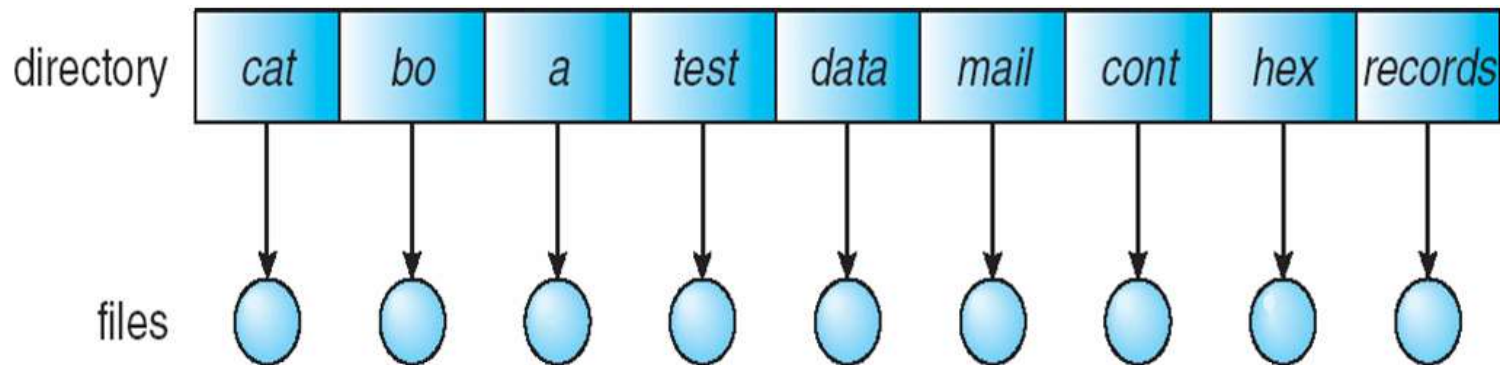- Rename a file
- Traverse the file system

# Organize the Directory (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, …)

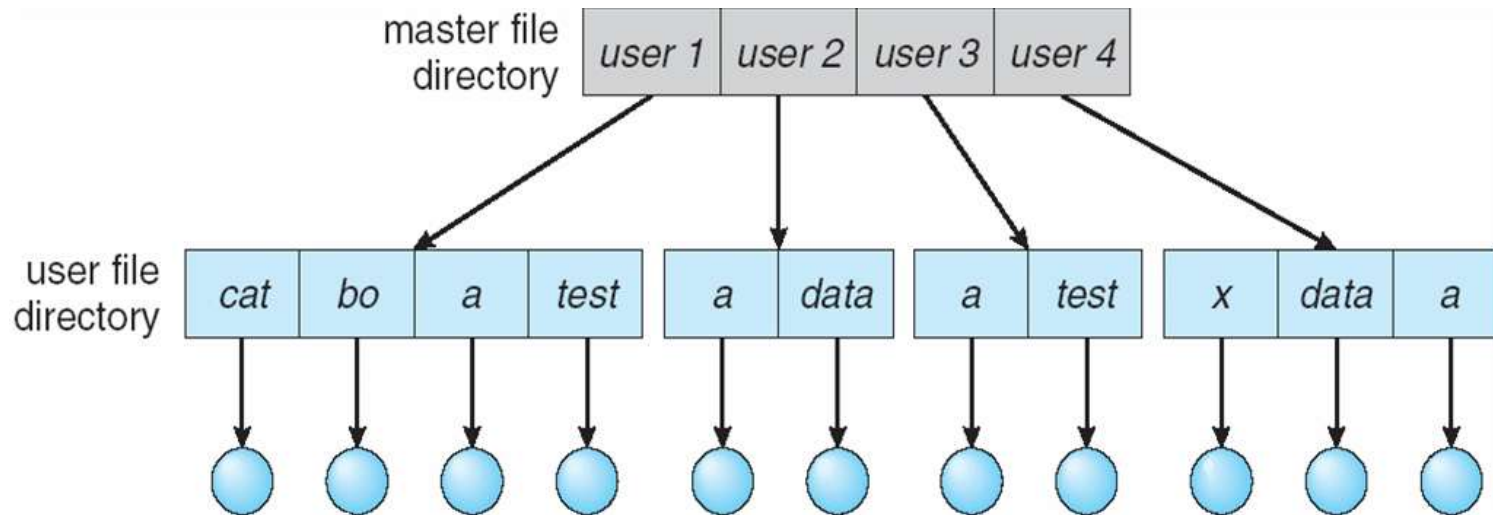■ A single directory for all users



Naming problem

Grouping problem
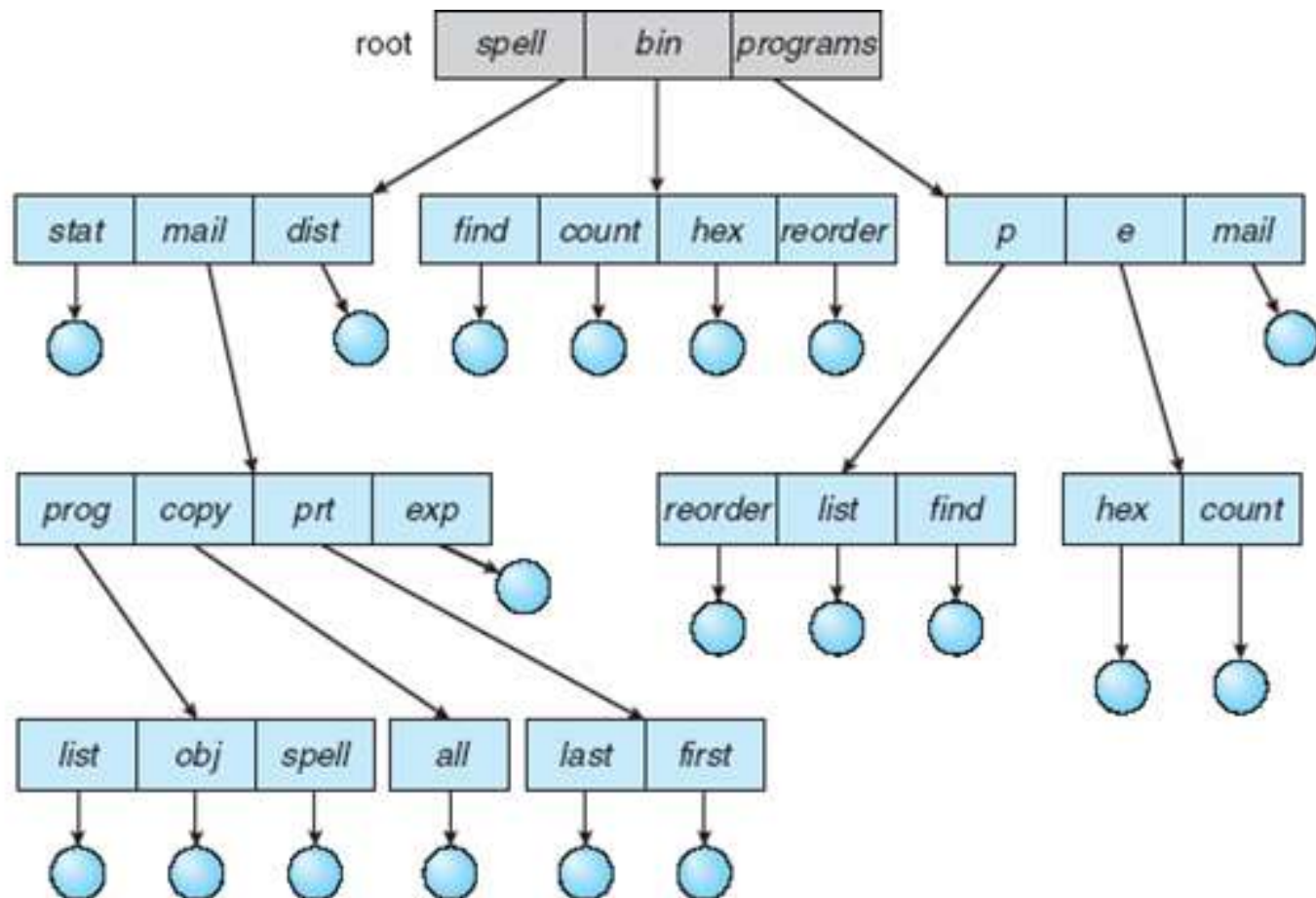
# Two-Level Directory

- Separate directory for each user



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

# Tree-Structured Directories

# Tree-Structured Directories (Cont)

- Efficient searching

- Grouping Capability

- Current directory (working directory)
  - cd /spell/mail/prog
  - type list

# Tree-Structured Directories (Cont)

- **Absolute** or **relative** path name
- Creating a new file is done in current directory
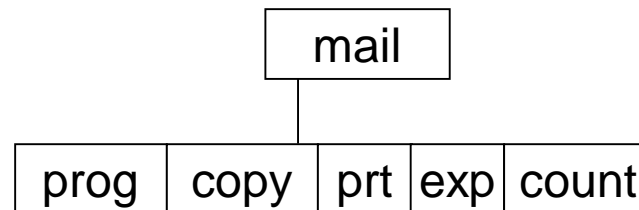- Delete a file

  rm <file-name>

- Creating a new subdirectory is done in current directory

  mkdir <dir-name>

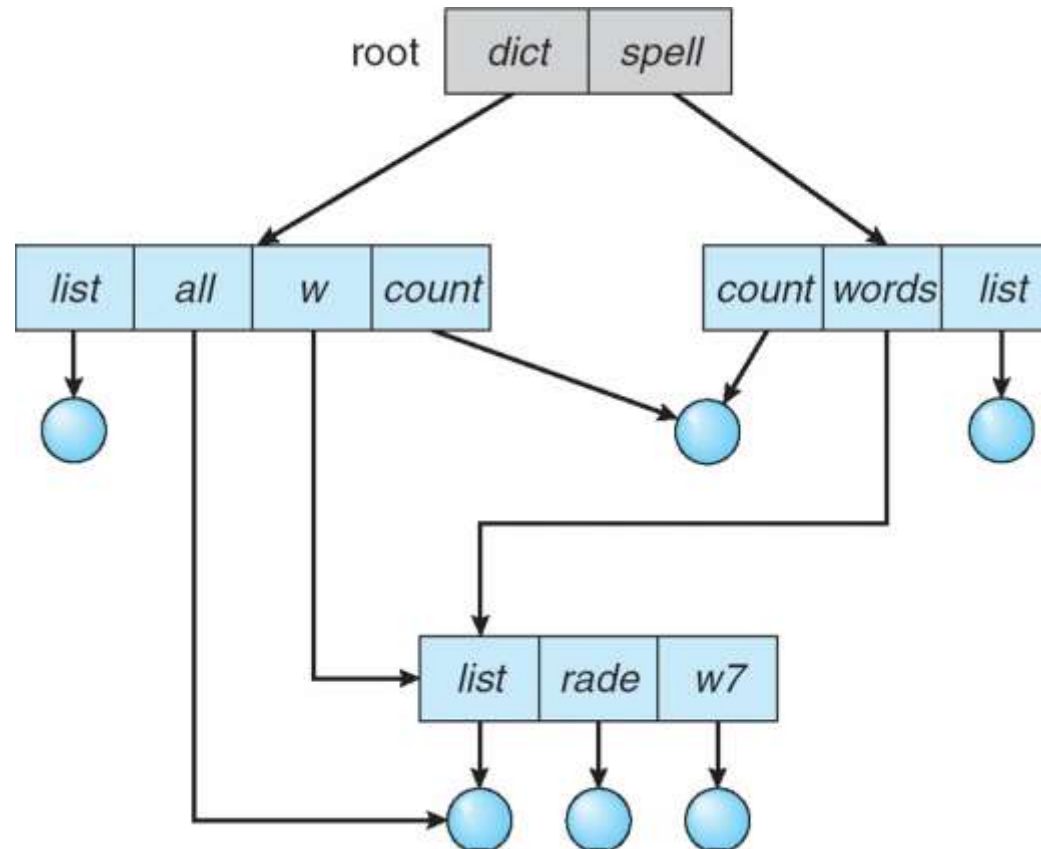  Example:  if in current directory   /mail

  mkdir count

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

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

# Acyclic-Graph Directories

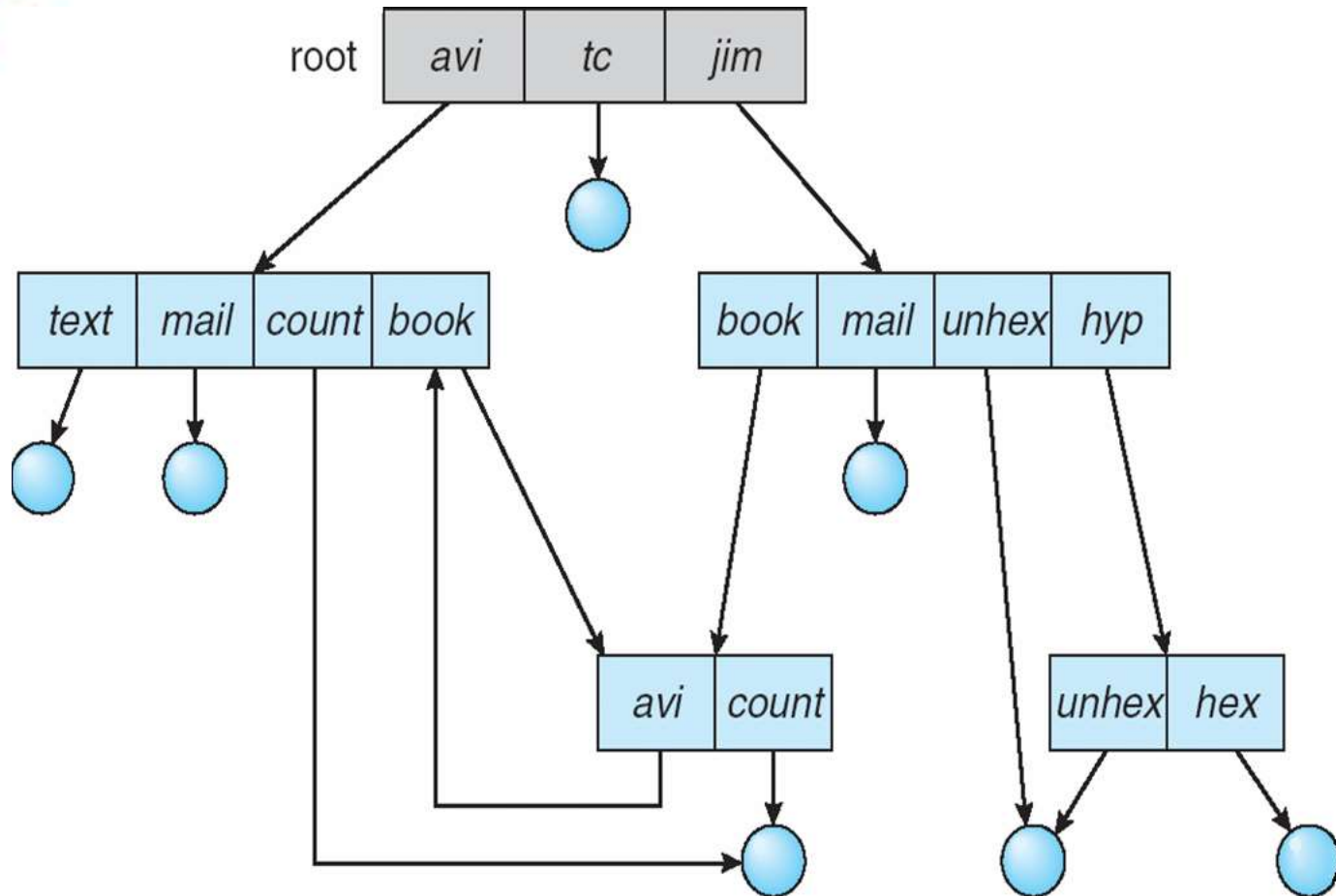- Have shared subdirectories and files

# Acyclic-Graph Directories (Cont.)

- ■ New directory entry type
  - ● **Link** – another name (pointer) to an existing file
  - ● **Resolve the link** – follow pointer to locate the file

# General Graph Directory
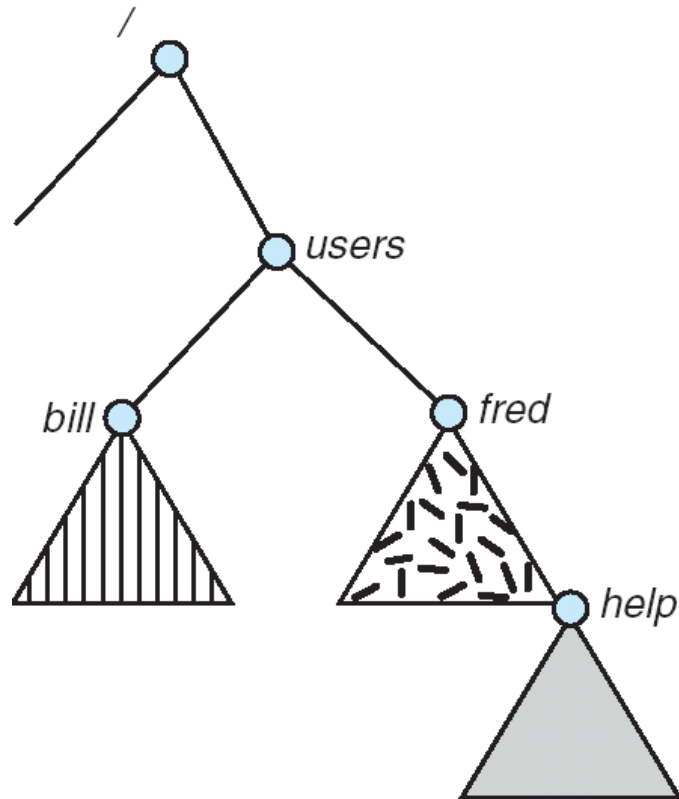
# General Graph Directory (Cont.)

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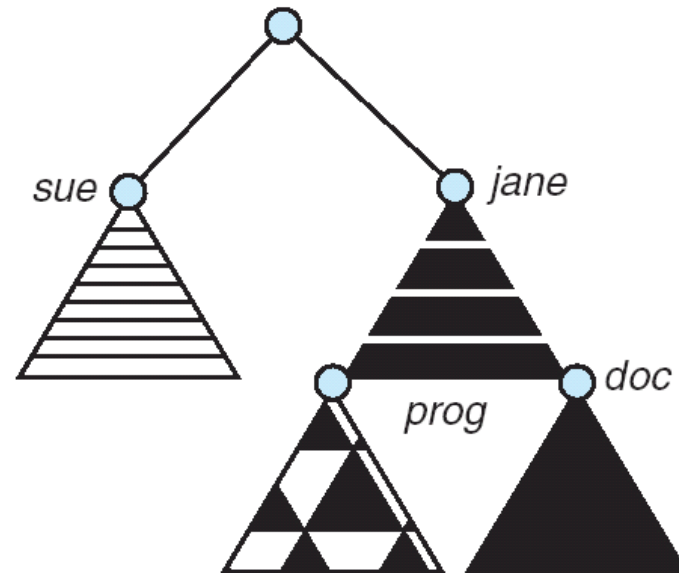  - Every time a new link is added use a cycle detection algorithm to determine whether it is OK

# File System Mounting

- A file system must be **mounted** before it can be accessed

- A unmounted file system (i.e. Fig. 11-11(b)) is mounted at a **mount point**

(a)                    (b)

# File Sharing

- Sharing of files on multi-user systems is desirable

- Sharing may be done through a **protection** scheme

- On distributed systems, files may be shared across a network

- Network File System (NFS) is a common distributed file-sharing method

# File Sharing – Multiple Users

- **User IDs** identify users, allowing permissions and protections to be per-user

- **Group IDs** allow users to be in groups, permitting group access rights

# File Sharing – Remote File Systems

- Uses networking to allow file system access between systems
  - Manually via programs like FTP
  - Automatically, seamlessly using **distributed file systems**
  - Semi automatically via the **world wide web**
- **Client-server** model allows clients to mount remote file systems from servers
  - Server can serve multiple clients
  - Client and user-on-client identification is insecure or complicated
  - **NFS** is standard UNIX client-server file sharing protocol
  - **CIFS** is standard Windows protocol
  - Standard operating system file calls are translated into remote calls
- Distributed Information Systems **(distributed naming services)** such as LDAP, DNS, NIS, Active Directory implement unified access to information needed for remote computing

# File Sharing – Failure Modes

- Remote file systems add new failure modes, due to network failure, server failure

- Recovery from failure can involve state information about status of each remote request

- Stateless protocols such as NFS include all information in each request, allowing easy recovery but less security

# File Sharing – Consistency Semantics

- **Consistency semantics** specify how multiple users are to access a shared file simultaneously
  - Similar to Ch 7 process synchronization algorithms
    - ▸ Tend to be less complex due to disk I/O and network latency (for remote file systems
  - Andrew File System (AFS) implemented complex remote file sharing semantics
  - Unix file system (UFS) implements:
    - ▸ Writes to an open file visible immediately to other users of the same open file
    - ▸ Sharing file pointer to allow multiple users to read and write concurrently
  - AFS has session semantics
    - ▸ Writes only visible to sessions starting after the file is closed

# Protection

- File owner/creator should be able to control:
  - what can be done
  - by whom

- Types of access
  - **Read**
  - **Write**
  - **Execute**
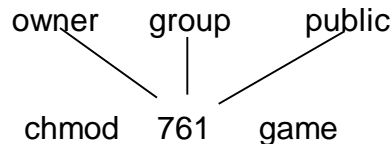  - **Append**
  - **Delete**
  - **List**

*Go, change the world*

- Mode of access:  read, write, execute
- Three classes of users

|  |  |  | RWX |
|---|---|---|---|
| a) **owner access** | 7 | $\Rightarrow$ | 1 1 1 |
|  |  |  | RWX |
| b) **group access** | 6 | $\Rightarrow$ | 1 1 0 |
|  |  |  | RWX |
| 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.
- For a particular file (say *game*) or subdirectory, define an appropriate access.

```
owner    group    public
      \      |      /
  chmod   761    game
```

Attach a group to a file

```
chgrp    G    game
```

# Windows XP Access-control List Management

# A Sample UNIX Directory Listing

| Permissions | Links | Owner | Group | Size | Date | Name |
|---|---|---|---|---|---|---|
| -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/ |

# Chapter 11:  File System Implementation

# Chapter 11: File System Implementation

- File-System Structure

- File-System Implementation

- Directory Implementation

- Allocation Methods

- Free-Space Management

- Efficiency and Performance

- Recovery

- NFS

- Example: WAFL File System

# Objectives

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

- To describe the implementation of remote file systems

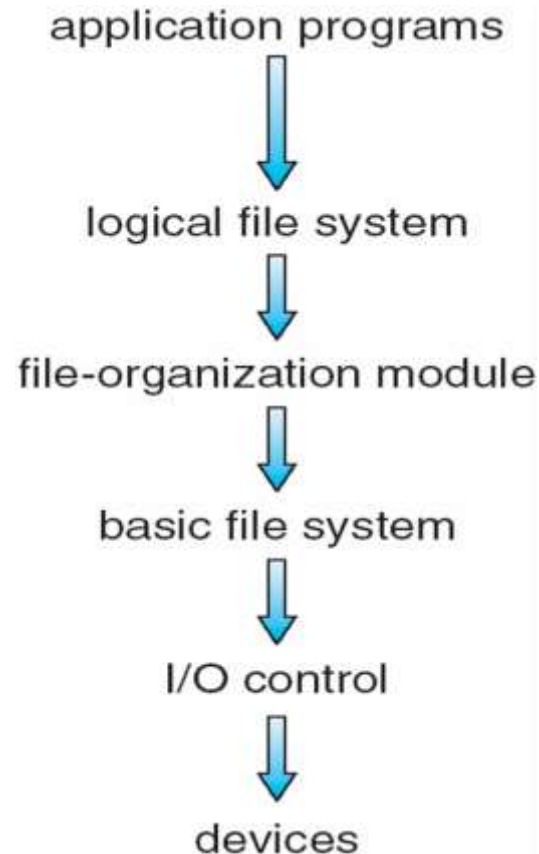- To discuss block allocation and free-block algorithms and trade-offs

# 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

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
  - Logical layers can be implemented by any coding method according to OS designer
- Many file systems, sometimes many within an operating system
  - Each with its own format (CD-ROM is ISO 9660; Unix has **UFS**, FFS; Windows has FAT, FAT32, NTFS as well as floppy, CD, DVD Blu-ray, Linux has more than 40 types, with **extended file system** ext2 and ext3 leading; plus distributed file systems, etc)
  - New ones still arriving – ZFS, GoogleFS, Oracle ASM, FUSE

# 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, master file table)** contains volume details

  - Total # of blocks, # of free blocks, block size, free block pointers or array

- Directory structure organizes the files

  - Names and inode numbers, master file table

- 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

# A Typical File Control Block

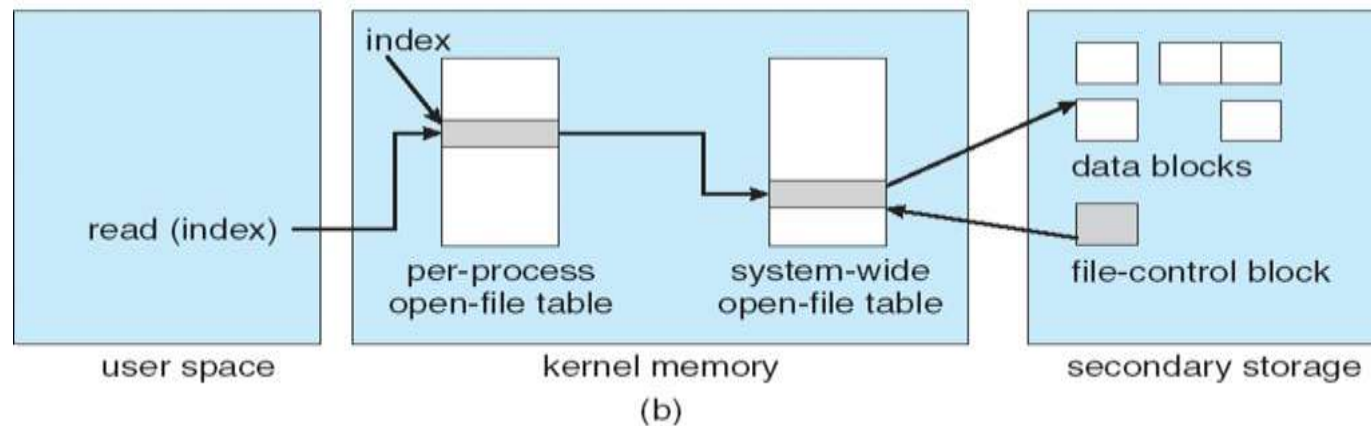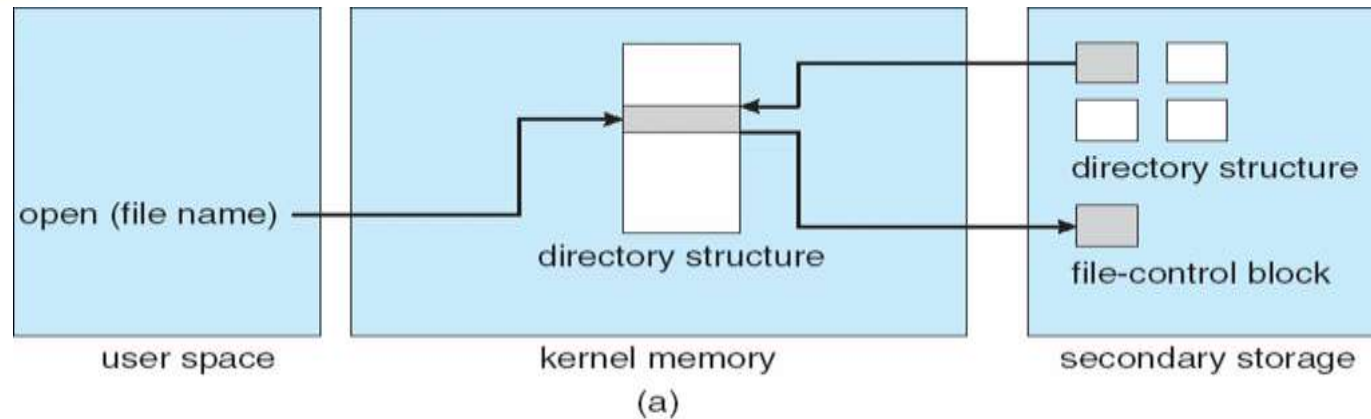| |
|---|
| file permissions |
| file dates (create, access, write) |
| file owner, group, ACL |
| file size |
| file data blocks or pointers to file data blocks |

# In-Memory File System Structures

- 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
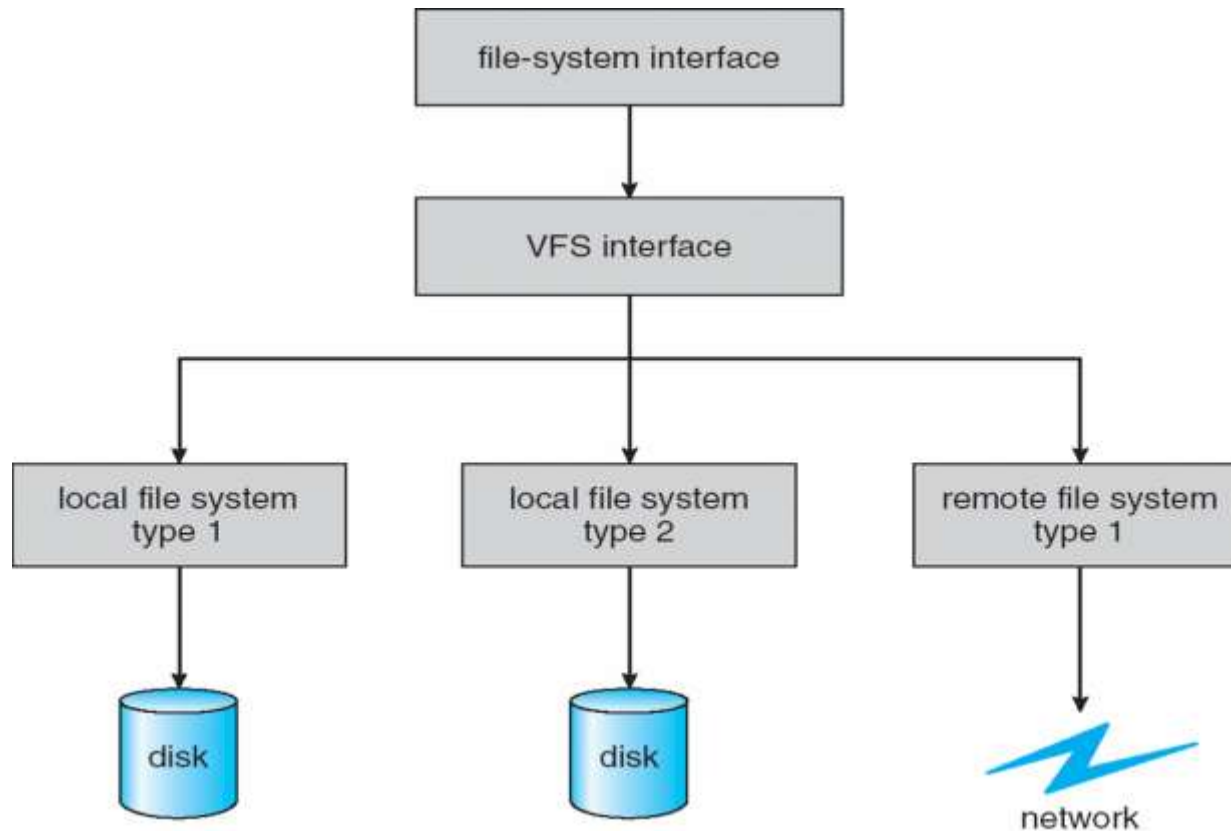    - If yes, add to mount table, allow access

# 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

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

# Schematic View of Virtual 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

# Directory Implementation

- **Linear list** of file names with pointer to the data blocks
  - Simple to program
  - Time-consuming to execute
    - Linear search time
    - Could keep ordered alphabetically via linked list or use B+ tree

- **Hash Table** – linear list with hash data structure
  - Decreases directory search time
  - **Collisions** – situations where two file names hash to the same location
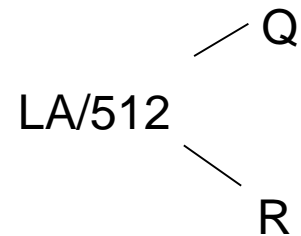  - Only good if entries are fixed size, or use chained-overflow method

# Allocation Methods - Contiguous

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

- **Contiguous allocation** – each file occupies set of contiguous blocks
  - Best performance in most cases
  - Simple – only starting location (block #) and length (number of blocks) are required
  - Problems include finding space for file, knowing file size, external fragmentation, need for **compaction off-line** (**downtime**) or **on-line**

# Contiguous Allocation
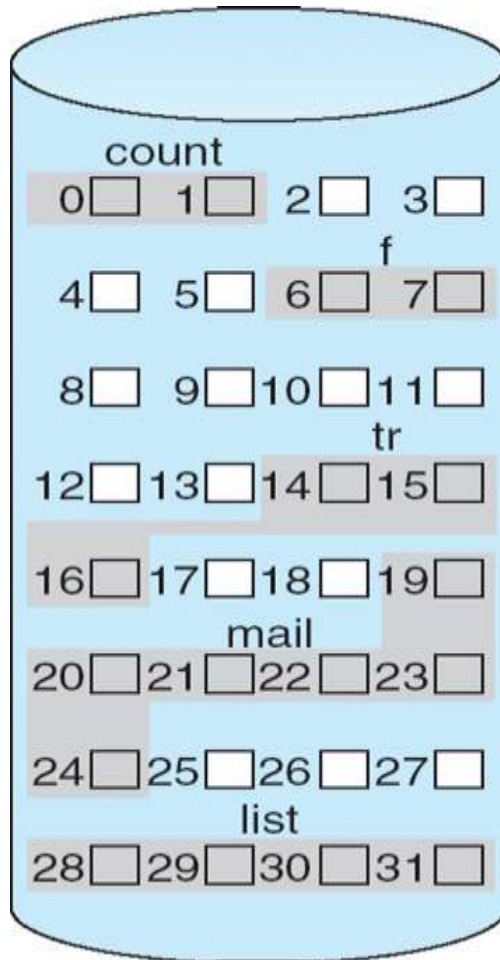
■ Mapping from logical to physical

$$LA/512 \begin{cases} Q \\ R \end{cases}$$

Block to be accessed = Q + starting address
Displacement into block = R

# Contiguous Allocation of Disk Space



| file | start | length |
|------|-------|--------|
| count | 0 | 2 |
| tr | 14 | 3 |
| mail | 19 | 6 |
| list | 28 | 4 |
| f | 6 | 2 |

# Extent-Based Systems

- Many newer file systems (i.e., Veritas File System) use a modified contiguous allocation scheme

- Extent-based file systems allocate disk blocks in extents

- An **extent** is a contiguous block of disks
  - Extents are allocated for file allocation
  - A file consists of one or more extents

# Allocation Methods - Linked

- **Linked allocation** – each file a linked list of blocks
  - File ends at nil pointer
  - No external fragmentation
  - Each block contains pointer to next block
  - No compaction, external fragmentation
  - Free space management system called when new block needed
  - Improve efficiency by clustering blocks into groups but increases internal fragmentation
  - Reliability can be a problem
  - Locating a block can take many I/Os and disk seeks
- FAT (File Allocation Table) variation
  - Beginning of volume has table, indexed by block number
  - Much like a linked list, but faster on disk and cacheable
  - New block allocation simple

# Linked Allocation

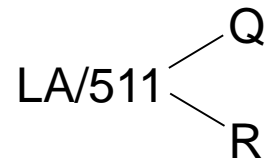- Each file is a linked list of disk blocks: blocks may be scattered anywhere on the disk

block  =  | pointer |
           |         |

# Linked Allocation

- Mapping

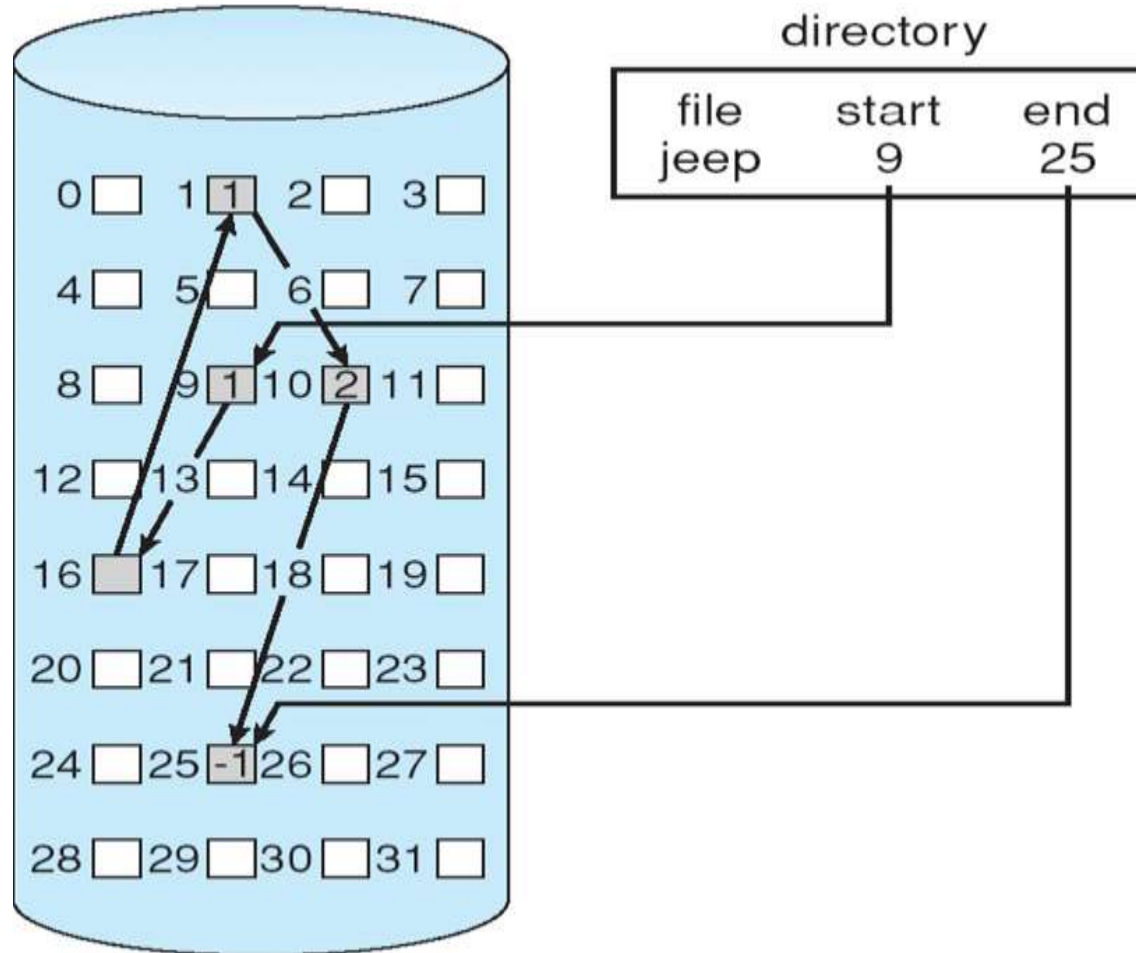$$LA/511 \begin{cases} Q \\ R \end{cases}$$

Block to be accessed is the Qth block in the linked chain of blocks representing the file.
Displacement into block = R + 1

# Linked Allocation

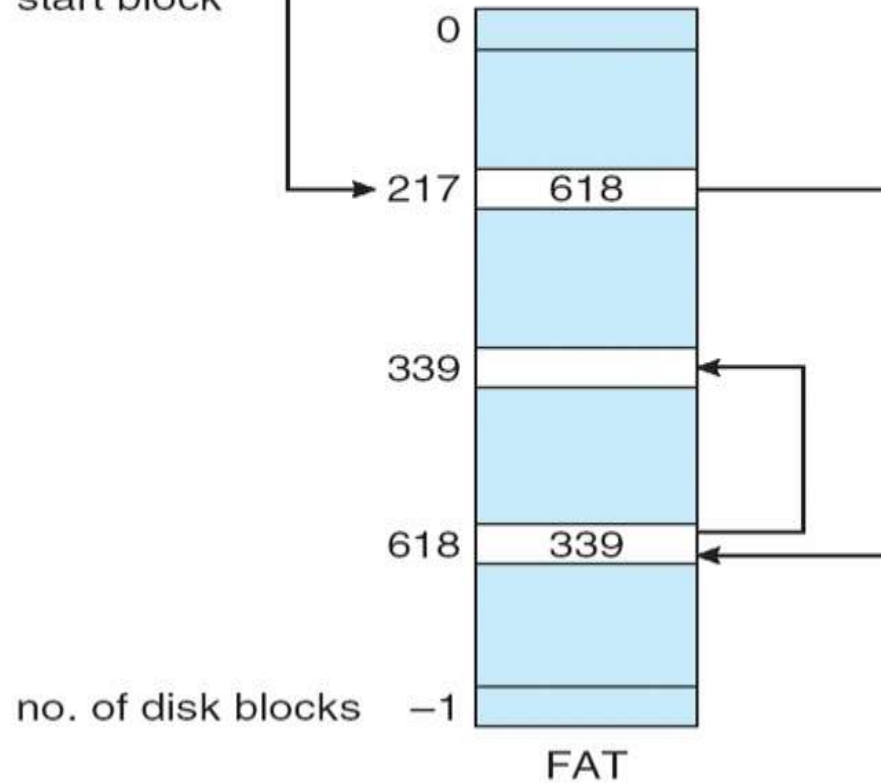# File-Allocation Table

directory entry

| test | • • • | 217 |
|------|-------|-----|

name             start block

| | |
|---|---|
| 0 | |
| 217 | 618 |
| 339 | |
| 618 | 339 |
| no. of disk blocks   −1 | |

FAT
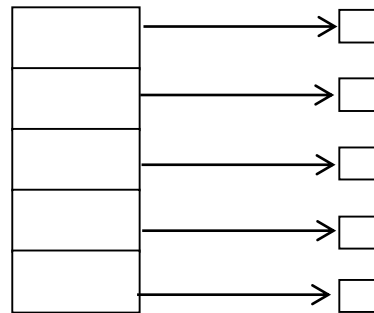
# Allocation Methods - Indexed

- **Indexed allocation**
    - Each file has its own **index block**(s) of pointers to its data blocks

- Logical view



index table

# Indexed Allocation (Cont.)

- Need index table

- Random access

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

- Mapping from logical to physical in a file of maximum size of 256K bytes and block size of 512 bytes. We need only 1 block for index table

$$LA/512 \begin{array}{c} Q \\ R \end{array}$$

      Q = displacement into index table
      R = displacement into block

- Mapping from logical to physical in a file of unbounded length (block size of 512 words)

- Linked scheme – Link blocks of index table (no limit on size)

$$LA / (512 \times 511) \begin{cases} Q_1 \\ R_1 \end{cases}$$

$Q_1$ = block of index table
$R_1$ is used as follows:

$$R_1 / 512 \begin{cases} Q_2 \\ R_2 \end{cases}$$

$Q_2$ = displacement into block of index table
$R_2$ displacement into block of file:

  

# Indexed Allocation – Mapping (Cont.)

- Two-level index (4K blocks could store 1,024 four-byte pointers in outer index -> 1,048,567 data blocks and file size of up to 4GB)

$$LA / (512 \times 512) \begin{cases} Q_1 \\ R_1 \end{cases}$$

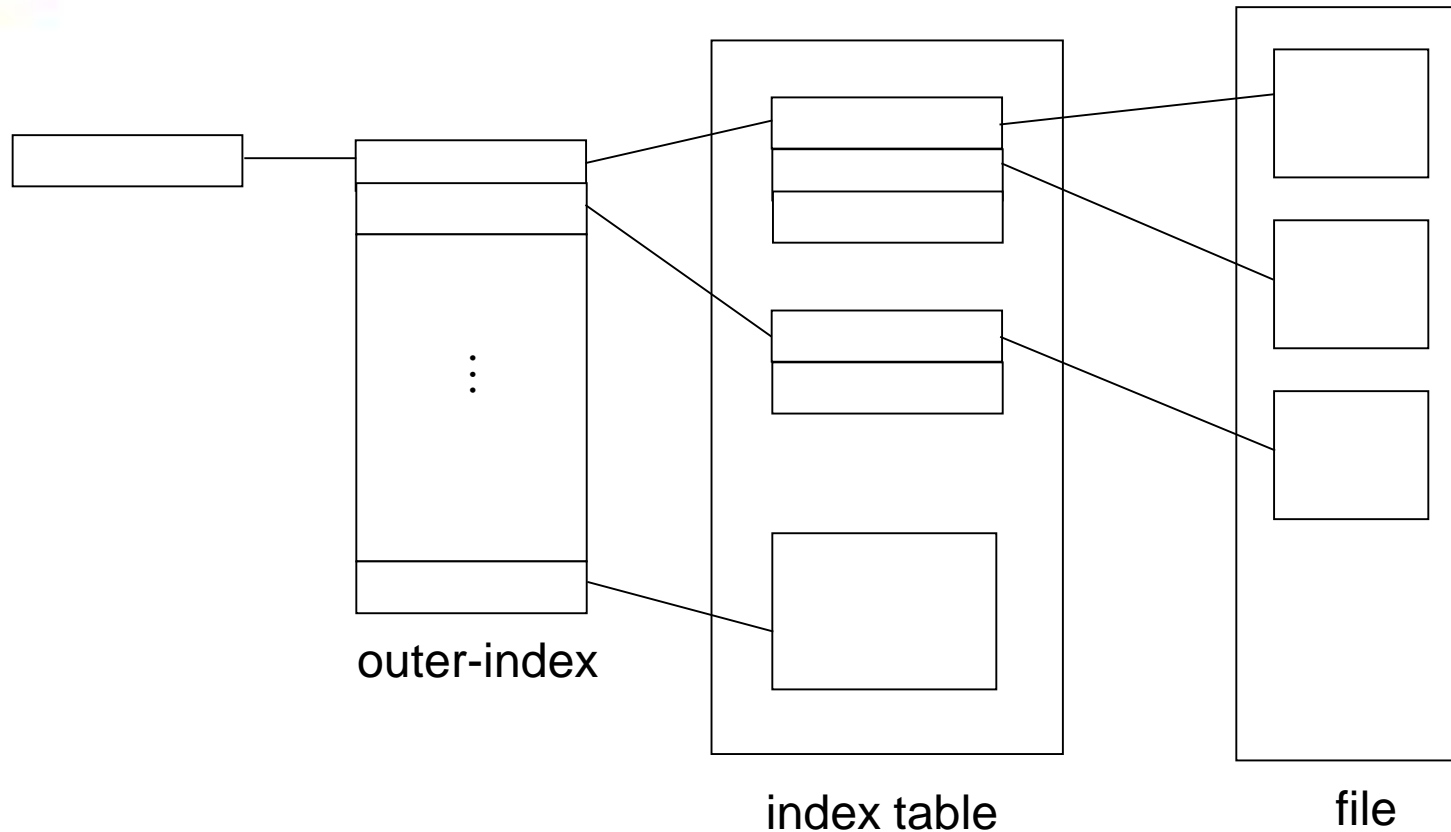$Q_1$ = displacement into outer-index
$R_1$ is used as follows:

$$R_1 / 512 \begin{cases} Q_2 \\ R_2 \end{cases}$$

$Q_2$ = displacement into block of index table
$R_2$ displacement into block of file:

# Indexed Allocation – Mapping (Cont.)



outer-index

index table

file

Note: More index blocks than can be addressed with 32-bit file pointer

# Performance

- Best method depends on file access type
  - Contiguous great for sequential and random
- Linked good for sequential, not random
- Declare access type at creation -> select either contiguous or linked
- Indexed more complex
  - Single block access could require 2 index block reads then data block read
  - Clustering can help improve throughput, reduce CPU overhead

# Performance (Cont.)

- Adding instructions to the execution path to save one disk I/O is reasonable
  - Intel Core i7 Extreme Edition 990x (2011) at 3.46Ghz = 159,000 MIPS
    - http://en.wikipedia.org/wiki/Instructions_per_second
  - Typical disk drive at 250 I/Os per second
    - 159,000 MIPS / 250 = 630 million instructions during one disk I/O
  - Fast SSD drives provide 60,000 IOPS
    - 159,000 MIPS / 60,000 = 2.65 millions instructions during one disk I/O

# Free-Space Management

- File system maintains **free-space list** to track available blocks/clusters
  - (Using term "block" for simplicity)
- **Bit vector** or **bit map** ($n$ blocks)

$$\begin{array}{cccccccc} 0 & 1 & 2 & & & & & n-1 \end{array}$$

| | | | | | | … | |
|---|---|---|---|---|---|---|---|

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

Block number calculation

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

CPUs have instructions to return offset within word of first "1" bit

- Bit map requires extra space
  - Example:

    block size = 4KB = $2^{12}$ bytes

    disk size = $2^{40}$ bytes (1 terabyte)

    $n = 2^{40}/2^{12} = 2^{28}$ bits (or 256 MB)

    if clusters of 4 blocks -> 64MB of memory

- Easy to get contiguous files

- Linked list (free list)
  - Cannot get contiguous space easily
  - No waste of space
  - No need to traverse the entire list (if # free blocks recorded)

# Linked Free Space List on Disk



free-space list head

# Free-Space Management (Cont.)

- Grouping
  - Modify linked list to store address of next *n-1* free blocks in first free block, plus a pointer to next block that contains free-block-pointers (like this one)

- Counting
  - Because space is frequently contiguously used and freed, with contiguous-allocation allocation, extents, or clustering
    - Keep address of first free block and count of following free blocks
    - Free space list then has entries containing addresses and counts

# Free-Space Management (Cont.)

- Space Maps
  - Used in ZFS
  - Consider meta-data I/O on very large file systems
    - ▸ Full data structures like bit maps couldn't fit in memory -> thousands of I/Os
  - Divides device space into **metaslab** units and manages metaslabs
    - ▸ Given volume can contain hundreds of metaslabs
  - Each metaslab has associated space map
    - ▸ Uses counting algorithm
  - But records to log file rather than file system
    - ▸ Log of all block activity, in time order, in counting format
  - Metaslab activity -> load space map into memory in balanced-tree structure, indexed by offset
    - ▸ Replay log into that structure
    - ▸ Combine contiguous free blocks into single entry

# End

# Secondary Storage Structure, Protection

# Chapter 12:  Mass-Storage Systems

- Overview of Mass Storage Structure

- Disk Structure

- Disk Attachment

- Disk Scheduling

- Disk Management

# 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
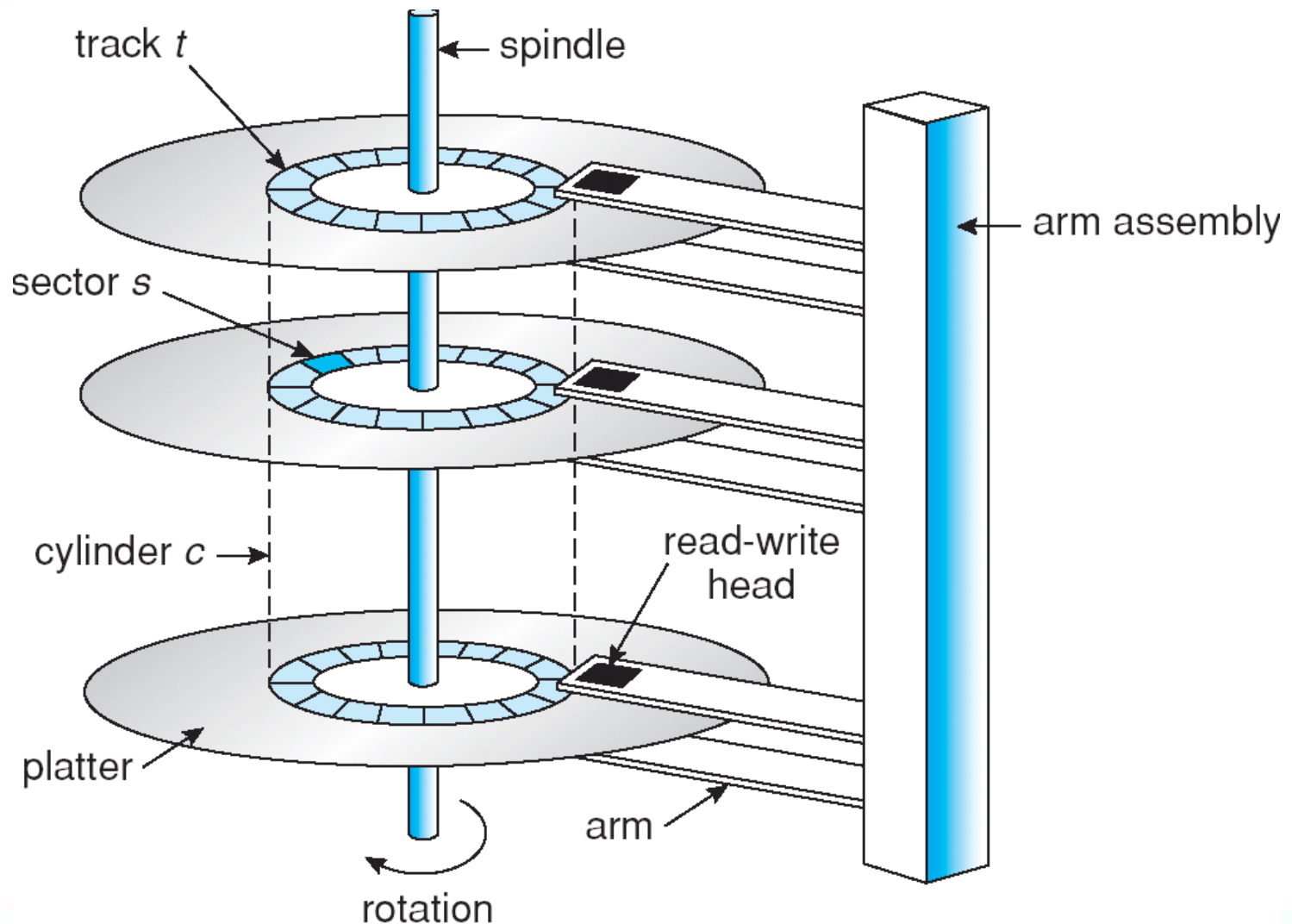
# Overview of Mass Storage Structure

- 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
- 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

# Moving-head Disk Mechanism

# 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

# 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.
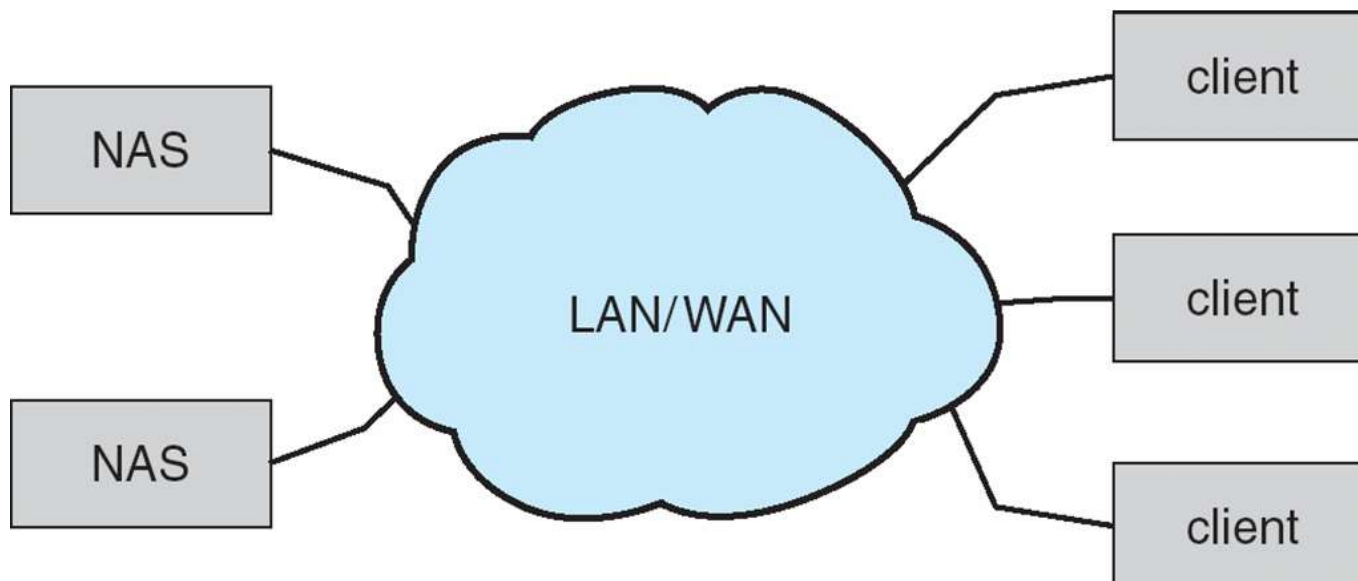
# Disk Attachment

- 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** (**SAN**s) in which many hosts attach to many storage units

  - Can be **arbitrated loop** (**FC-AL**) of 126 devices

# 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

# 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 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

Illustration shows total head movement of 640 cylinders.

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

# 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.

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

head starts at 53

- 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.
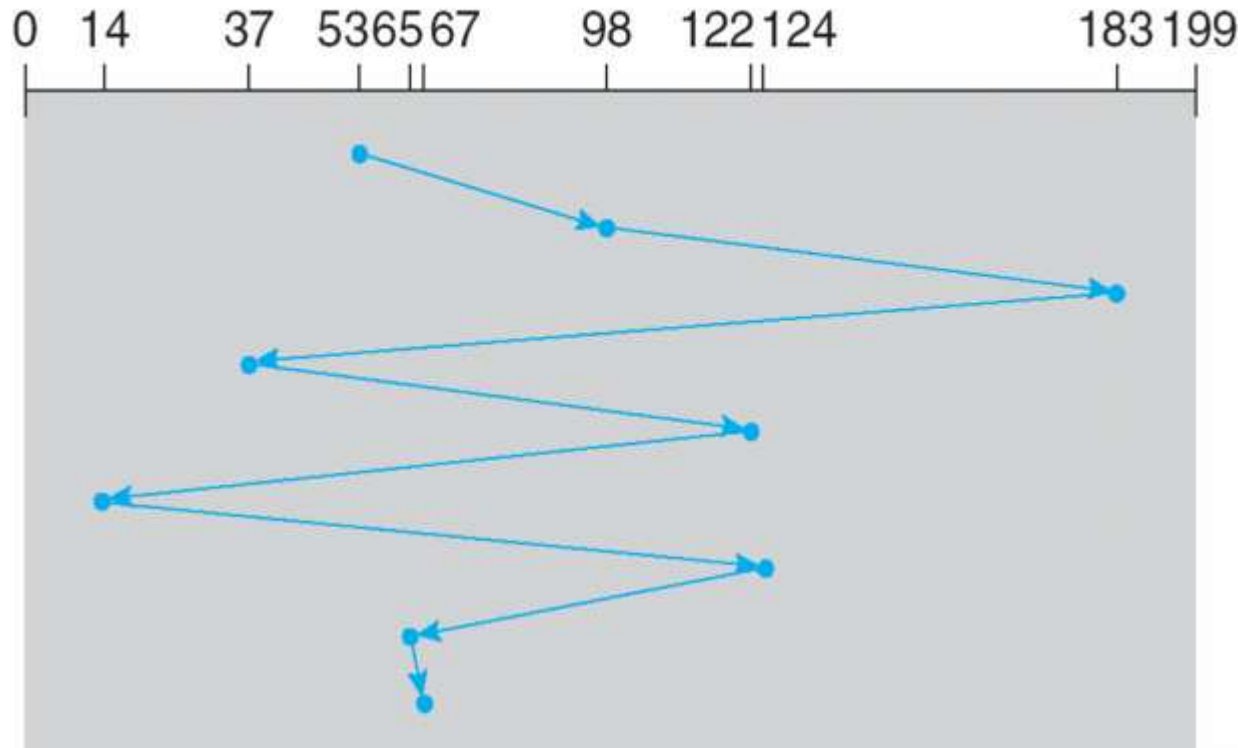
- Sometimes called the *elevator algorithm*.

- Illustration shows total head movement of 208 cylinders.

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

head starts at 53

# C-SCAN

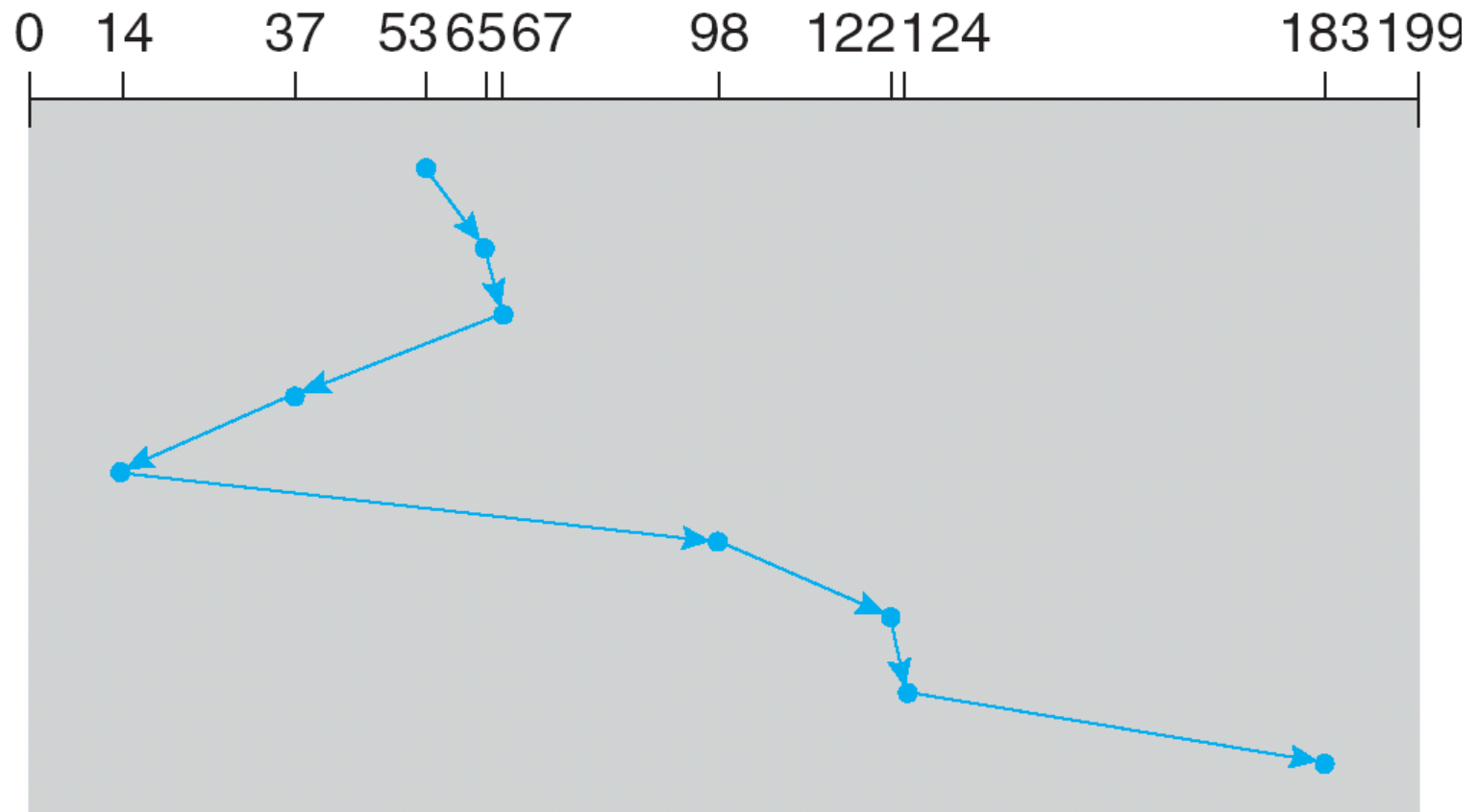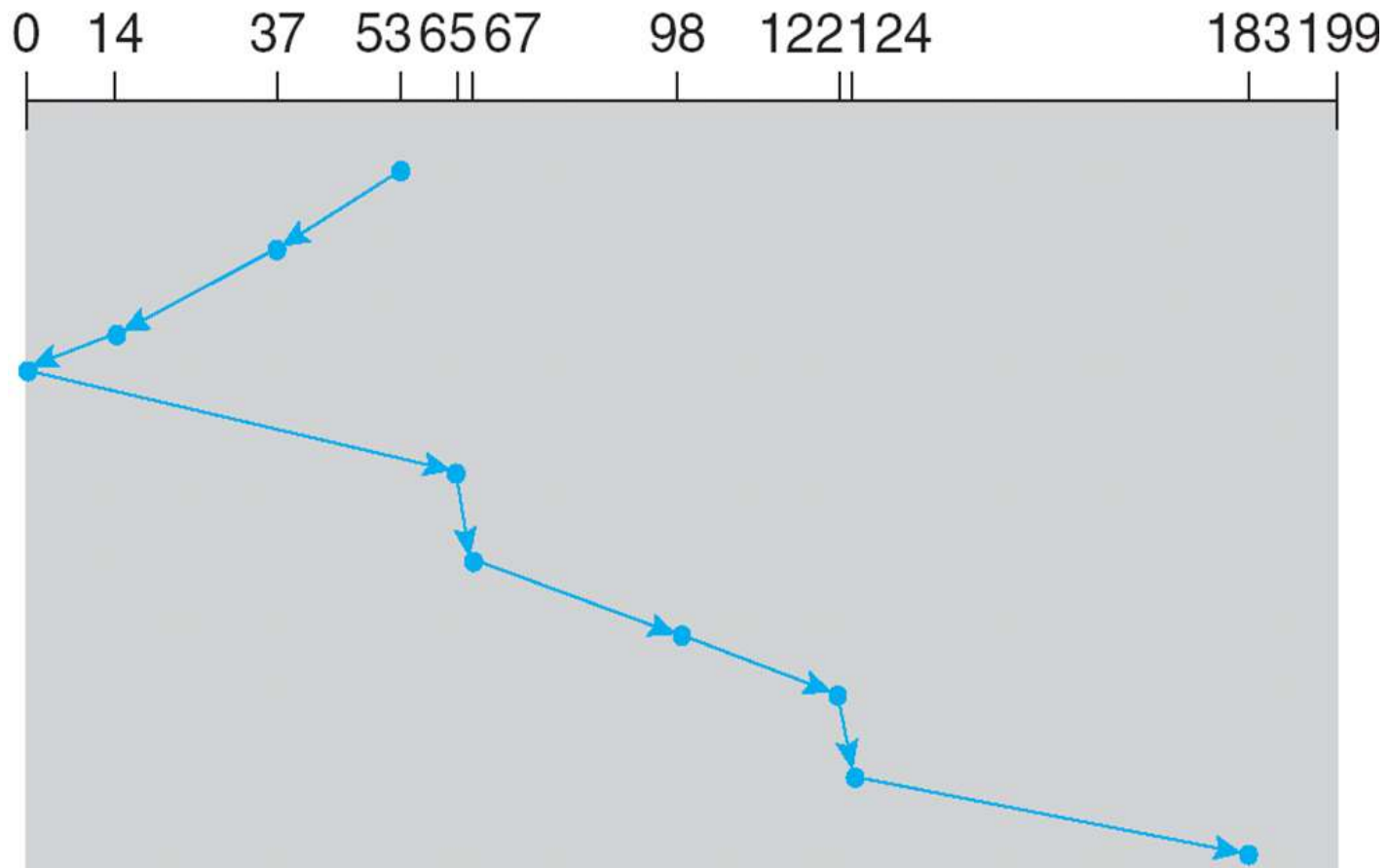- Provides a more uniform wait time than SCAN.

- 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.

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.

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".

- 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

# Data Structures for Swapping on Linux Systems

**End**

# Protection

# Protection

- Goals of Protection
- Principles of Protection
- Domain of Protection
- Access Matrix

# Objectives

- Discuss the goals and principles of protection in a modern computer system

- Explain how protection domains combined with an access matrix are used to specify the resources a process may access

- Examine capability and language-based protection systems

# Goals of Protection

- In one protection model, computer consists of a collection of objects, hardware or software

- Each object has a unique name and can be accessed through a well-defined set of operations

- Protection problem - ensure that each object is accessed correctly and only by those processes that are allowed to do so

# Principles of Protection

*Go, change the world*

- Guiding principle – **principle of least privilege**
  - Programs, users and systems should be given just enough **privileges** to perform their tasks
  - Limits damage if entity has a bug, gets abused
  - Can be static (during life of system, during life of process)
  - Or dynamic (changed by process as needed) – **domain switching**, **privilege escalation**
  - "Need to know" a similar concept regarding access to data

- Must consider "grain" aspect
  - Rough-grained  privilege management easier, simpler, but least privilege now done in large chunks
    - For example, traditional Unix processes either have abilities of the associated user, or of root
  - Fine-grained management more complex, more overhead, but more protective
    - File ACL lists, RBAC
- Domain can be user, process, procedure

# Domain Structure

- Access-right = *<object-name, rights-set>*
  where *rights-set* is a subset of all valid operations that can
  be performed on the object

- Domain = set of access-rights

$D_1$

< $O_3$, {read, write} >
< $O_1$, {read, write} >
< $O_2$, {execute} >

$D_2$    $D_3$

< $O_2$, {write} >    < $O_4$, {print} >    < $O_1$, {execute} >
< $O_3$, {read} >

# Domain Implementation (UNIX)

- Domain = user-id

- Domain switch accomplished via file system
  - ▸ Each file has associated with it a domain bit (setuid bit)
  - ▸ When file is executed and setuid = on, then user-id is set to owner of the file being executed
  - ▸ When execution completes user-id is reset

- Domain switch accomplished via passwords
  - ● su command temporarily switches to another user's domain when other domain's password provided

- Domain switching via commands
  - ● sudo command prefix executes specified command in another domain (if original domain has privilege or password given)

# Domain Implementation (MULTICS)

- Let $D_i$ and $D_j$ be any two domain rings
- If $j < I \Rightarrow D_i \subseteq D_j$

# Multics Benefits and Limits

- Ring / hierarchical structure provided more than the basic kernel / user or root / normal user design

- Fairly complex -> more overhead

- But does not allow strict need-to-know
  - Object accessible in $D_j$ but not in $D_i$, then *j* must be < *i*
  - But then every segment accessible in $D_i$ also accessible in $D_j$

# Access Matrix

- View protection as a matrix (**access matrix**)
- Rows represent domains
- Columns represent objects
- `Access(i, j)` is the set of operations that a process executing in Domain$_i$ can invoke on Object$_j$

| object / domain | $F_1$ | $F_2$ | $F_3$ | printer |
|---|---|---|---|---|
| $D_1$ | read | | read | |
| $D_2$ | | | | print |
| $D_3$ | | read | execute | |
| $D_4$ | read write | | read write | |

- If a process in Domain $D_i$ tries to do "op" on object $O_j$, then "op" must be in the access matrix

- User who creates object can define access column for that object

- Can be expanded to dynamic protection
  - Operations to add, delete access rights
  - Special access rights:
    - *owner of $O_i$*
    - *copy op from $O_i$ to $O_j$ (denoted by "*")*
    - *control – $D_i$ can modify $D_j$ access rights*
    - *transfer – switch from domain $D_i$ to $D_j$*
  - *Copy* and *Owner* applicable to an object
  - *Control* applicable to domain object

- **Access matrix** design separates mechanism from policy
  - Mechanism
    - ‣ Operating system provides access-matrix + rules
    - ‣ If ensures that the matrix is only manipulated by authorized agents and that rules are strictly enforced
  - Policy
    - ‣ User dictates policy
    - ‣ Who can access what object and in what mode
- But doesn't solve the general confinement problem

# Access Matrix of Figure A with Domains as Objects

| object<br>domain | $F_1$ | $F_2$ | $F_3$ | laser<br>printer | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|---|---|---|---|---|---|---|---|---|
| $D_1$ | read | | read | | | switch | | |
| $D_2$ | | | | print | | | switch | switch |
| $D_3$ | | read | execute | | | | | |
| $D_4$ | read<br>write | | read<br>write | | switch | | | |

# Access Matrix with *Copy* Rights

| object<br>domain | $F_1$ | $F_2$ | $F_3$ |
|---|---|---|---|
| $D_1$ | execute | | write* |
| $D_2$ | execute | read* | execute |
| $D_3$ | execute | | |

(a)

| object<br>domain | $F_1$ | $F_2$ | $F_3$ |
|---|---|---|---|
| $D_1$ | execute | | write* |
| $D_2$ | execute | read* | execute |
| $D_3$ | execute | read | |

(b)

# Access Matrix With *Owner* Rights

| domain \ object | $F_1$ | $F_2$ | $F_3$ |
|---|---|---|---|
| $D_1$ | owner execute | | write |
| $D_2$ | | read* owner | read* owner write |
| $D_3$ | execute | | |

(a)

| domain \ object | $F_1$ | $F_2$ | $F_3$ |
|---|---|---|---|
| $D_1$ | owner execute | | write |
| $D_2$ | | owner read* write* | read* owner write |
| $D_3$ | | write | write |

(b)

# Modified Access Matrix of Figure B

| object<br>domain | $F_1$ | $F_2$ | $F_3$ | laser<br>printer | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|---|---|---|---|---|---|---|---|---|
| $D_1$ | read | | read | | | switch | | |
| $D_2$ | | | | print | | | switch | switch<br>control |
| $D_3$ | | read | execute | | | | | |
| $D_4$ | write | | write | | switch | | | |

**End**