# Operating Systems

## File-System Interface & Implementation

# File-System Interface

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

# File System Implementation

- File-System Structure

- File-System Implementation

- Directory Implementation

- Allocation Methods

- Free-Space Management

- Efficiency and Performance

# 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 Locking Example – Java API

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

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

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

# File Types – Name, Extension

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

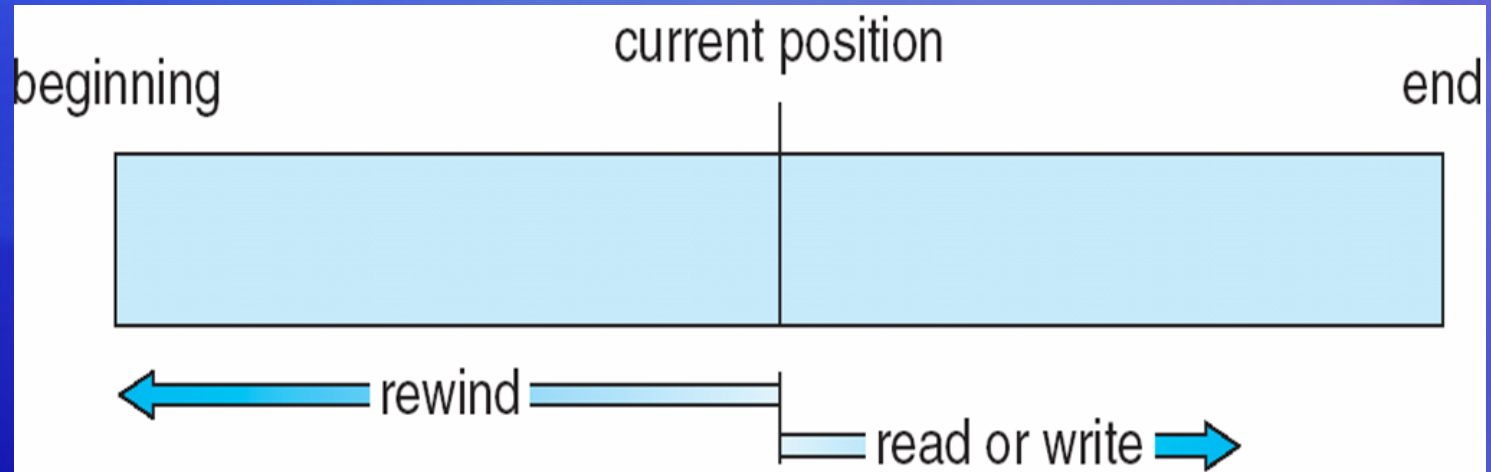# Access Methods

**Sequential Access**

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

**Direct Access**

read *n*
write *n*
position to *n*
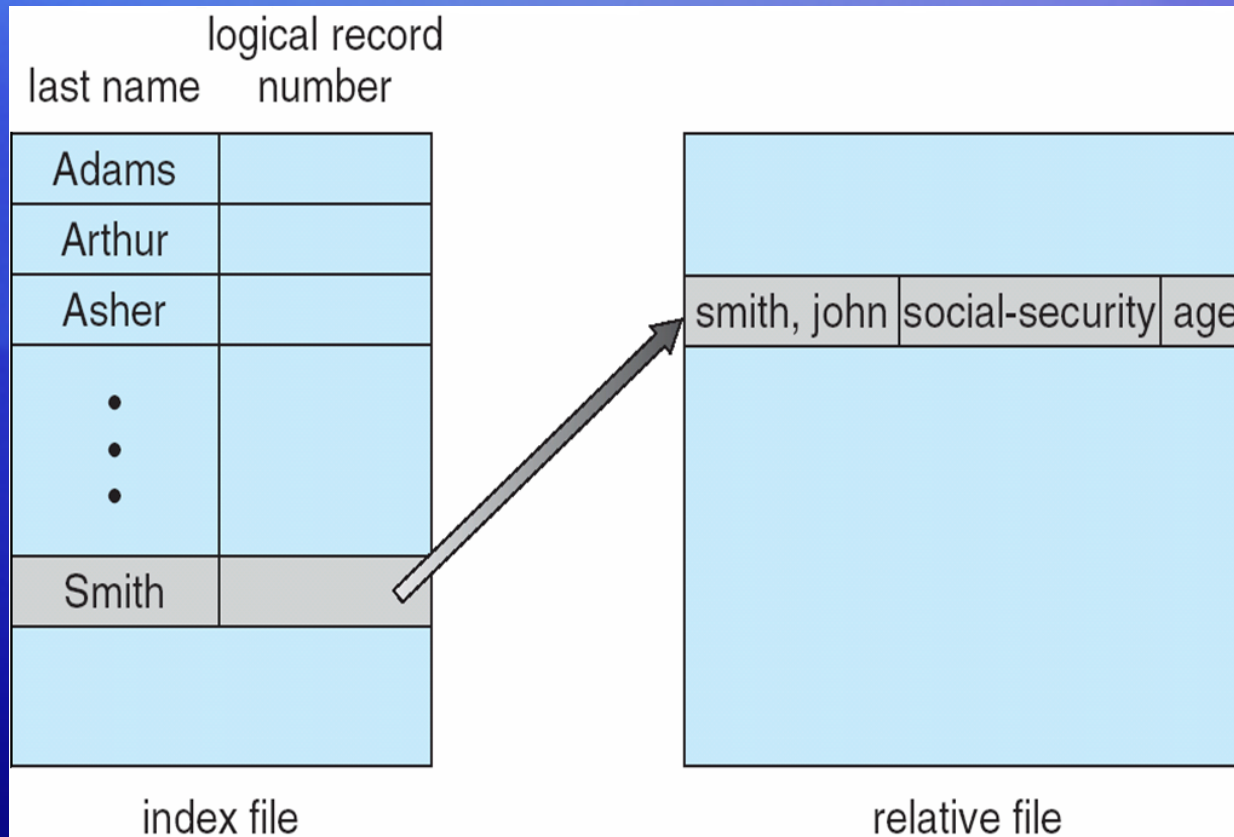read next
write next
rewrite *n*

*n* = relative block number

# Sequential-access File
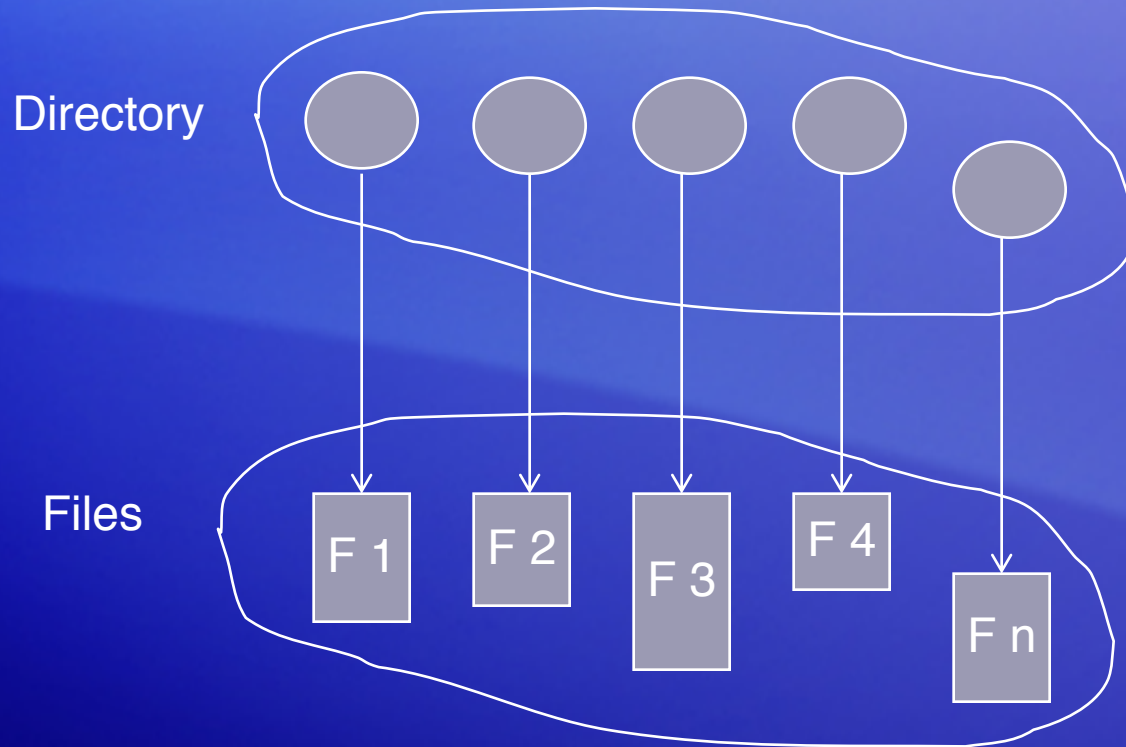
# Simulation of Sequential Access on Direct-access File

| 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

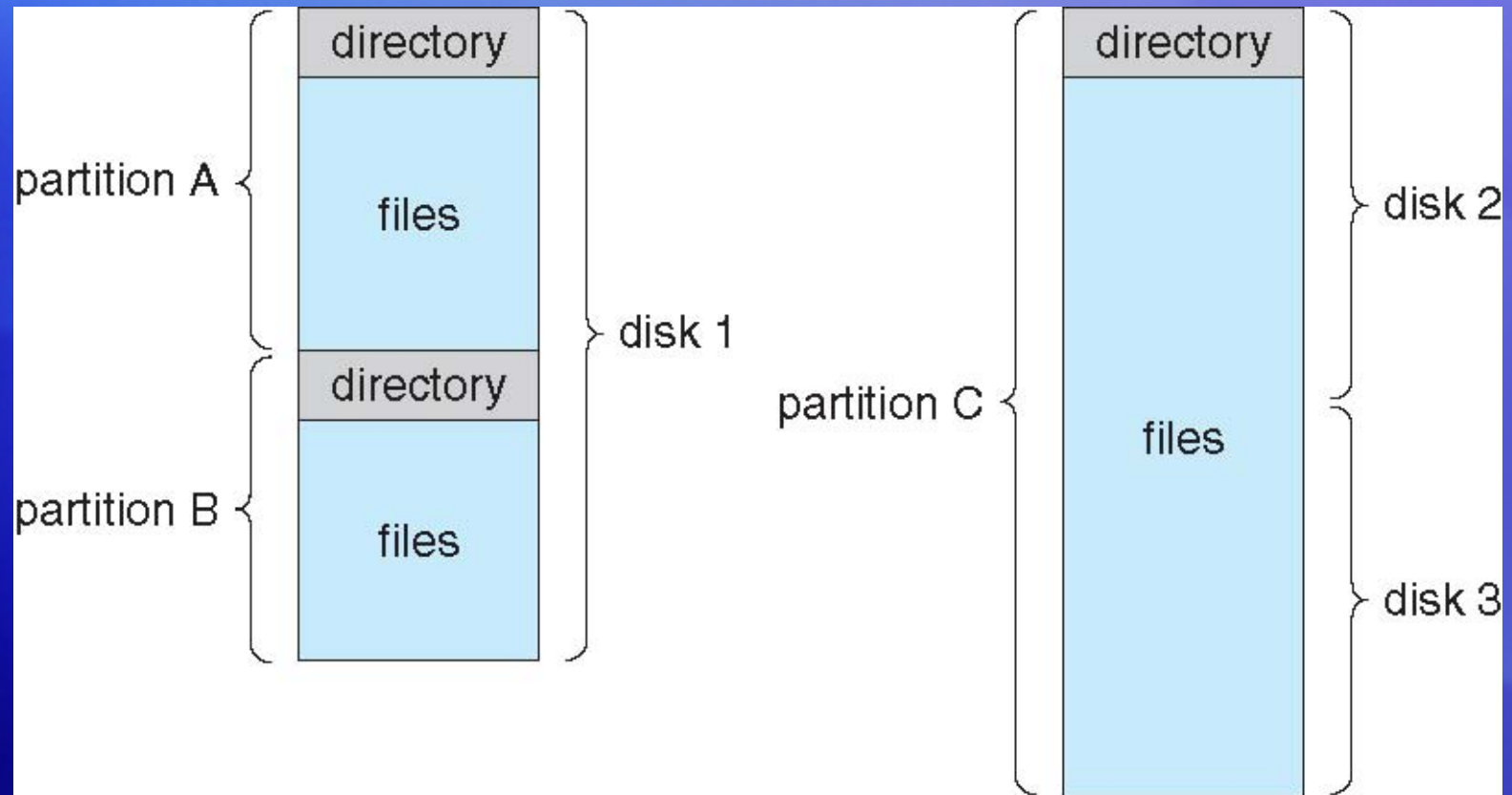A collection of nodes containing information about all files



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

- 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

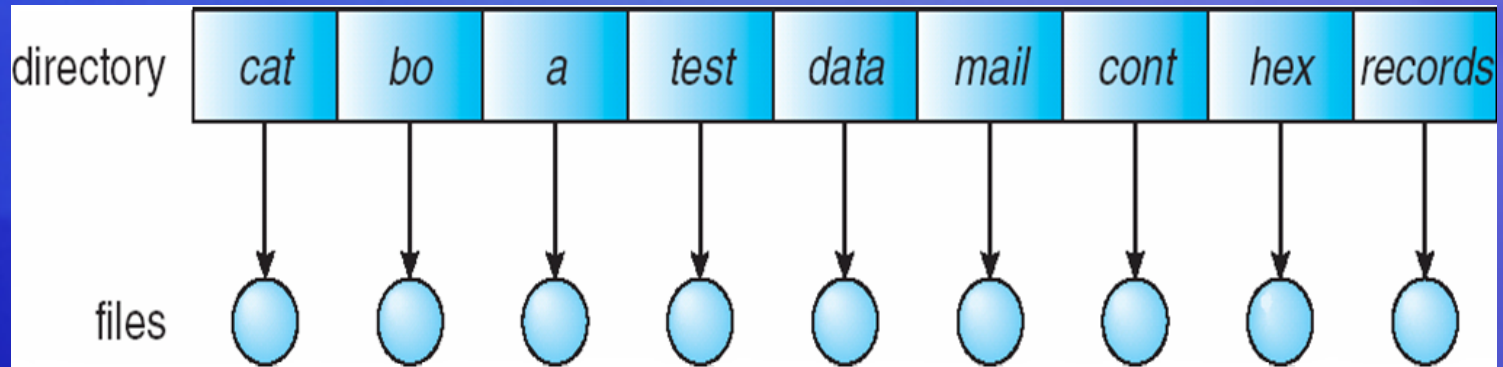# Operations Performed on Directory

- 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, …)

# Single-Level Directory
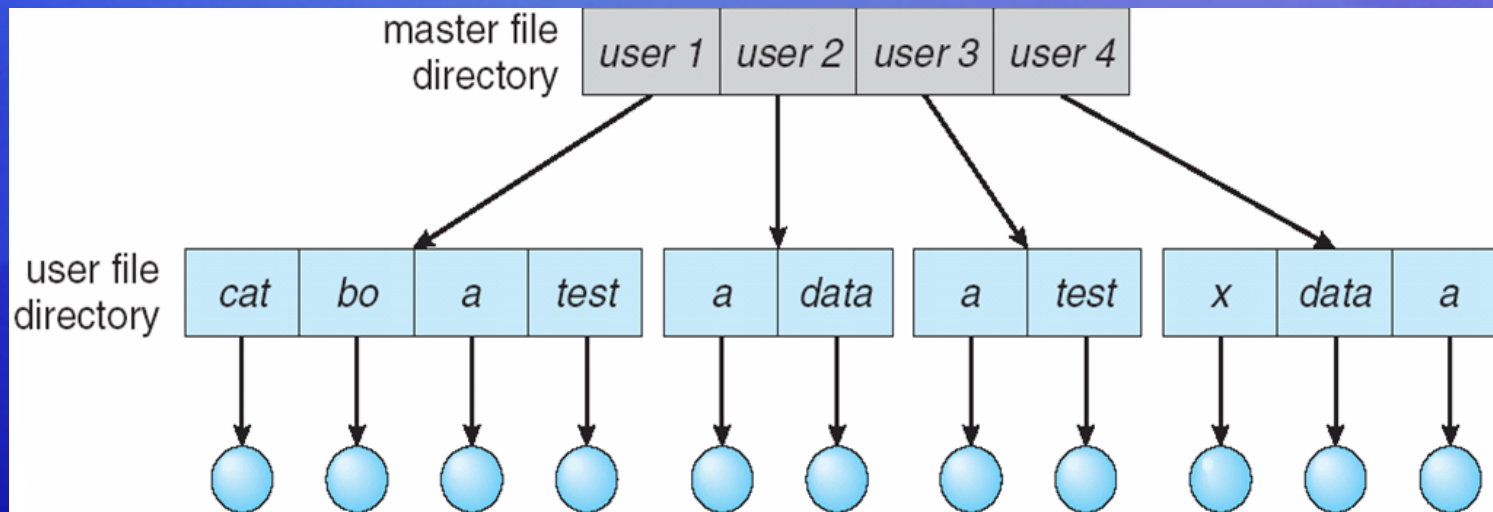
A single directory for all users

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

files
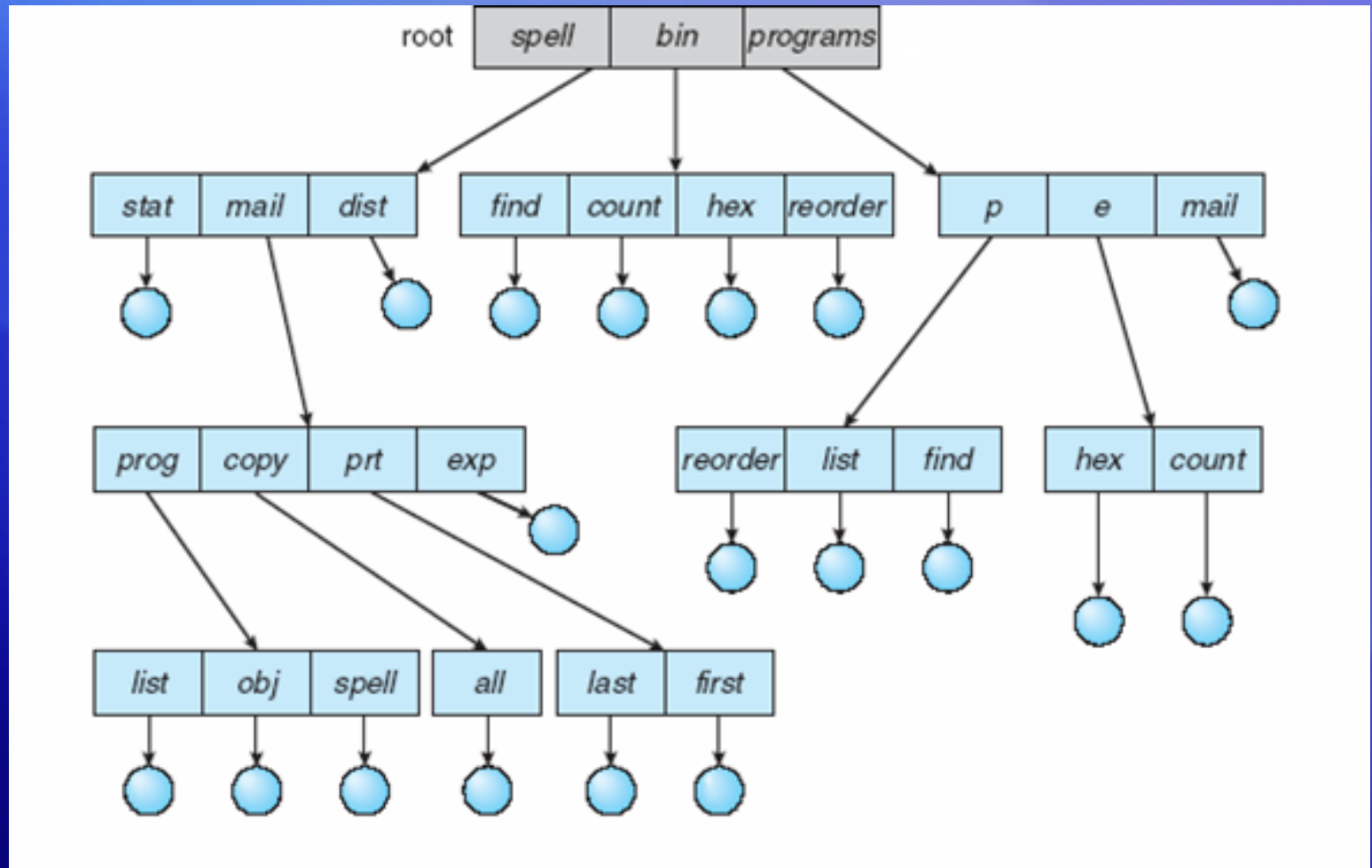
Naming problem

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

# 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
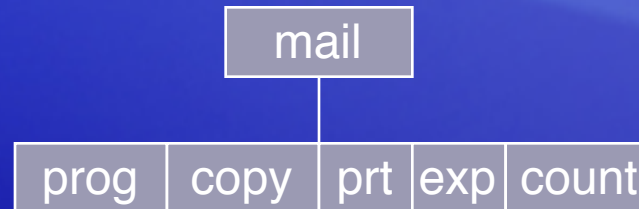- Creating a new file is done in current directory
- Delete a file

  **rm <file-name>**

- Creating a new subdirectory is done in current directory

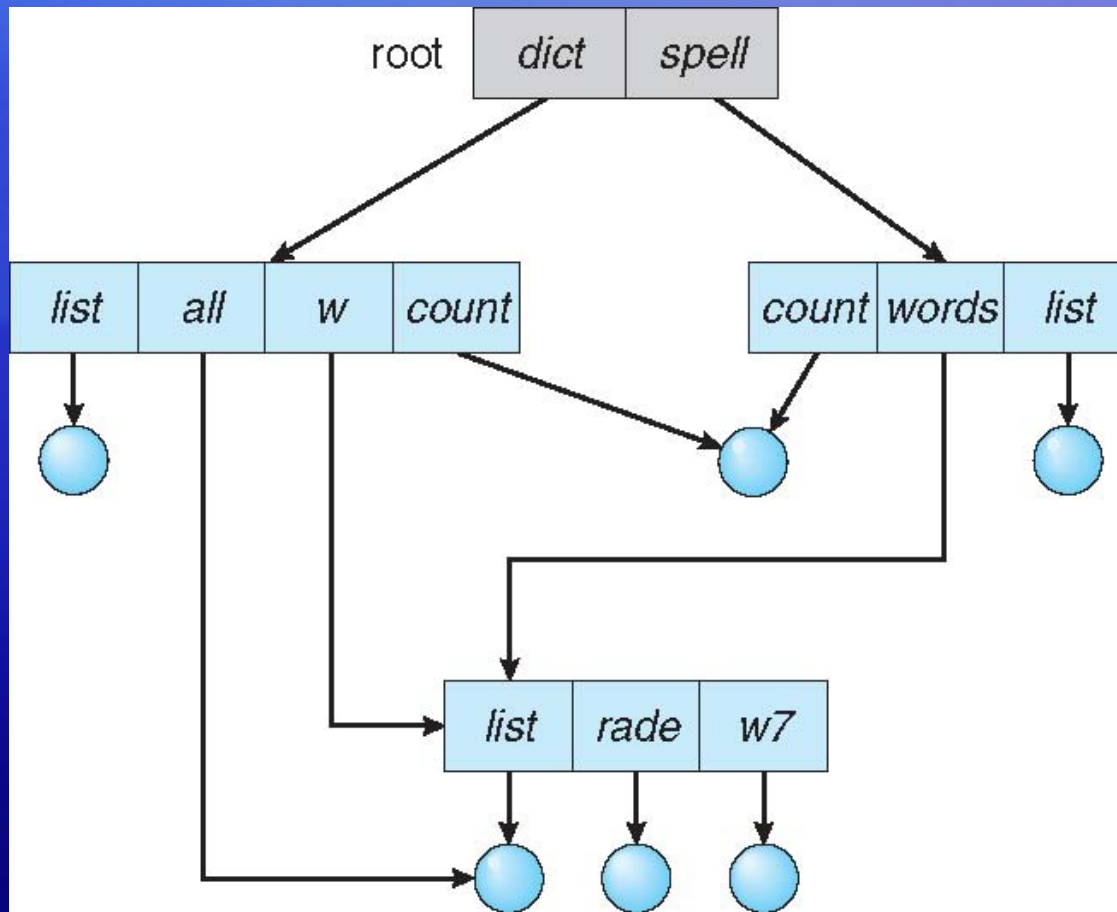  **mkdir <dir-name>**

- Example:  if in current directory  **/mail**

  **mkdir count**

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

Deleting "mail" $\Rightarrow$ deleting the entire subtree rooted by "mail"
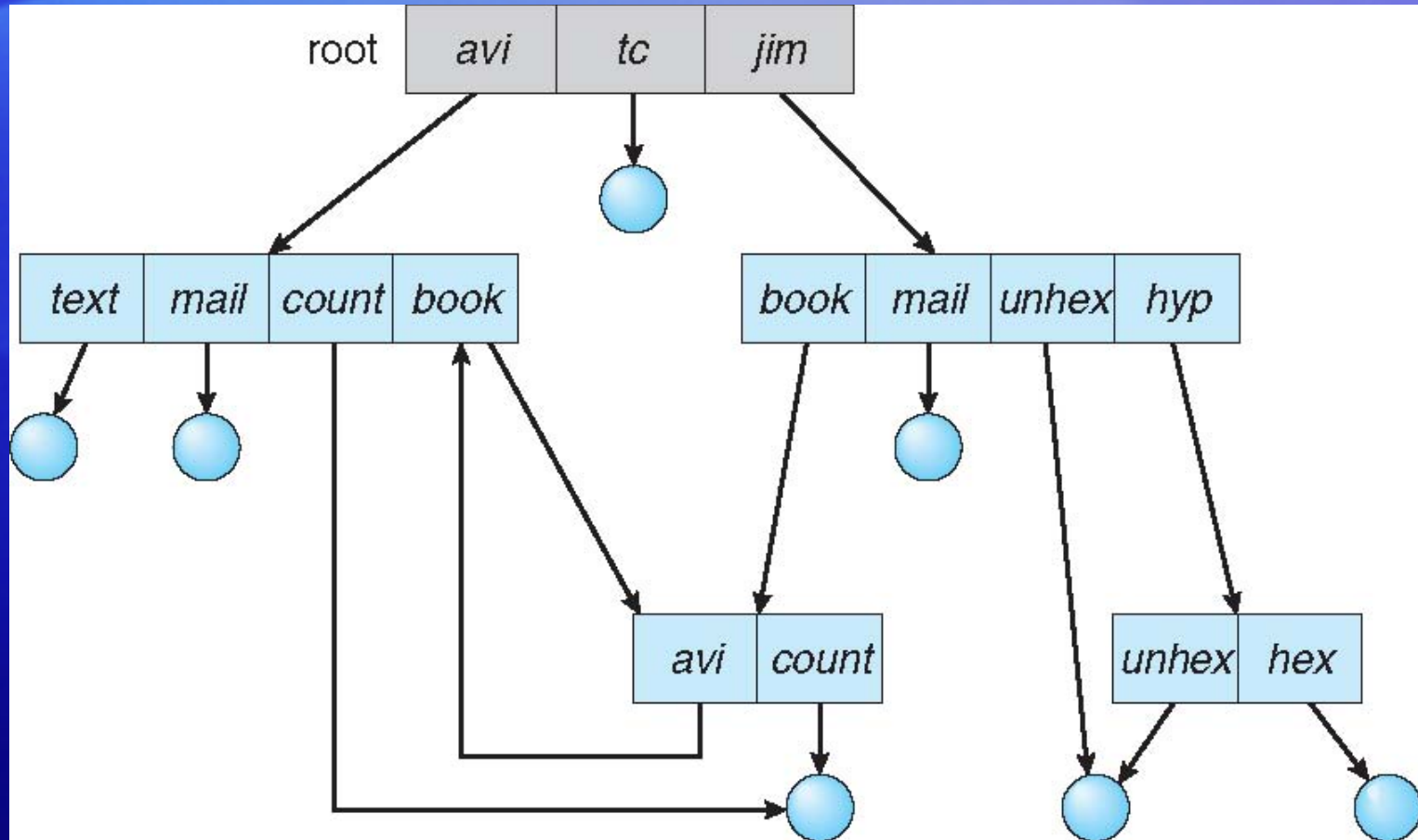
# Acyclic-Graph Directories

Have shared subdirectories and files

# Acyclic-Graph Directories (Cont.)

- Two different names (aliasing)

- If *dict* deletes *list* $\Rightarrow$ dangling pointer

- Solutions:
  - Backpointers, so we can delete all pointers
    Variable size records a problem
  - Backpointers using a daisy chain organization
  - Entry-hold-count solution

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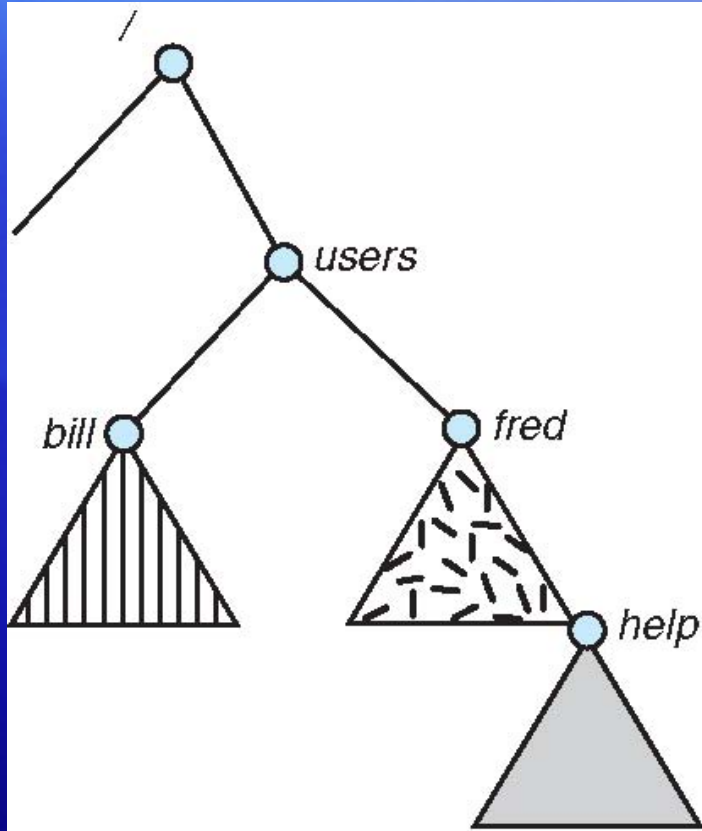
# General Graph Directory

# General Graph Directory (Cont.)

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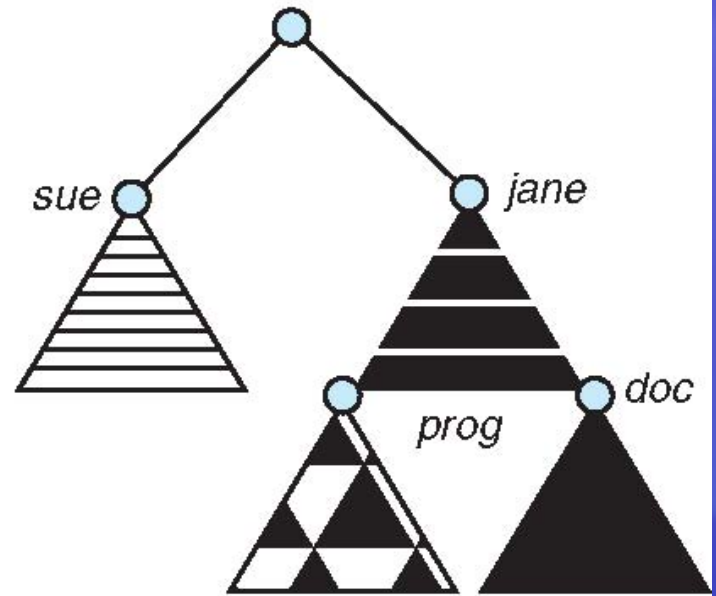
# File System Mounting

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

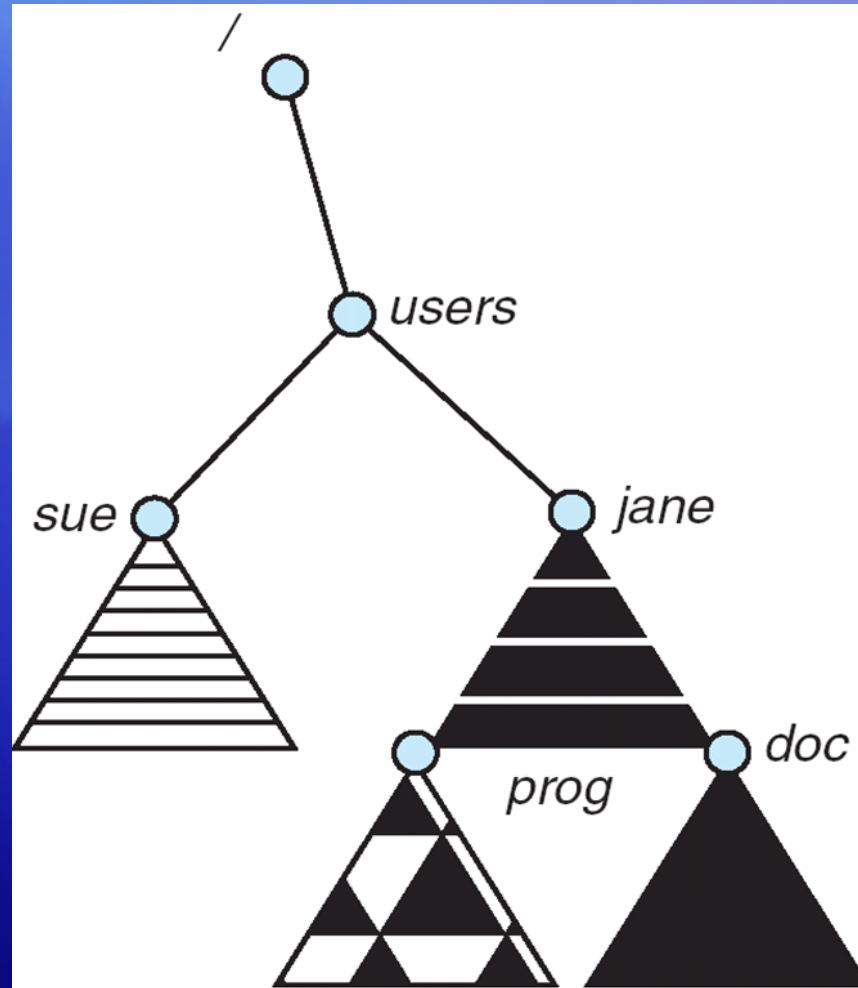- A unmounted file system is mounted at a **mount point**

(a)

(b)

# Mount Point

# 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

# Protection

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

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

# Access Lists and Groups

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

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

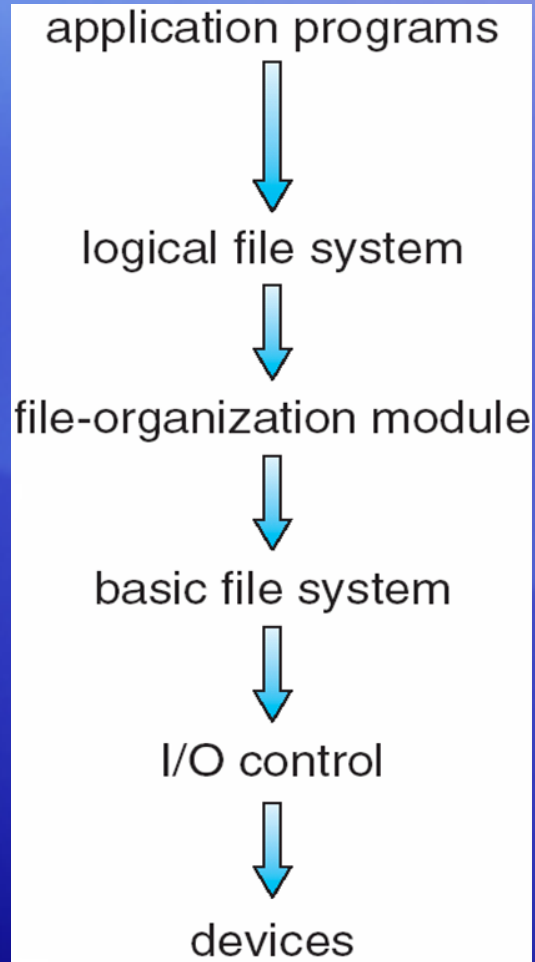Attach a group to a file

chgrp    G    game

# File-System Structure

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

# Layered File System

application programs

↓

logical file system

↓

file-organization module

↓

basic file system

↓

I/O control

↓

devices

# File System Layers

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

# File System Layers (Cont.)

- **Logical file system** manages metadata information
  - Translates file name into file number, file handle, location by maintaining file control blocks (**inodes** in Unix)
  - Directory management
  - Protection
- Layering useful for reducing complexity and redundancy, but adds overhead and can decrease performance
  - 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
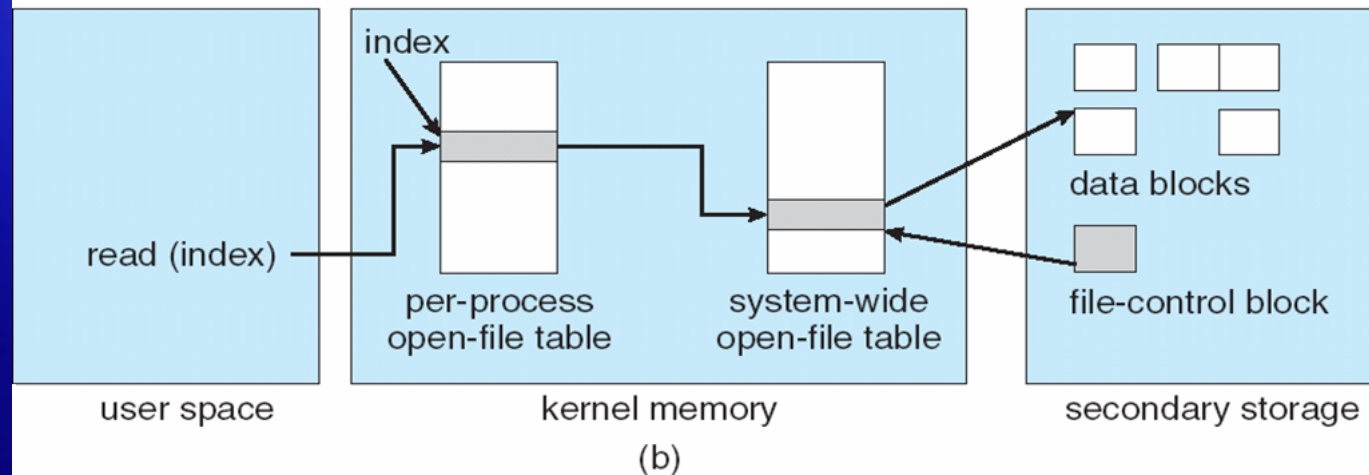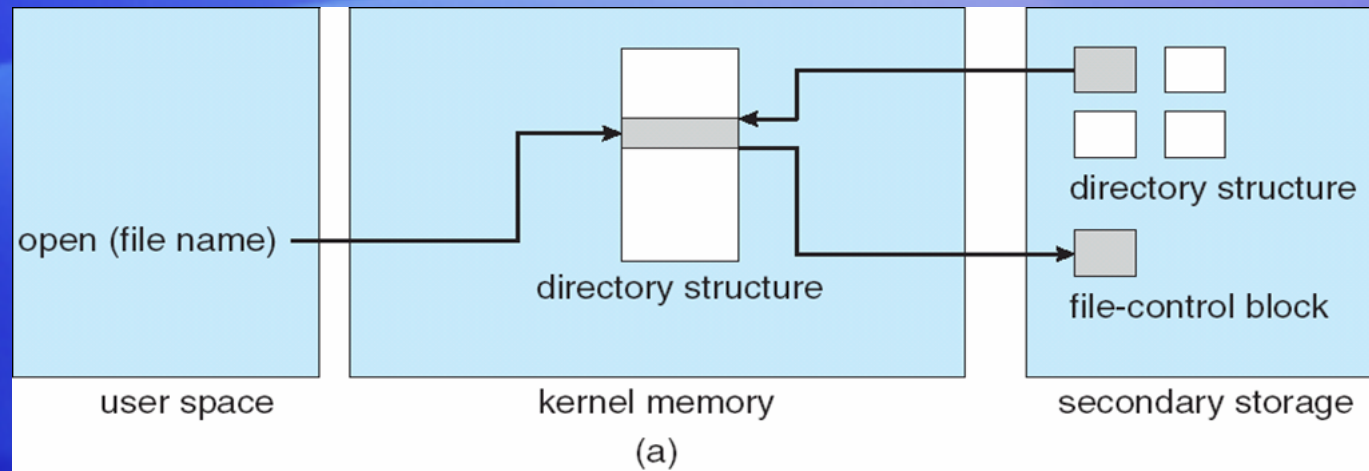
# A Typical File Control Block

| |
|---|
| 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

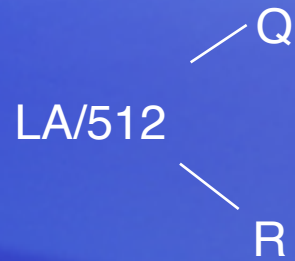# Directory Implementation

- **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**
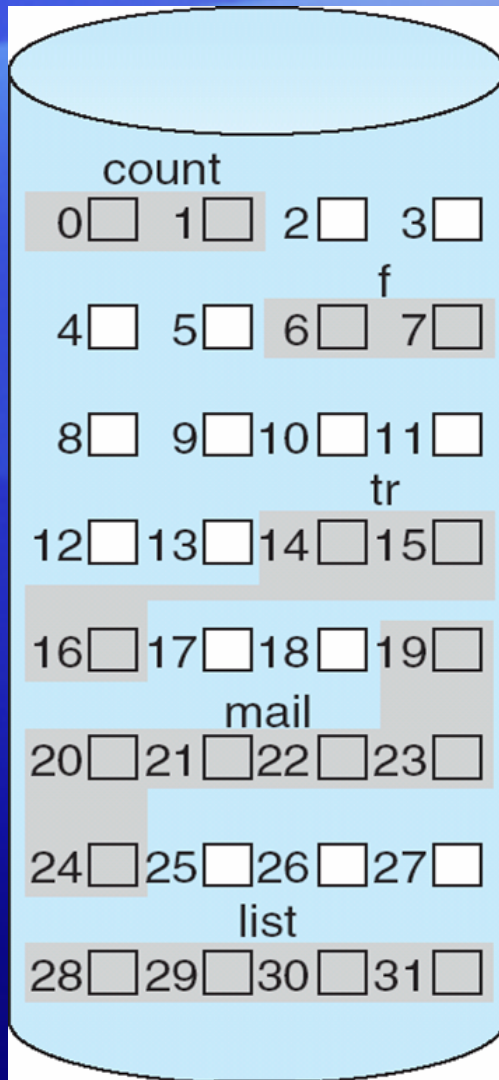
# Contiguous Allocation

Mapping from logical to physical

```
                    ╱ Q
            LA/512
                    ╲
                      R
```

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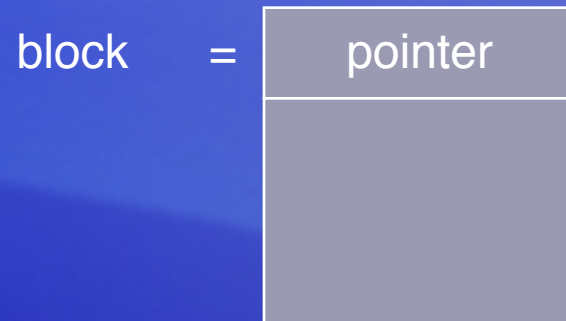Displacement into block = R

# Contiguous Allocation of Disk Space

# 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 |
             |---------|
             |         |

Mapping

$$LA/511 \begin{cases} Q \\ R \end{cases}$$
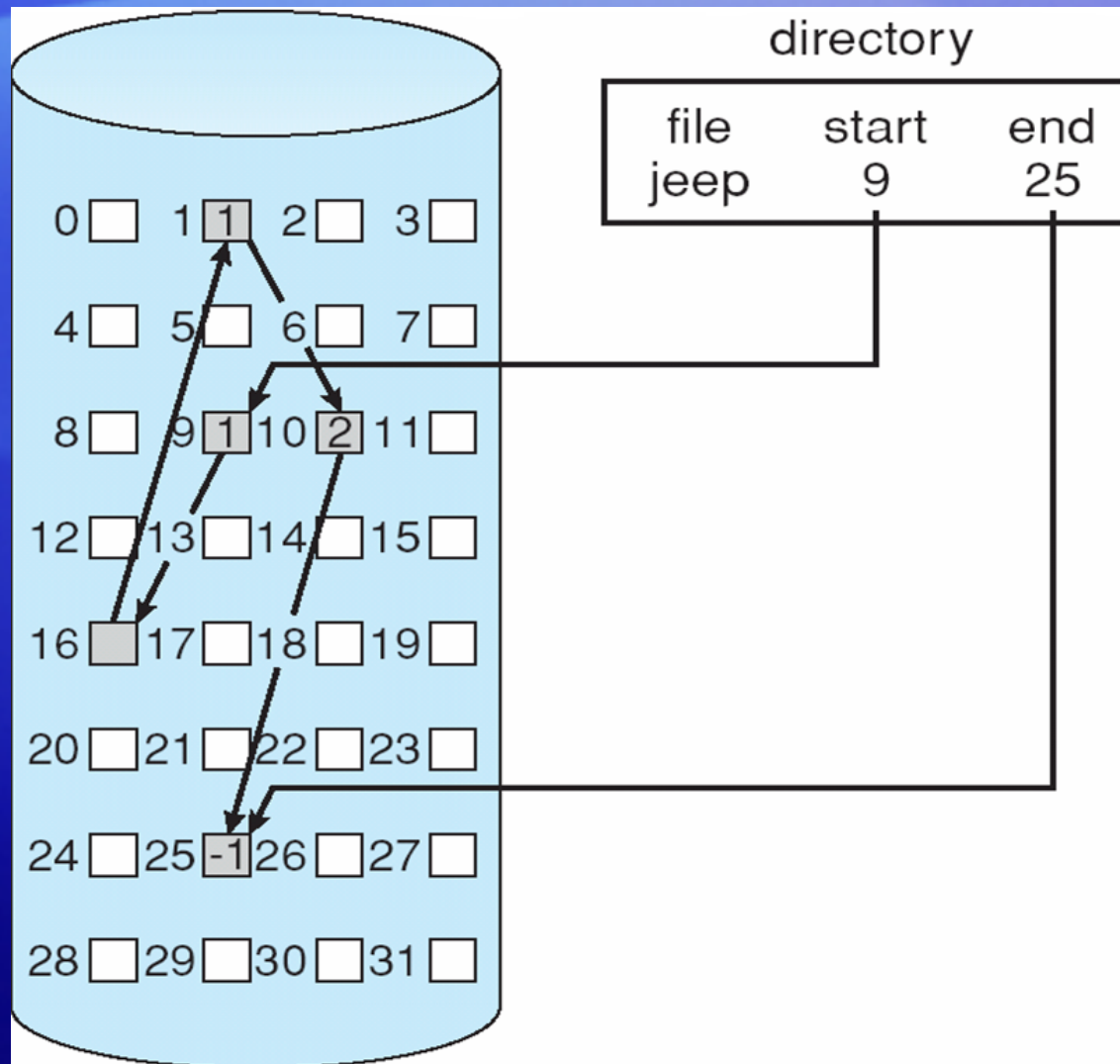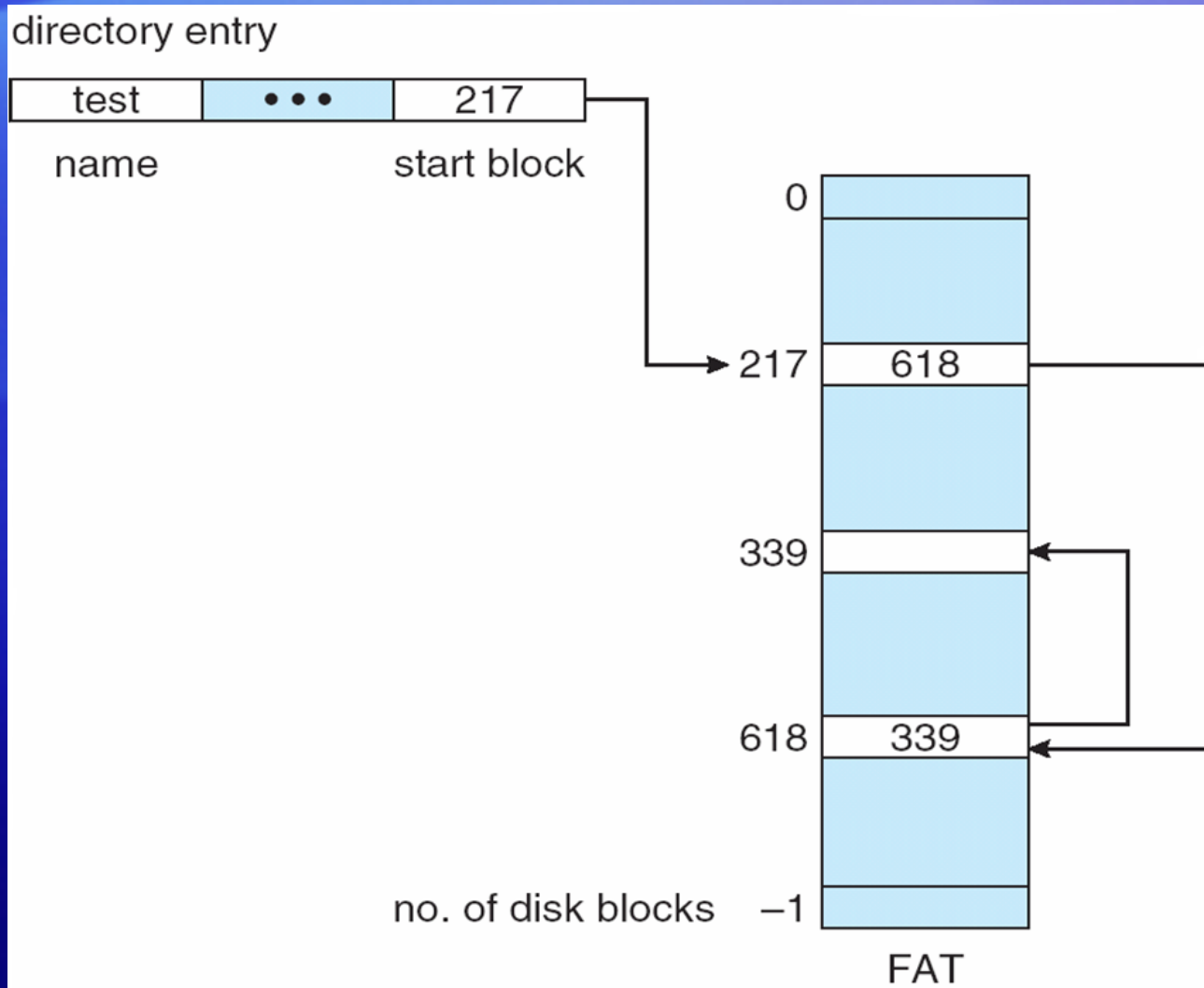
Block to be accessed is the Qth block in the linked chain of blocks representing the file.
Displacement into block = R + 1
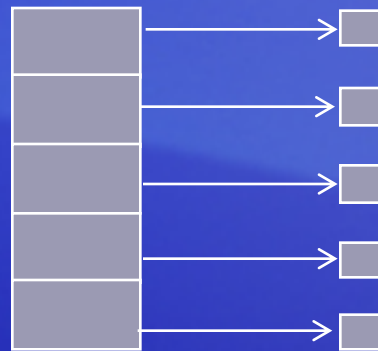
# Linked Allocation

# File-Allocation Table

# Allocation Methods - Indexed
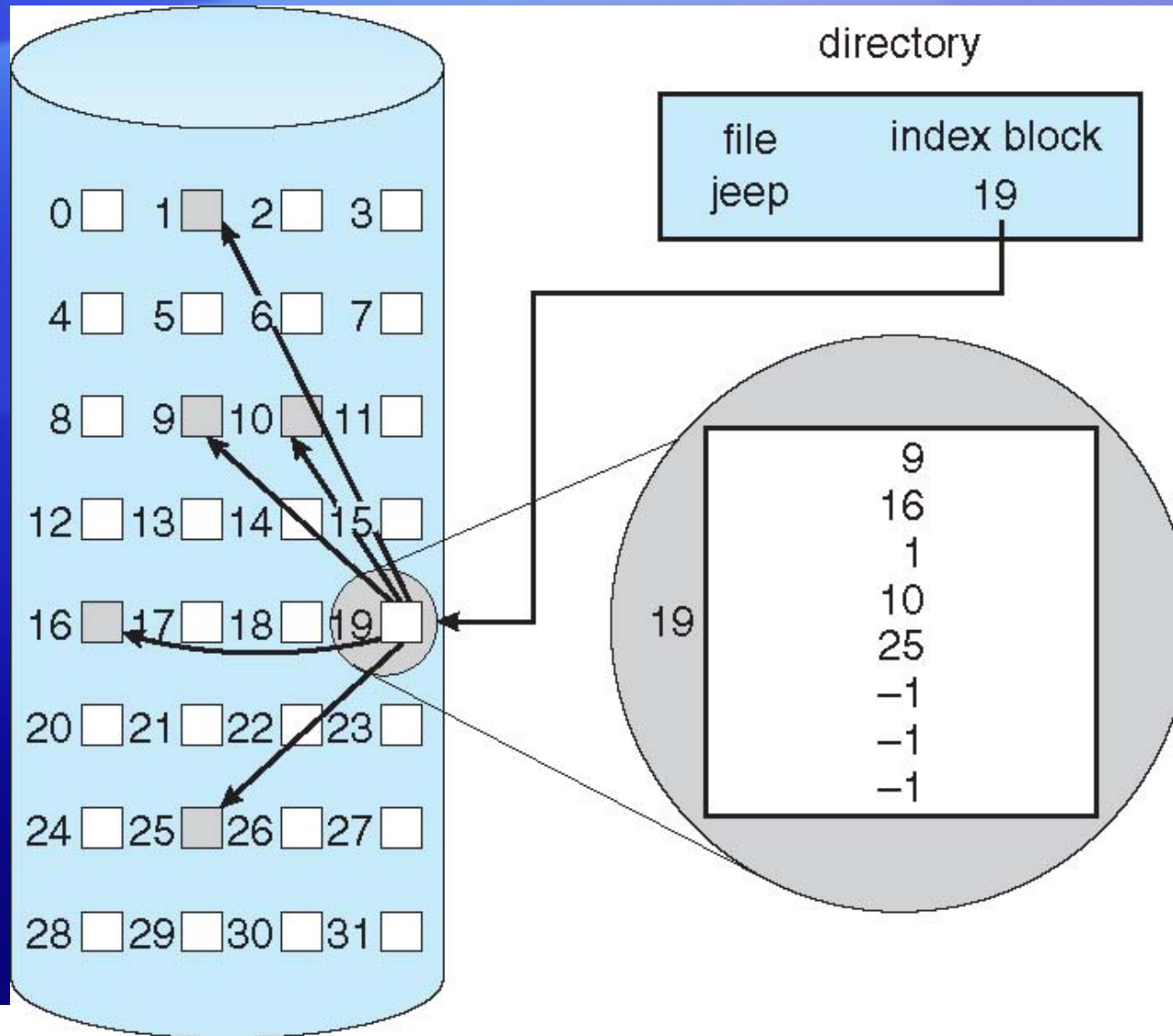
**Indexed allocation**

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

Logical view



index table
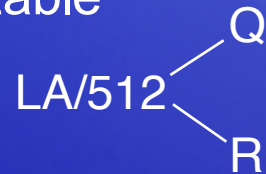
# Example of Indexed Allocation

# 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 $\diagdown$ Q

R

      Q = displacement into index table
      R = displacement into block

# Indexed Allocation – Mapping (Cont.)

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:

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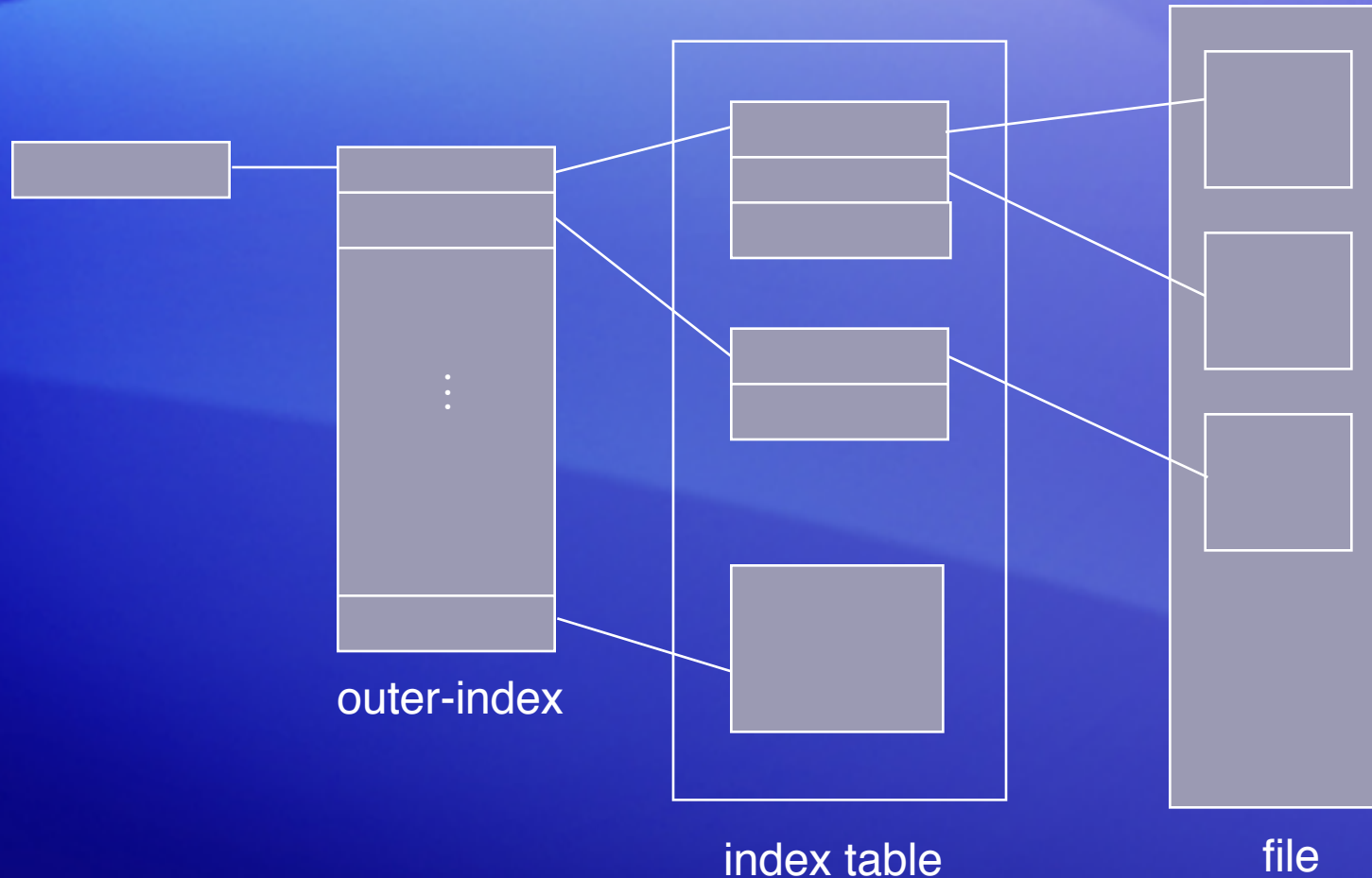
$$LA / (512 \times 512) \begin{cases} Q_1 \\ R_1 \end{cases}$$

$Q_1$ = 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:

outer-index

index table

file

# 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)



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

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
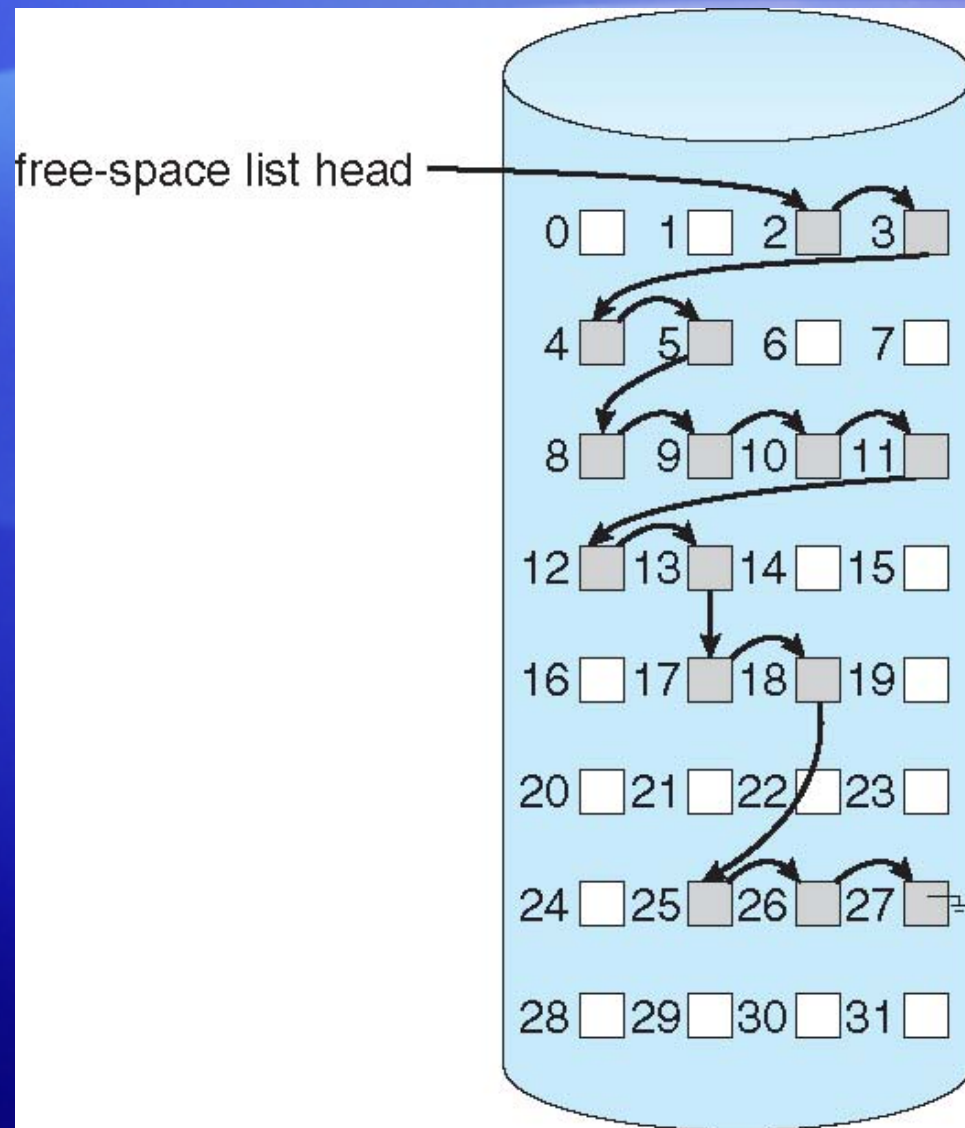
Linked list (free list)

      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 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

# Efficiency and Performance

- Efficiency dependent on:
    - Disk allocation and directory algorithms
    - Types of data kept in file's directory entry
    - Pre-allocation or as-needed allocation of metadata structures
    - Fixed-size or varying-size data structures

# Efficiency and Performance (Cont.)

- Performance
  - Keeping data and metadata close together
  - **Buffer cache** – separate section of main memory for frequently used blocks
  - **Synchronous** writes sometimes requested by apps or needed by OS
    - No buffering / caching – writes must hit disk before acknowledgement
    - **Asynchronous** writes more common, buffer-able, faster
  - **Free-behind** and **read-ahead** – techniques to optimize sequential access
  - Reads frequently slower than writes

# Reference Book

"Operating System Concepts" by Silberchartz, Galvin, Gagne, Wiley India Publications.