# File Systems
# Interface and Implementation

# Chapter 11: File-System Interface

- File and Directory
- File System Structure
- Access Methods
- File-System Mounting
- File Sharing
- Protection
- File-System Implementation
- Space Allocation Methods
- Recovery

# Objectives

- To explain the function of file systems

- To describe the interfaces to file systems

- To discuss file-system protection

- 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 Concept

- Contiguous logical address space

- Types:

  - Data

    - numeric

    - character

    - binary

  - Program

- Contents defined by file's creator

  - Many types

    - Consider **text file, source file, executable file**
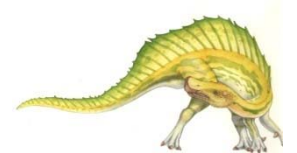
# File Attributes

- **Name** – only information kept in human-readable form

- **Identifier** – unique tag (number) identifies file within file system

- **Type** – needed for systems that support different types

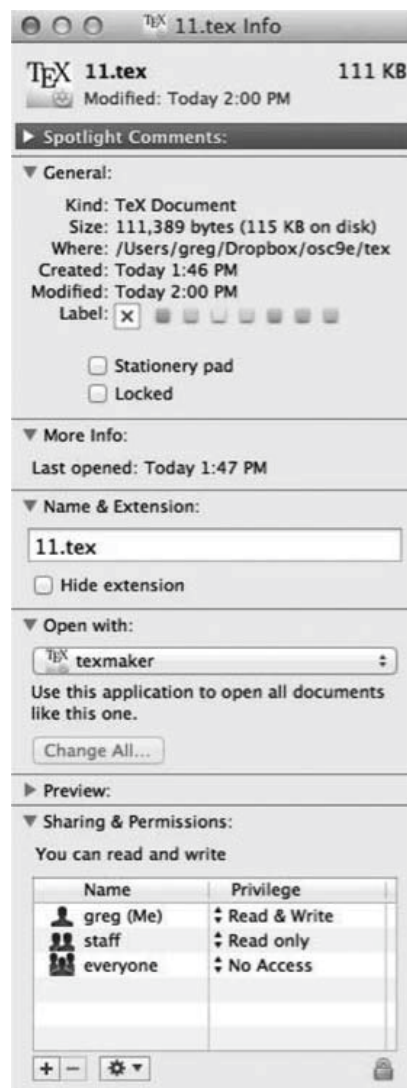- **Location** – pointer to file location on device

# File Attributes

- **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

- Many variations, including extended file attributes such as file checksum

- Information kept in the directory structure

# File info Window on Mac OS X

# File Operations

- File is an **abstract data type**

- **Create**

- **Write** – at **write pointer** location

- **Read** – at **read pointer** location

- **Reposition within file** - **seek**

- **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:

  - **Open-file table**: tracks open files

  - File pointer:  pointer to last read/write location, per process that has the file open

  - **File-open count**: counter of number of times a file is open – to allow removal of data from open-file table when last processes closes it

  - Disk location of the file: cache of data access information

  - Access rights: per-process access mode information

# Open File Locking

- Provided by some operating systems and file systems

    - Similar to reader-writer locks

    - **Shared lock** similar to reader lock – several processes can acquire concurrently

    - **Exclusive lock** similar to writer lock

- Mediates access to a file

- Mandatory or advisory:

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

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

# File Types – Name, Extension

- Extensions are not mandatory in some Oses

- Extension usage and association and automatically open the software

- UNIX uses magic number to indicate type of file

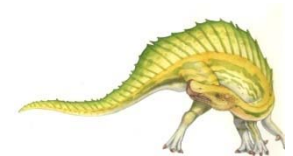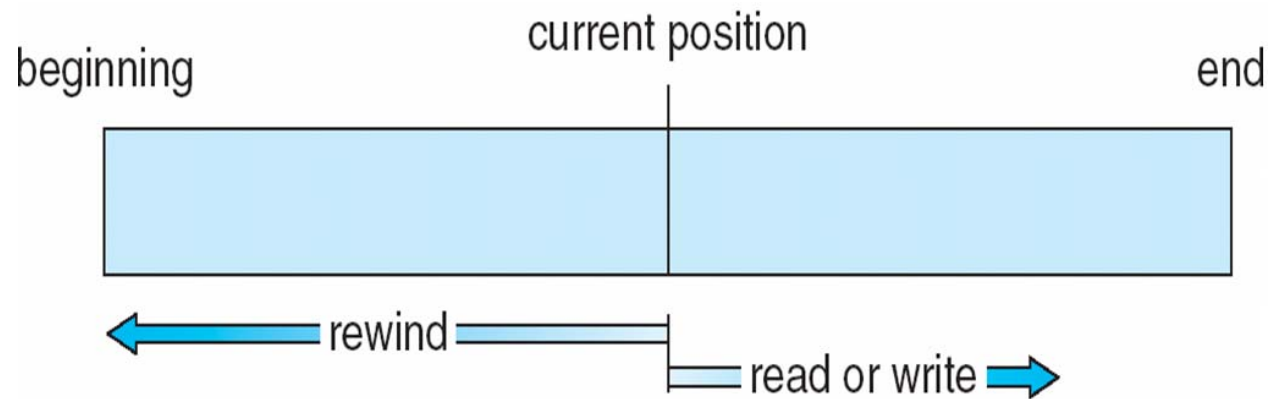| file type | usual extension | function |
|---|---|---|
| executable | exe, com, bin or none | ready-to-run machine-language program |
| object | obj, o | compiled, machine language, not linked |
| source code | c, cc, java, pas, asm, a | source code in various languages |
| batch | bat, sh | commands 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 |

# 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
- Executable file is mandatory structure for any OS

# Sequential-access File

# Access Methods

- **Sequential Access**

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

- **Direct Access –** file is fixed length logical records

  ```
  read n
  write n
  position to n
          read next
          write next
  rewrite n
  ```

  $n$ = relative block number

- Relative block numbers allow OS to decide where file should be placed

# Simulation of Sequential Access on Direct-access File

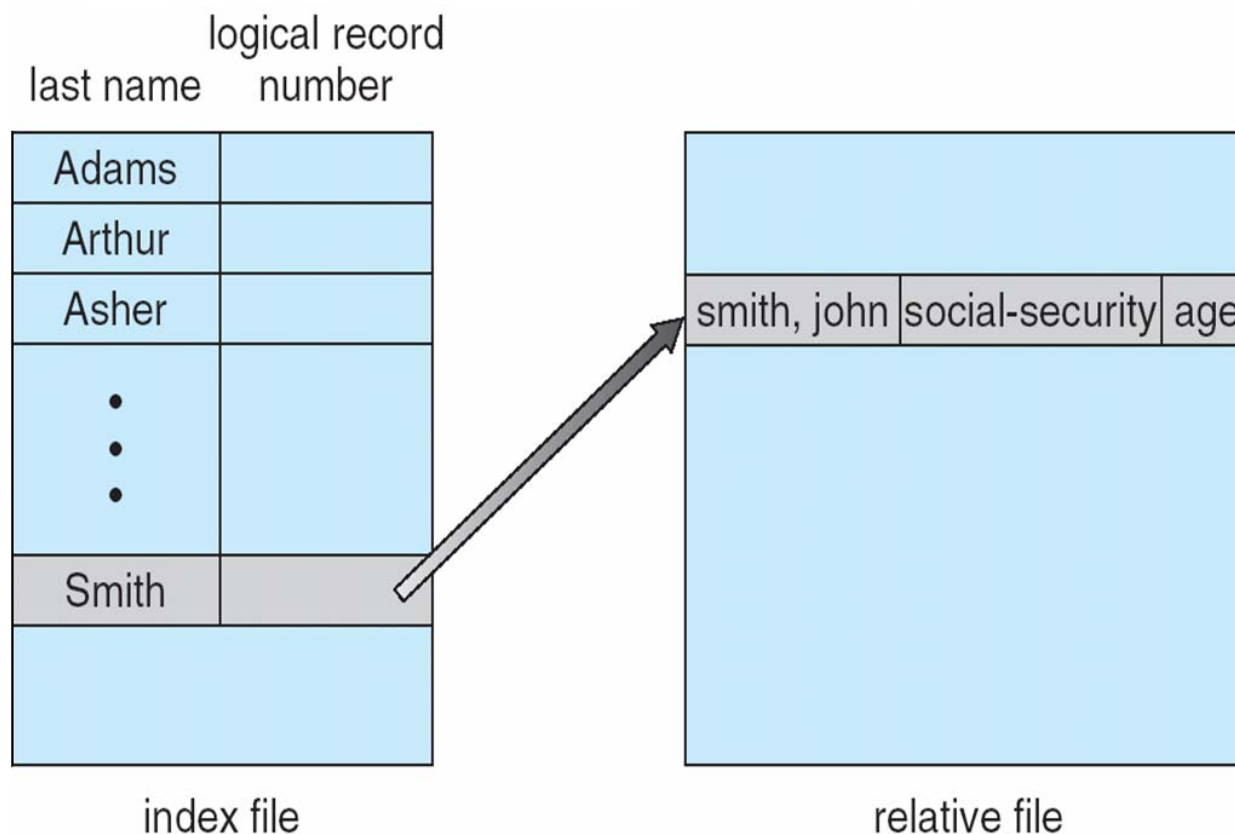| sequential access | implementation for direct access |
|---|---|
| reset | $cp = 0;$ |
| read next | read $cp;$<br>$cp = cp + 1;$ |
| write next | write $cp;$<br>$cp = cp + 1;$ |

# Other Access Methods

- Can be built on top of base methods

- General involve creation of an index for the file

- Keep index in memory for fast determination of location of data to be operated on (consider UPC code plus record of data about that item)

- If too large, index (in memory) of the index (on disk)

- IBM indexed sequential-access method (ISAM)

  - Small master index, points to disk blocks of secondary index

  - File kept sorted on a defined key

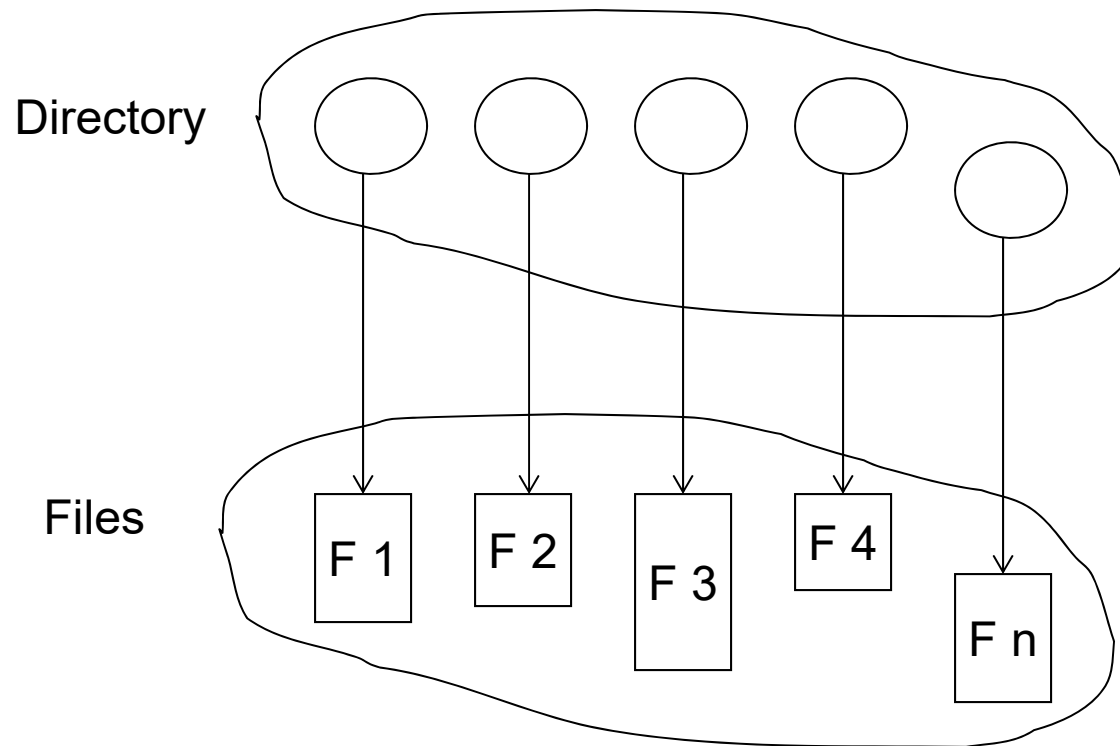  - All done by the OS

# Example of Index and Relative Files



VMS operating system provides 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

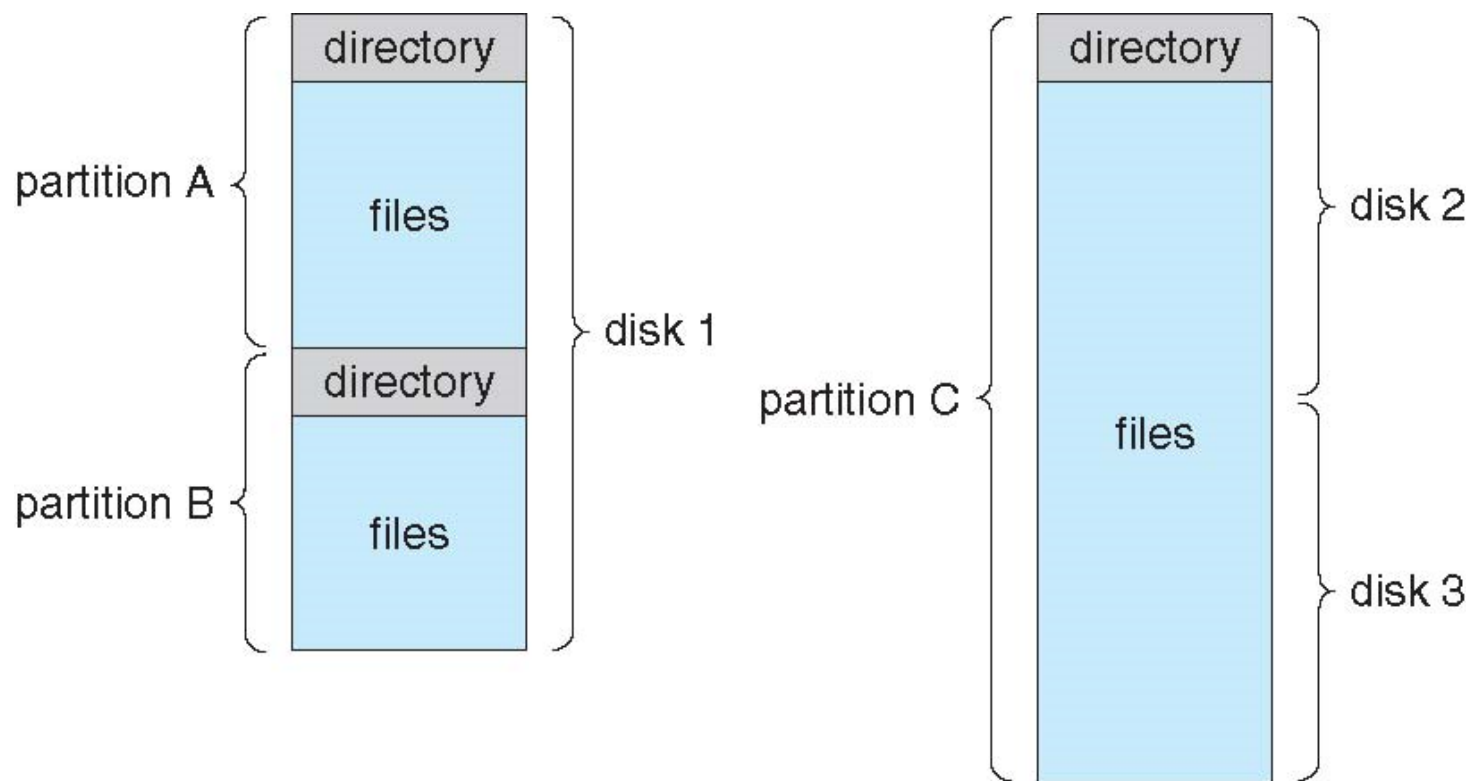Both the directory structure and the files reside on disk

# 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**

- There are many **special-purpose file systems**, frequently all within the same operating system or computer

# A Typical File-system Organization

# Types of File Systems

- Special purpose file systems in Solaris

    - tmpfs – memory-based volatile FS for fast, temporary I/O

    - objfs – interface into kernel memory to get kernel symbols for debugging

    - ctfs – contract file system for managing daemons

    - lofs – loopback file system allows one FS to be accessed in place of another

    - procfs – kernel interface to process structures

    - ufs, zfs – general purpose file systems

# Operations Performed on Directory

- Search for a file

- Create a file

- Delete a file

- List a directory

- Rename a file

- Traverse the file system

# Directory Organization

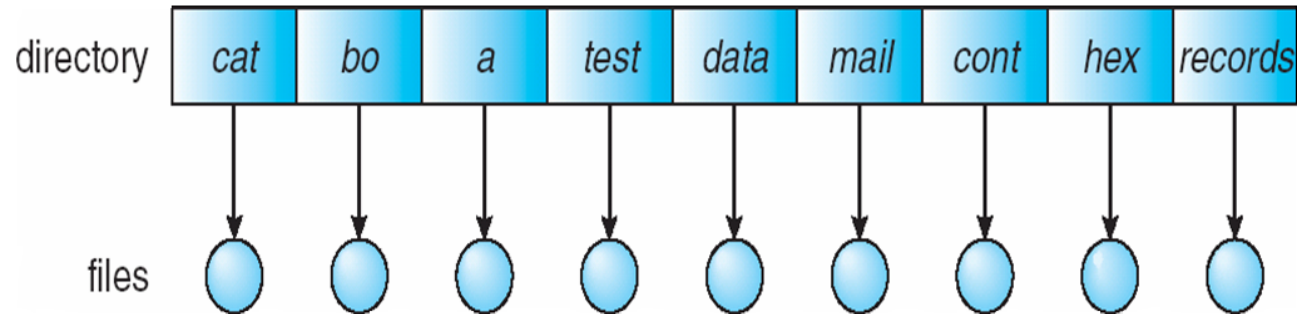The directory is organized logically to obtain

- Efficiency – locating a file quickly

- Naming – convenient to users

  - Two users can have same name for different files

  - The same file can have several different names

- Grouping – logical grouping of files by properties, (e.g., all Java programs, all games, …)

# Single-Level Directory

- A single directory for all users



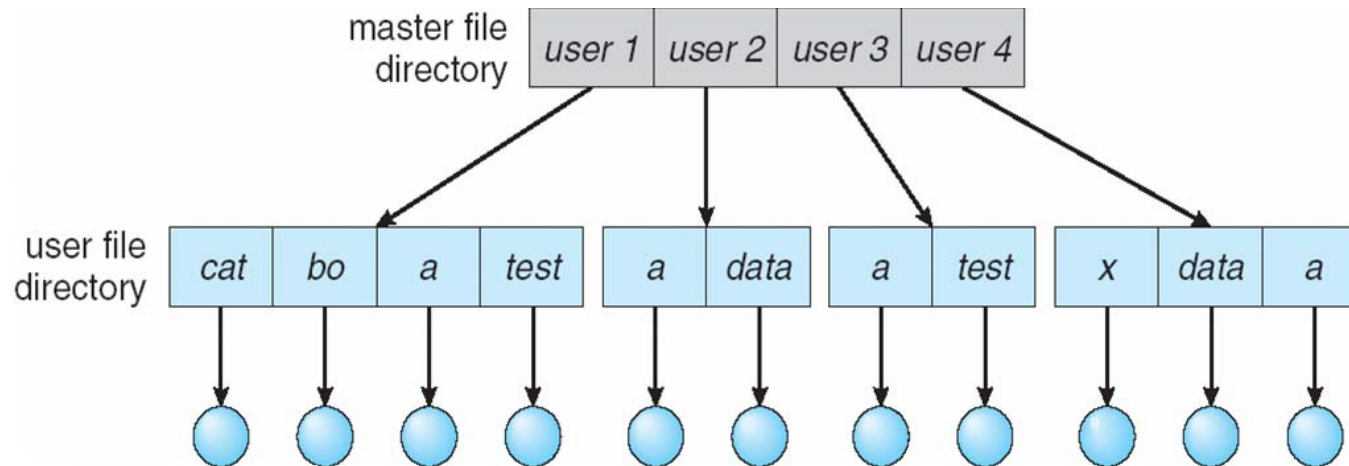| directory | cat | bo | a | test | data | mail | cont | hex | records |

- Naming problem
- Grouping problem

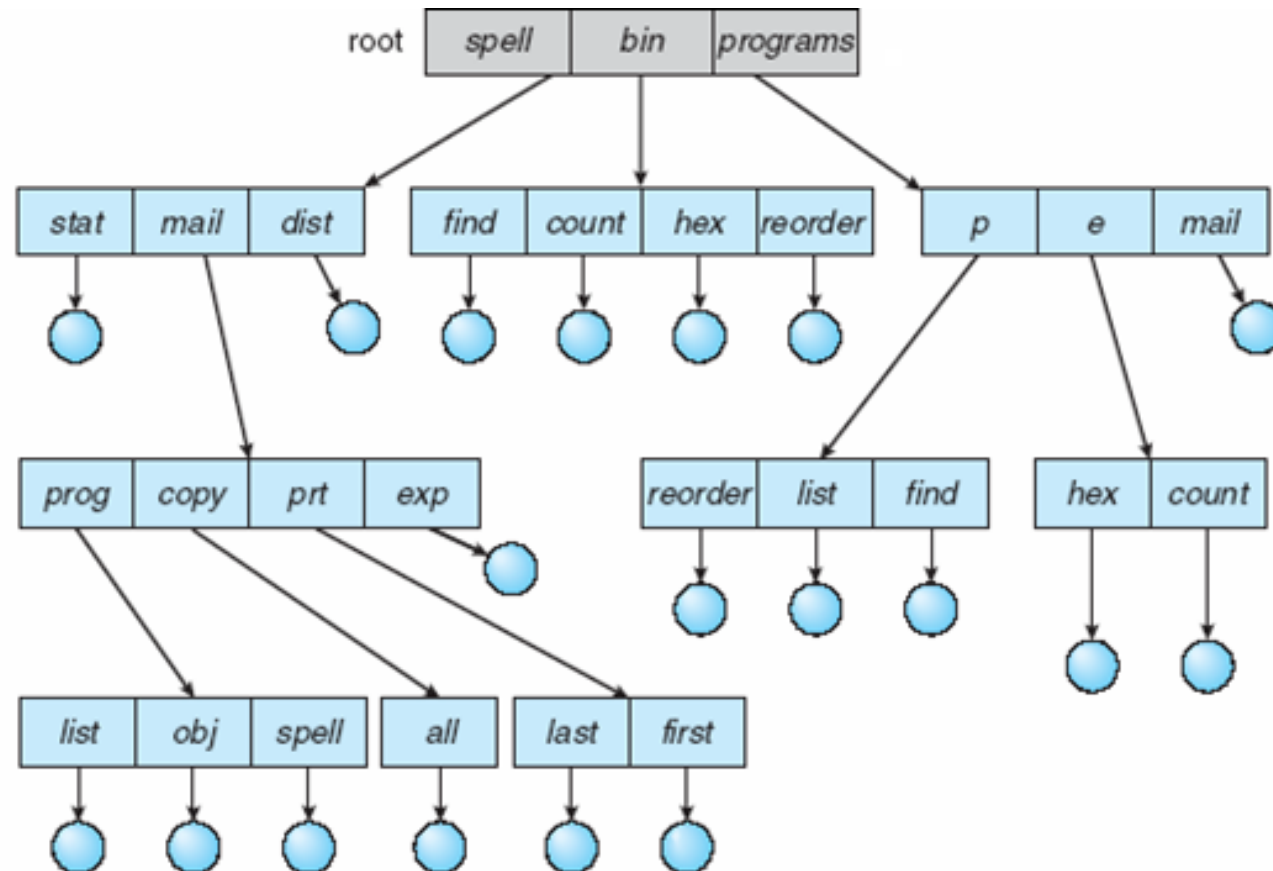# Two-Level Directory

- 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

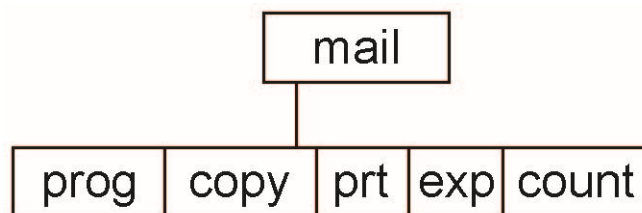      `rm <file-name>`

- Creating a new subdirectory is done in current directory

      `mkdir <dir-name>`

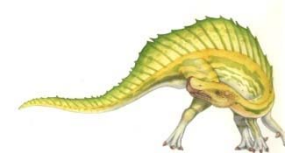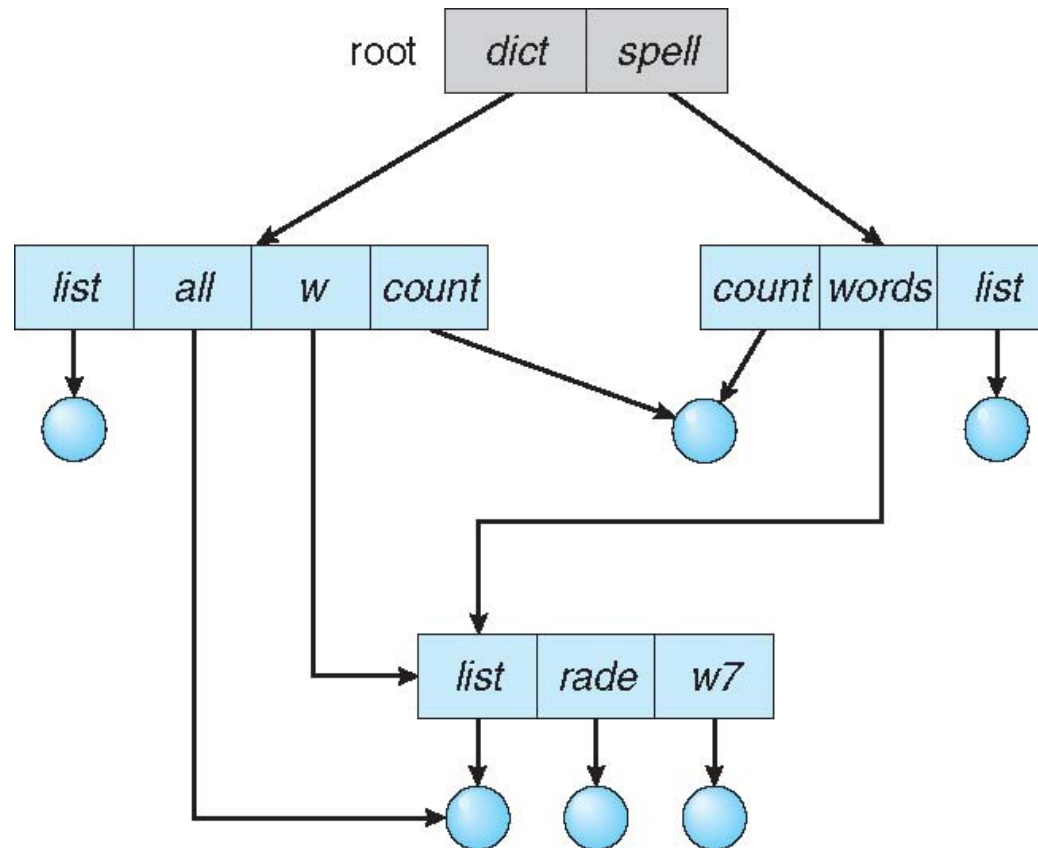    Example: if in current directory `/mail`

      `mkdir count`



    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
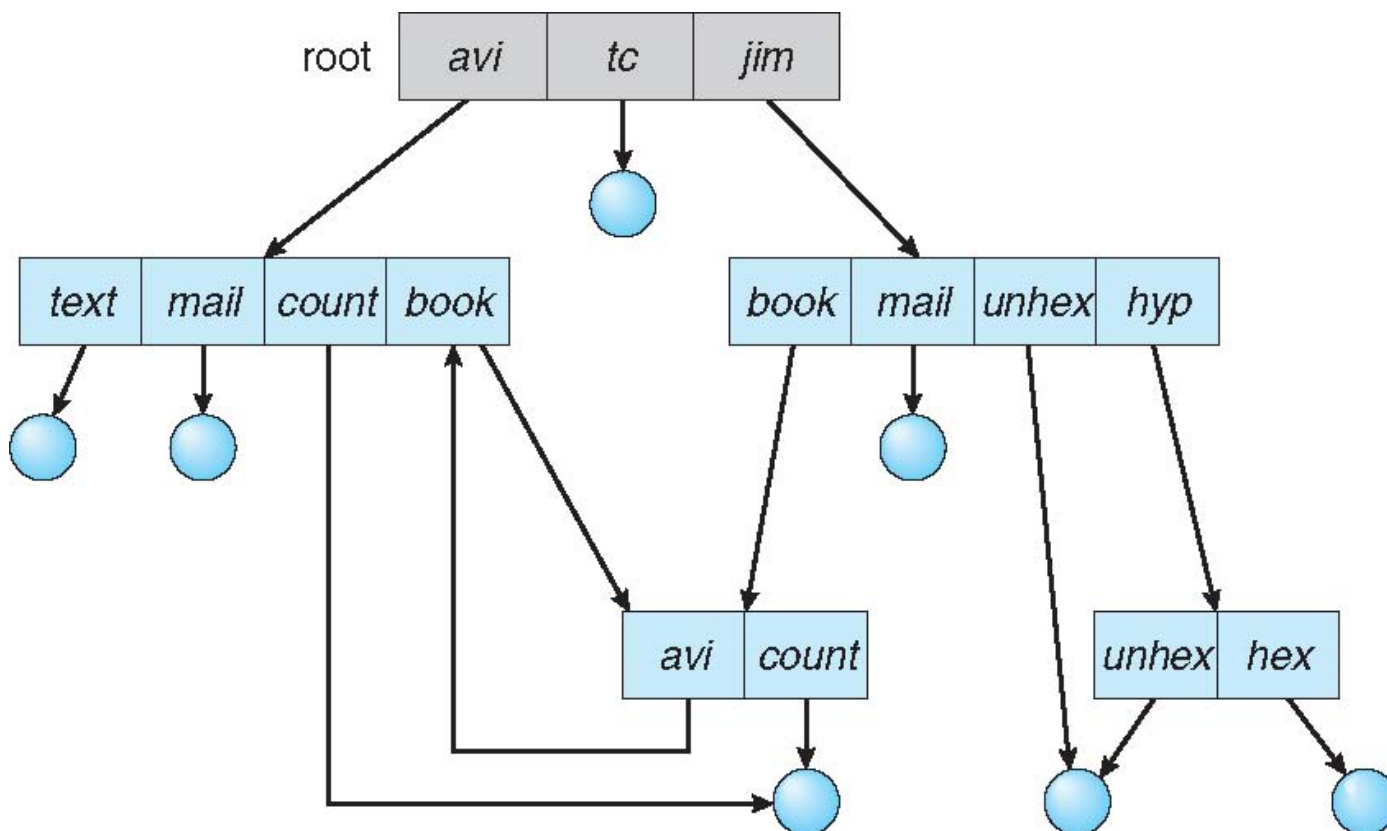
  Solutions:

  - 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
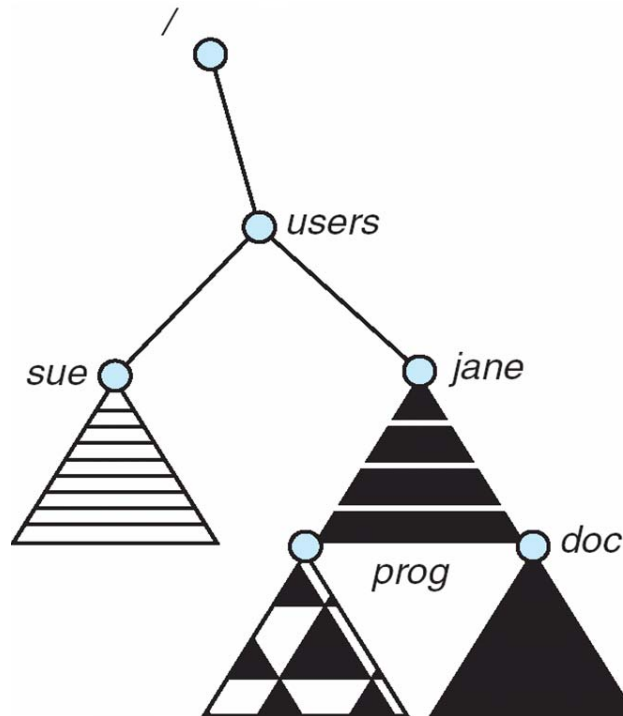
# File System Mounting

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

- A unmounted file system (i.e., Fig. 11-11(b)) 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

- If multi-user system

  - **User IDs** identify users, allowing permissions and protections to be per-user
    **Group IDs** allow users to be in groups, permitting group access rights

  - Owner of a file / directory

  - Group of a file / directory

# 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

- All file systems have failure modes

  - For example corruption of directory structures or other non-user data, called **metadata**

- 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 v3 include all information in each request, allowing easy recovery but less security

# Protection

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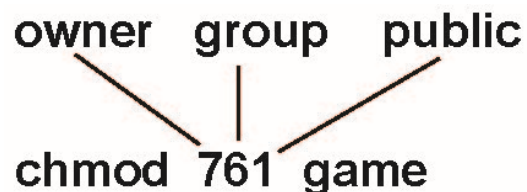    - **Append**
    - **Delete**
    - **List**

# Access Lists and Groups

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

|  |  |  | | 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.

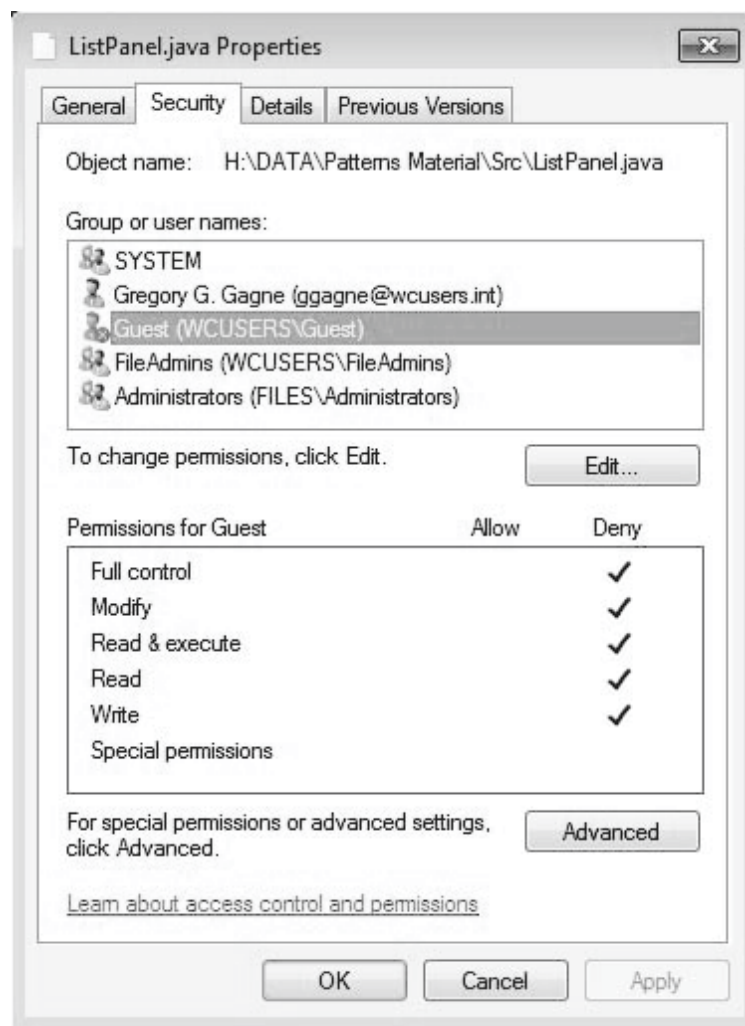- For a particular file (say *game*) or subdirectory, define an appropriate access.

owner   group   public

chmod  761  game

Attach a group to a file

chgrp     G     game

# Windows 7 Access-Control List Management

# A Sample UNIX Directory Listing

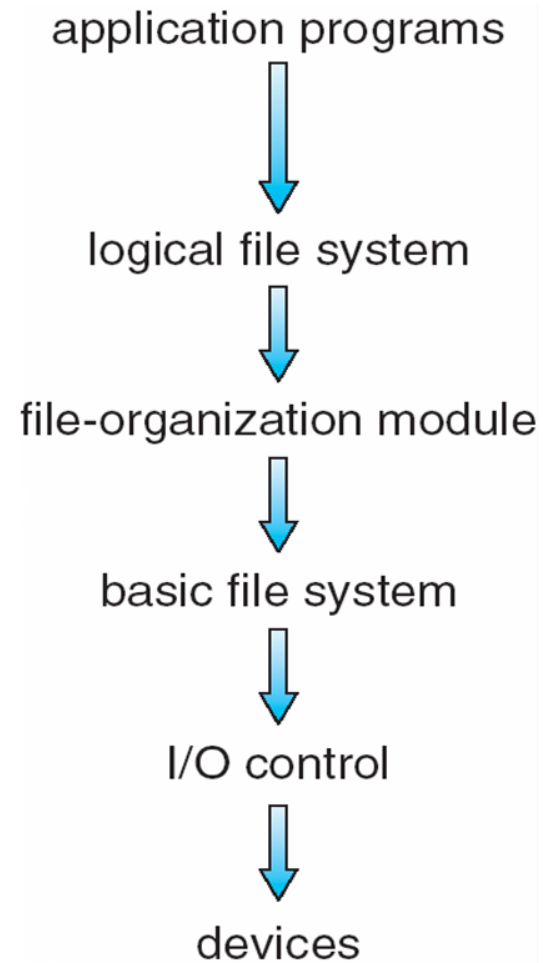| | | | | | |
|---|---|---|---|---|---|
| -rw-rw-r-- | 1 pbg | staff | 31200 | Sep 3 08:30 | intro.ps |
| drwx------ | 5 pbg | staff | 512 | Jul 8 09.33 | private/ |
| drwxrwxr-x | 2 pbg | staff | 512 | Jul 8 09:35 | doc/ |
| drwxrwx--- | 2 pbg | student | 512 | Aug 3 14:13 | student-proj/ |
| -rw-r--r--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

# File-System Structure

- File structure
    - Logical storage unit
    - Collection of related information
- **File system** resides on secondary storage (disks)
    - Provides 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 (called **inode** in Unix)
- **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 commands
- 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

# File System Layers (Cont.)

- 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

# File-System Implementation (Cont.)

- Per-file **File Control Block** (**FCB**) contains many details about the file

  - inode number, permissions, size, dates

  - NFTS stores it in master file table using relational DB structures

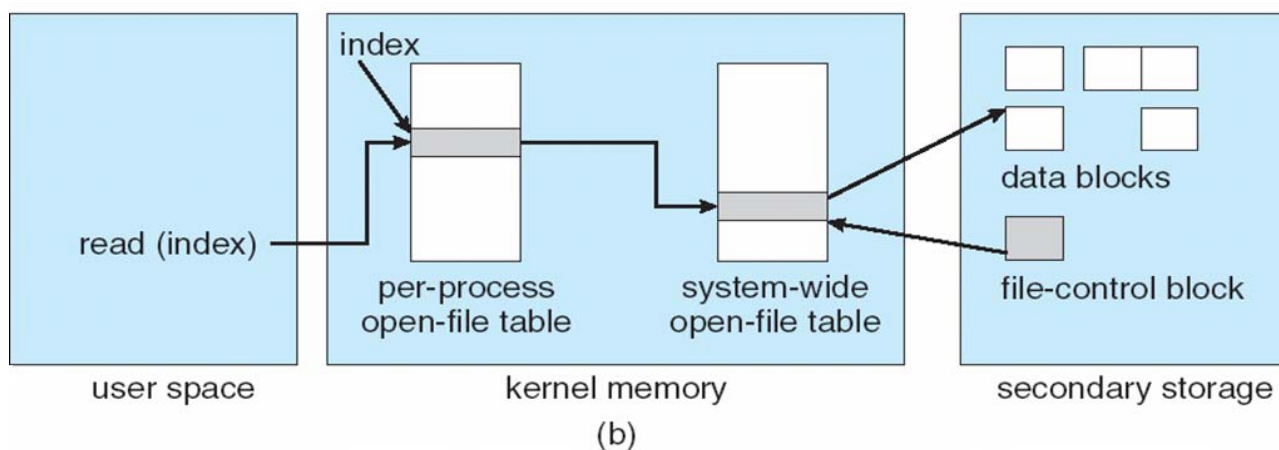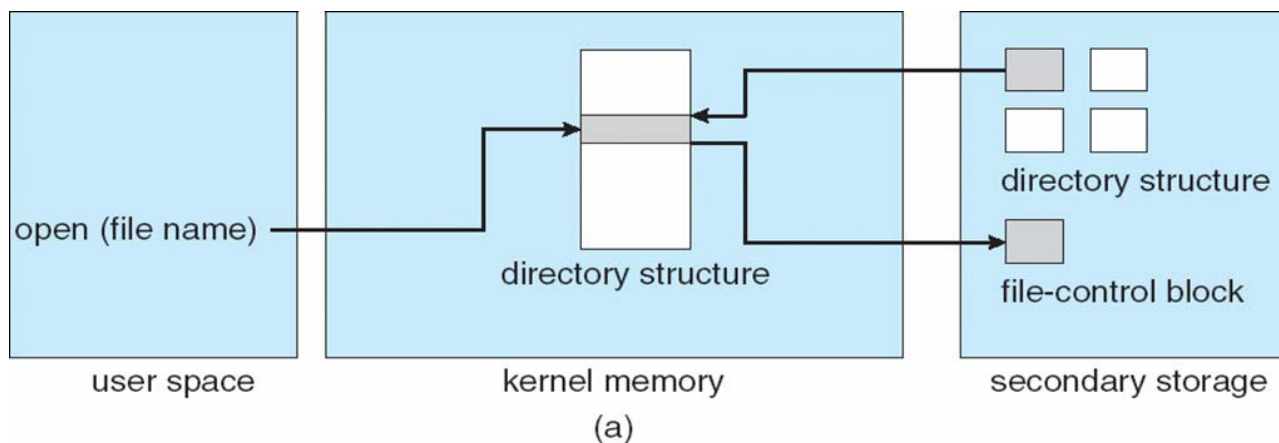| |
|---|
| file permissions |
| file dates (create, access, write) |
| file owner, group, ACL |
| file size |
| file data blocks or pointers to file data blocks |

# In-Memory File System Structures

- 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



open (file name)

directory structure

directory structure

file-control block

user space    kernel memory    secondary storage

(a)

index

read (index)

per-process
open-file table

system-wide
open-file table

data blocks

file-control block

user space    kernel memory    secondary storage
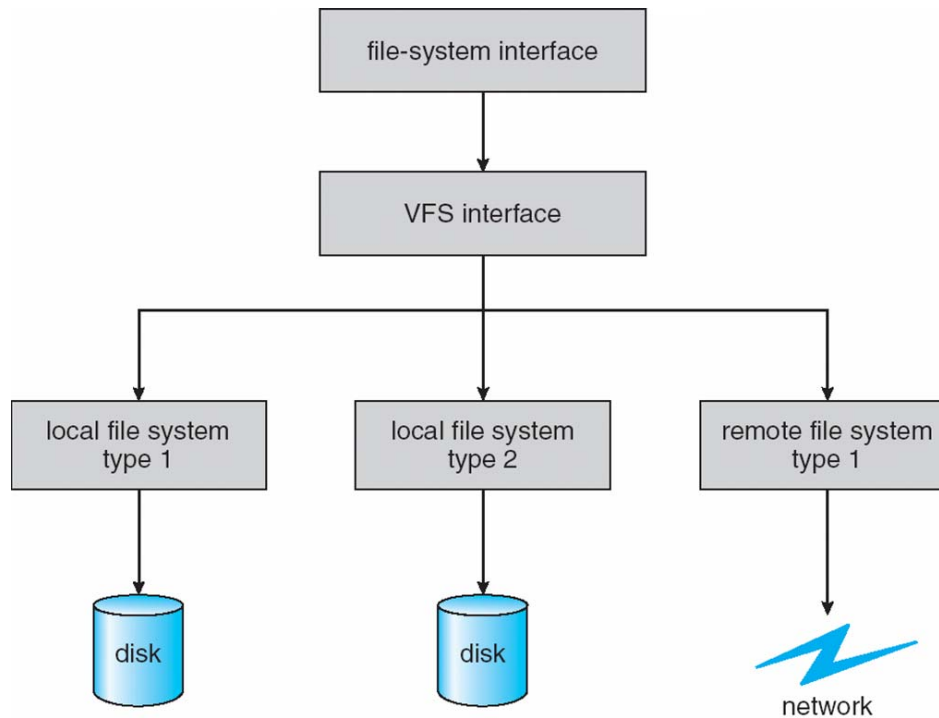
(b)

# Virtual File Systems

- **Virtual File Systems** (**VFS**) on Unix provide an object-oriented way of implementing file systems

- VFS allows the same system call interface (the API) to be used for different types of file systems

  - Separates file-system generic operations from implementation details

  - Implementation can be one of many file systems types, or network file system

    ▸ Implements **vnodes** which hold inodes or network file details

  - Then dispatches operation to appropriate file system implementation routines

# Virtual File Systems (Cont.)

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

# Virtual File System Implementation

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

# 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
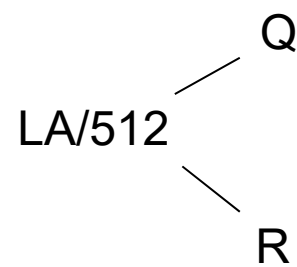
# Allocation Methods - Contiguous

- 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**
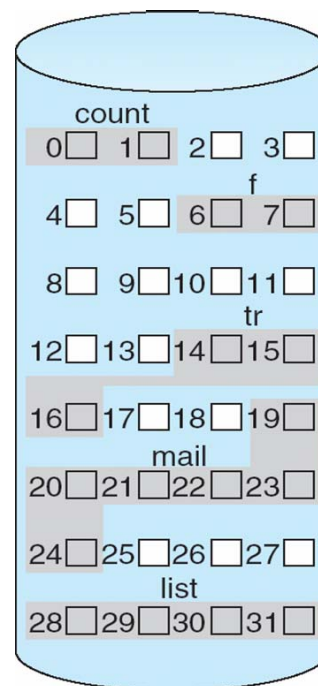
# Contiguous Allocation

☐ Mapping from logical to physical

$LA/512$

$$\begin{array}{c} Q \\ \diagdown \\ \diagup \\ R \end{array}$$

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

**directory**

| 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
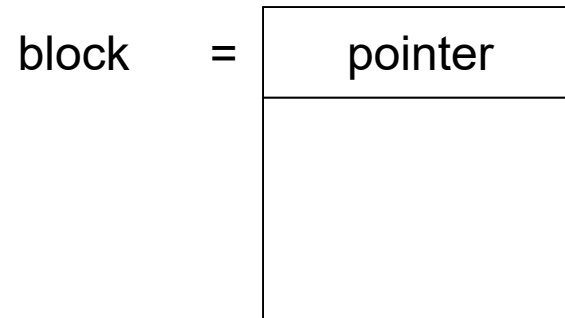
# Allocation Methods - Linked

- **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

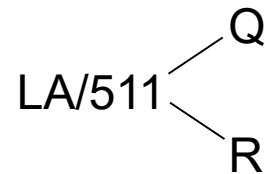  - Locating a block can take many I/Os and disk seeks

# Linked Allocation

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

block  =  | pointer |
          |         |

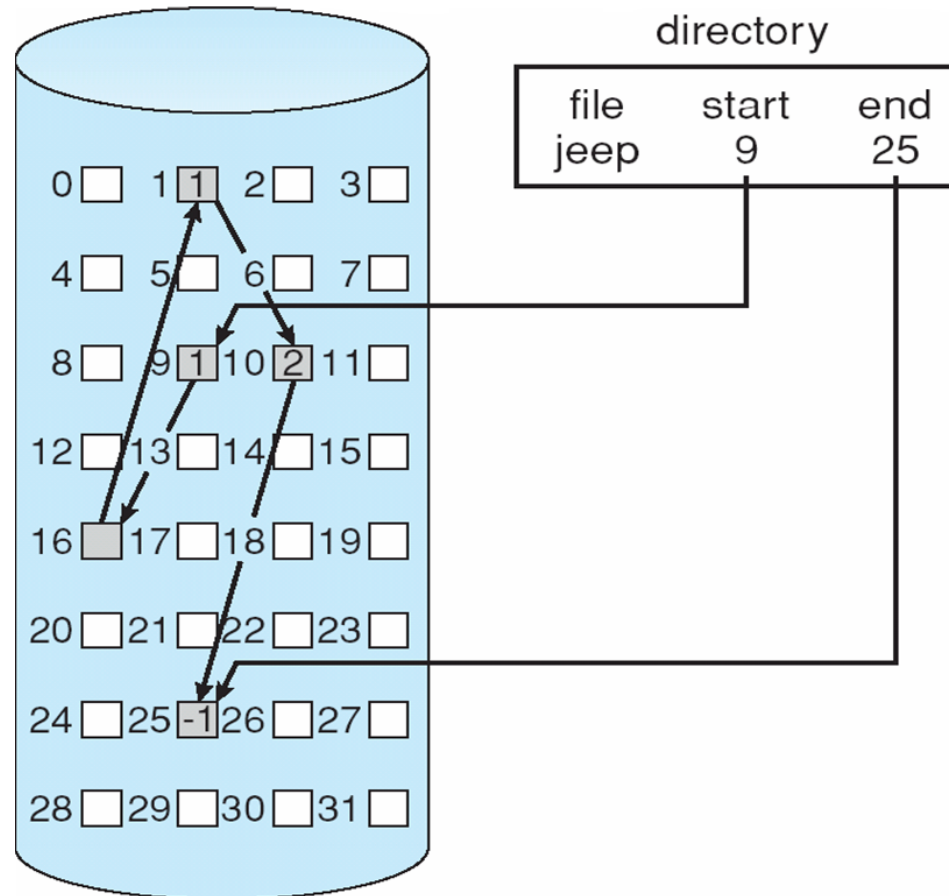- Mapping

$$LA/511 \begin{array}{c} Q \\ R \end{array}$$

Block to be accessed is the Qth block in the linked chain of blocks representing the file.

Displacement into block = R + 1

# Linked Allocation
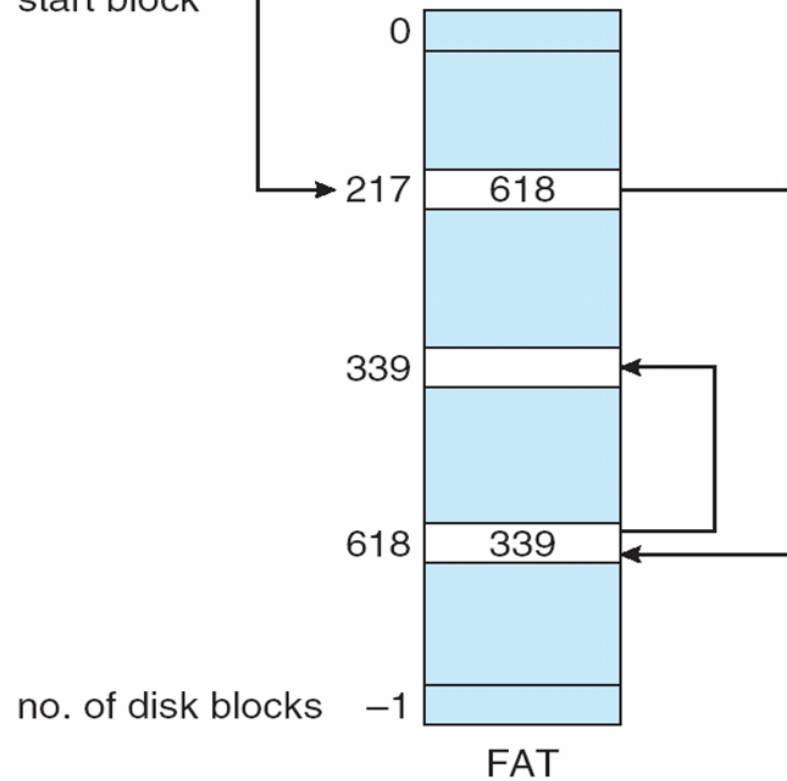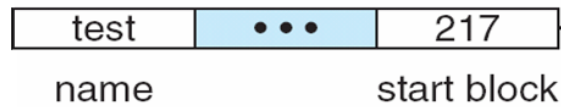
# Allocation Methods – Linked (Cont.)

- 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
    - Directory entry has first block that points to next block
    - New block allocation requires finding first empty entry and extend EOF to point to this new block
    - FAT is cached to avoid disk seeks; random access is faster

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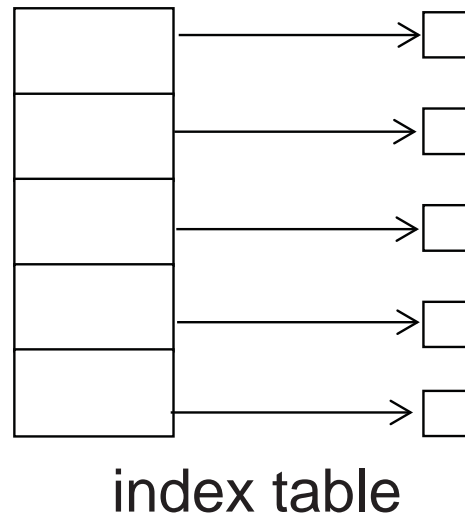
- **Indexed allocation**
  - Each file has its own **index block**(s) of pointers to its data blocks
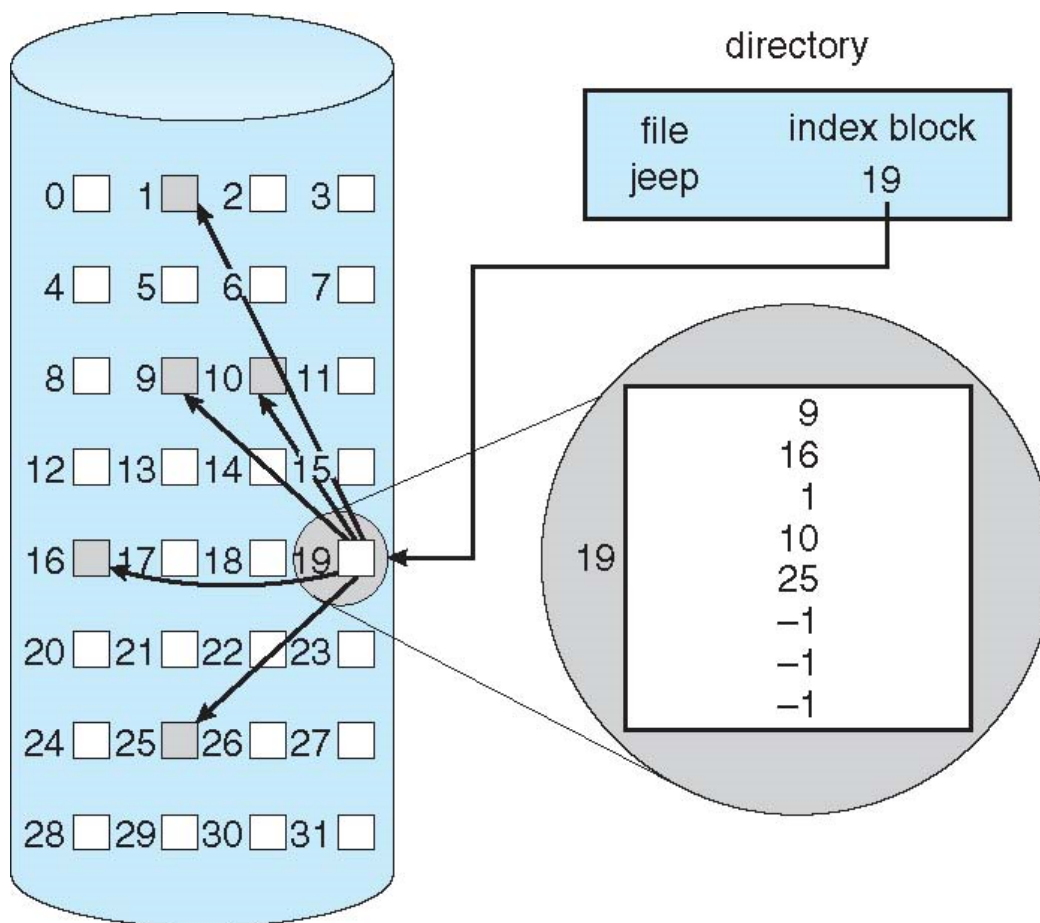- Need index table

- Random access

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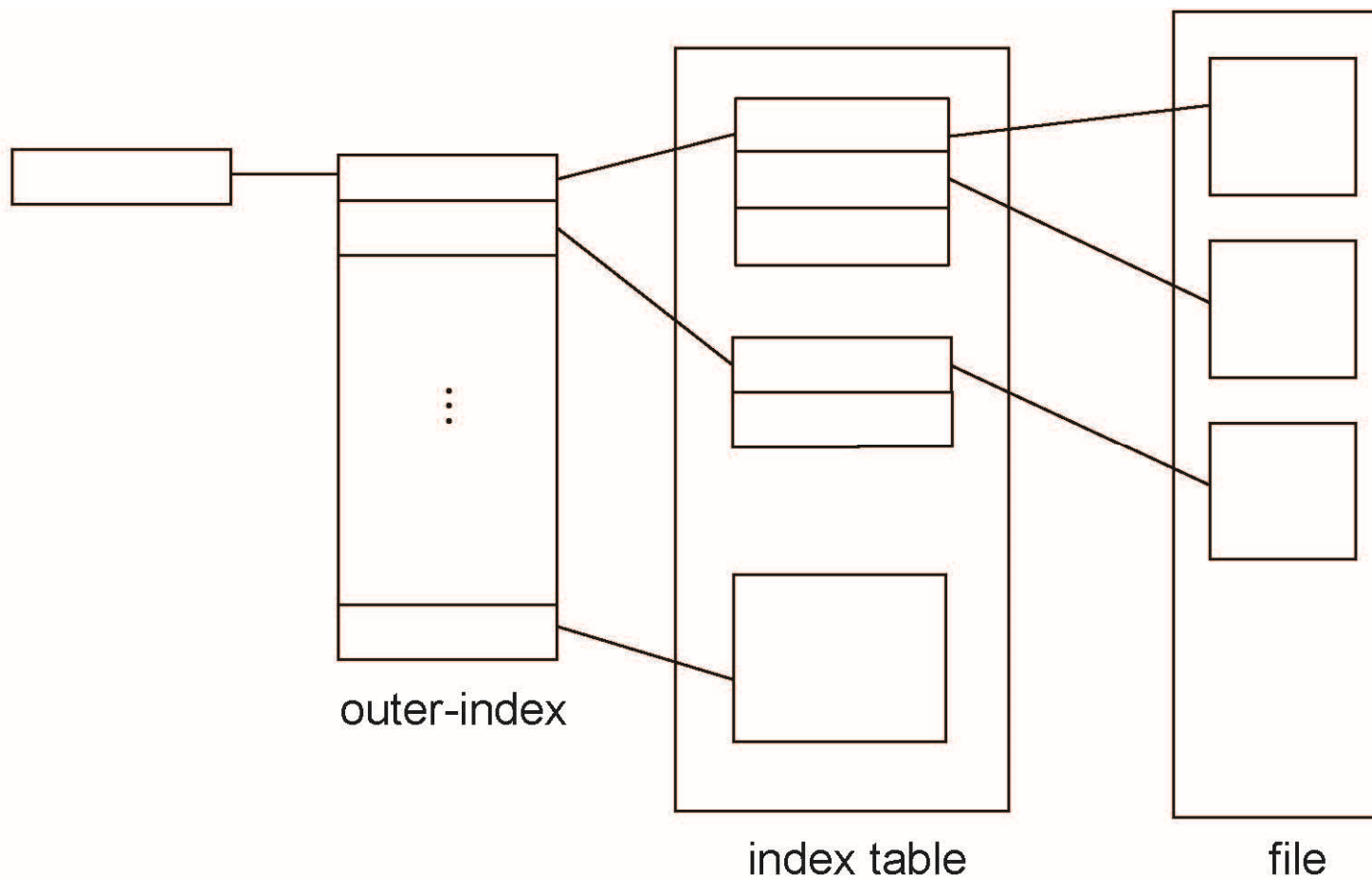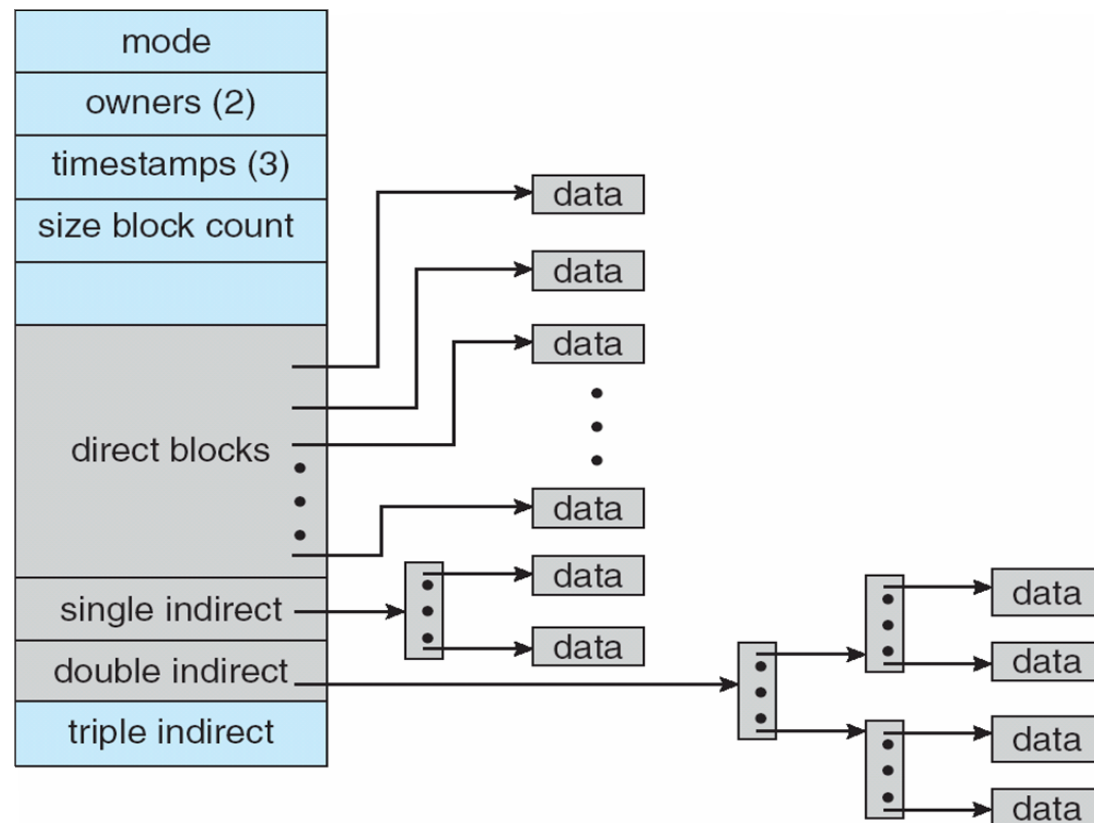
- Logical view

index table

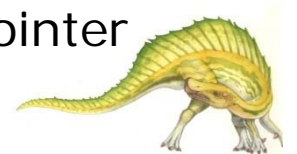# Example of Indexed Allocation

outer-index

index table

file

# Combined Scheme: UNIX UFS

4K bytes per block, 32-bit addresses



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

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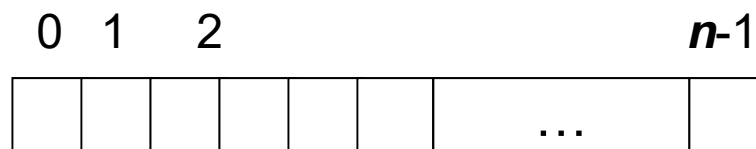
- Best method depends on file access type
    - Contiguous great for sequential and random
- Linked good for sequential, not random
- Declare access type at creation -> select either contiguous or linked
- Indexed more complex
    - Single block access could require 2 index block reads then data block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

# 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)

$$\underset{0 \quad 1 \quad 2 \qquad\qquad\qquad\qquad\qquad\qquad n\text{-}1}{\boxed{\phantom{xx}\phantom{xx}\phantom{xx}\phantom{xx}\phantom{xx}\phantom{xx}\cdots\phantom{xx}}}$$

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

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
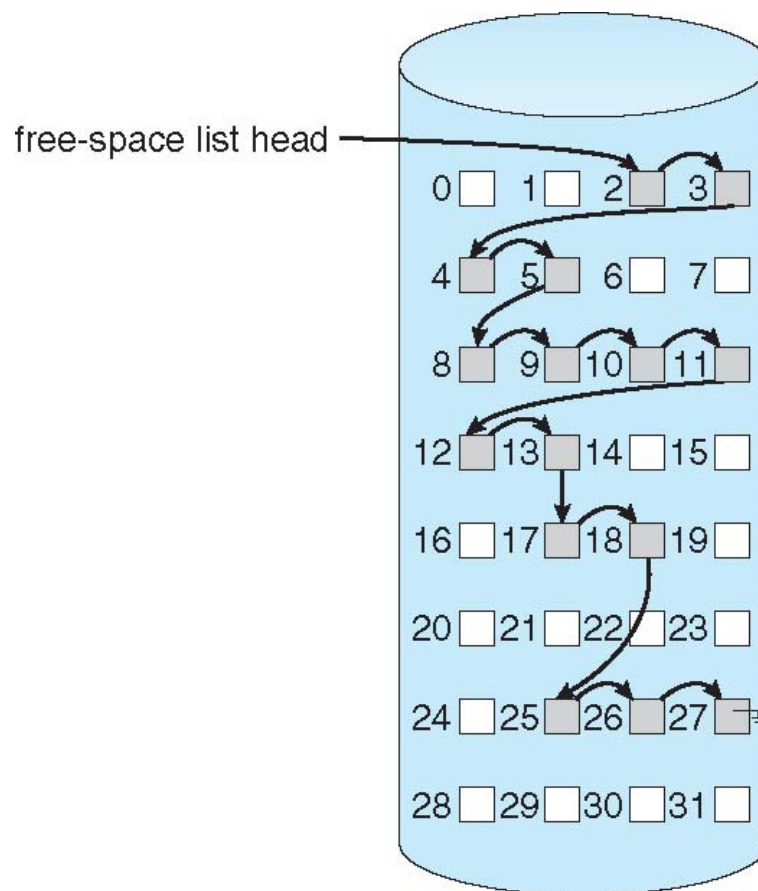
# Free-Space Management (Cont.)

- 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 32MB)

    if clusters of 4 blocks -> 8MB of memory

- Easy to get contiguous files

# Linked Free Space List on Disk

- Linked list (free list)

  - Cannot get contiguous space easily

  - No waste of space

  - No need to traverse the entire list (if # free blocks recorded)



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

# 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

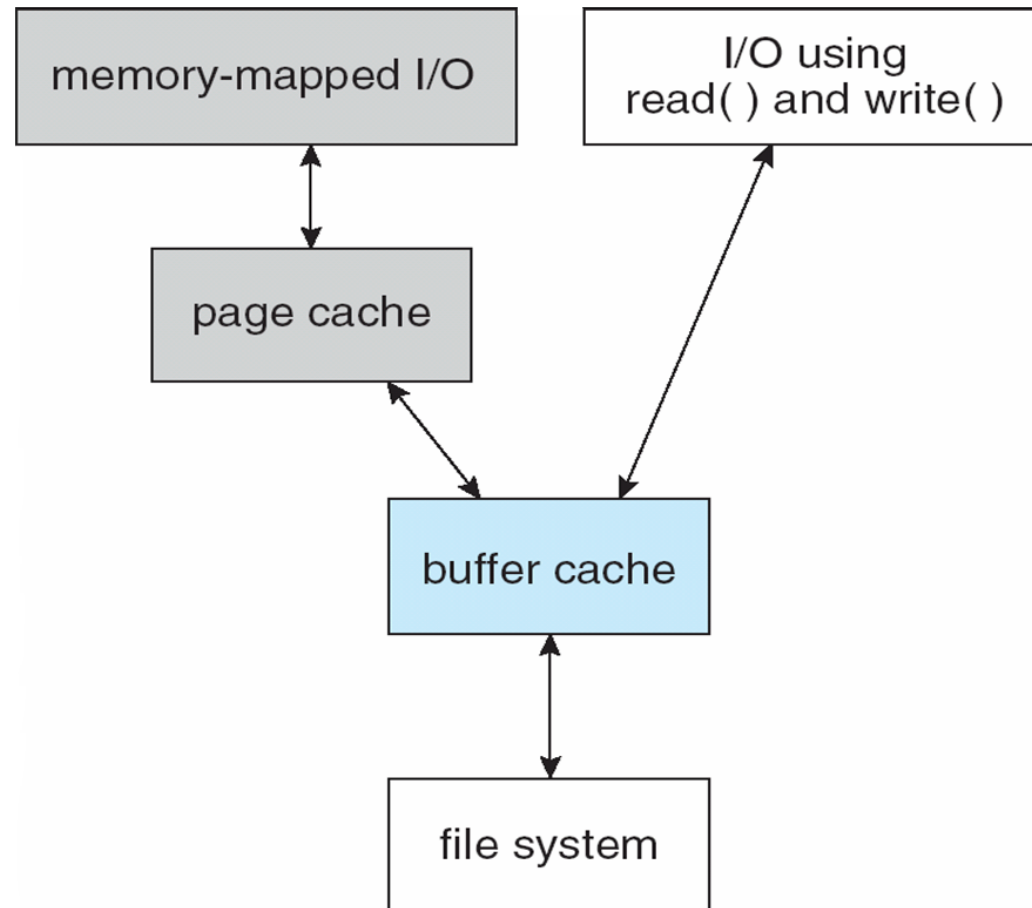  - Reads frequently slower than writes

# Page Cache

- A **page cache** caches pages rather than disk blocks using virtual memory techniques and addresses

- Memory-mapped I/O uses a page cache

- Routine I/O through the file system uses the buffer (disk) cache

- This leads to the following figure

# I/O Without a Unified Buffer Cache
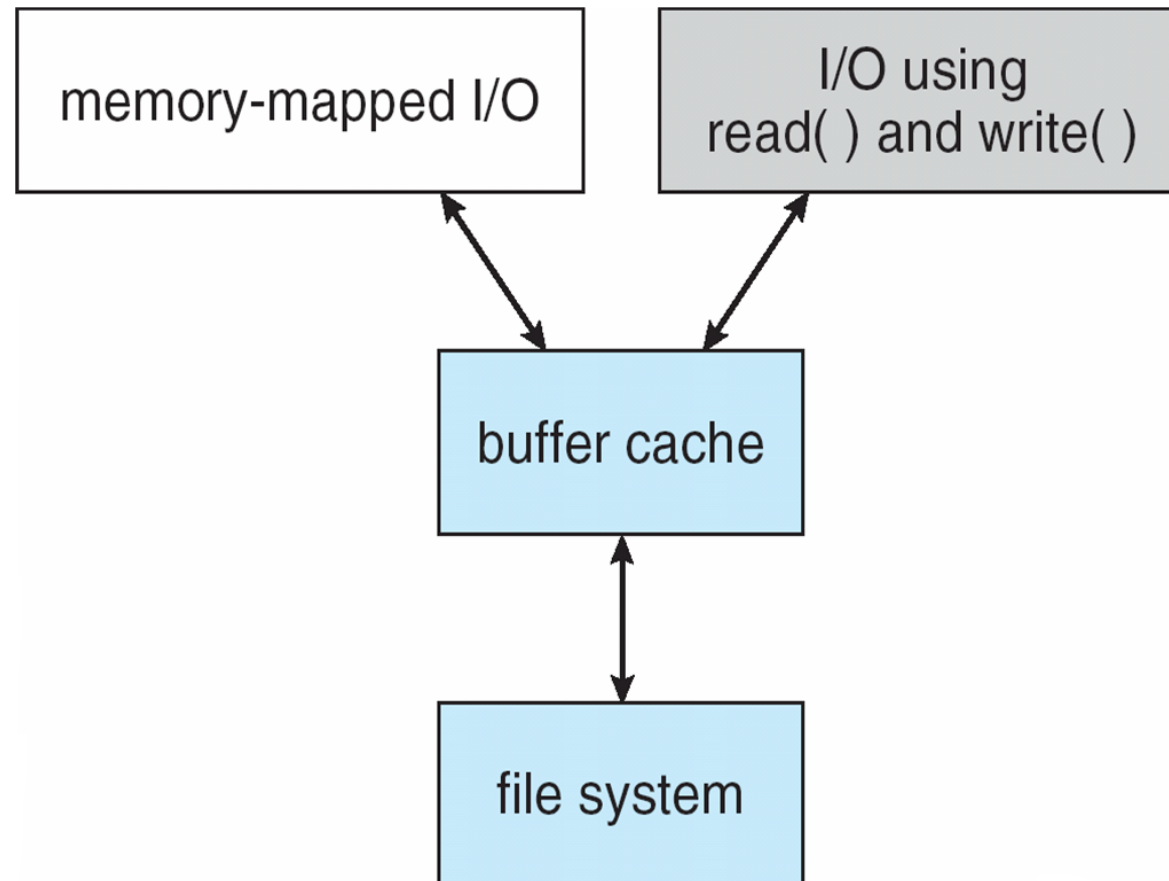
# Unified Buffer Cache

- A **unified buffer cache** uses the same page cache to cache both memory-mapped pages and ordinary file system I/O to avoid **double caching**

- But which caches get priority, and what replacement algorithms to use?

# I/O Using a Unified Buffer Cache

# Recovery

- **Consistency checking** – compares data in directory structure with data blocks on disk, and tries to fix inconsistencies (fsck in Linux)
    - Can be slow and sometimes fails

- Use system programs to **back up** data from disk to another storage device (magnetic tape, other magnetic disk, optical)

- Recover lost file or disk by **restoring** data from backup

# Log Structured File Systems

- **Log structured** (or **journaling**) file systems record each metadata update to the file system as a **transaction**

- All transactions are written to a log

  - A transaction is considered committed once it is written to the log (sequentially)

  - Sometimes to a separate device or section of disk

  - However, the file system may not yet be updated

- The transactions in the log are asynchronously written to the file system structures

  - When the file system structures are modified, the transaction is removed from the log

- If the file system crashes, all remaining transactions in the log must still be performed

- Faster recovery from crash, removes chance of inconsistency of metadata