

11. File-System Interface

Sungyoung Lee

*College of Engineering
KyungHee University*

Contents

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

n Requirements for long-term information storage

- ü It must be possible to store a very large amount of information
- ü The information must survive the termination of the process using it
- ü Multiple processes must be able to access the information concurrently

n File system

- ü Implement an abstraction for secondary storage (files)
- ü Organize files logically (directories)
- ü Permit sharing of data between processes, people, and machines
- ü Protect data from unwanted access (security)

n File

- ü A named collection of related information that is recorded on secondary storage
 - § persistent through power failures and system reboots
- ü OS provides a uniform logical view of information storage via files

n File structures

- ü Flat: byte sequence
- ü Structured:
 - § Lines
 - § Fixed length records
 - § Variable length records

File Concept

n Contiguous logical address space

n Types:

ü Data

§ numeric

§ character

§ binary

ü Program

File Structure

n None

- ü Sequence of words, bytes

n Simple record structure

- ü Lines
- ü Fixed length
- ü Variable length

n Complex Structures

- ü Formatted document
- ü Relocatable load file

n Can simulate last two with first method by inserting appropriate control characters

n Who decides:

- ü Operating system
- ü Program

File Attributes

n Name

- ü only information kept in human-readable form

n Type

- ü needed for systems that support different types

n Location

- ü pointer to file location on device

n Size

- ü current file size

n Protection

- ü controls who can do reading, writing, executing

n Time, date, and user identification

- ü data for protection, security, and usage monitoring

n Information about files are kept in the directory structure, which is maintained on the disk

File Attributes

n Attributes or metadata

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

File Operations

- n Create
- n Write
- n Read
- n Reposition within file – file seek
- n Delete
- n Truncate
- n Open(F_i)
 - ü search the directory structure on disk for entry F_i , and move the content of entry to memory
- n Close (F_i)
 - ü move the content of entry F_i in memory to directory structure on disk

File Operations

n Unix operations

```
int creat (const char *pathname, mode_t mode);  
int open (const char *pathname, int flags, mode_t mode);  
int close (int fd);  
ssize_t read (int fd, void *buf, size_t count);  
ssize_t write (int fd, const void *buf, size_t count);  
off_t lseek (int fd, off_t offset, int whence);  
int stat (const char *pathname, struct stat *buf);  
int chmod (const char *pathname, mode_t mode);  
int chown (const char *pathname, uid_t owner, gid_t grp);  
int flock (int fd, int operation);  
int fcntl (int fd, int cmd, long arg);
```

n Files may have types

- ü Understood by file systems
 - § device, directory, symbolic link, etc.
- ü Understood by other parts of OS or runtime libraries
 - § executable, dll, source code, object code, text, etc.
- ü Understood by application programs
 - § jpg, mpg, avi, mp3, etc.

n Encoding file types

- ü Windows encodes type in name
 - § .com, .exe, .bat, .dll, .jpg, .avi, .mp3, etc.
- ü Unix encodes type in contents
 - § magic numbers (e.g., 0xcafebabe for Java class files)
 - § initial characters (e.g., #! for shell scripts)

File Types – Name, Extension

file type	usual extension	function
executable	exe, com, bin or none	read 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, rrf, doc	various word-processor formats
library	lib, a, so, dll, mpeg, mov, rm	libraries of routines for programmers
print or view	arc, zip, tar	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	binary file containing audio or A/V information

- n Some file systems provide different access methods that specify different ways for accessing data in a file.

- n **Sequential access**
 - ü read bytes one at a time, in order
- n **Direct access**
 - ü random access given block/byte number
- n **Record access**
 - ü File is an array of fixed- or variable-length records, read/written sequentially or randomly by record #
- n **Index access**
 - ü File system contains an index to a particular field of each record in a file, reads specify a value for that field and the system finds the records via the index (DBs)

Access Methods

n Sequential Access

read next

write next

reset

no read after last write
(rewrite)

n Direct Access

read n

write n

position to n

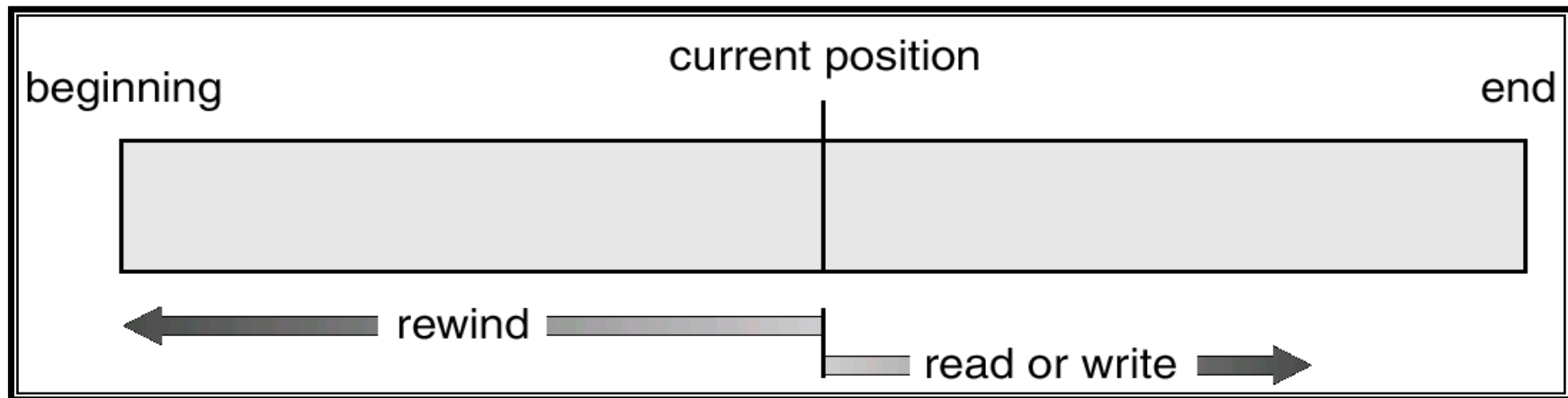
read next

write next

rewrite n

n = relative block number

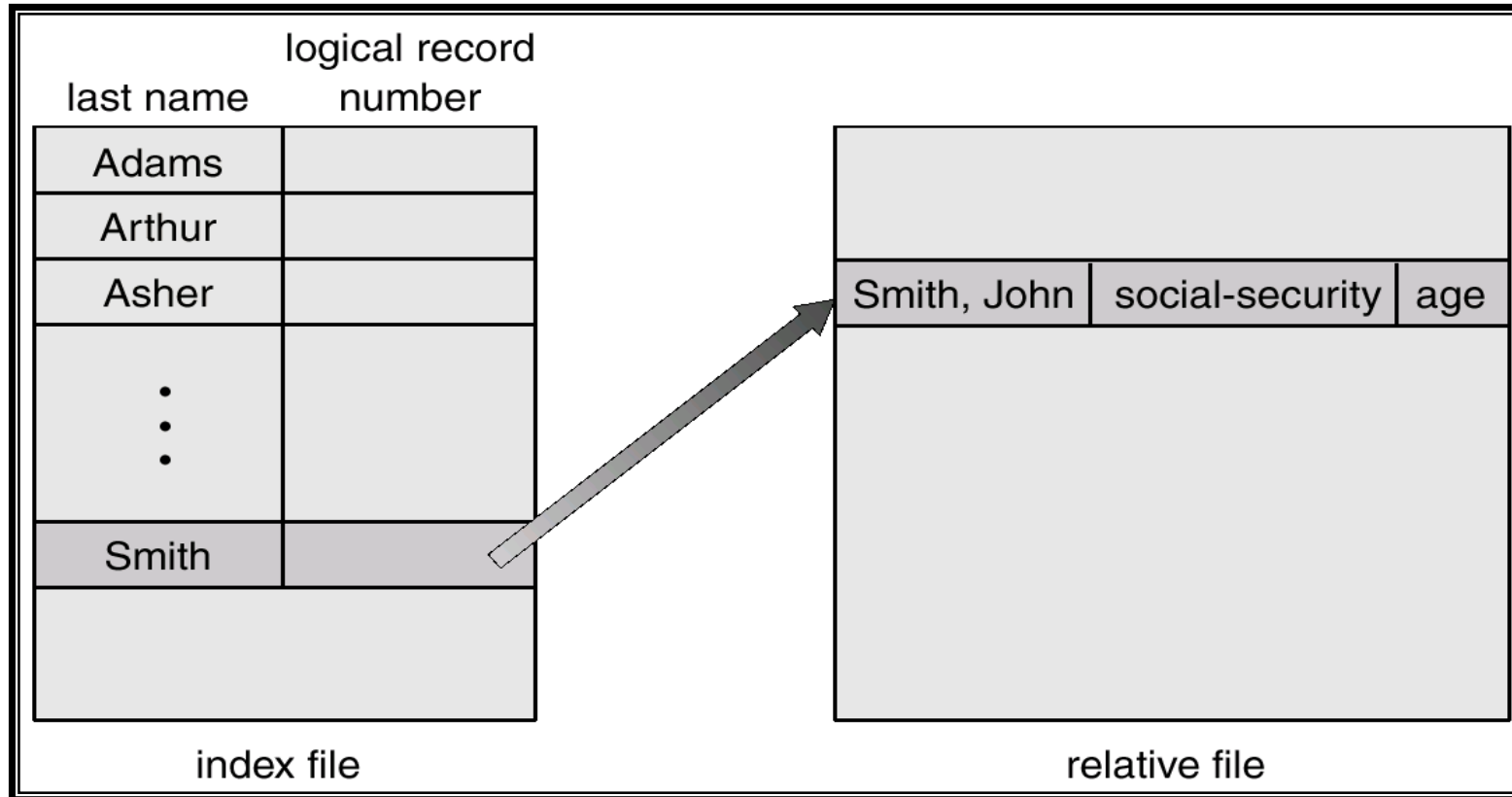
Sequential-access File



Simulation of Sequential Access on a Direct-access File

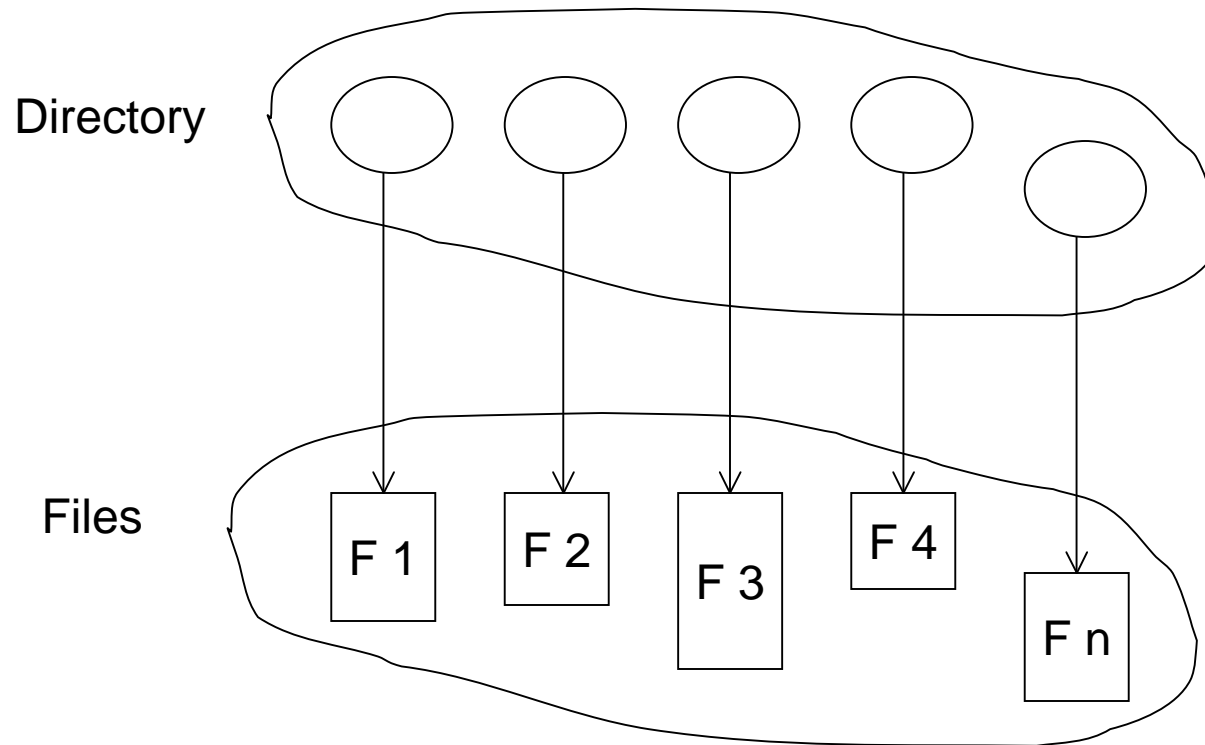
sequential access	implementation for direct access
<i>reset</i>	<i>cp = 0;</i>
<i>read next</i>	<i>read cp;</i> <i>cp = cp+1;</i>
<i>write next</i>	<i>write cp;</i> <i>cp = cp+1;</i>

Example of Index and Relative Files



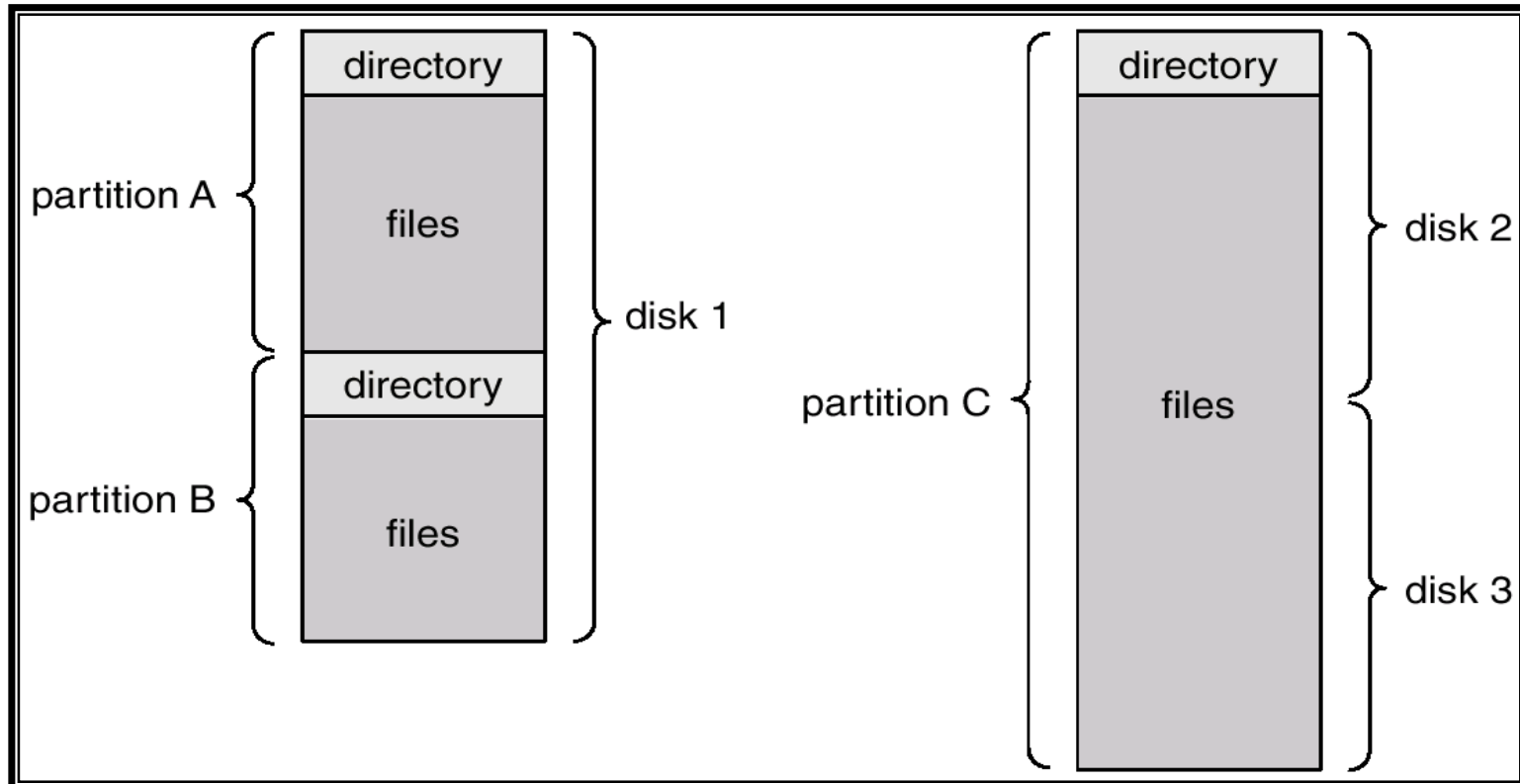
Directory Structure

n A collection of nodes containing information about all files



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

A Typical File-system Organization



n Directories

- ü For users, they provide a structured way to organize files
- ü For the file system, they provide a convenient naming interface that allows the implementation to separate logical file organization from physical file placement on the disk

n A hierarchical directory system

- ü Most file systems support multi-level directories
- ü Most file systems support the notion of a current directory (or working directory)
 - § Relative names specified with respect to current directory
 - § Absolute names start from the root of directory tree

n A directory is ...

- ü typically just a file that happens to contain special metadata
 - § Only need to manage one kind of secondary storage unit
- ü directory = list of (file name, file attributes)
- ü attributes include such things as:
 - § size, protection, creation time, access time,
 - § location on disk, etc.
- ü usually unordered (effectively random)
 - § Entries usually sorted by program that reads directory

Information in a Device Directory

- n Name
- n Type
- n Address
- n Current length
- n Maximum length
- n Date last accessed (for archival)
- n Date last updated (for dump)
- n Owner ID (who pays)
- n Protection information (discuss later)

Operations Performed on Directory

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

Directory Operations

n Unix operations

ü Directories implemented in files

§ Use file operations to manipulate directories

ü C runtime libraries provides a higher-level abstraction for reading directories

§ `DIR *opendir (const char *name);`

§ `struct dirent *readdir (DIR *dir);`

§ `void seekdir (DIR *dir, off_t offset);`

§ `int closedir (DIR *dir);`

ü Other directory-related system calls

§ `int rename (const char *oldpath, const char *newpath);`

§ `int link (const char *oldpath, const char *newpath);`

§ `int unlink (const char *pathname);`

Organize the Directory (Logically) to Obtain

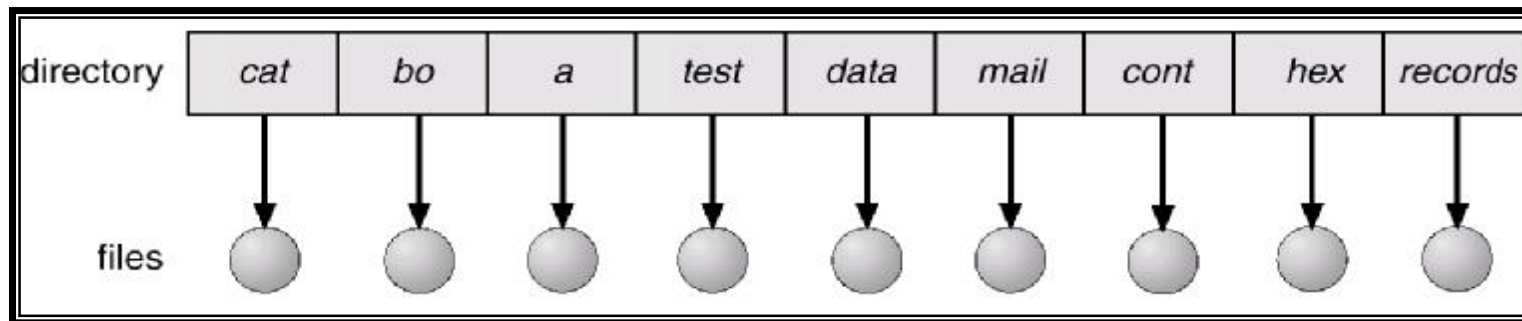
- n **Efficiency** – locating a file quickly

- n **Naming** – convenient to users
 - ü Two users can have same name for different files
 - ü The same file can have several different names

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

Single-Level Directory

- n A single directory for all users

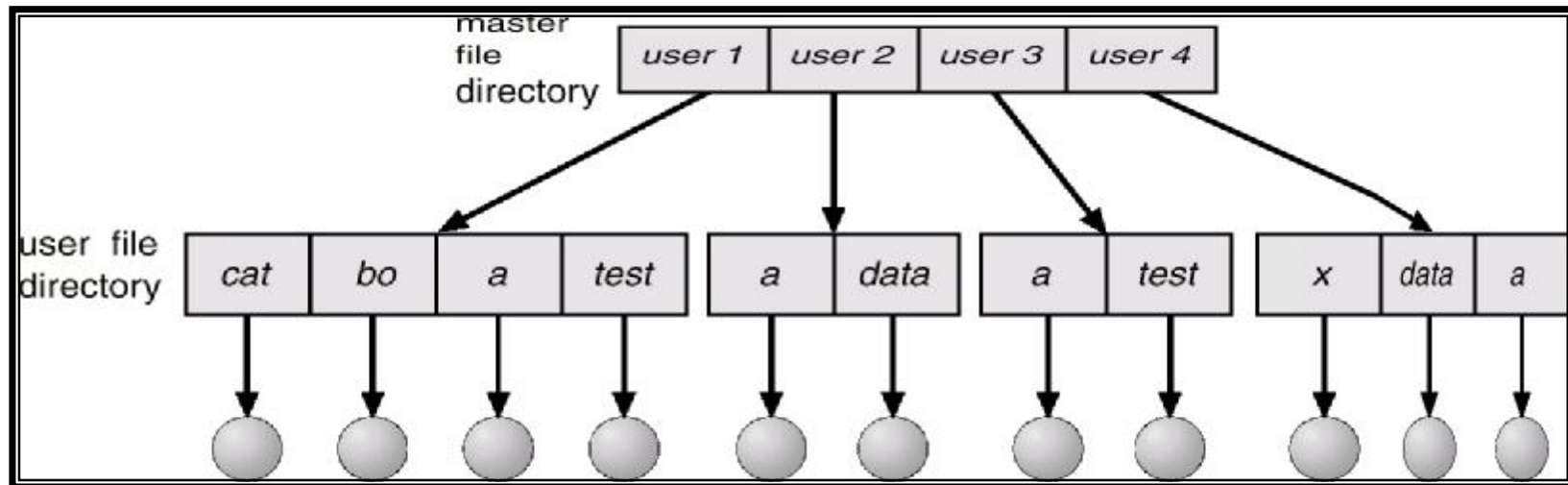


Naming problem

Grouping problem

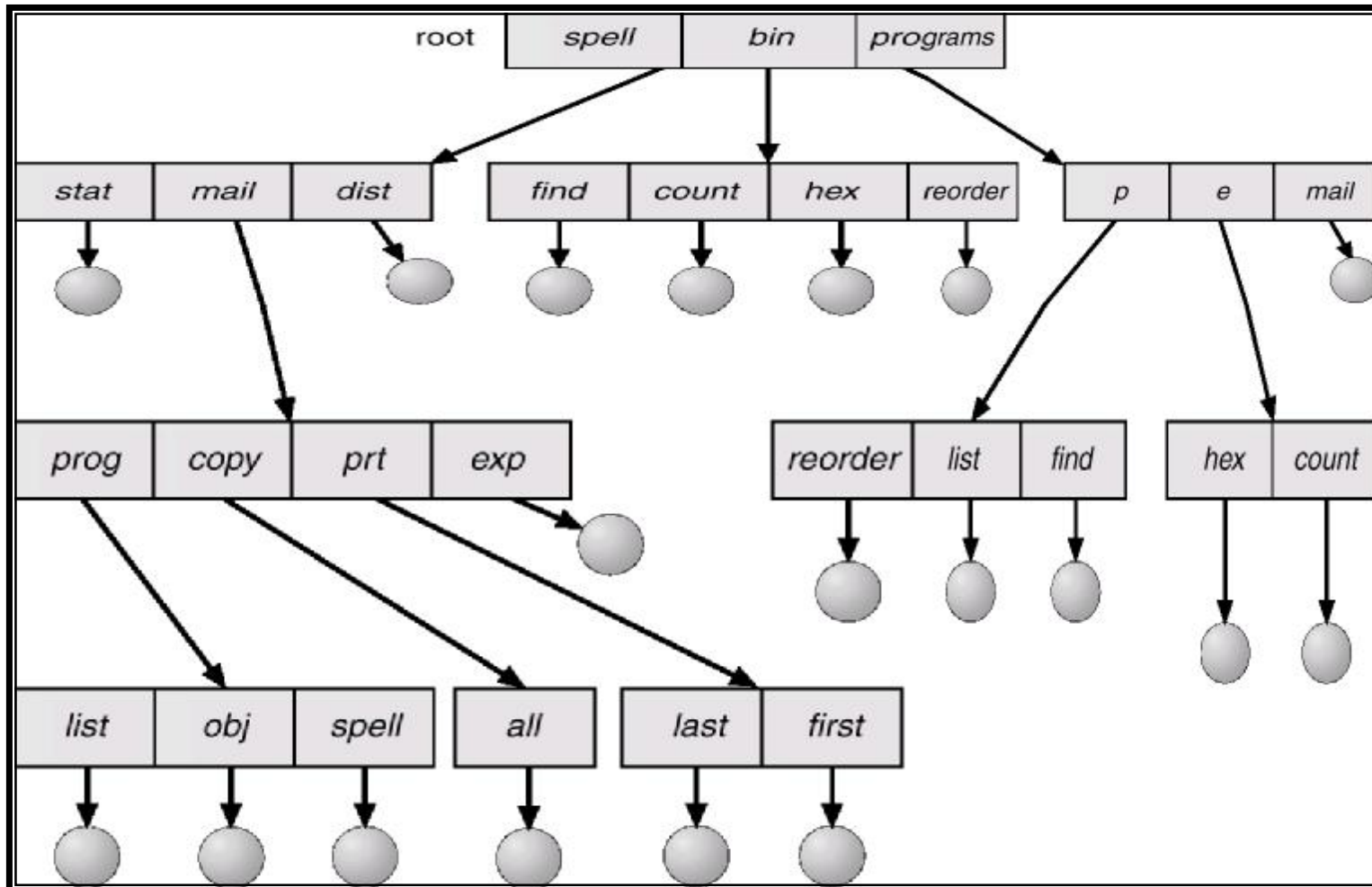
Two-Level Directory

n 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'd)

- n Efficient searching
- n Grouping Capability
- n Current directory (working directory)
 - ü **cd** /spell/mail/prog
 - ü **type** list

Tree-Structured Directories (Cont'd)

n **Absolute** or **relative** path name

n Creating a new file is done in current directory

n Delete a file

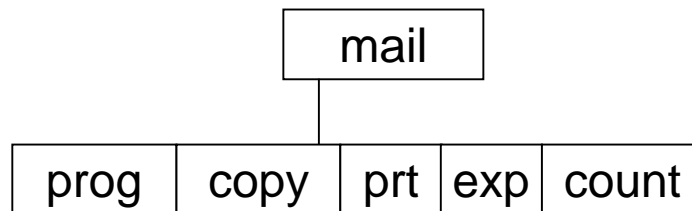
rm <file-name>

n Creating a new subdirectory is done in current directory

mkdir <dir-name>

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

mkdir count



Deleting “mail” ⇒ deleting the entire subtree rooted by “mail”

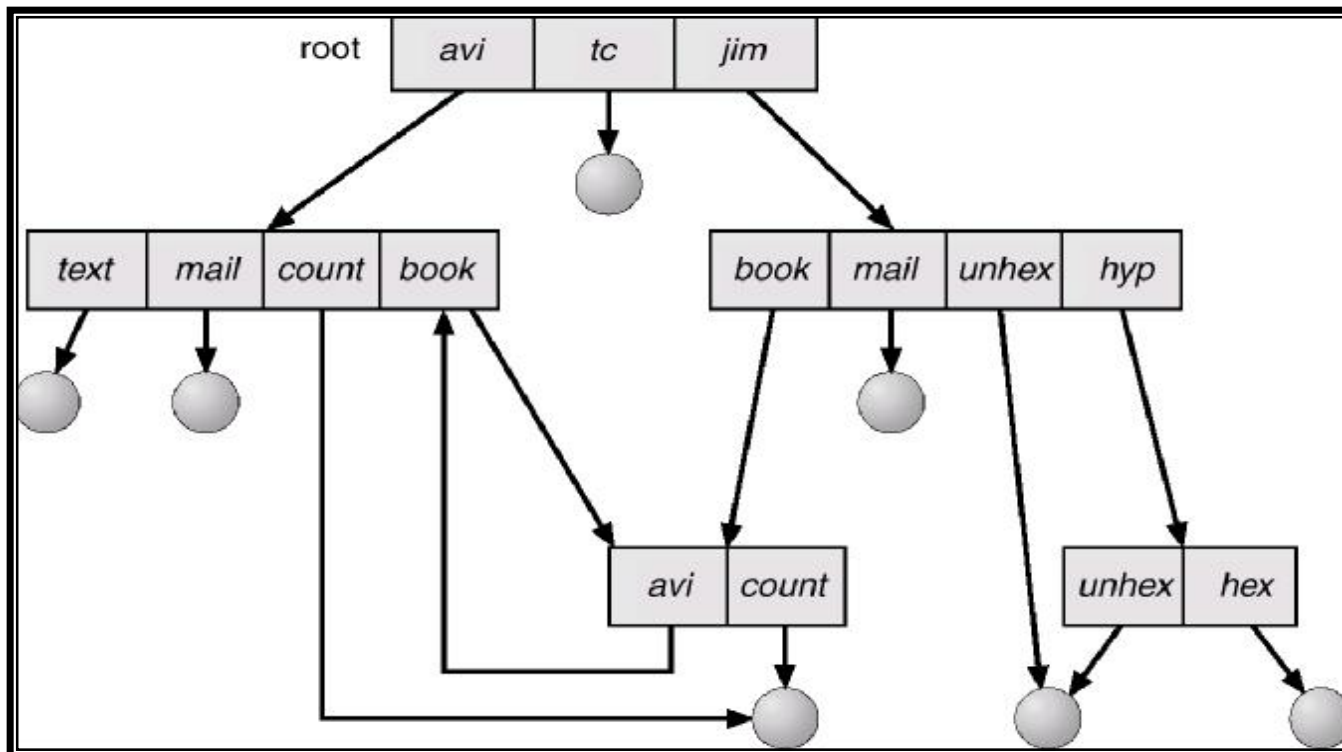
Acyclic-Graph Directories (Cont'd)

- n Two different names (aliasing)
- n If *dict* deletes *list* \Rightarrow dangling pointer

Solutions:

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

General Graph Directory



General Graph Directory (Cont'd)

n How do we guarantee no cycles?

- ü Allow only links to file not subdirectories
- ü Garbage collection
- ü Every time a new link is added use a cycle detection algorithm to determine whether it is OK

Pathname Translation

n `open("/a/b/c", ...)`

- ü Open directory "/" (well known, can always find)
- ü Search the directory for "a", get location of "a"
- ü Open directory "a", search for "b", get location of "b"
- ü Open directory "b", search for "c", get location of "c"
- ü Open file "c"
- ü (Of course, permissions are checked at each step)

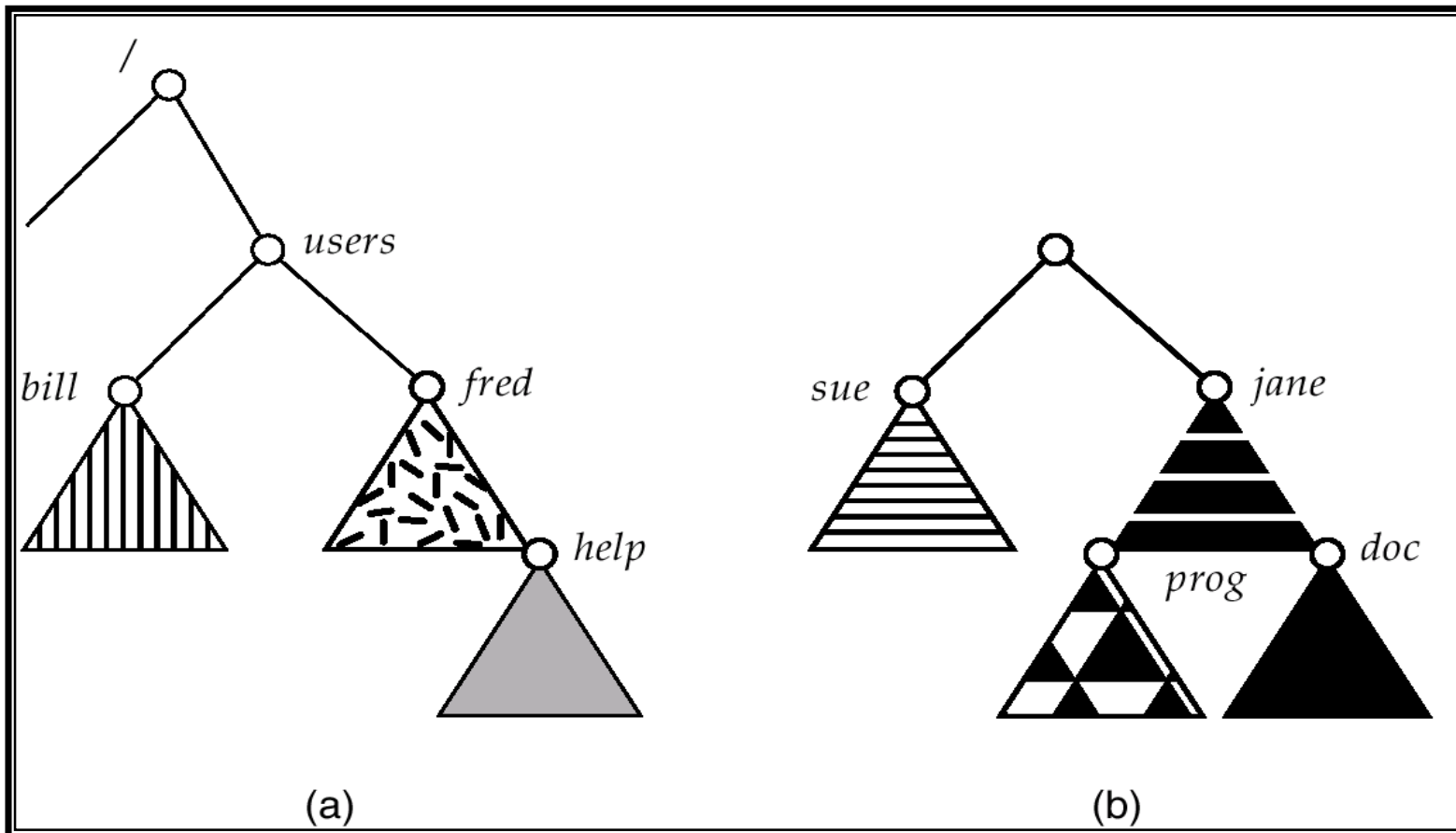
n System spends a lot of time walking down directory paths

- ü This is why open is separate from read/write
- ü OS will cache prefix lookups to enhance performance
 - § /a/b, /a/bb, /a/bbb, etc. all share the "/a" prefix

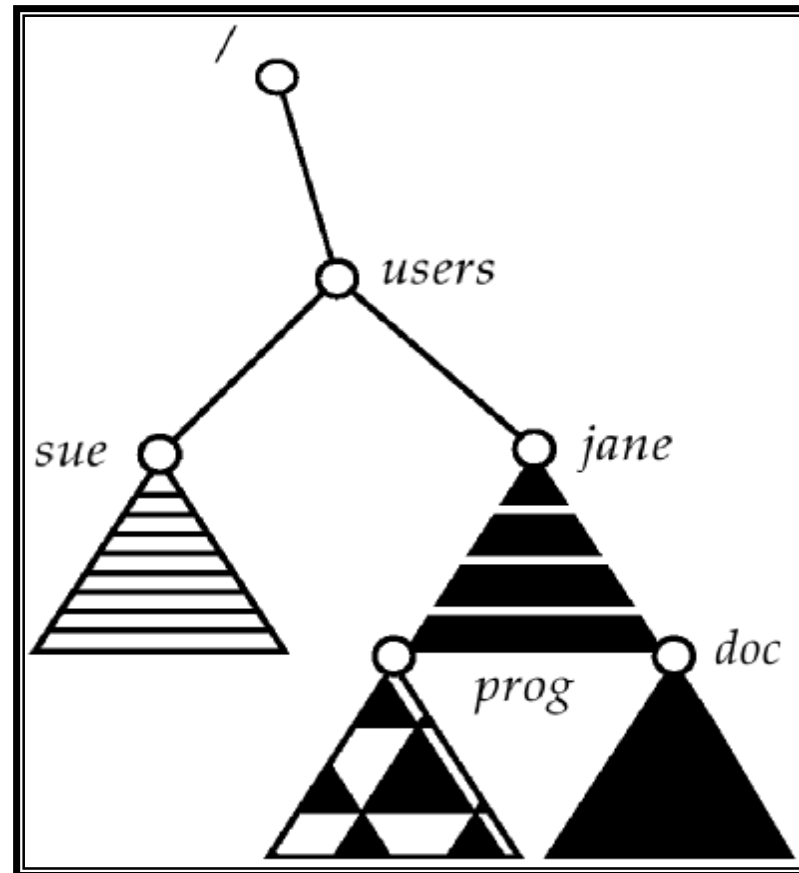
File System Mounting

- n A file system must be **mounted** before it can be accessed
- n A unmounted file system (i.e. Fig. 11-11(b)) is mounted at a **mount point**
- n *Example*
 - ü Windows: to drive letters (e.g., C:\, D:\, ...)
 - ü Unix: to an existing empty directory (= mount point)

(a) Existing (b) Unmounted Partition



Mount Point



File Sharing

- n Sharing of files on multi-user systems is desirable
- n Sharing may be done through a *protection* scheme
- n On distributed systems, files may be shared across a network
- n Network File System (NFS) is a common distributed file-sharing method

Consistency Semantics

n UNIX semantics

- ü Writes to an open file are visible immediately to other users that have this file open at the same time.
- ü One mode of sharing allows users to share the pointer of current location into the file.
 - § via `fork()` or `dup()`

n AFS session semantics

- ü Writes to an open file are not visible immediately
- ü Once a file is closed, the changes made to it are visible only in sessions starting later

n Immutable-shared-files semantics

- ü Once a file is declared as shared by its creator, it cannot be modified

Protection

n File owner/creator should be able to control:

- ü what can be done
- ü by whom

n Types of access

- ü Read
- ü Write
- ü Execute
- ü Append
- ü Delete
- ü List

n Representing protection

ü Access control lists (ACLs)

§ For each object, keep list of subjects and their allowed actions

ü Capabilities

§ For each subject, keep list of objects and their allowed actions

	objects		
	/etc/passwd	/home/hong	/home/guest
subjects	root	rw	rw
	hong	r	r
Capability	guest	-	-
	ACL		

n ACLs vs. Capabilities

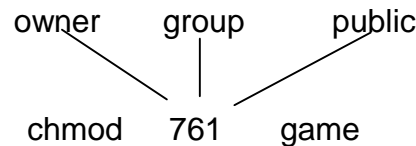
- ü Two approaches differ only in how the table is represented
- ü Capabilities are easy to transfer
 - § They are like keys; can hand them off
 - § They make sharing easy
- ü In practice, ACLs are easier to manage
 - § Object-centric, easy to grant and revoke
 - § To revoke capabilities, need to keep track of all subjects that have the capability – hard to do, given that subjects can hand off capabilities
- ü ACLs grow large when object is heavily shared
 - § Can simplify by using “groups”
 - § Additional benefit: change group membership affects all objects that have this group in its ACL

Access Lists and Groups

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

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

- n Ask manager to create a group (unique name), say G, and add some users to the group
- n For a particular file (say *game*) or subdirectory, define an appropriate access



Attach a group to a file

chgrp G game

n Advisory lock on a whole file

- ü int **flock** (int fd, int operation)

- § LOCK_SH(shared), LOCK_EX(exclusive), LOCK_UN(unlock)

n POSIX record lock

- ü discretionary lock: can lock portions of a file

- ü If a process dies, its locks are automatically removed

- ü int **fcntl** (int fd, int cmd, struct flock *lock);

- § cmd: F_GETLK, F_SETLK, F_SETLKW

- § struct flock { type, whence, start, len, pid };

n System V mandatory lock

- ü A file is marked as a candidate by setting the setgid bit and removing the group execute bit

- ü Must mount the file system to permit mandatory file locks

- ü Use the existing flock()/fcntl() to apply locks

- ü Every read() and write() is checked for locking