

*Operating
Systems:
Internals
and Design
Principles*

Chapter 12

File Management

Ninth Edition
By William Stallings

Files

- Data collections created by users
- The File System is one of the most important parts of the OS to a user
- Desirable properties of files:

Long-term existence

- Files are stored on disk or other secondary storage and do not disappear when a user logs off

Sharable between processes

- Files have names and can have associated access permissions that permit controlled sharing

Structure

- Files can be organized into hierarchical or more complex structure to reflect the relationships among files

File Systems

- Provide a means to store data organized as files as well as a collection of functions that can be performed on files
- Maintain a set of attributes associated with the file
- Typical operations include:
 - Create
 - Delete
 - Open
 - Close
 - Read
 - Write

File Structure

Four terms are commonly used when discussing files:

Field

Record

File

Database

Structure Terms

Field

- Basic element of data
- Contains a single value
- Fixed or variable length

Database

- Collection of related data
- Relationships among elements of data are explicit
- Designed for use by a number of different applications
- Consists of one or more types of files

File

- Collection of similar records
- Treated as a single entity
- May be referenced by name
- Access control restrictions usually apply at the file level

Record

- Collection of related fields that can be treated as a unit by some application program
- Fixed or variable length

File Management System

- Meet the data management needs of the user
- Guarantee that the data in the file are valid
- Optimize performance
- Provide I/O support for a variety of storage device types
- Minimize the potential for lost or destroyed data
- Provide a standardized set of I/O interface routines to user processes
- Provide I/O support for multiple users in the case of multiple-user systems

Minimal User Requirements

■ Each user:

1

- Should be able to create, delete, read, write and modify files

2

- May have controlled access to other users' files

3

- May control what type of accesses are allowed to the users' files

4

- Should be able to move data between files

5

- Should be able to back up and recover files in case of damage

6

- Should be able to access his or her files by name rather than by numeric identifier

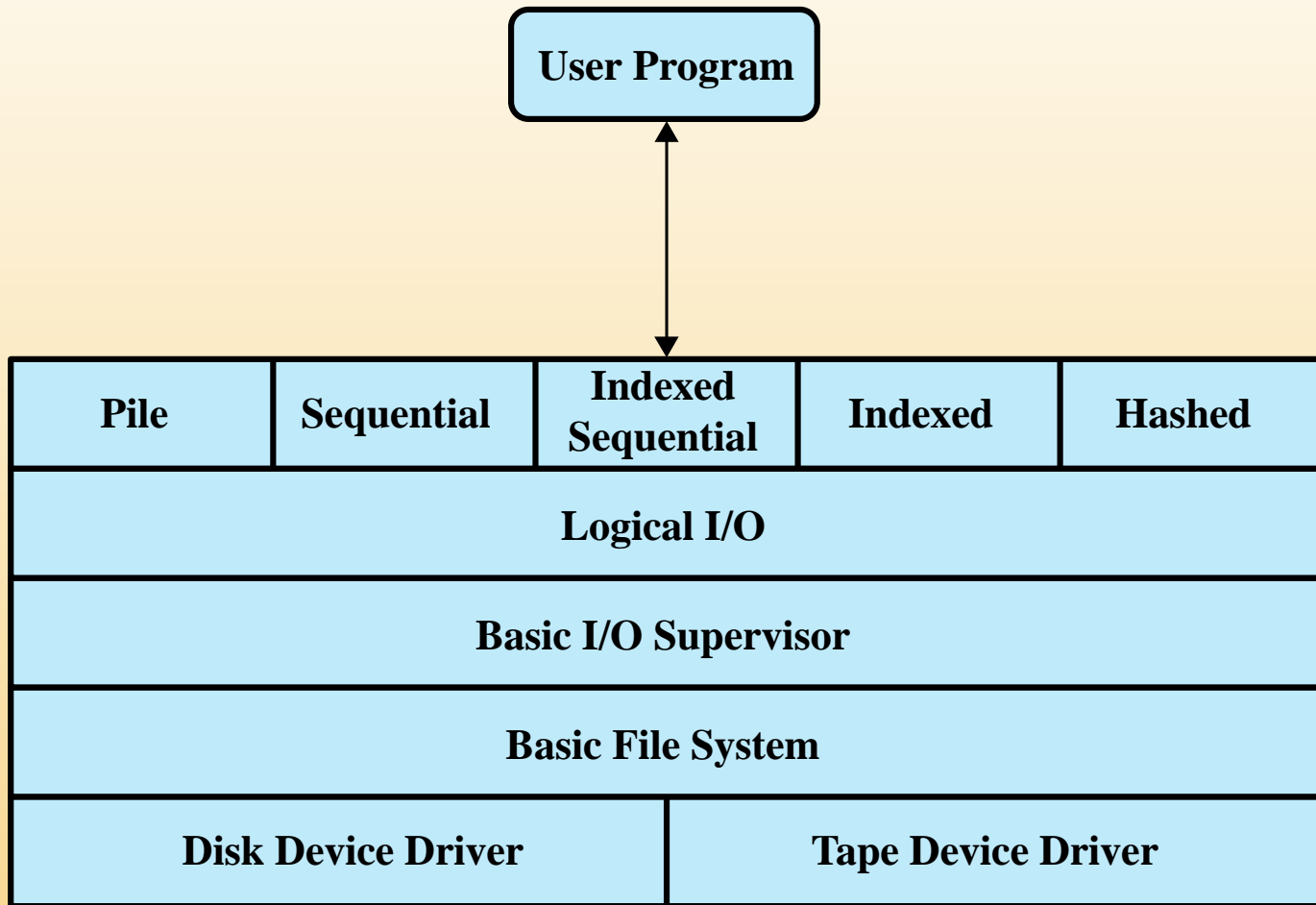


Figure 12.1 File System Software Architecture

Device Drivers

- Lowest level
- Communicates directly with peripheral devices
- Responsible for starting I/O operations on a device
- Processes the completion of an I/O request
- Considered to be part of the operating system

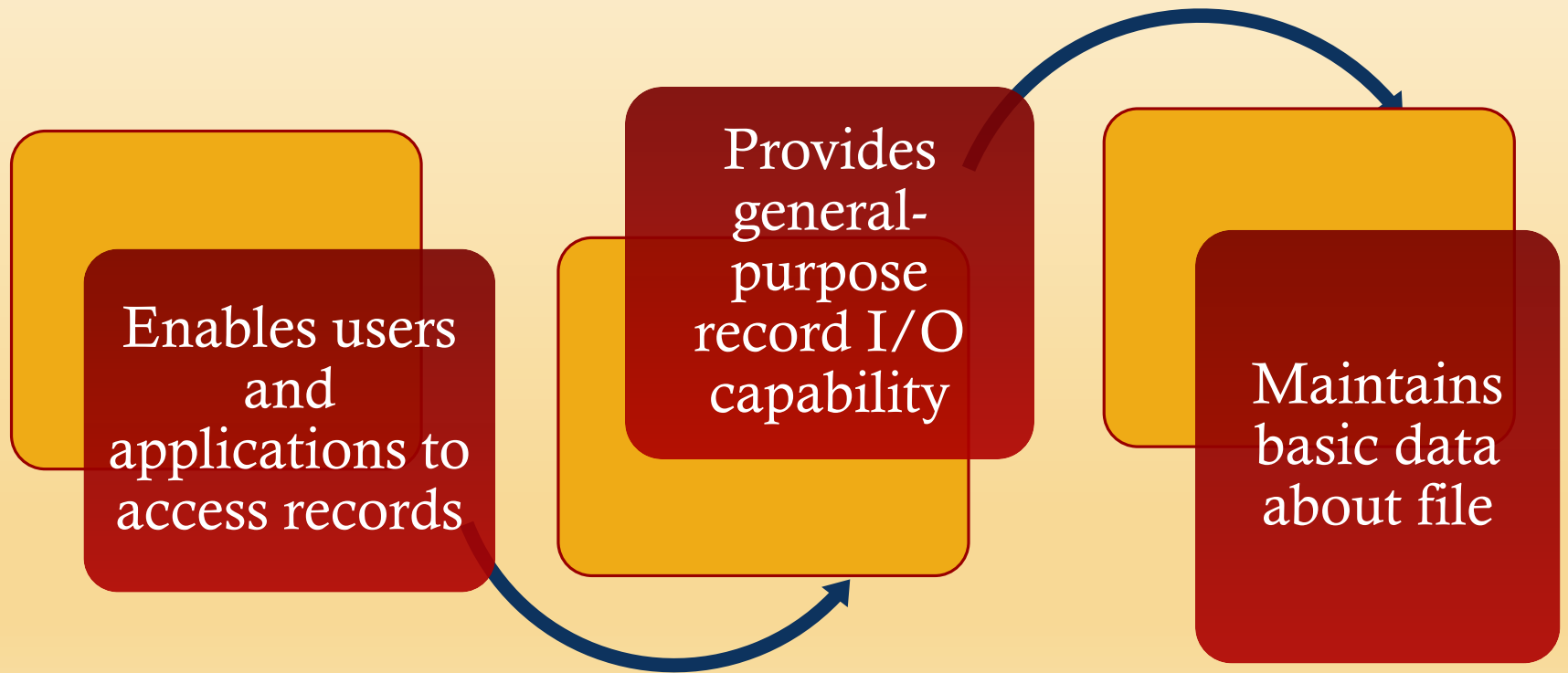
Basic File System

- Also referred to as the physical I/O level
- Primary interface with the environment outside the computer system
- Deals with blocks of data that are exchanged with disk or tape systems
- Concerned with the placement of blocks on the secondary storage device
- Concerned with buffering blocks in main memory
- Does not understand the content of the data or the structure of the files involved
- Considered part of the operating system

Basic I/O Supervisor

- Responsible for all file I/O initiation and termination
- At this level, control structures are maintained that deal with device I/O, scheduling, and file status
- Selects the device on which I/O is to be performed
- Concerned with scheduling disk and tape accesses to optimize performance
- I/O buffers are assigned and secondary memory is allocated at this level
- Part of the operating system

Logical I/O



Access Method

- Level of the file system closest to the user
- Provides a standard interface between applications and the file systems and devices that hold the data
- Different access methods reflect different file structures and different ways of accessing and processing the data

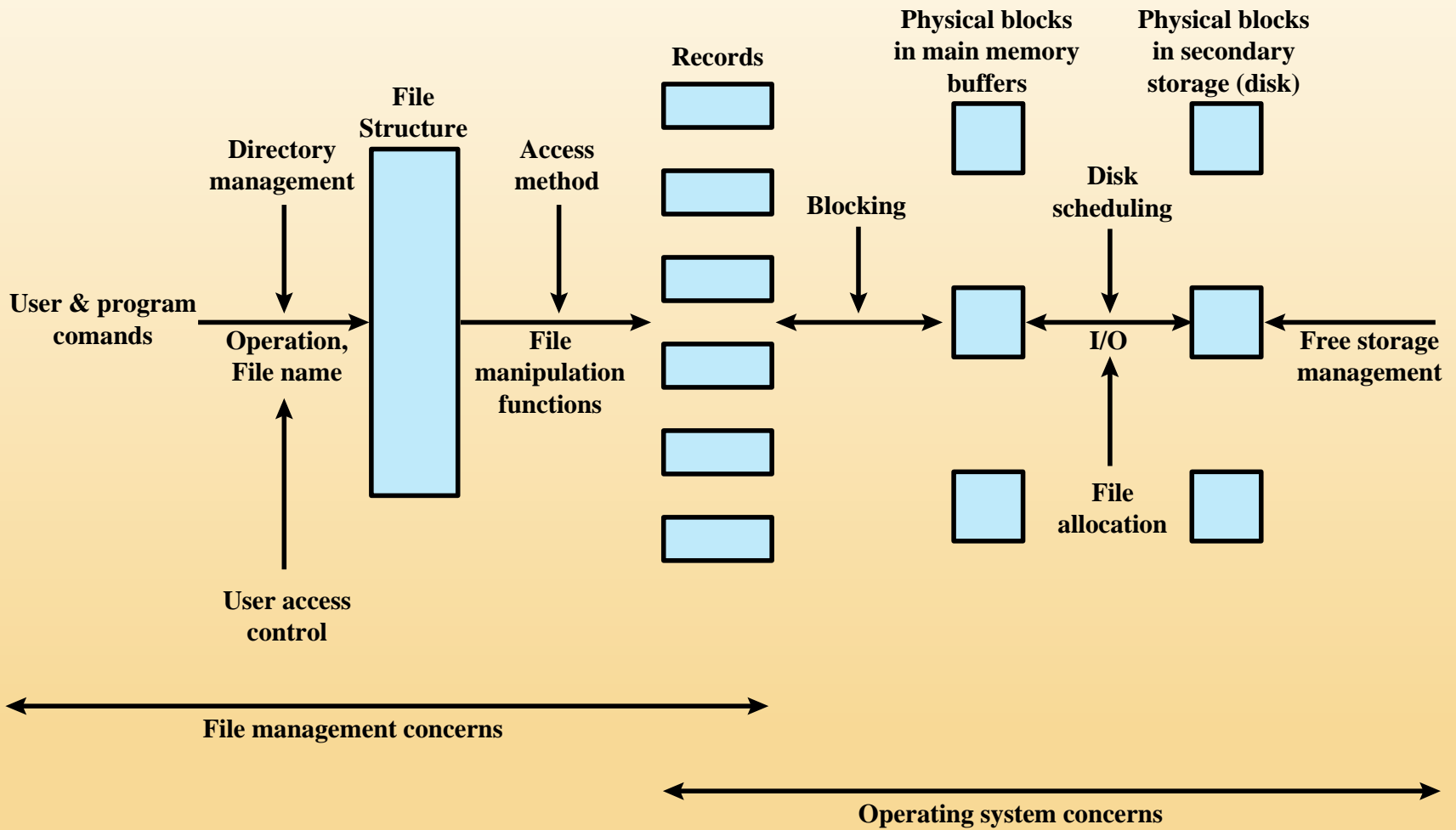
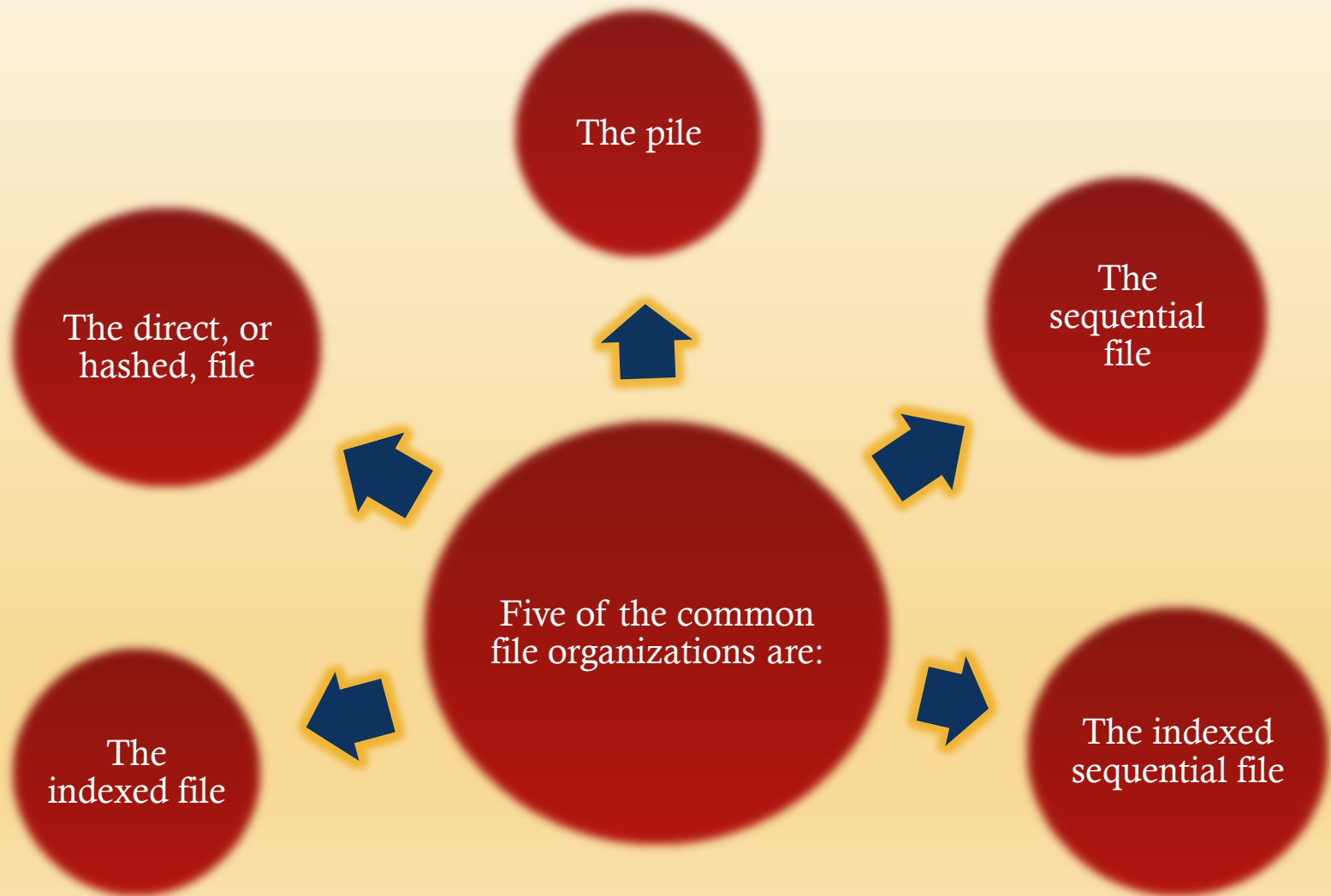


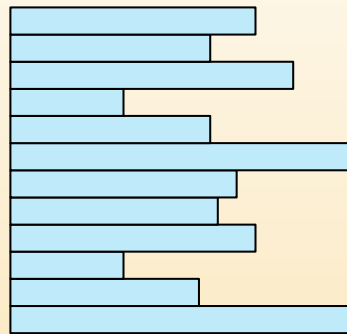
Figure 12.2 Elements of File Management

File Organization and Access

- *File organization* is the logical structuring of the records as determined by the way in which they are accessed
- In choosing a file organization, several criteria are important:
 - Short access time
 - Ease of update
 - Economy of storage
 - Simple maintenance
 - Reliability
- Priority of criteria depends on the application that will use the file

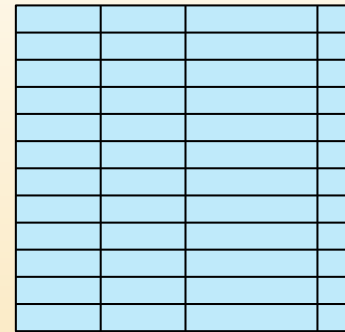
File Organization Types





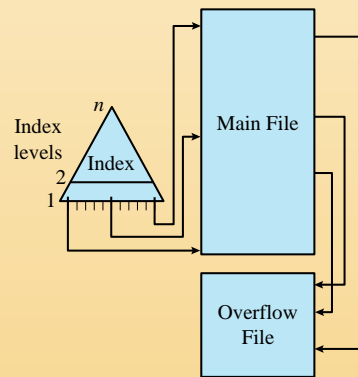
Variable-length records
 Variable set of fields
 Chronological order

(a) Pile File

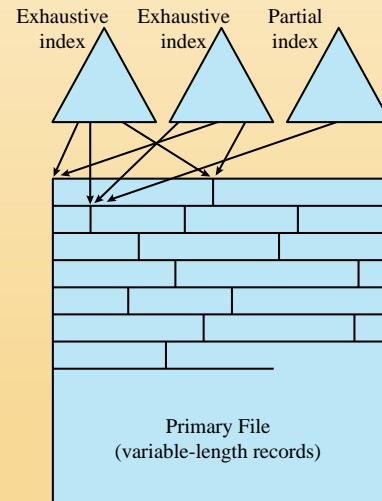


Fixed-length records
 Fixed set of fields in fixed order
 Sequential order based on key field

(b) Sequential File



(c) Indexed Sequential File

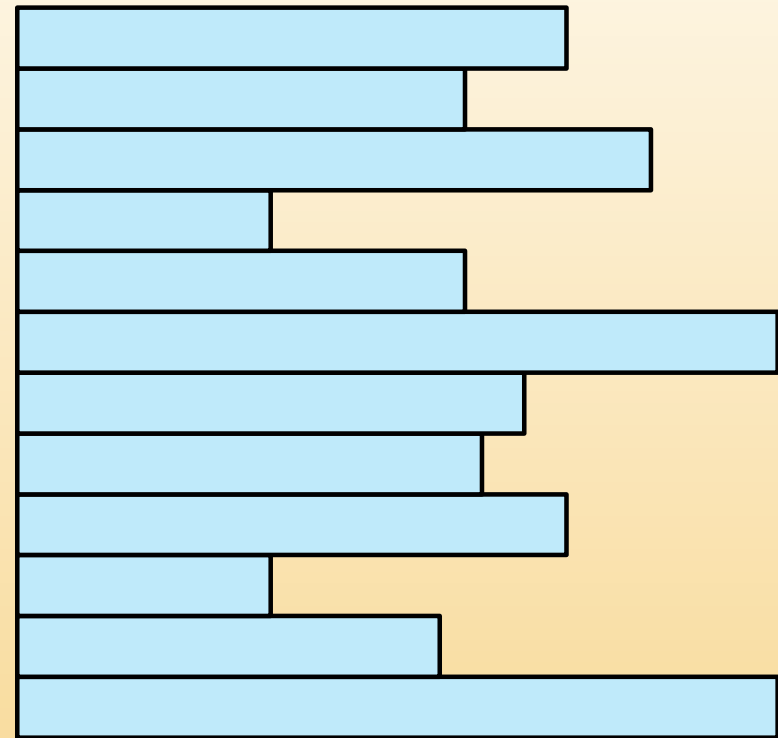


(d) Indexed File

Figure 12.3 Common File Organizations

The Pile

- Least complicated form of file organization
- Data are collected in the order they arrive
- Each record consists of one burst of data
- Purpose is simply to accumulate the mass of data and save it
- Record access is by exhaustive search



Variable-length records
Variable set of fields
Chronological order

(a) Pile File

The Sequential File

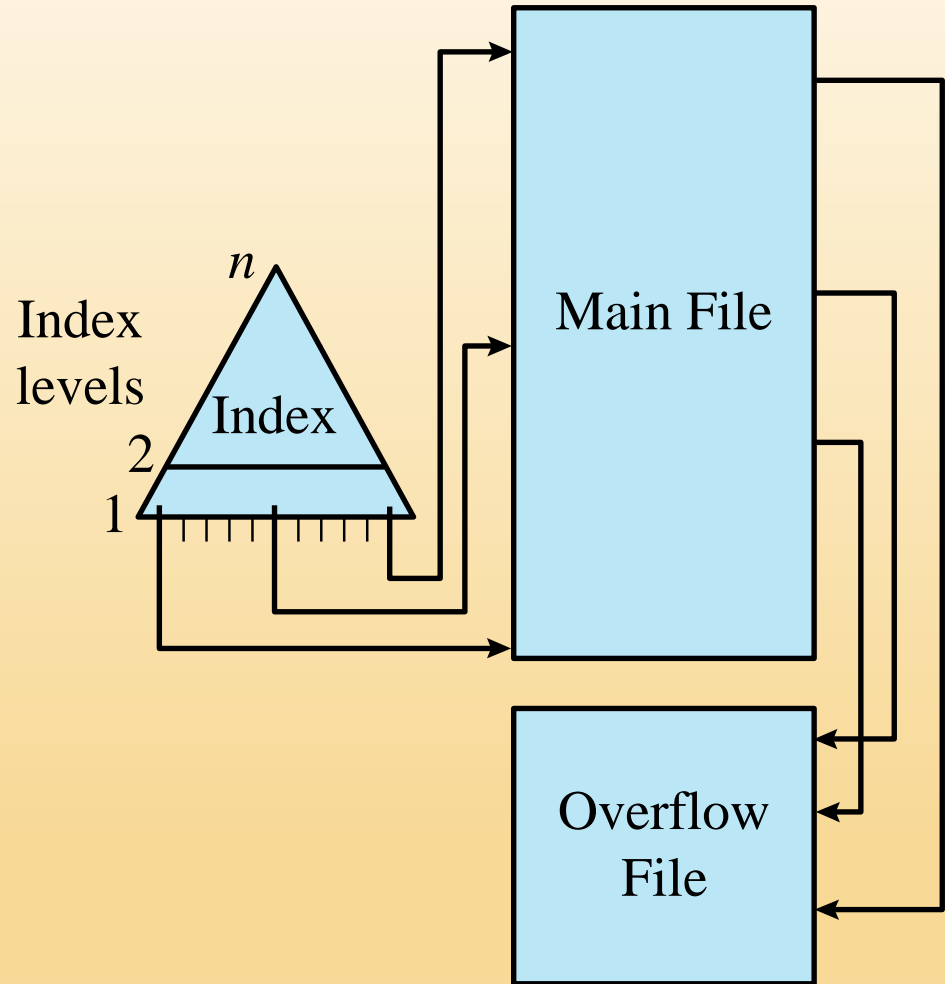
- Most common form of file structure
- A fixed format is used for records
- Key field uniquely identifies the record
- Typically used in batch applications
- Only organization that is easily stored on tape as well as disk

Fixed-length records
Fixed set of fields in fixed order
Sequential order based on key field

(b) Sequential File

Indexed Sequential File

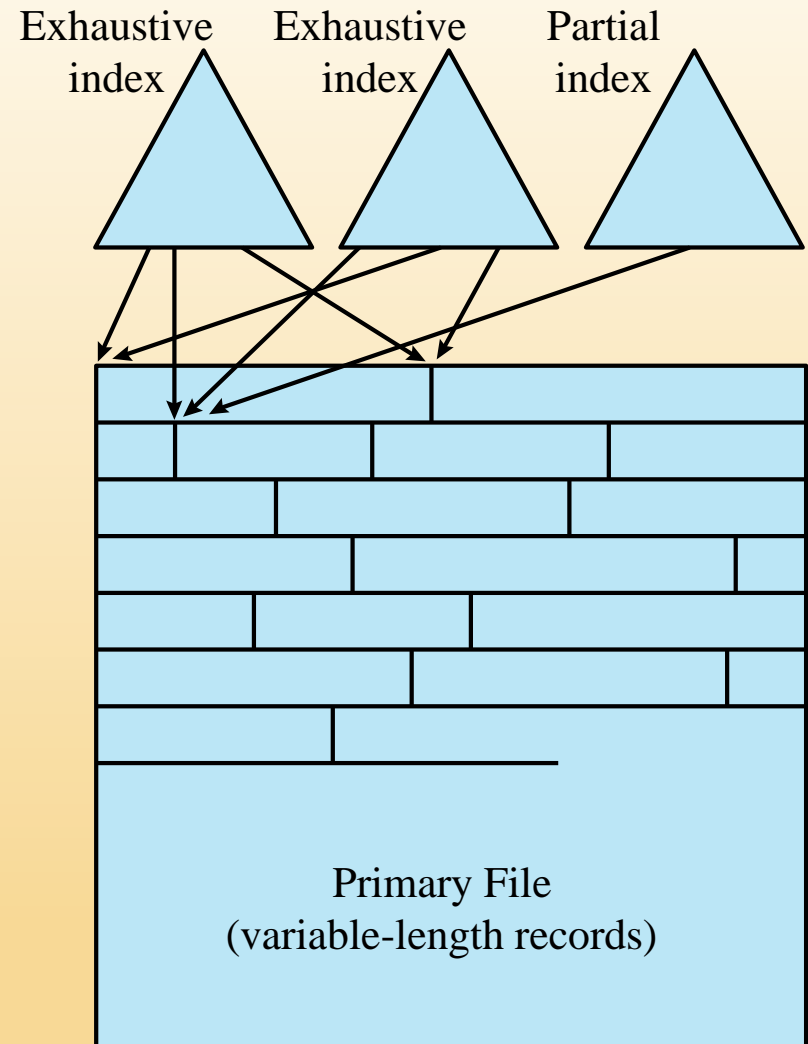
- Adds an index to the file to support random access
- Adds an overflow file
- Greatly reduces the time required to access a single record
- Multiple levels of indexing can be used to provide greater efficiency in access



(c) Indexed Sequential File

Indexed File

- Records are accessed only through their indexes
- Variable-length records can be employed
- Exhaustive index contains one entry for every record in the main file
- Partial index contains entries to records where the field of interest exists
- Used mostly in applications where timeliness of information is critical
- Examples would be airline reservation systems and inventory control systems



(d) Indexed File

Direct or Hashed File

- Access directly any block of a known address
- Makes use of hashing on the key value
- Often used where:
 - Very rapid access is required
 - Fixed-length records are used
 - Records are always accessed one at a time

Examples are:

- Directories
- Pricing tables
- Schedules
- Name lists

B-Trees

- A balanced tree structure with all branches of equal length
- Standard method of organizing indexes for databases
- Commonly used in OS file systems
- Provides for efficient searching, adding, and deleting of items

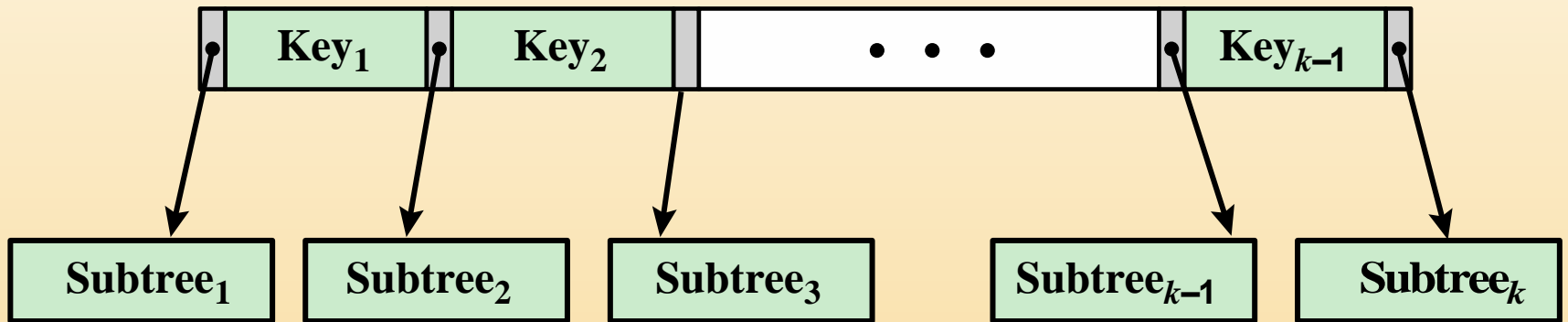


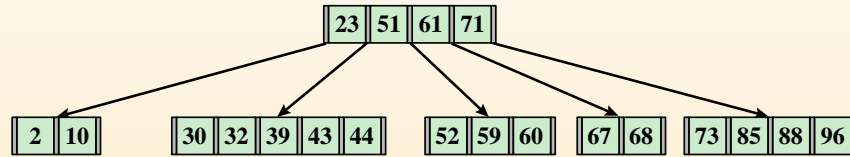
Figure 12.4 A B-tree Node with k Children

B-Tree

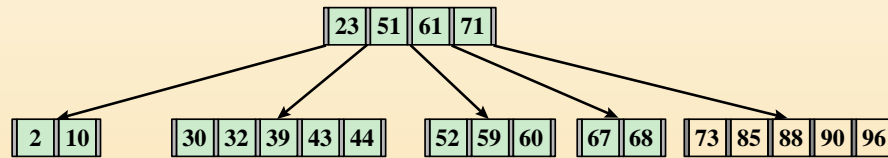
Characteristics

A B-tree is characterized by its minimum degree d and satisfies the following properties:

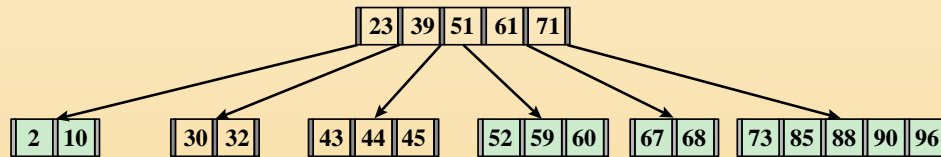
- Every node has at most $2d - 1$ keys and $2d$ children or, equivalently, $2d$ pointers
- Every node, except for the root, has at least $d - 1$ keys and d pointers, as a result, each internal node, except the root, is at least half full and has at least d children
- The root has at least 1 key and 2 children
- All leaves appear on the same level and contain no information. This is a logical construct to terminate the tree; the actual implementation may differ
- A nonleaf node with k pointers contains $k - 1$ keys



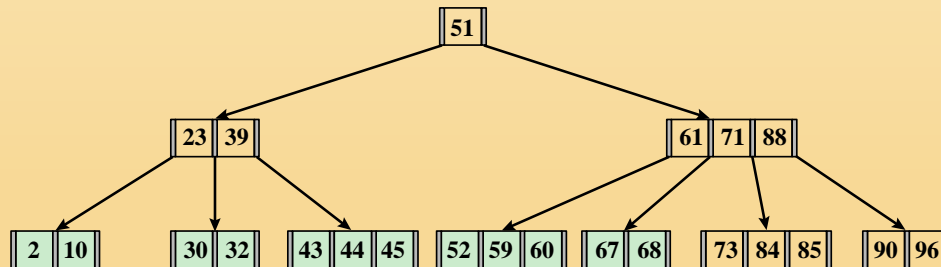
(a) B-tree of minimum degree $d = 3$.



(b) Key = 90 inserted. This is a simple insertion into a node.



(c) Key = 45 inserted. This requires splitting a node into two parts and promoting one key to the root node.



(d) Key = 84 inserted. This requires splitting a node into two parts and promoting one key to the root node. This then requires the root node to be split and a new root created.

Figure 12.5 Inserting Nodes into a B-tree

Table 12.1

Information Elements of a File Directory

(Table can be found on page 537 in textbook)

Basic Information	
File Name	Name as chosen by creator (user or program). Must be unique within a specific directory.
File Type	For example: text, binary, load module, etc.
File Organization	For systems that support different organizations
Address Information	
Volume	Indicates device on which file is stored
Starting Address	Starting physical address on secondary storage (e.g., cylinder, track, and block number on disk)
Size Used	Current size of the file in bytes, words, or blocks
Size Allocated	The maximum size of the file
Access Control Information	
Owner	User who is assigned control of this file. The owner may be able to grant/deny access to other users and to change these privileges.
Access Information	A simple version of this element would include the user's name and password for each authorized user.
Permitted Actions	Controls reading, writing, executing, transmitting over a network
Usage Information	
Date Created	When file was first placed in directory
Identity of Creator	Usually but not necessarily the current owner
Date Last Read Access	Date of the last time a record was read
Identity of Last Reader	User who did the reading
Date Last Modified	Date of the last update, insertion, or deletion
Identity of Last Modifier	User who did the modifying
Date of Last Backup	Date of the last time the file was backed up on another storage medium
Current Usage	Information about current activity on the file, such as process or processes that have the file open, whether it is locked by a process, and whether the file has been updated in main memory but not yet on disk

Operations Performed on a Directory

- To understand the requirements for a file structure, it is helpful to consider the types of operations that may be performed on the directory:



Two-Level Scheme

There is one directory for each user and a master directory

Master directory has an entry for each user directory providing address and access control information

Each user directory is a simple list of the files of that user

Names must be unique only within the collection of files of a single user

File system can easily enforce access restriction on directories

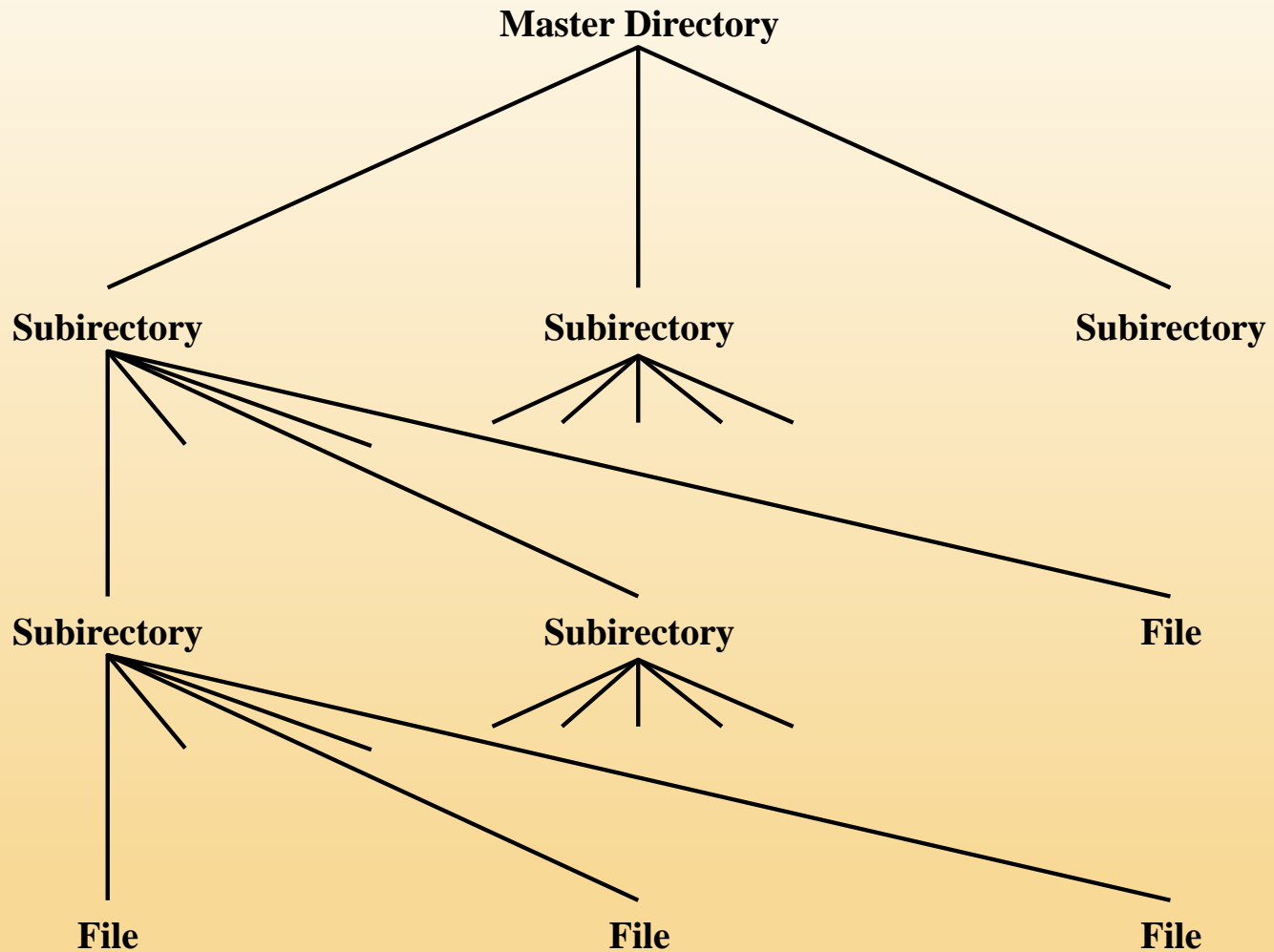


Figure 12.6 Tree-Structured Directory

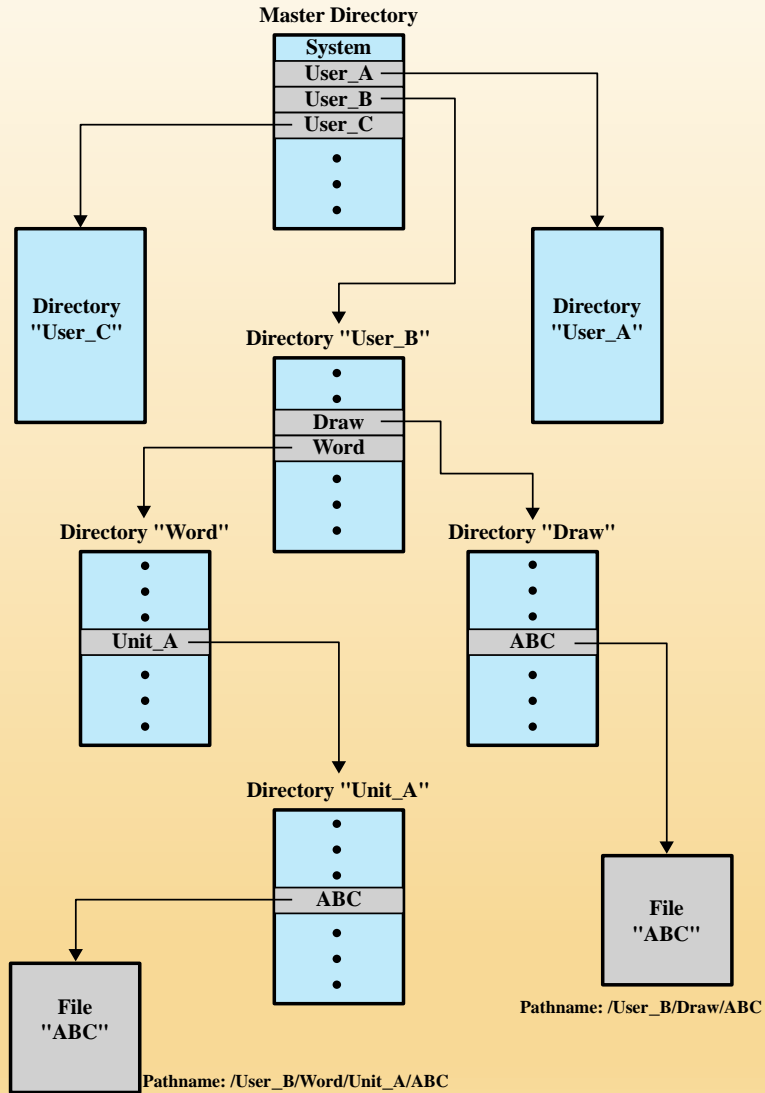


Figure 12.7 Example of Tree-Structured Directory

File Sharing

Two issues arise
when allowing files
to be shared among
a number of users:

Access rights

Management of
simultaneous
access

Access Rights

- *None*

- The user would not be allowed to read the user directory that includes the file

- *Knowledge*

- The user can determine that the file exists and who its owner is and can then petition the owner for additional access rights

- *Execution*

- The user can load and execute a program but cannot copy it

- *Reading*

- The user can read the file for any purpose, including copying and execution

- *Appending*

- The user can add data to the file but cannot modify or delete any of the file's contents

- *Updating*

- The user can modify, delete, and add to the file's data

- *Changing protection*

- The user can change the access rights granted to other users

- *Deletion*

- The user can delete the file from the file system

User Access Rights

Owner

Usually the initial creator of the file

Has full rights

May grant rights to others

Specific Users

Individual users who are designated by user ID

User Groups

A set of users who are not individually defined

All

All users who have access to this system

These are public files

Record Blocking

- Blocks are the unit of I/O with secondary storage
 - For I/O to be performed records must be organized as blocks
- Given the size of a block, three methods of blocking can be used:

1) Fixed-Length Blocking – fixed-length records are used, and an integral number of records are stored in a block

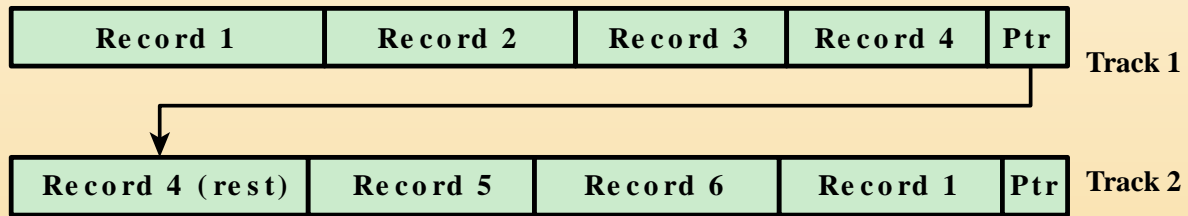
Internal fragmentation – unused space at the end of each block

2) Variable-Length Spanned Blocking – variable-length records are used and are packed into blocks with no unused space

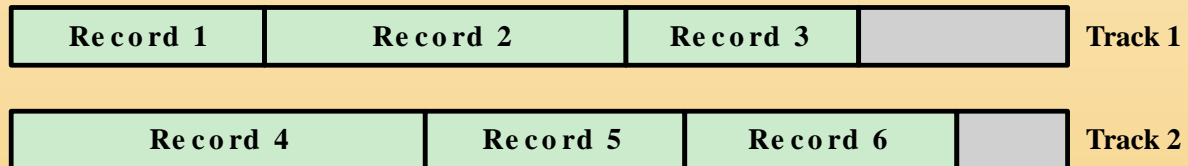
3) Variable-Length Unspanned Blocking – variable-length records are used, but spanning is not employed



(a) Fixed Blocking



(b) Variable Blocking: Spanned



(c) Variable Blocking: Unspanned

Figure 12.8 Record Blocking Methods

File Allocation

- On secondary storage, a file consists of a collection of blocks
- The operating system or file management system is responsible for allocating blocks to files
- The approach taken for file allocation may influence the approach taken for free space management
- Space is allocated to a file as one or more *portions* (contiguous set of allocated blocks)
- *File allocation table (FAT)*
 - Data structure used to keep track of the portions assigned to a file

Preallocation vs Dynamic Allocation

- A preallocation policy requires that the maximum size of a file be declared at the time of the file creation request
- For many applications it is difficult to estimate reliably the maximum potential size of the file
 - Tends to be wasteful because users and application programmers tend to overestimate size
- Dynamic allocation allocates space to a file in portions as needed

Portion Size

- In choosing a portion size there is a trade-off between efficiency from the point of view of a single file versus overall system efficiency
- Items to be considered:
 - 1) Contiguity of space increases performance, especially for `Retrieve_Next` operations, and greatly for transactions running in a transaction-oriented operating system
 - 2) Having a large number of small portions increases the size of tables needed to manage the allocation information
 - 3) Having fixed-size portions simplifies the reallocation of space
 - 4) Having variable-size or small fixed-size portions minimizes waste of unused storage due to overallocation

Alternatives

Two major alternatives:

Variable, large contiguous portions

- Provides better performance
- The variable size avoids waste
- The file allocation tables are small

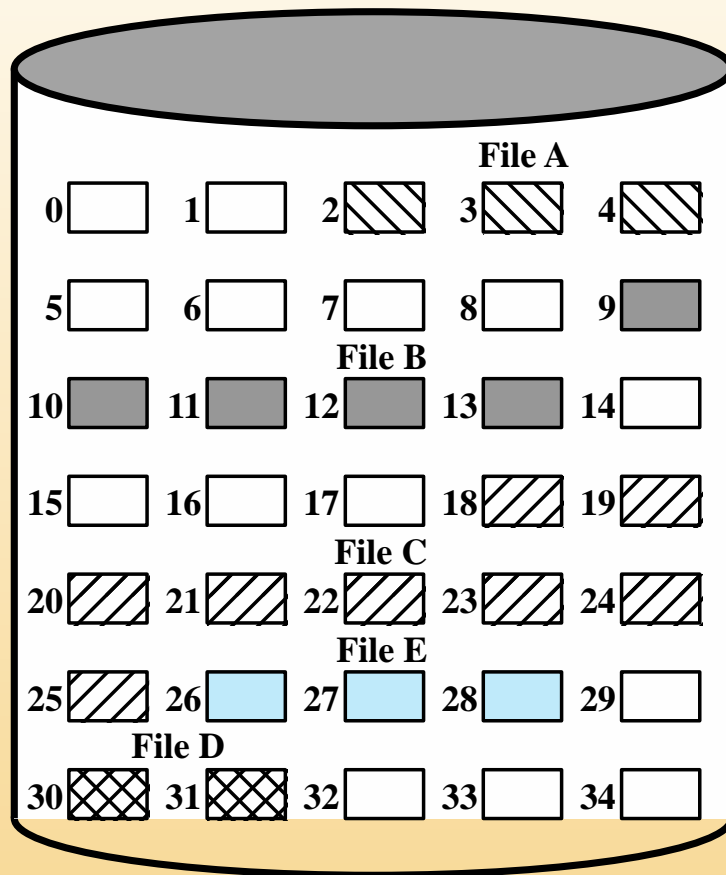
Blocks

- Small fixed portions provide greater flexibility
- They may require large tables or complex structures for their allocation
- Contiguity has been abandoned as a primary goal
- Blocks are allocated as needed

Table 12.2

File Allocation Methods

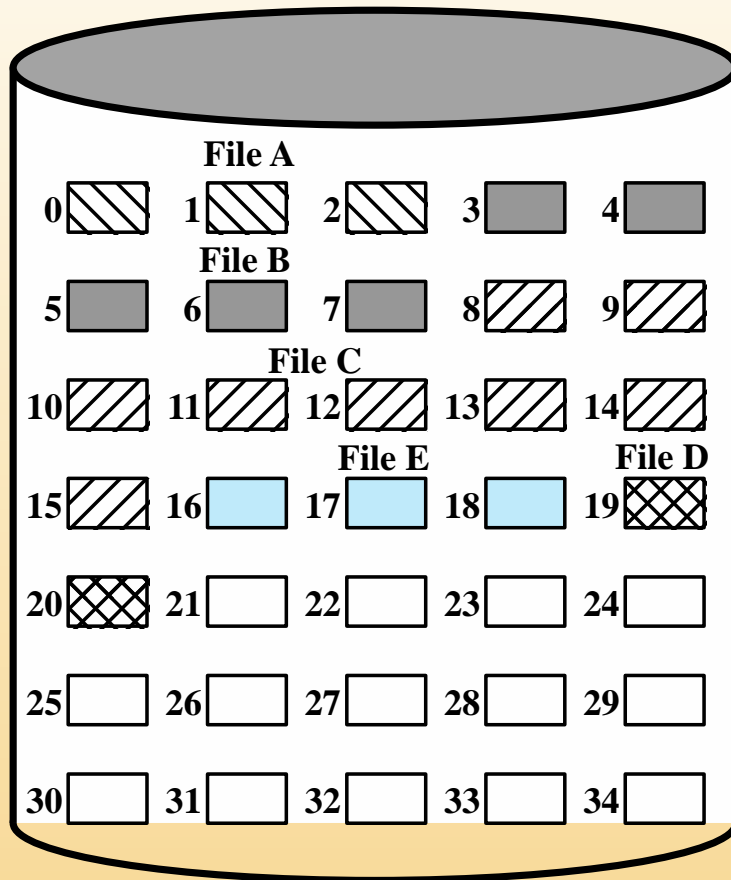
	Contiguous	Chained	Indexed	
Preallocation?	Necessary	Possible	Possible	
Fixed or variable size portions?	Variable	Fixed blocks	Fixed blocks	Variable
Portion size	Large	Small	Small	Medium
Allocation frequency	Once	Low to high	High	Low
Time to allocate	Medium	Long	Short	Medium
File allocation table size	One entry	One entry	Large	Medium



File Allocation Table

File Name	Start Block	Length
File A	2	3
File B	9	5
File C	18	8
File D	30	2
File E	26	3

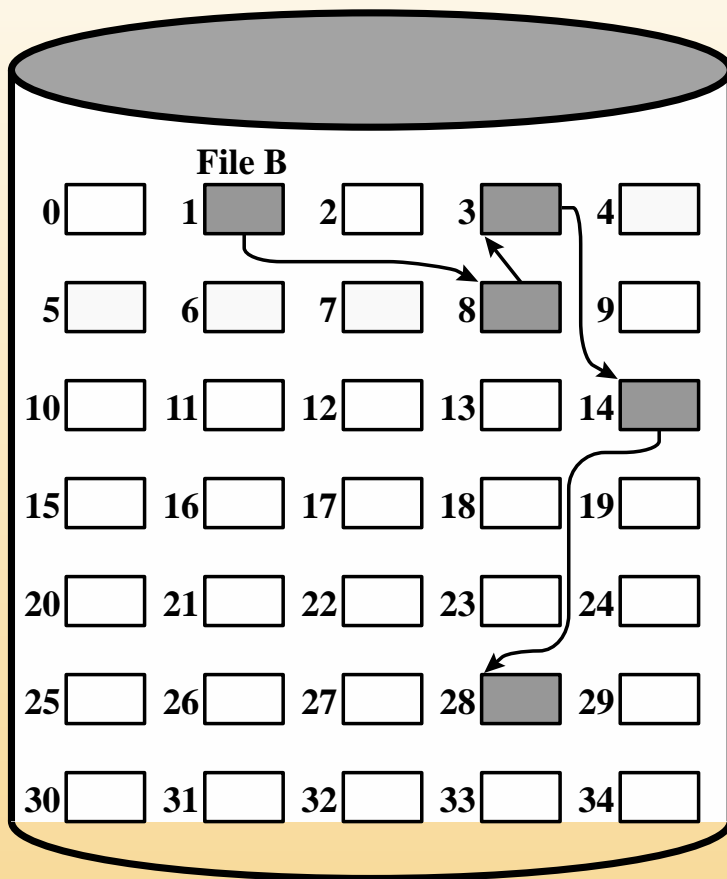
Figure 12.9 Contiguous File Allocation



File Allocation Table

File Name	Start Block	Length
File A	0	3
File B	3	5
File C	8	8
File D	19	2
File E	16	3

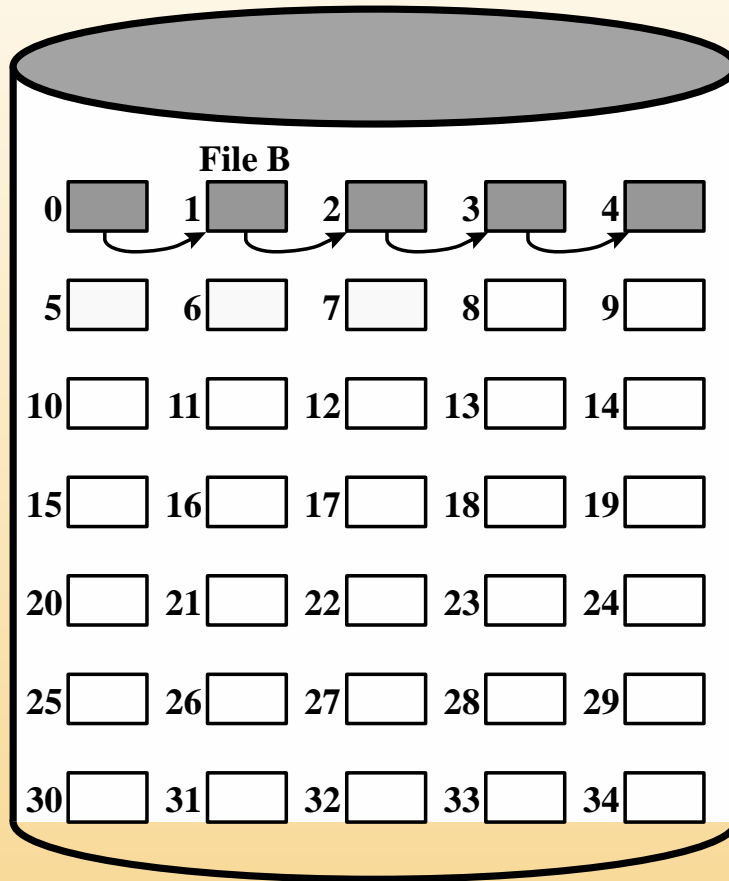
Figure 12.10 Contiguous File Allocation (After Compaction)



File Allocation Table

File Name	Start Block	Length
...
File B	1	5
...

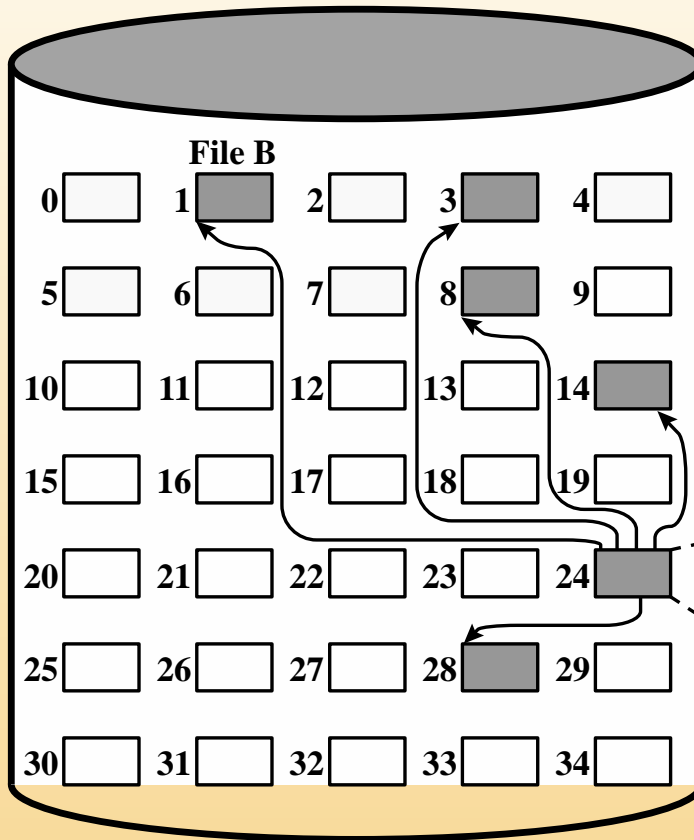
Figure 12.11 Chained Allocation



File Allocation Table

File Name	Start Block	Length
...
File B	0	5
...

Figure 12.12 Chained Allocation (After Consolidation)

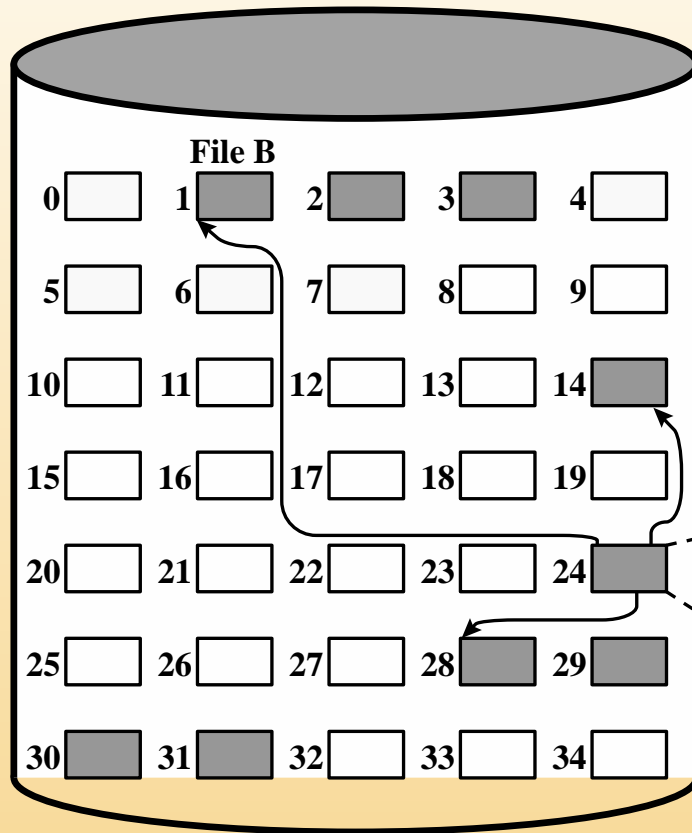


File Allocation Table

File Name	Index Block
...	...
File B	24
...	...

- 1
 - 8
 - 3
 - 14
 - 28

Figure 12.13 Indexed Allocation with Block Portions



File Allocation Table

File Name	Index Block
...	...
File B	24
...	...

Start Block	Length
1	3
28	4
14	1

Figure 12.14 Indexed Allocation with Variable-Length Portions

Free Space Management

- Just as allocated space must be managed, so must the unallocated space
- To perform file allocation, it is necessary to know which blocks are available
- A *disk allocation table* is needed in addition to a file allocation table

Bit Tables

- This method uses a vector containing one bit for each block on the disk
- Each entry of a 0 corresponds to a free block, and each 1 corresponds to a block in use

Advantages:

- Works well with any file allocation method
- It is as small as possible

Chained Free Portions

- The free portions may be chained together by using a pointer and length value in each free portion
- Negligible space overhead because there is no need for a disk allocation table
- Suited to all file allocation methods

Disadvantages:

- Leads to fragmentation
- Every time you allocate a block you need to read the block first to recover the pointer to the new first free block before writing data to that block

Indexing

- Treats free space as a file and uses an index table as it would for file allocation
- For efficiency, the index should be on the basis of variable-size portions rather than blocks
- This approach provides efficient support for all of the file allocation methods

Free Block List

Each block is assigned a number sequentially

The list of the numbers of all free blocks is maintained in a reserved portion of the disk

Depending on the size of the disk, either 24 or 32 bits will be needed to store a single block number

The size of the free block list is 24 or 32 times the size of the corresponding bit table and must be stored on disk

There are two effective techniques for storing a small part of the free block list in main memory:

The list can be treated as a push-down stack with the first few thousand elements of the stack kept in main memory

The list can be treated as a FIFO queue, with a few thousand entries from both the head and the tail of the queue in main memory

Volumes

- A collection of addressable sectors in secondary memory that an OS or application can use for data storage
- The sectors in a volume need not be consecutive on a physical storage device
 - They need only appear that way to the OS or application
- A volume may be the result of assembling and merging smaller volumes

UNIX File Management

- In the UNIX file system, six types of files are distinguished:

Regular, or ordinary

- Contains arbitrary data in zero or more data blocks

Directory

- Contains a list of file names plus pointers to associated inodes (index nodes)

Special

- Contains no data but provides a mechanism to map physical devices to file names

Named pipes

- An interprocess communications facility

Links

- An alternative file name for an existing file

Symbolic links

- A data file that contains the name of the file to which it is linked

Inodes

- All types of UNIX files are administered by the OS by means of inodes
- An inode (index node) is a control structure that contains the key information needed by the operating system for a particular file
- Several file names may be associated with a single inode
 - An active inode is associated with exactly one file
 - Each file is controlled by exactly one inode

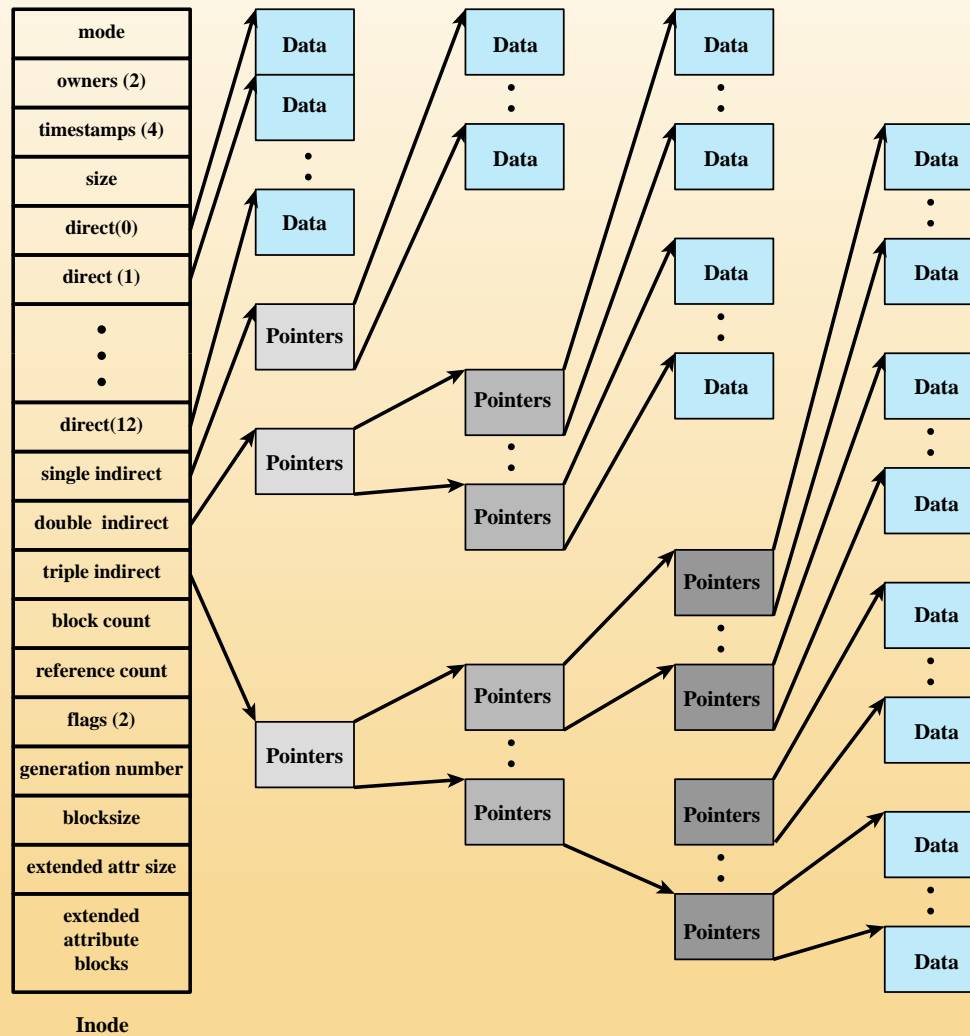


Figure 12.15 Structure of FreeBSD inode and File

File Allocation

- File allocation is done on a block basis
- Allocation is dynamic, as needed, rather than using preallocation
- An indexed method is used to keep track of each file, with part of the index stored in the inode for the file
- In all UNIX implementations the inode includes a number of direct pointers and three indirect pointers (single, double, triple)

Table 12.3

Capacity of a FreeBSD File with 4-Kbyte Block Size

Level	Number of Blocks	Number of Bytes
Direct	12	48K
Single Indirect	512	2M
Double Indirect	$512 \times 512 = 256K$	1G
Triple Indirect	$512 \times 256K = 128M$	512G

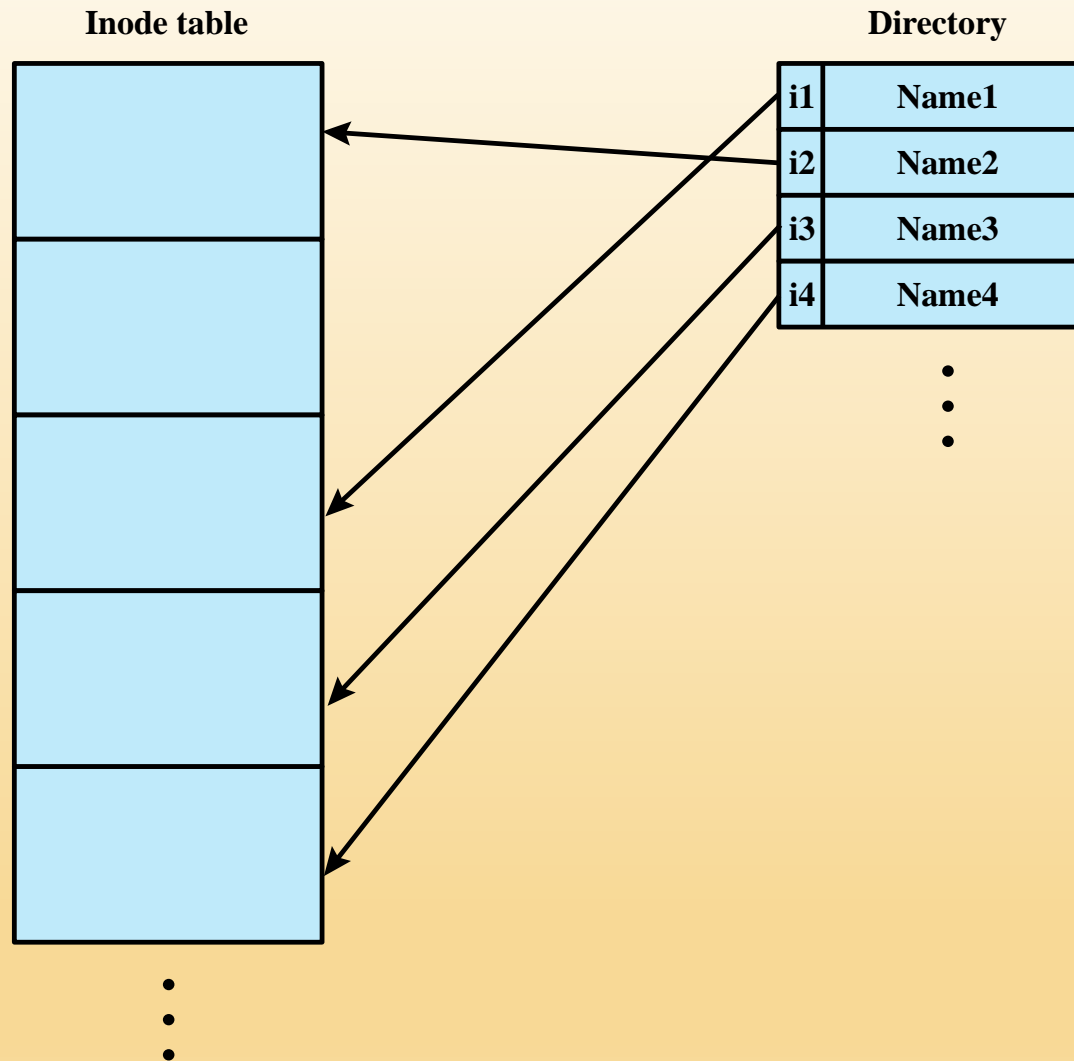


Figure 12.16 UNIX Directories and Inodes

Volume Structure

- A UNIX file system resides on a single logical disk or disk partition and is laid out with the following elements:

Boot block

Contains code required to boot the operating system

Superblock

Contains attributes and information about the file system

Inode table

Collection of inodes for each file

Data blocks

Storage space available for data files and subdirectories

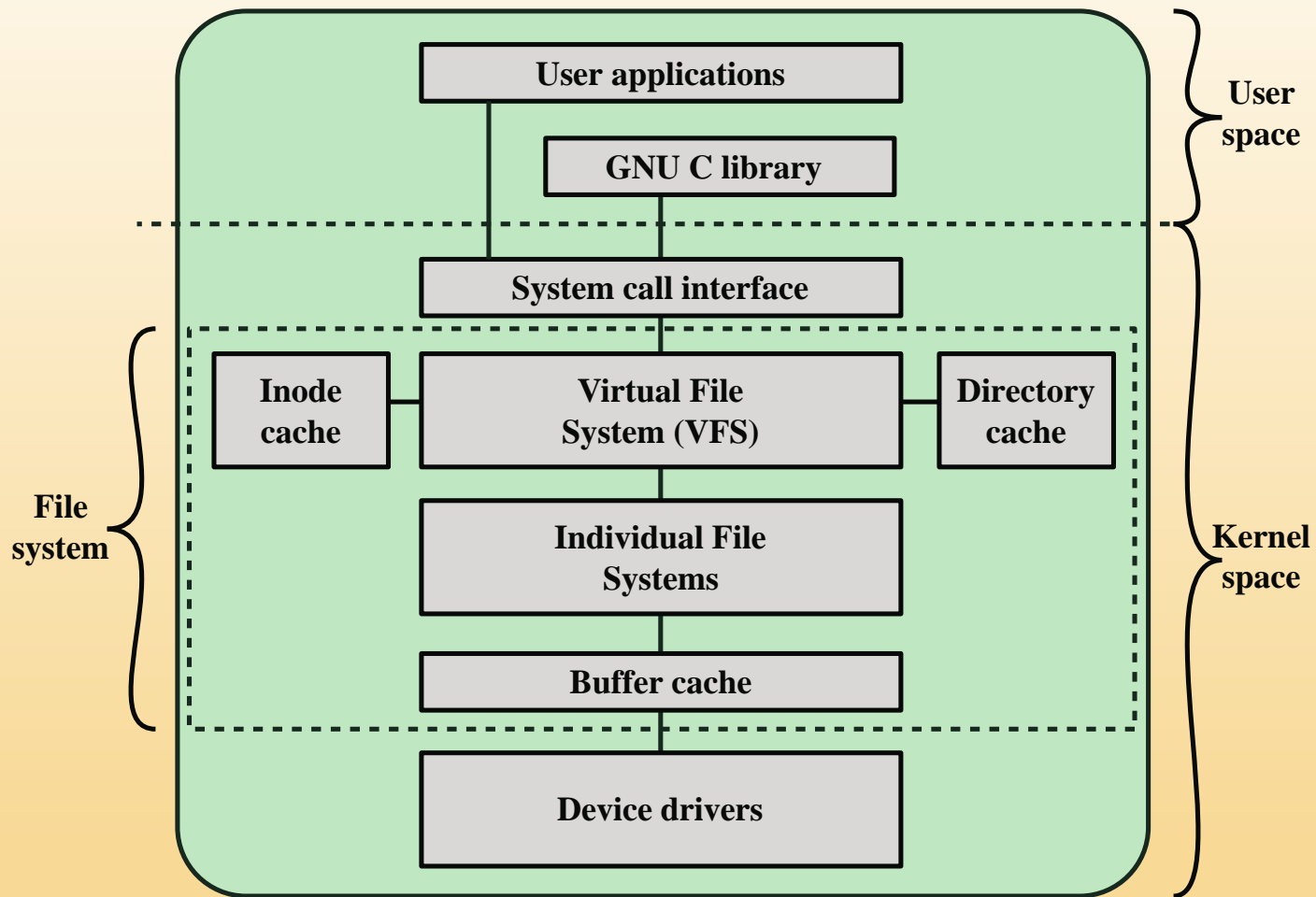


Figure 12.17 Linux Virtual File System Context

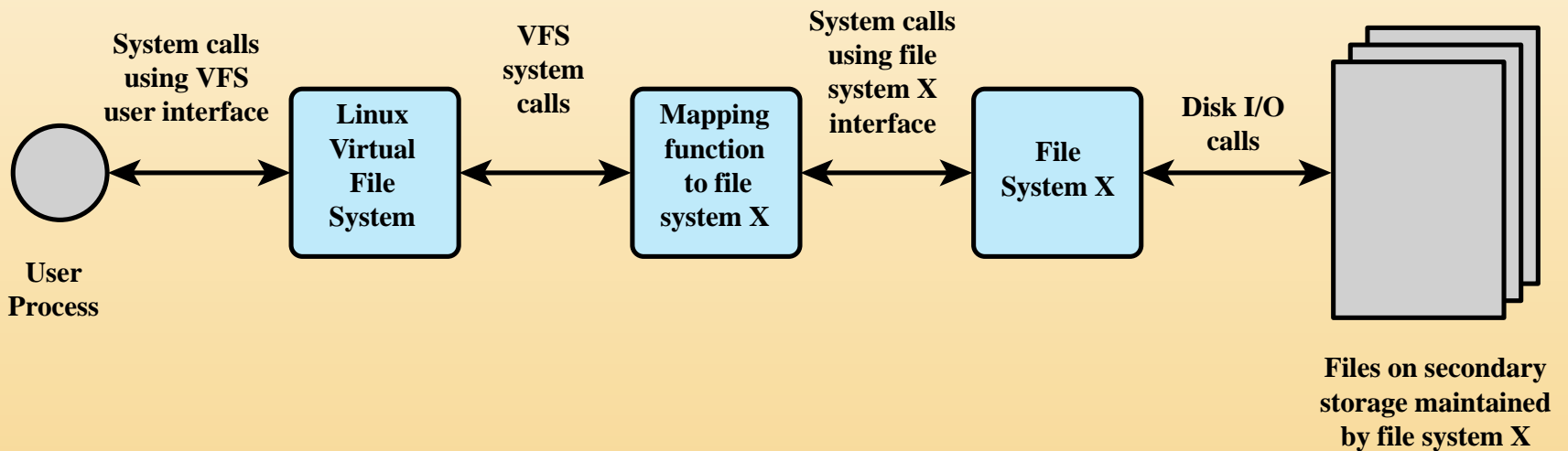


Figure 12.18 Linux Virtual File System Concept

Primary Object Types in VFS

Superblock Object

- Represents a specific mounted file system

Dentry Object

- Represents a specific directory entry

Inode Object

- Represents a specific file

File Object

- Represents an open file associated with a process

The Superblock Object

- Stores information describing a specific file system
- Typically, the superblock corresponds to the file system superblock or file system control block, which is stored in a special sector on disk
- The superblock object consists of a number of data items
 - Examples include:
 - The device this file system is mounted on
 - The basic block size of the file system
 - Dirty flag, to indicate that the superblock has been changed but not written back to disk
 - File system type
 - Flags, such as a read-only flag
 - Pointer to the root of the file system directory
 - List of open files
 - Semaphore for controlling access to the file system
 - List of superblock operations

The Inode Object

- An inode is associated with each file
- The inode object holds all the information about a named file except its name and the actual data contents of the file
- Items contained in an inode object include owner, group, permissions, access times for a file, size of data it holds, and number of links
- The inode object also includes an inode operations object that describes the file system's implemented functions that the VFS can invoke on an inode

The Dentry Object

- A dentry (directory entry) is a specific component in a path
- The component may be either a directory name or a file name
- Dentry objects facilitate quick lookups to files and directories, and are used in a dentry cache for that purpose
- The dentry object includes a pointer to the inode and superblock
- It also includes a pointer to the parent dentry and pointers to any subordinate dentries

The File Object

- The file object is used to represent a file opened by a process
- The object is created in response to the `open ()` system call, and destroyed in response to the `close ()` system call
- The file object also includes an inode operations object that describes the file system's implemented functions that the VFS can invoke on a file object
- The methods defined for the file object include read, write, open, release, and lock
- The file object consists of a number of items, including:
 - Dentry object associated with the file
 - File system containing the file
 - File objects usage counter
 - User's user ID
 - User's group ID
 - File pointer, which is the current position in the file from which the next operation will take place

Caches

- The VFS employs three caches to improve performance:
 - Inode cache
 - Because every file and directory is represented by a VFS inode, a directory listing command or a file access command causes a number of inodes to be accessed
 - The inode cache stores recently visited inodes to make access quicker
 - Directory cache
 - The directory cache stores the mapping between the full directory names and their inode numbers
 - This speeds up the process of listing a directory
 - Buffer cache
 - The buffer cache is independent of the file systems and is integrated into the mechanisms that the Linux kernel uses to allocate and read and write data buffers
 - As the real file systems read data from the underlying physical disks, this results in requests to the block device drivers to read physical blocks from the device that they control
 - If the same data is needed often, it will be retrieved from the buffer cache rather than read from the disk

Windows File System

- The developers of Windows NT designed a new file system, the New Technology File System (NTFS) which is intended to meet high-end requirements for workstations and servers
- Key features of NTFS:
 - Recoverability
 - Security
 - Large disks and large files
 - Multiple data streams
 - Journaling
 - Compression and encryption
 - Hard and symbolic links

NTFS Volume and File Structure

- NTFS makes use of the following disk storage concepts:

Sector

- The smallest physical storage unit on the disk
- The data size in bytes is a power of 2 and is almost always 512 bytes

Cluster

- One or more contiguous sectors
- The cluster size in sectors is a power of 2

Volume

- A logical partition on a disk, consisting of one or more clusters and used by a file system to allocate space
- Can be all or a portion of a single disk or it can extend across multiple disks
- The maximum volume size for NTFS is 2^{64} clusters

Table 12.4

Windows NTFS Partition and Cluster Sizes

Volume Size	Sectors per Cluster	Cluster Size
≤ 512 Mbyte	1	512 bytes
512 Mbyte - 1 Gbyte	2	1K
1 Gbyte - 2 Gbyte	4	2K
2 Gbyte - 4 Gbyte	8	4K
4 Gbyte - 8 Gbyte	16	8K
8 Gbyte - 16 Gbyte	32	16K
16 Gbyte - 32 Gbyte	64	32K
> 32 Gbyte	128	64K

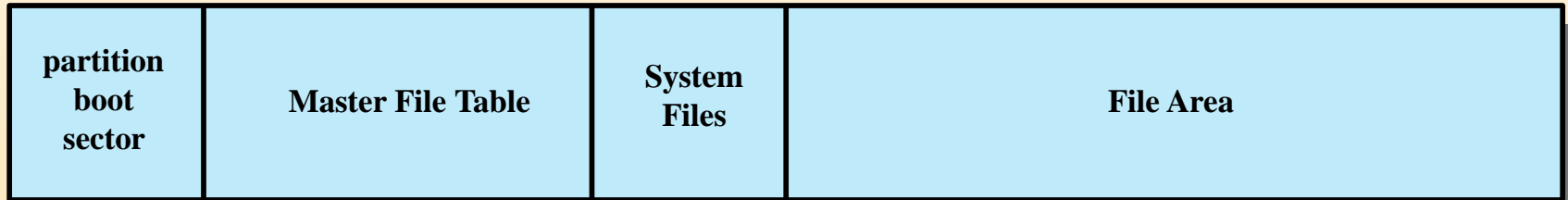


Figure 12.19 NTFS Volume Layout

Master File Table (MFT)

- The heart of the Windows file system is the MFT
- The MFT is organized as a table of 1,024-byte rows, called records
- Each row describes a file on this volume, including the MFT itself, which is treated as a file
- Each record in the MFT consists of a set of attributes that serve to define the file (or folder) characteristics and the file contents

Table 12.5

Windows NTFS File and Directory Attribute Types

Attribute Type	Description
Standard information	Includes access attributes (read-only, read/write, etc.); time stamps, including when the file was created or last modified; and how many directories point to the file (link count).
Attribute list	A list of attributes that make up the file and the file reference of the MFT file record in which each attribute is located. Used when all attributes do not fit into a single MFT file record.
File name	A file or directory must have one or more names.
Security descriptor	Specifies who owns the file and who can access it.
Data	The contents of the file. A file has one default unnamed data attribute and may have one or more named data attributes.
Index root	Used to implement folders.
Index allocation	Used to implement folders.
Volume information	Includes volume-related information, such as the version and name of the volume.
Bitmap	Provides a map representing records in use on the MFT or folder.

Note: Colored rows refer to required file attributes; the other attributes are optional.

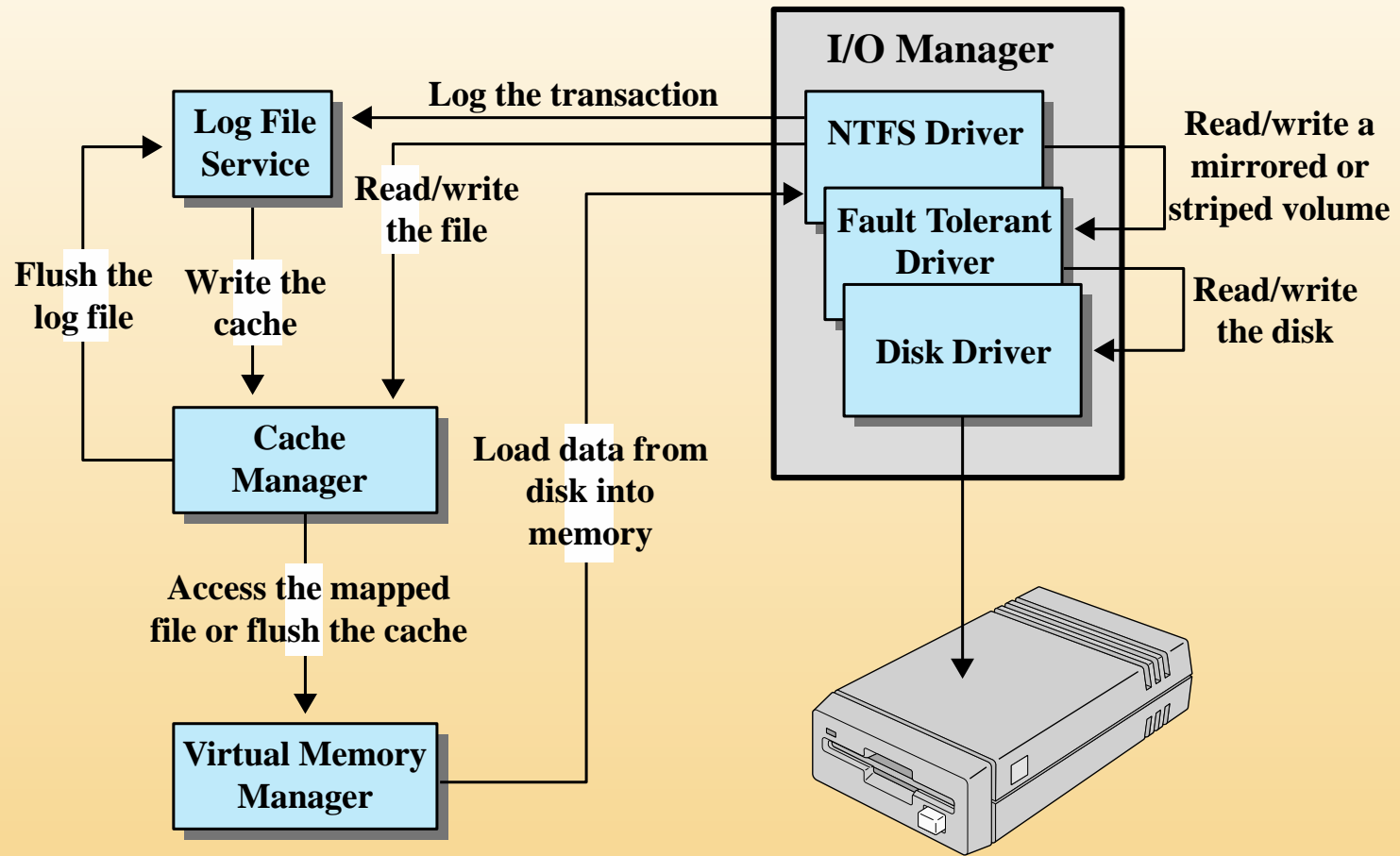
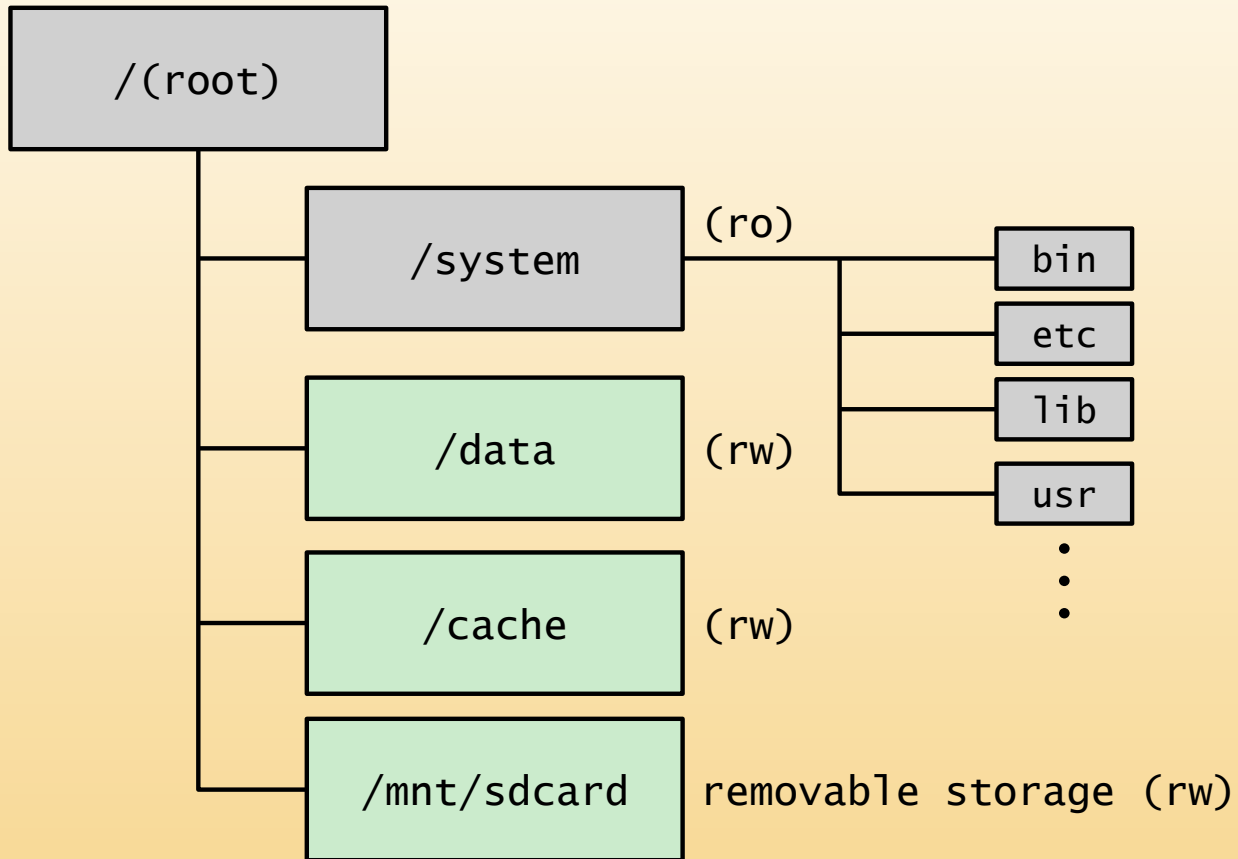


Figure 12.20 Windows NTFS Components



ro: mounted as read only
rw: mounted as read and write

Figure 12.21 Typical Directory Tree of Android

SQLite

- Most widely deployed SQL database engine in the world
- Based on the Structured Query Language (SQL)
- Designed to provide a streamlined SQL-based database management system suitable for embedded systems and other limited memory systems
- The full SQLite library can be implemented in under 400 KB
- In contrast to other database management systems, SQLite is not a separate process that is accessed from the client application
 - The SQLite library is linked in, and thus becomes an integral part of the application program

Summary

- File structure
- File management systems
- File organization and access
 - The pile
 - The sequential file
 - The indexed sequential file
 - The indexed file
 - The direct or hashed file
- B-Trees
- File directories
 - Contents
 - Structure
 - Naming
- File sharing
 - Access rights
 - Simultaneous access
- Record blocking
- Android file management
 - File system
 - SQLite
- Secondary storage management
 - File allocation
 - Free space management
 - Volumes
 - Reliability
- UNIX file management
 - Inodes
 - File allocation
 - Directories
 - Volume structure
- Linux virtual file system
 - Superblock object
 - Inode object
 - Dentry object
 - File object
 - Caches
- Windows file system
 - Key features of NTFS
 - NTFS volume and file structure
 - Recoverability