

*Operating
Systems:
Internals
and Design
Principles*

Chapter 12

File Management

Eighth 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

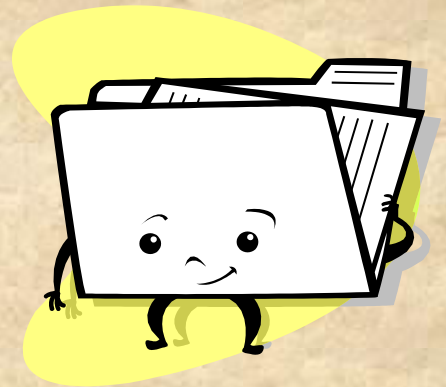
- files have names and can have associated access permissions that permit controlled sharing

Structure

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

File Systems

- File systems provide a means to store data organized as files as well as operations that can be performed on files
- Typical operations include:
 - Create, Delete, Open, Close, Read, Write
- Maintain a set of attributes associated with the file
 - Name, Extension, Size, Date created, Date last modified



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 (files or also called tables)
- relationships among elements of data are explicit
- designed for use by a number of different applications
- consists of one or more types of files

Record

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

File

- collection of similar records
- treated as a single entity
- may be referenced by name
- access control restrictions usually apply at the file level

File Management System Objectives

- 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 for Files

■ 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 files

4

- should be able to restructure the files in a form appropriate to the problem

5

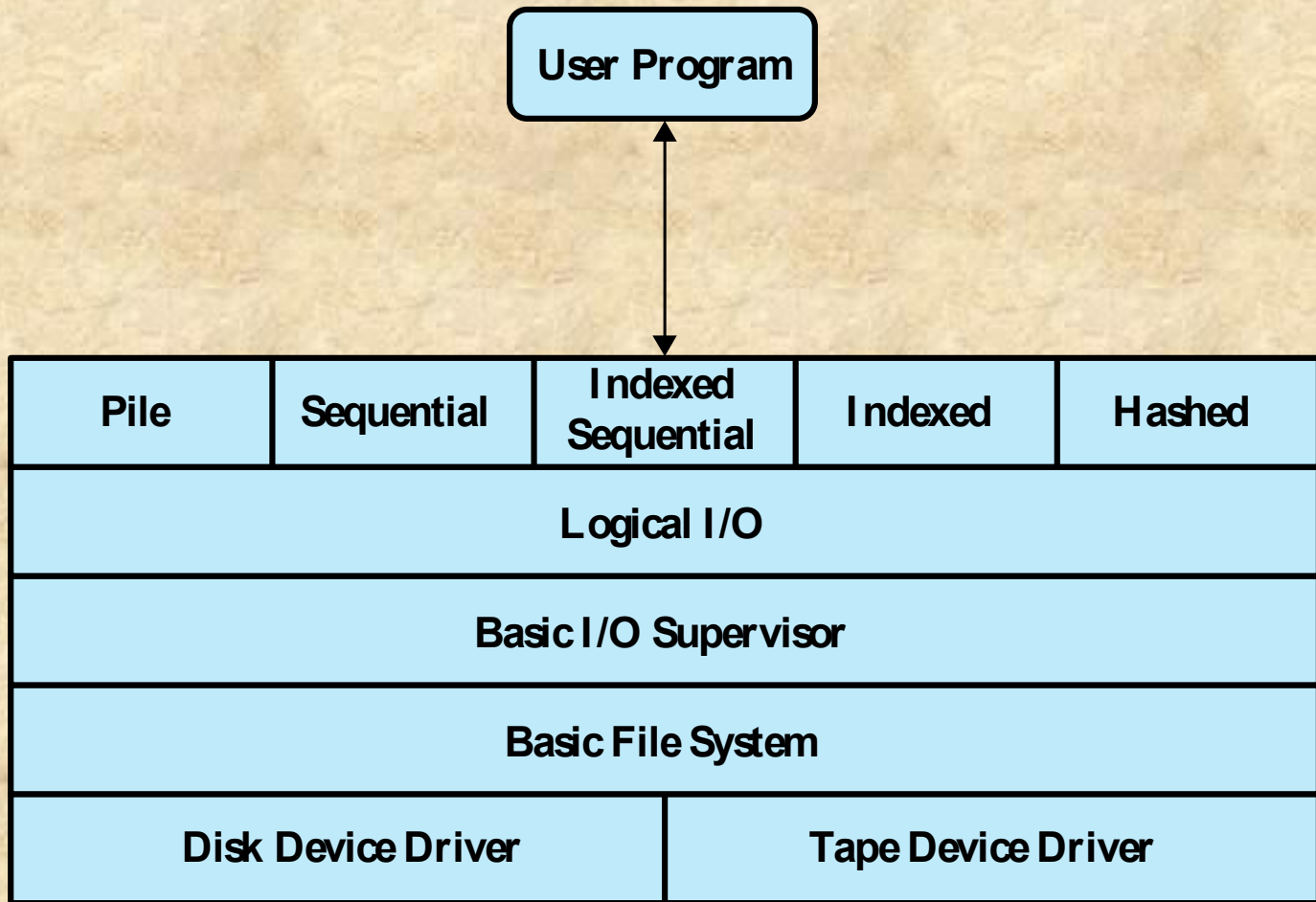
- should be able to move data between files

6

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

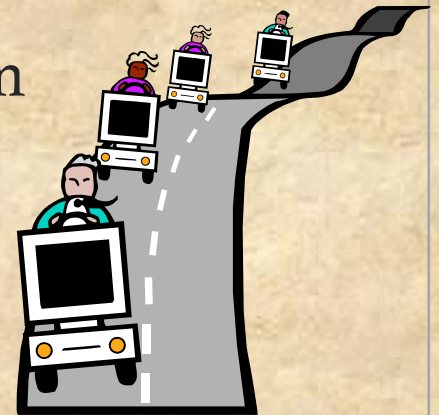
7

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



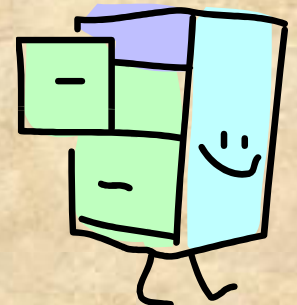
Device Drivers

- 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
- Considered part of the operating system

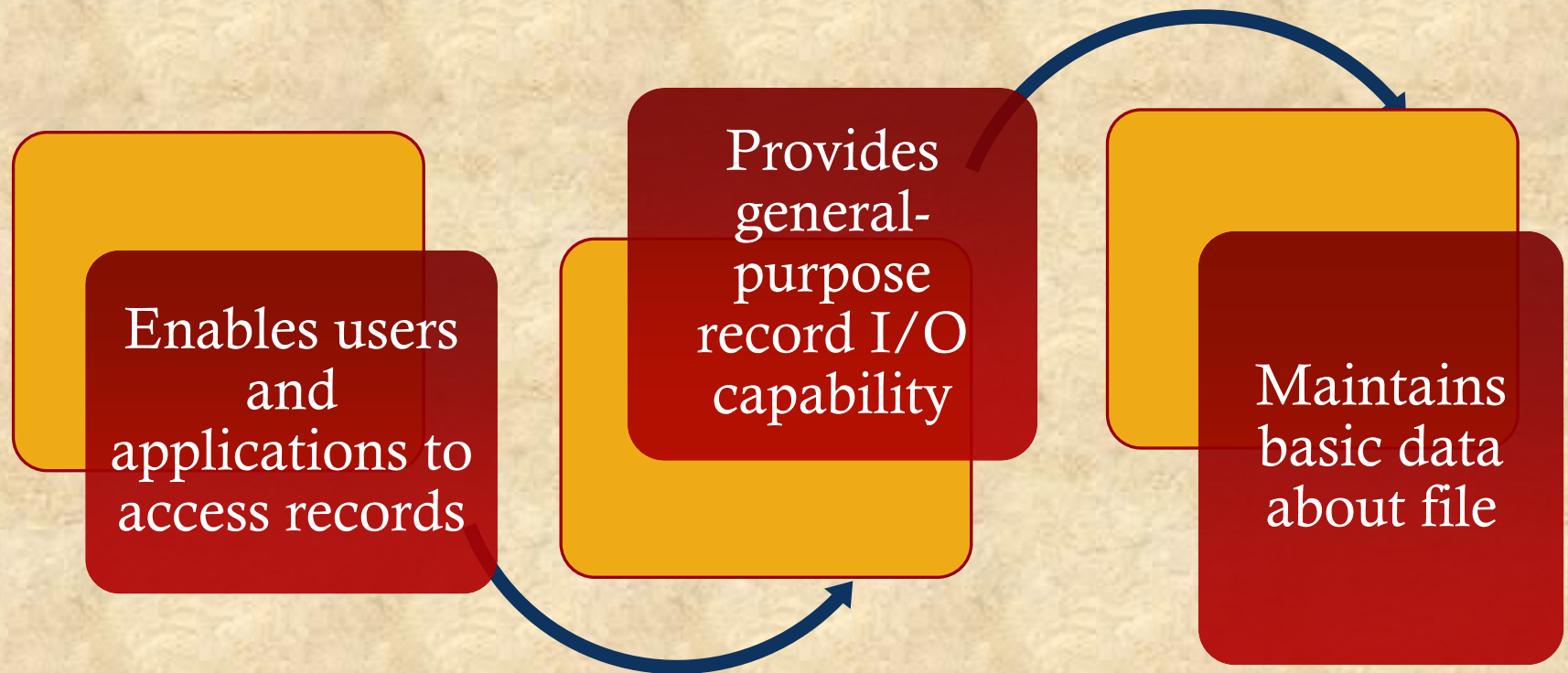




Basic I/O Supervisor

- Responsible for all file I/O initiation and termination
- Control structures that deal with device I/O, scheduling, and file status are maintained
- 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

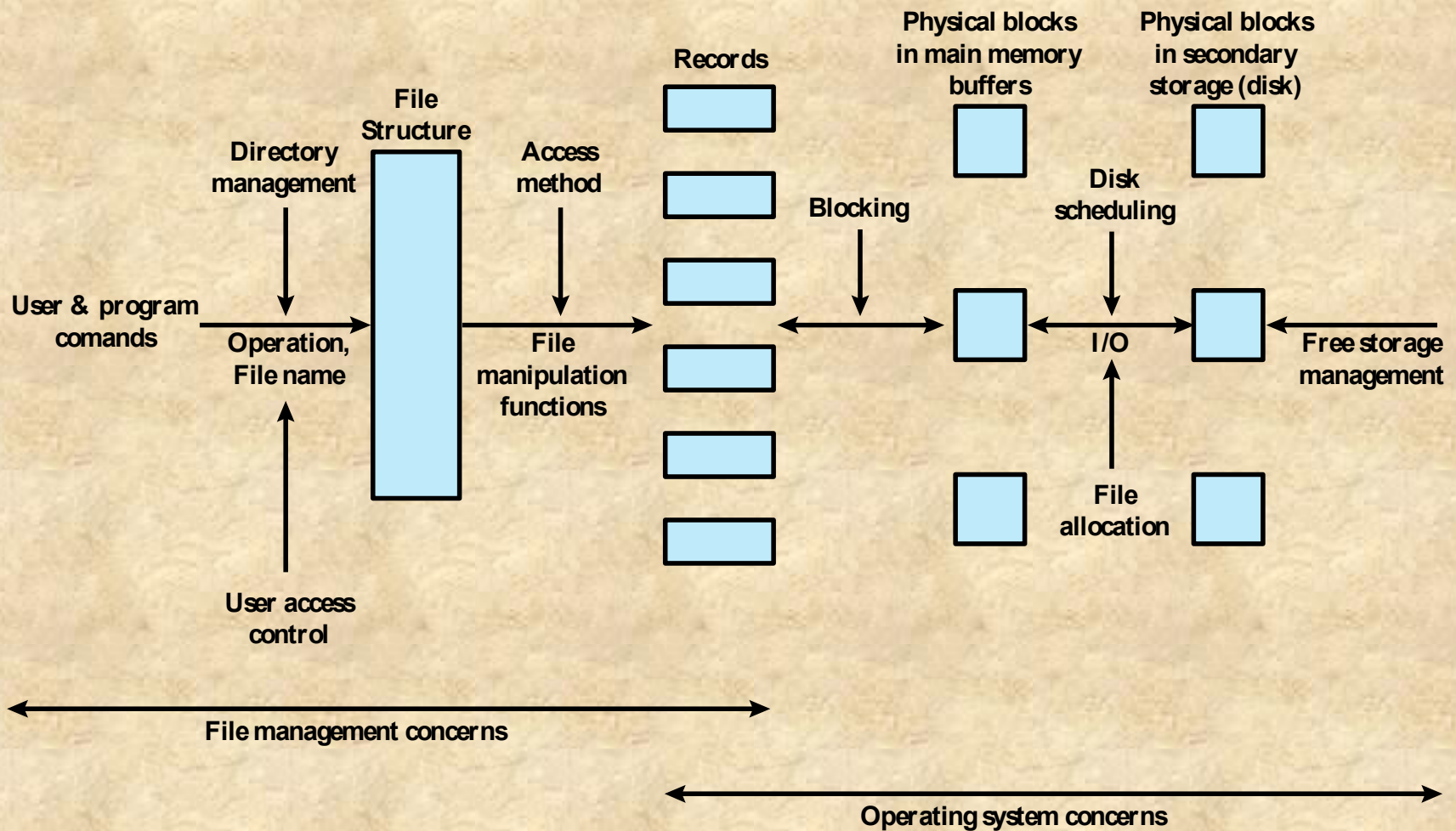
Logical I/O



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
- Sequential and indexed access methods are the most common



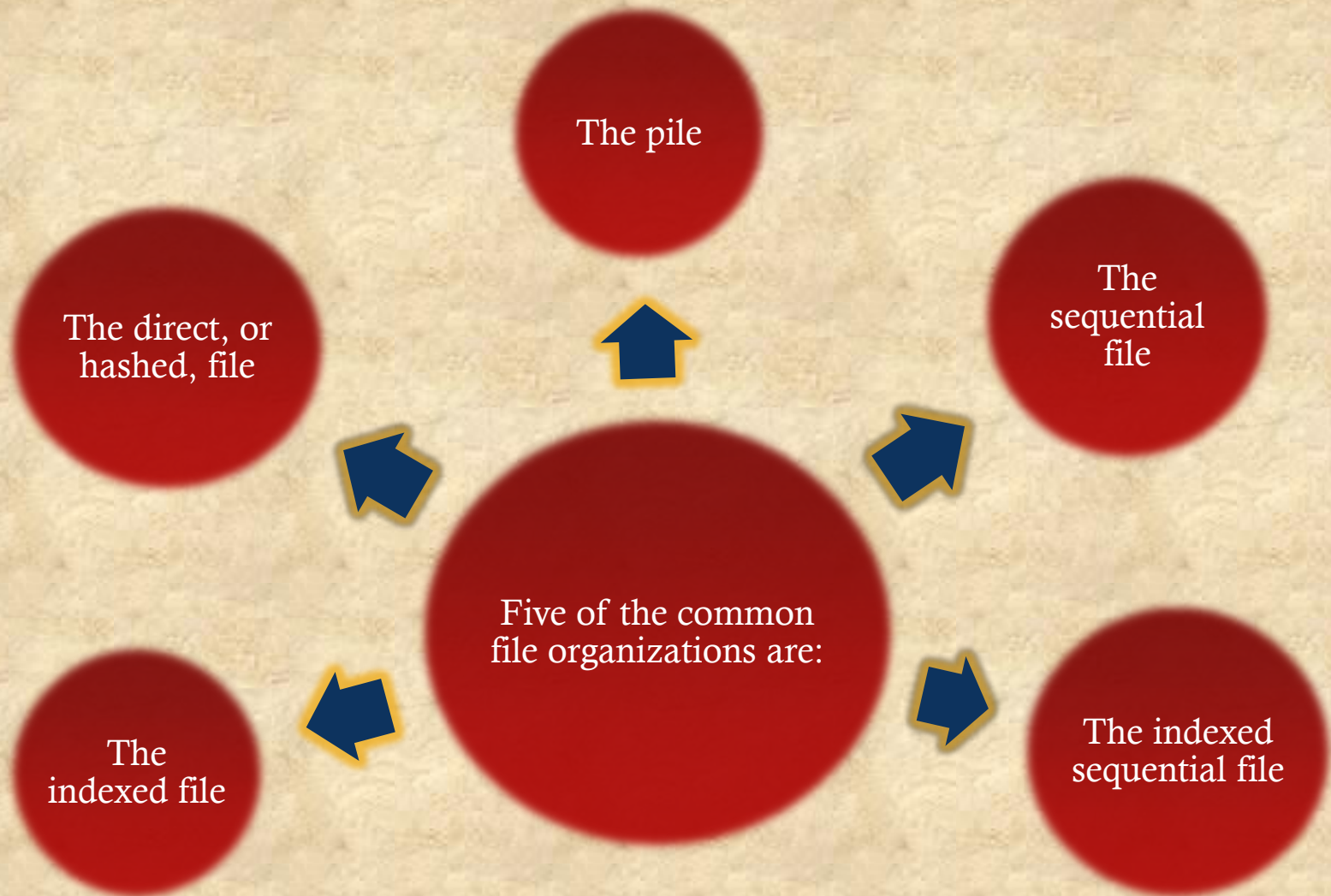


File Organization and Access

- ***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

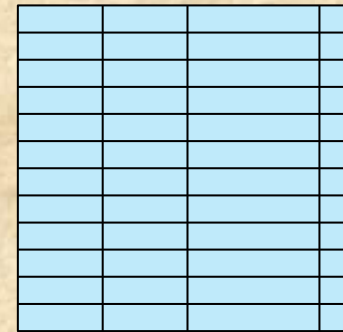


File Organization Types





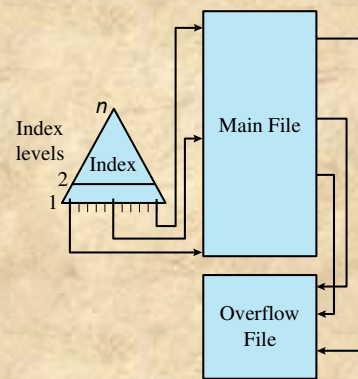
Variable-length records
Variable set of fields
Chronological order



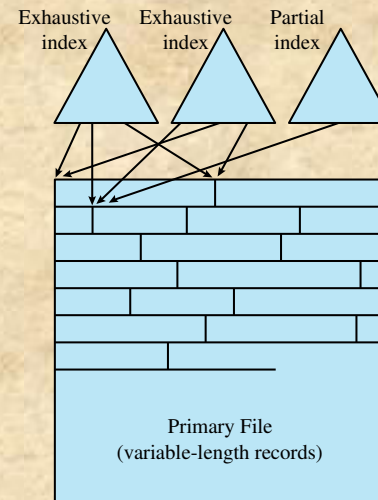
Fixed-length records
Fixed set of fields in fixed order
Sequential order based on key field

(a) Pile File

(b) Sequential File



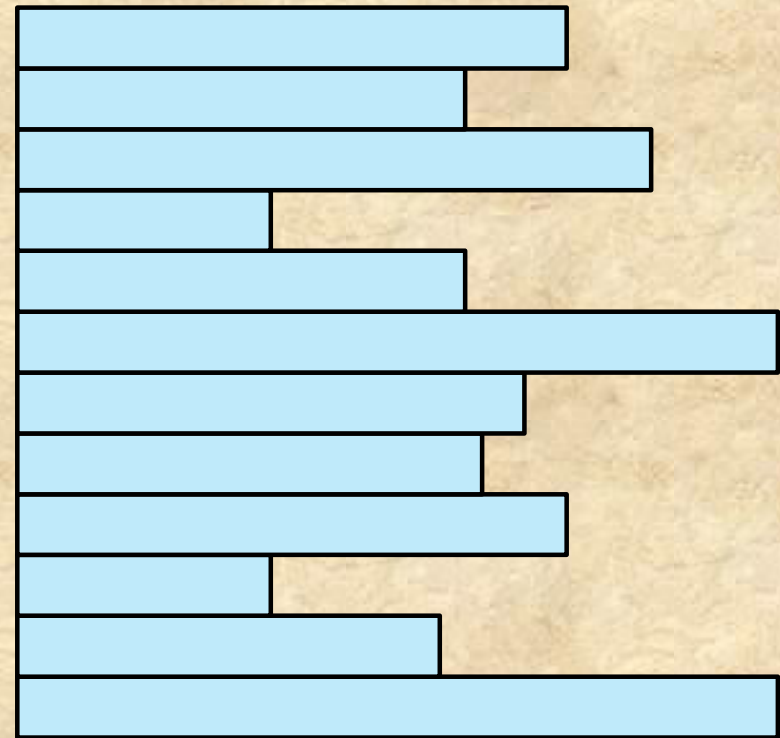
(c) Indexed Sequential File



(d) Indexed File

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

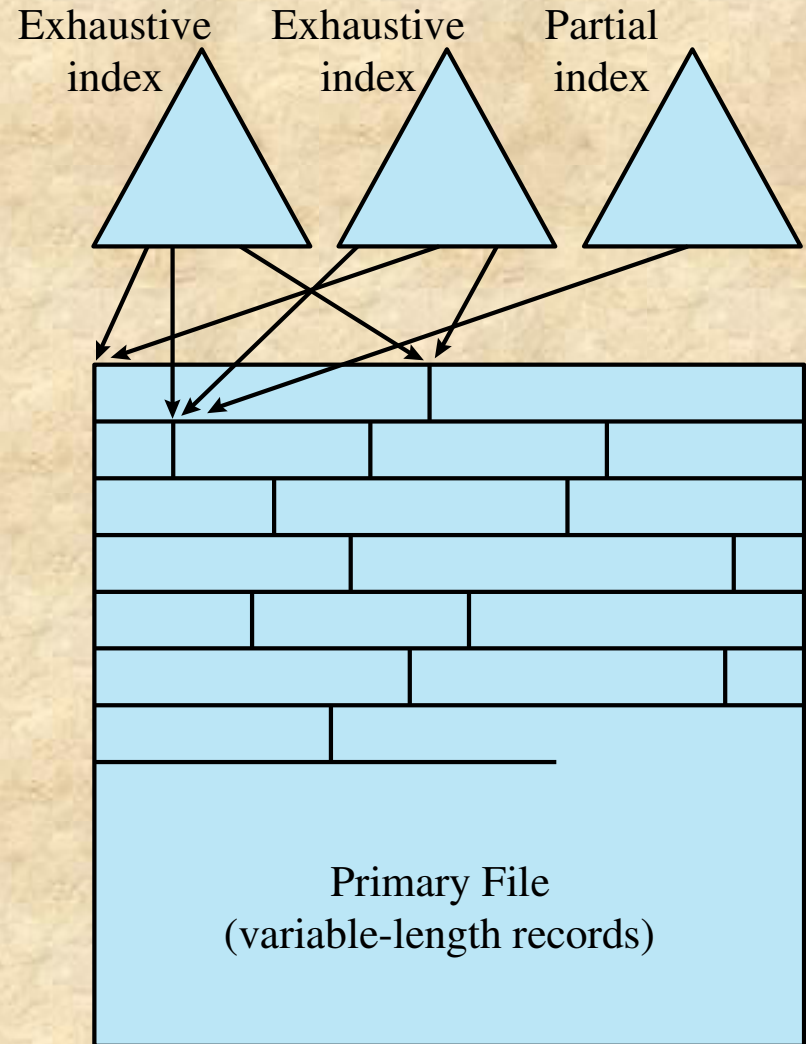
Fixed set of fields in fixed order

Sequential order based on key field

(b) Sequential File

Indexed File

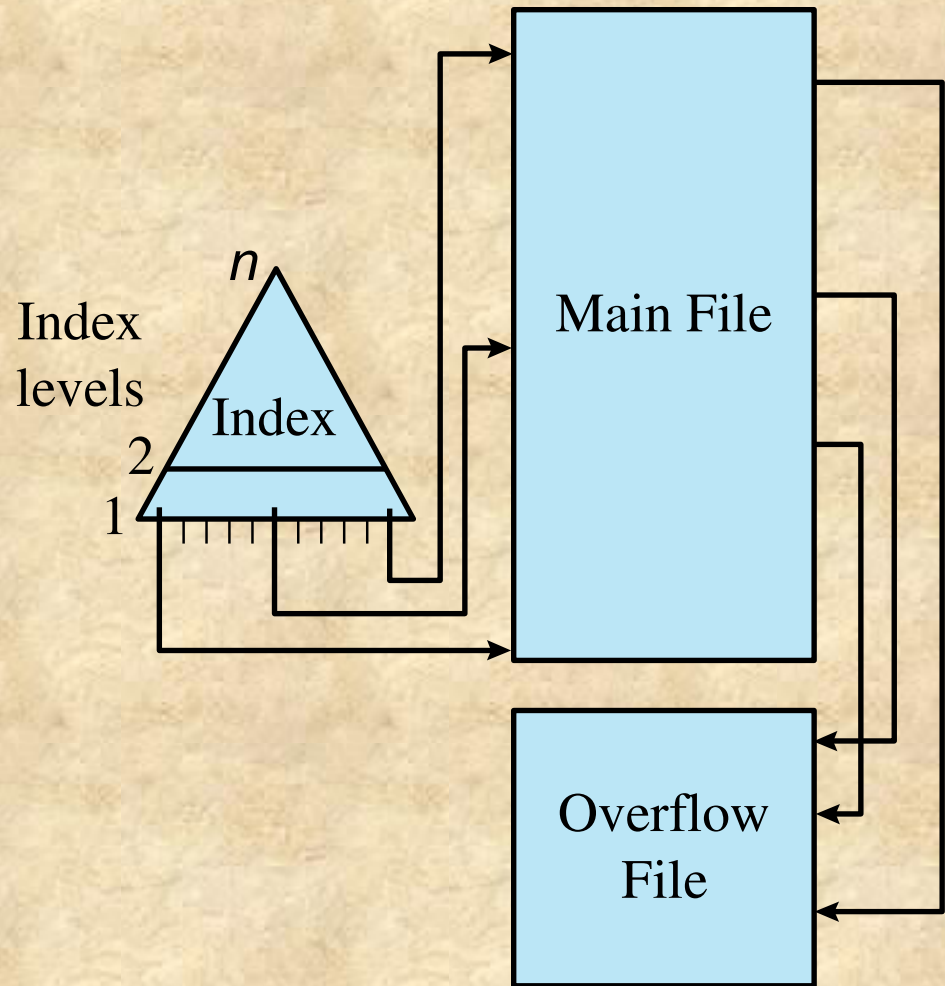
- Records are accessed only through their indexes
- Variable-length records can be employed
- Exhaustive index contains one entry for every record in the main file
- Partial index contains entries to records where the field of interest exists
- Used mostly in applications where timeliness of information is critical
- Examples would be airline reservation systems and inventory control systems
- See the video on how an index works:
<https://www.youtube.com/watch?v=zDzu6vka0rQ>



(d) Indexed 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

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
- Watch this YouTube video:
 - <https://www.youtube.com/watch?v=MfhjkfocRR0>

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
- First see how a B-Tree works in this video:
 - https://www.youtube.com/watch?v=k5J9M5_IMzg



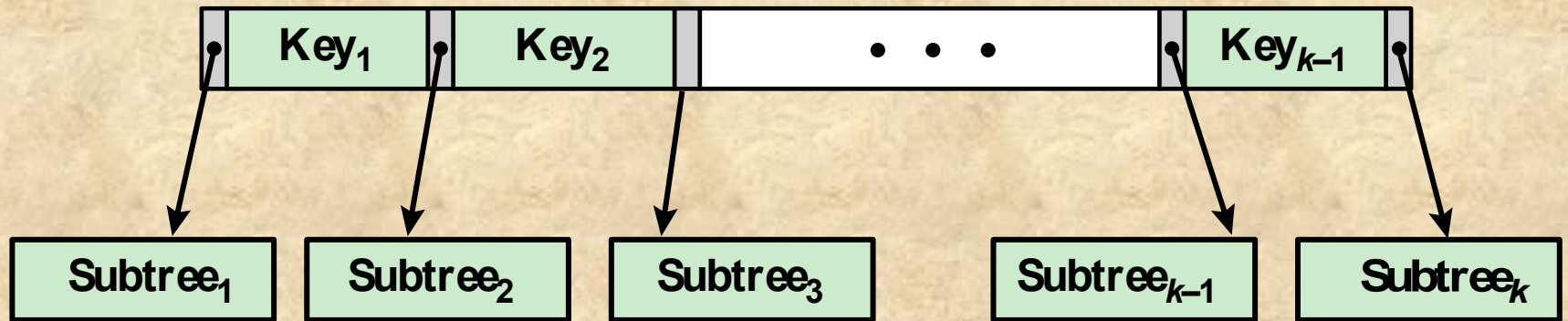
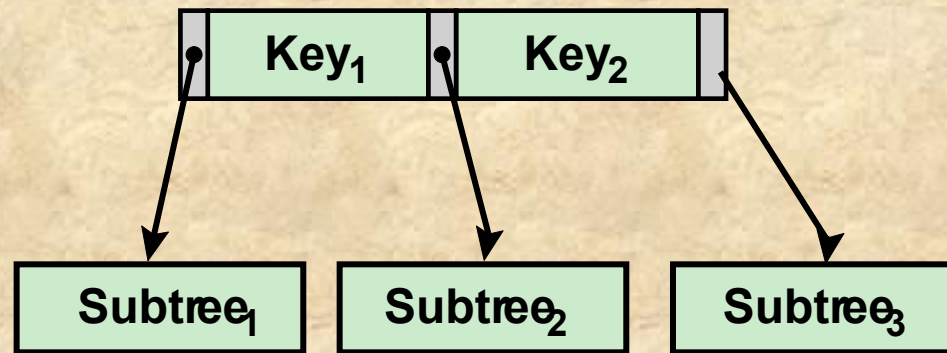


Figure 12.4 A B-tree Node with k Children



A B-tree Node with 3 Children

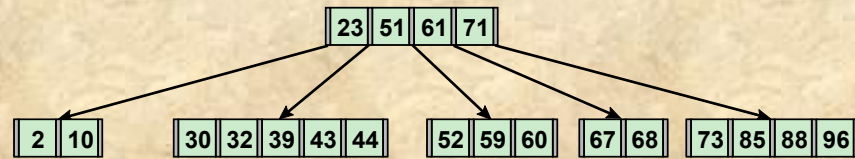
B-Tree

Characteristics

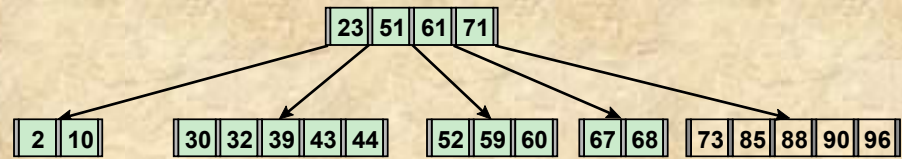
A B-tree is characterized by its minimum degree d and satisfies the following properties:

- 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

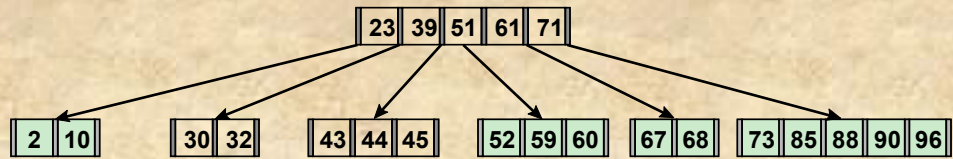




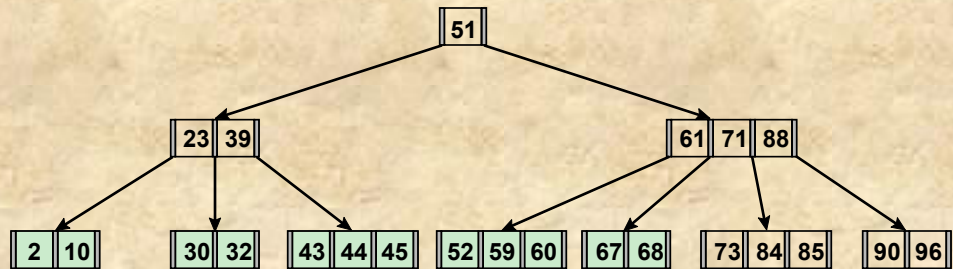
(a) B-tree of minimum degree $d = 3$.



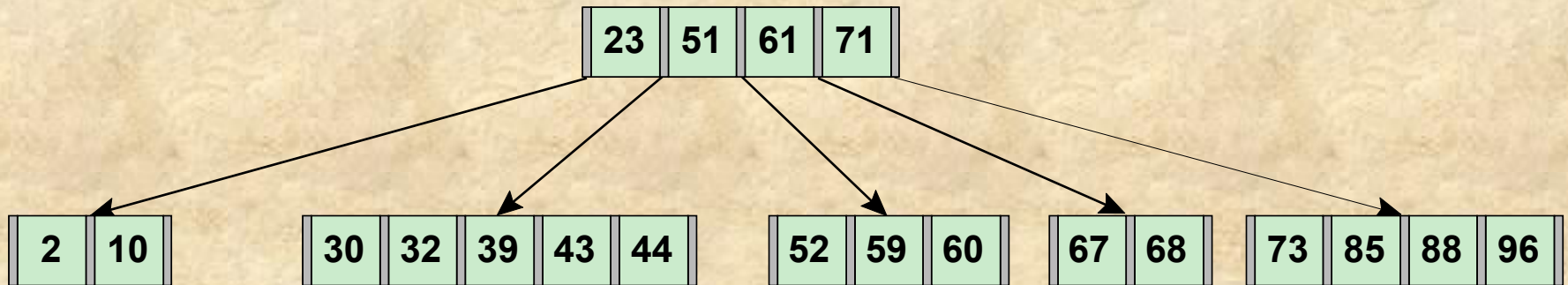
(b) Key = 90 inserted. This is a simple insertion into a node.



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

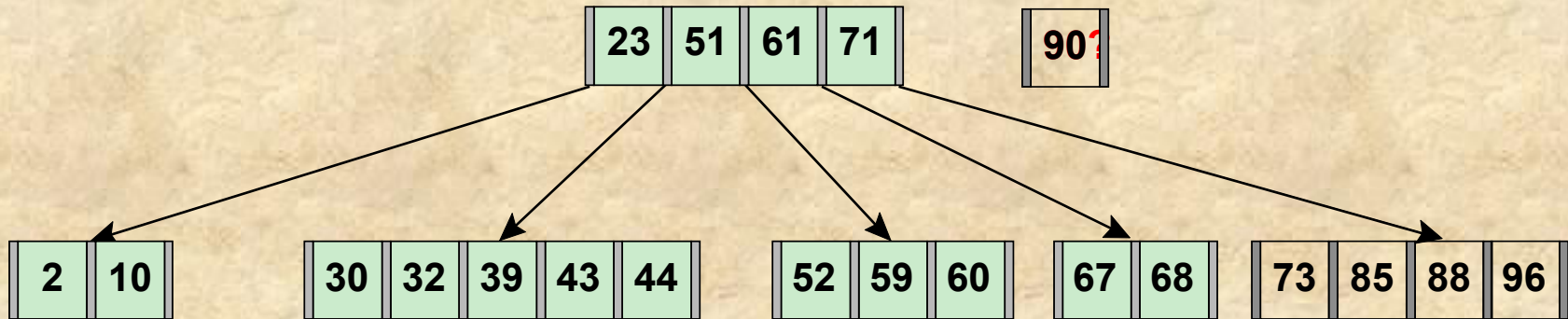


(d) Key = 84 inserted. This requires splitting a node into two parts and promoting one key to the root node. This then requires the root node to be split and a new root created.



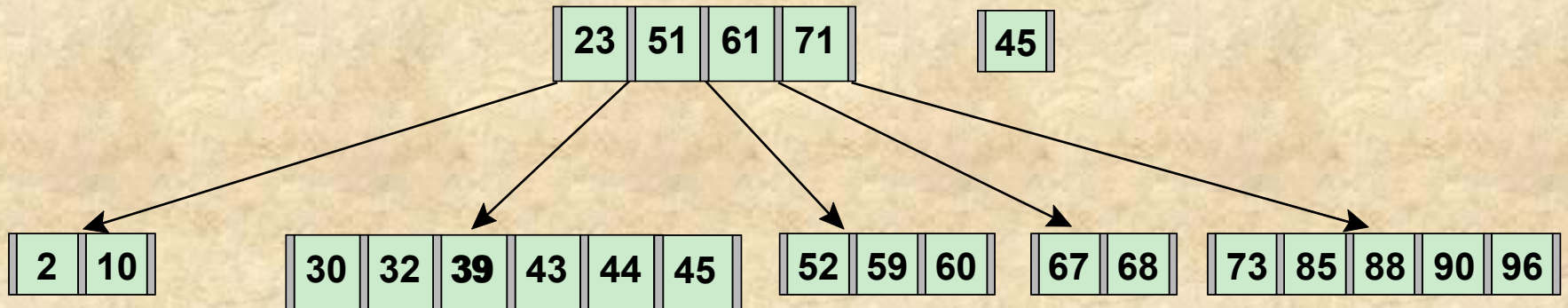
(a) B-tree of minimum degree $d = 3$.

- To search for a key, you start at the root node. Let us search for 90.
- If the key you want is in the node, you're done. If not, you go down one level. There are three cases:
 1. The key you want is less than the smallest key in this node. Take the leftmost pointer down to the next level. Not our case for 90.
 2. The value of the key is between the values of two adjacent keys in this node. Take the pointer between these keys down to the next level. Not our case for 90.
 3. The key you want is greater than the largest key in this node. Take the rightmost pointer down to the next level. This is our case, so take the rightmost node.



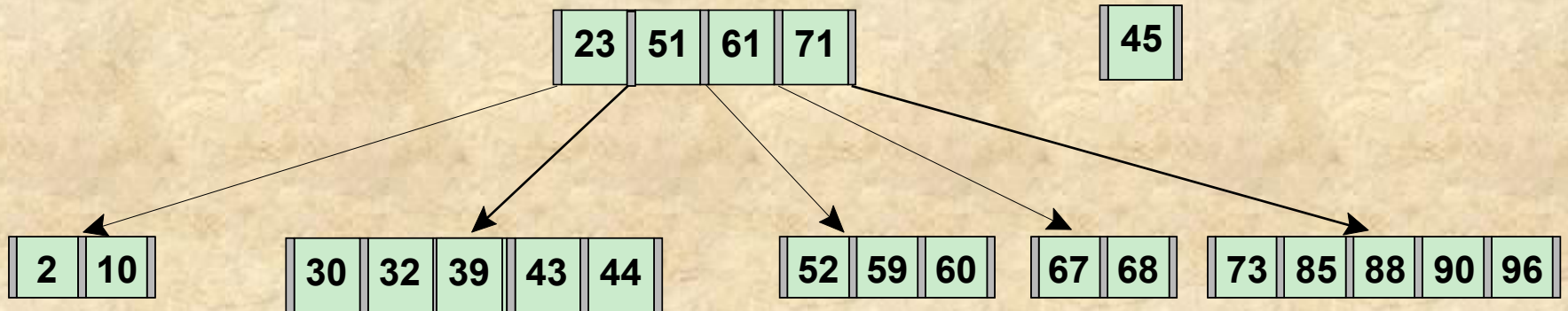
(b) Key = 90 inserted. This is a simple insertion into a node.

- If this node has fewer than $2d - 1 = 2(3) - 1 = 5$ keys, then insert the key into this node in the proper sequence.
- So 90 goes in-between 88 and 96.
- Now, insert 45...



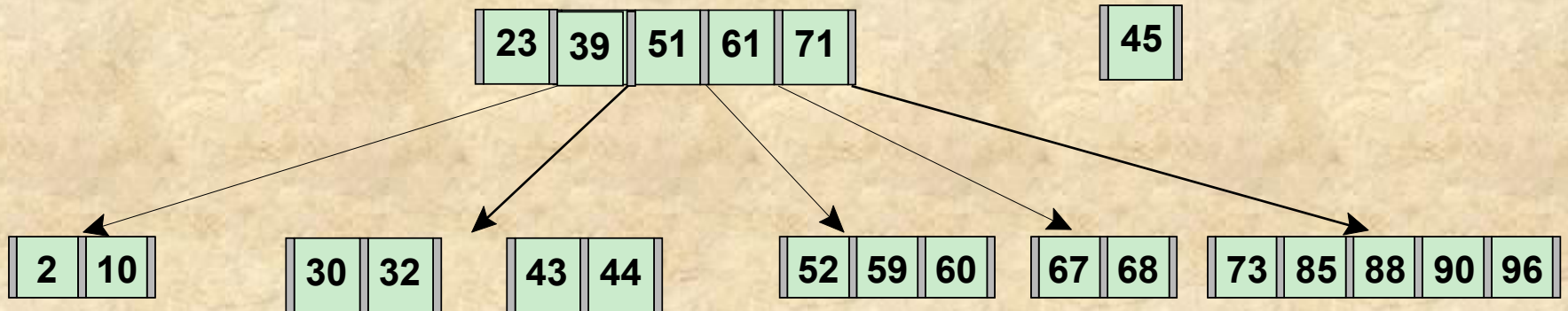
(c) Key = 45 inserted. This requires splitting a node into two parts and promoting one key to the root node.

- To search for a key, you start at the root node. Let us insert 45.
- The key we want (45) is NOT in the root node. $45 < 51$, so follow left pointer of 51 down to that child node.
- The value 45 in this node is greater than the largest key (44) in this node, so it has to be added to the right of key 44. But now the node is over full!
- If the node is full (having $2d - 1 = 5$ keys, which is the case) when a key needs to be added, then identify the median key. The median key is the middle key if there is an odd number of keys (which is key 39 in this case) or if an equal number of keys, the median becomes the average of the “two middle” keys.



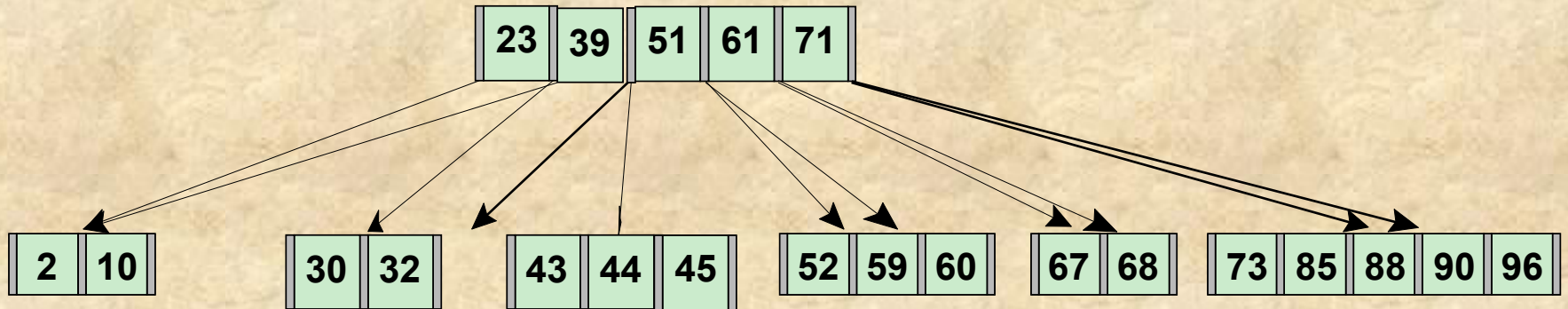
(c) Key = 45 inserted. This requires splitting a node into two parts and promoting one key to the root node.

- Split this node around its median key into two new nodes with $d - 1 = 2$ keys each
- Promote the median key to the next higher level, as described in step 4. The promoted node is inserted into the parent node as follows: if the parent node is already full, the parent node must be split and its median key promoted to the next highest layer. However, in this example, it is not the case.



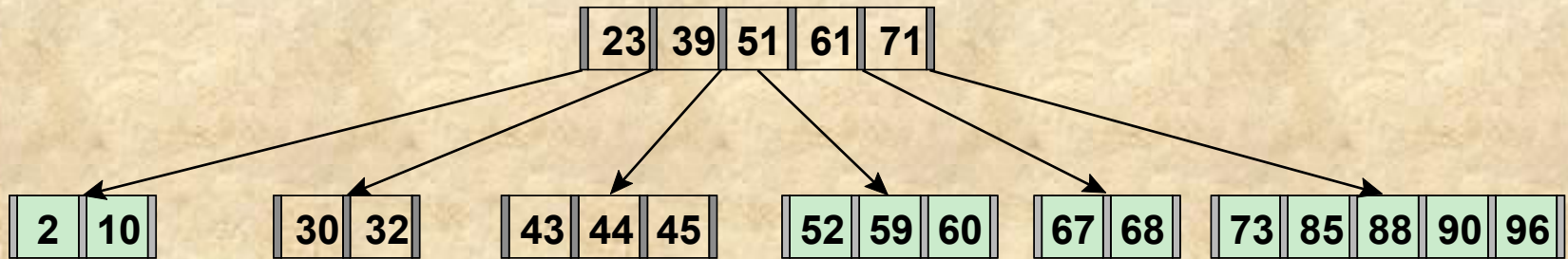
(c) Key = 45 inserted. This requires splitting a node into two parts and promoting one key to the root node.

- Split this node around its median key into two new nodes with $d - 1 = 2$ keys each
- Promote the median key to the next higher level, as described in step 4. The promoted node is inserted into the parent node as follows: if the parent node is already full, the parent node must be split and its median key promoted to the next highest layer. However, in this example, it is not the case.
- If the new key (45) has a value less than the median key, insert it into the left-hand new node; otherwise insert it into the right-hand new node; this is the current case. The result is that the original node has been split into two nodes, one with $d - 1$ keys and one with d keys.

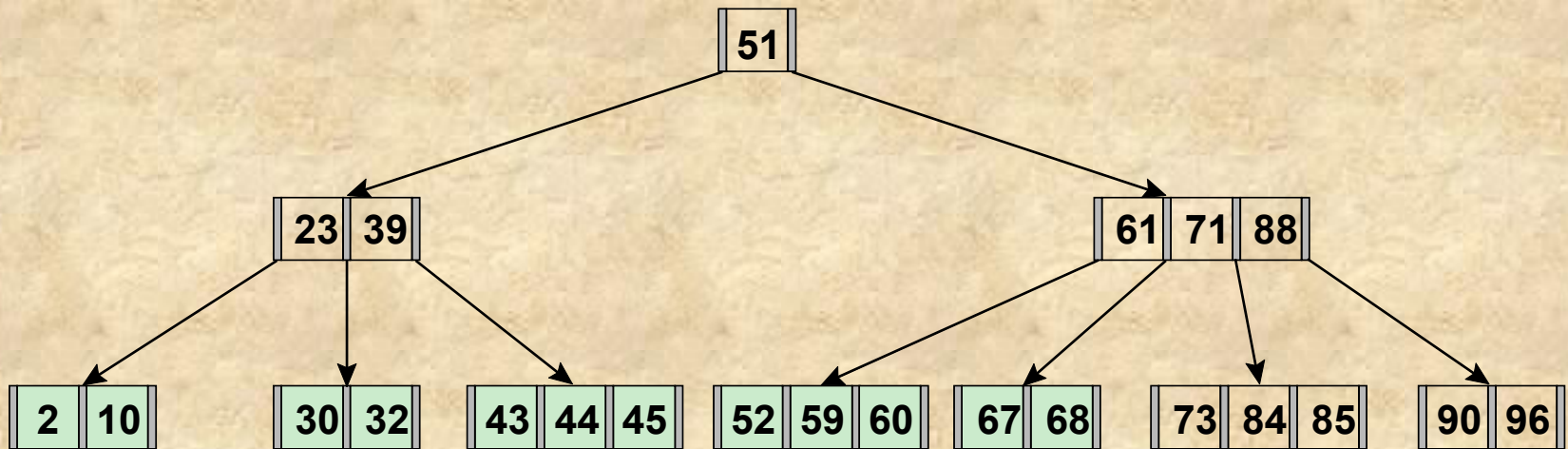


(c) Key = 45 inserted. This requires splitting a node into two parts and promoting one key to the root node.

- Split this node around its median key into two new nodes with $d - 1 = 2$ keys each
- Promote the median key to the next higher level, as described in step 4. The promoted node is inserted into the parent node as follows: if the parent node is already full, the parent node must be split and its median key promoted to the next highest layer. However, in this example, it is not the case.
- If the new key (45) has a value less than the median key, insert it into the left-hand new node; otherwise insert it into the right-hand new node; this is the current case. The result is that the original node has been split into two nodes, one with $d - 1$ keys and one with d keys.



(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.

■ Simulators for B-Trees

- <https://s3.amazonaws.com/learneroo/visual-algorithms/BTree.html>
- <http://ats.oka.nu/b-tree/b-tree.html>
- Note that the implementations of these simulators differ. B-trees may be implemented with slight variations. These variations are not for exam purposes, for example B+ Trees, however, if you are interested, you are welcome to go and Google them.
- You should be able to, in a test or exam, traverse a B-Tree as explained in this lecture and be able to insert a node and hence show the tree structure after inserting the node.
- Remember that B-trees are used in file systems to point to information or records in a file and be able to access that information very quickly by traversing the tree of pointers in order to very quickly get access to a certain record in a file.

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:

Search

Create
files

Delete
files

List
directory

Update
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

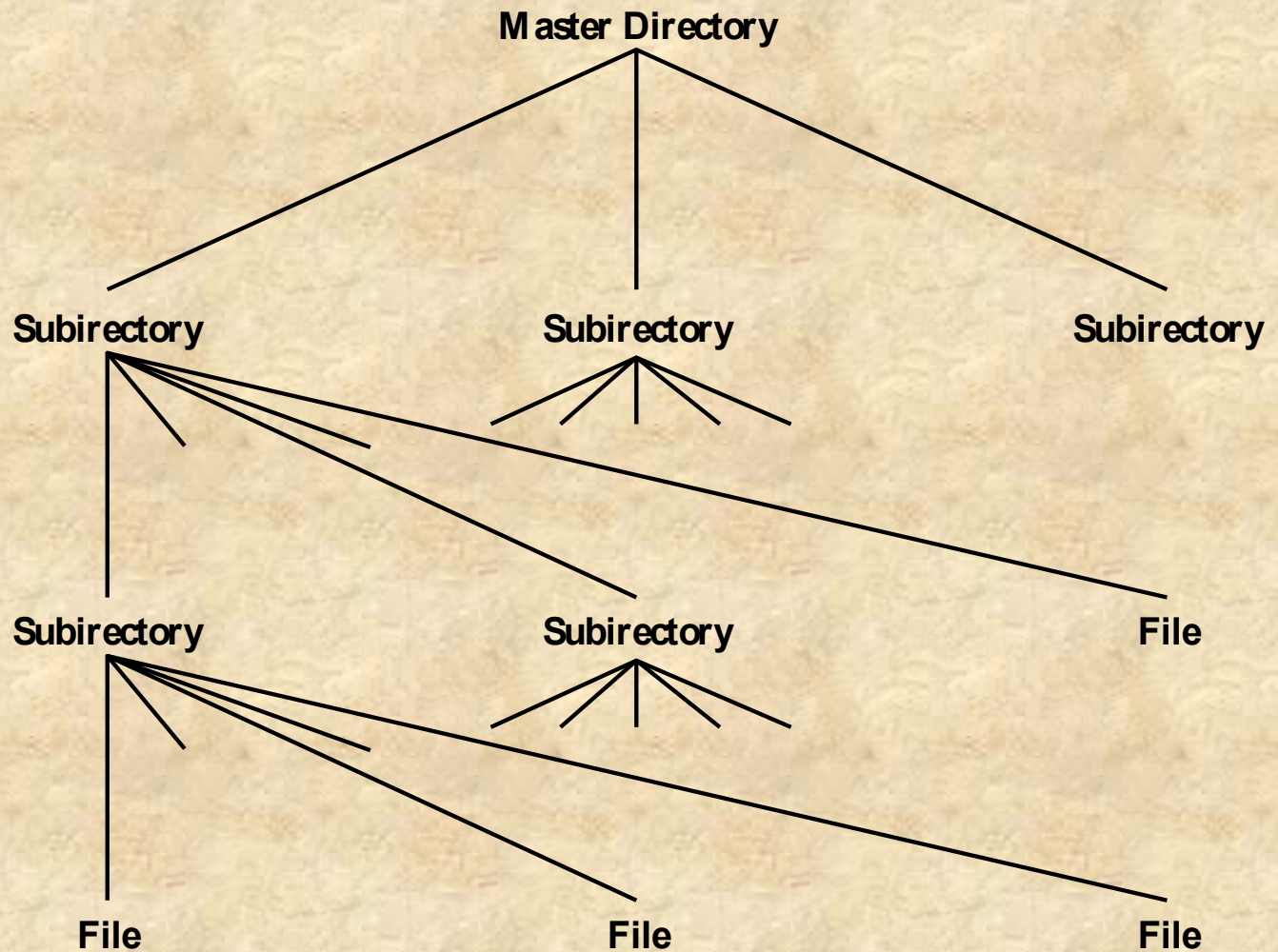


Figure 12.6 Tree-Structured Directory

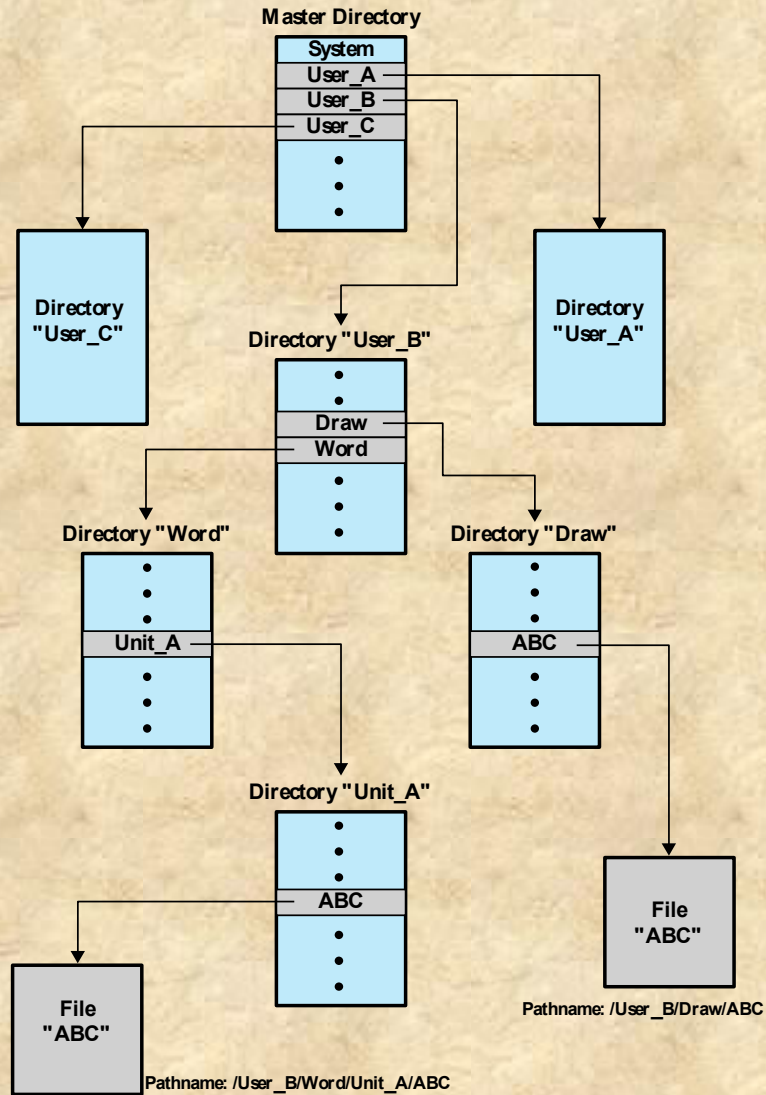


Figure 12.7 Example of Tree-Structured Directory

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 original 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
their specific
user IDs

User Groups

a set of users
who are not
individually
defined, e.g.
database
administrators

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



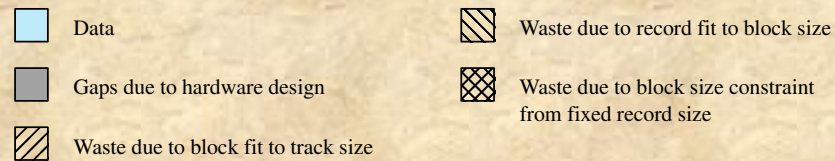
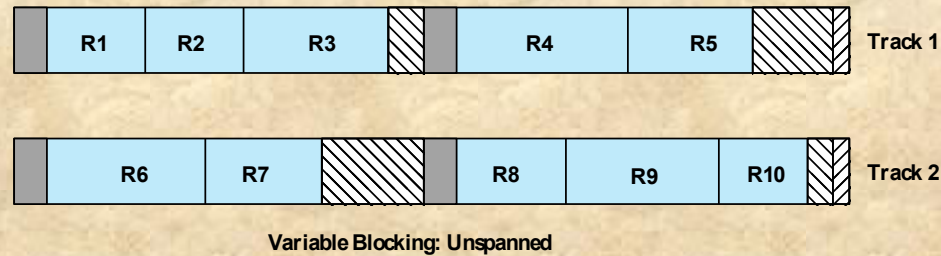
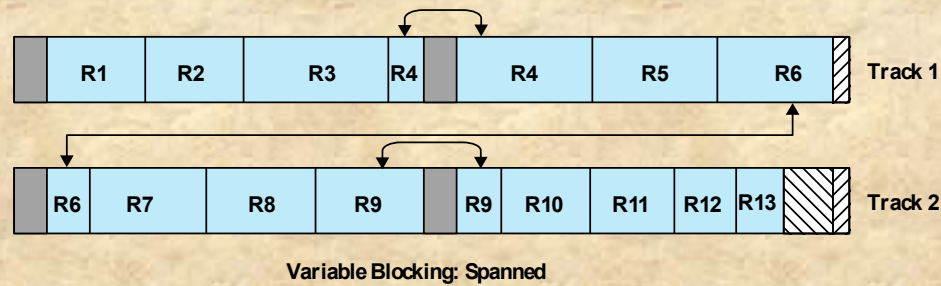
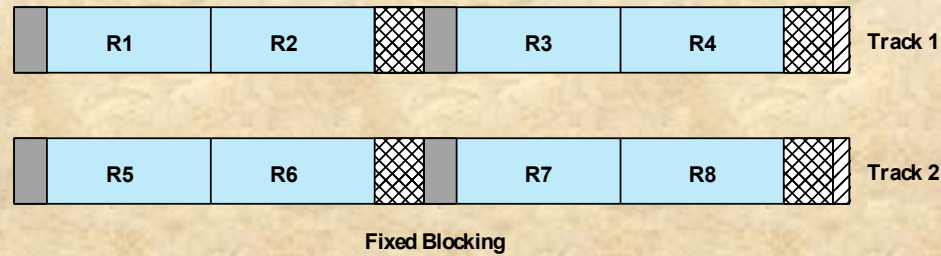
- Given the size of a block, three methods of blocking can be used:

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
Internal fragmentation may happen again



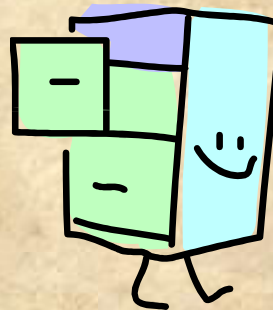


File Allocation

- 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
- Such trade-off items to be considered:
 - 1) contiguity of space (meaning bordering of storage 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

Alternatives

Two major alternatives:

Variable, large contiguous portions

- provides better performance
- the variable size avoids waste
- the file allocation tables are small



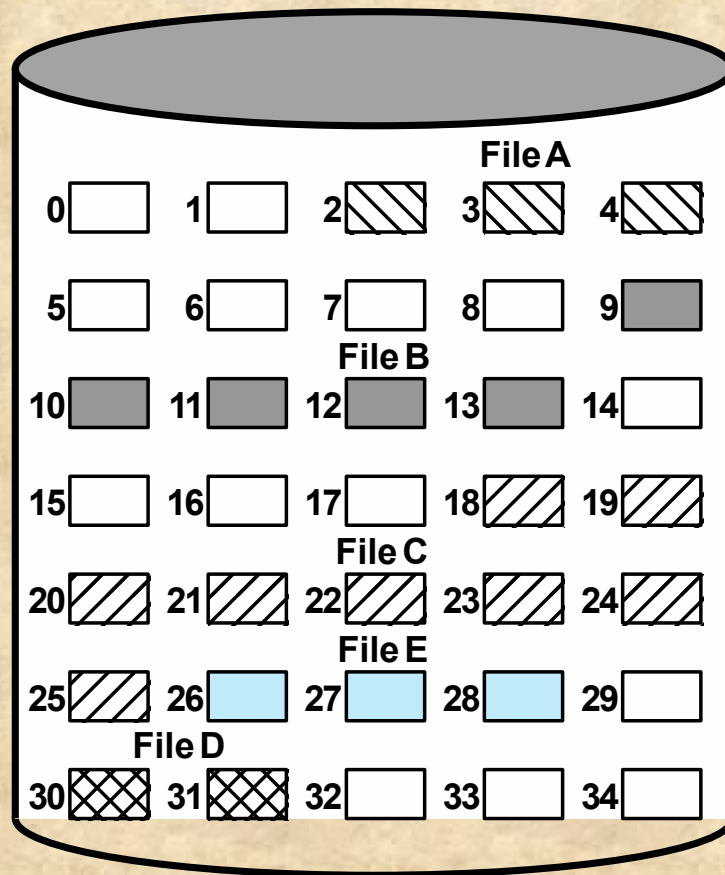
Blocks

- 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

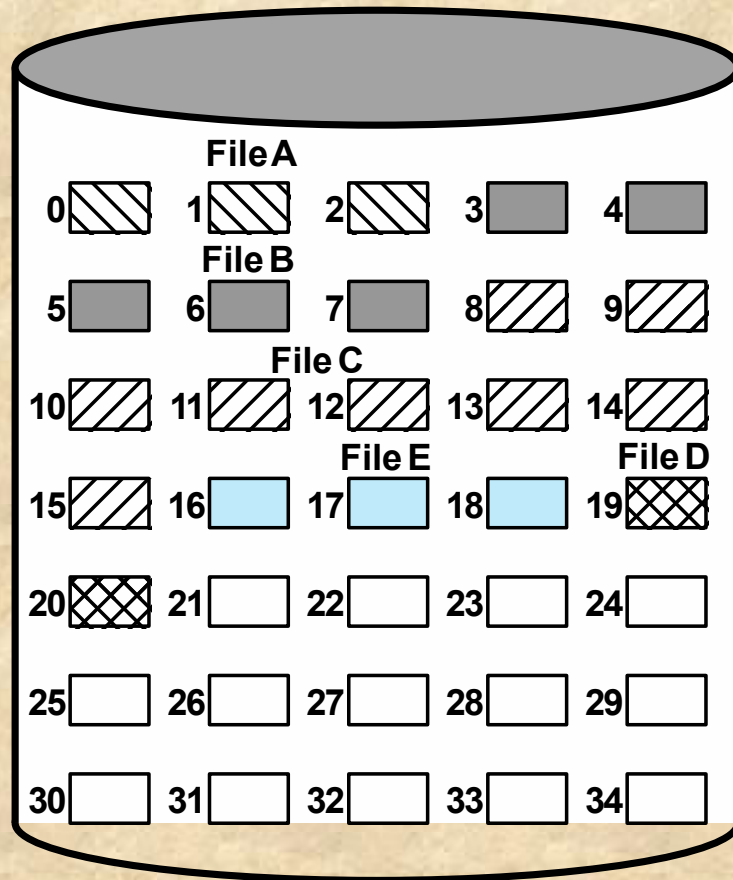
File Allocation Methods

	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

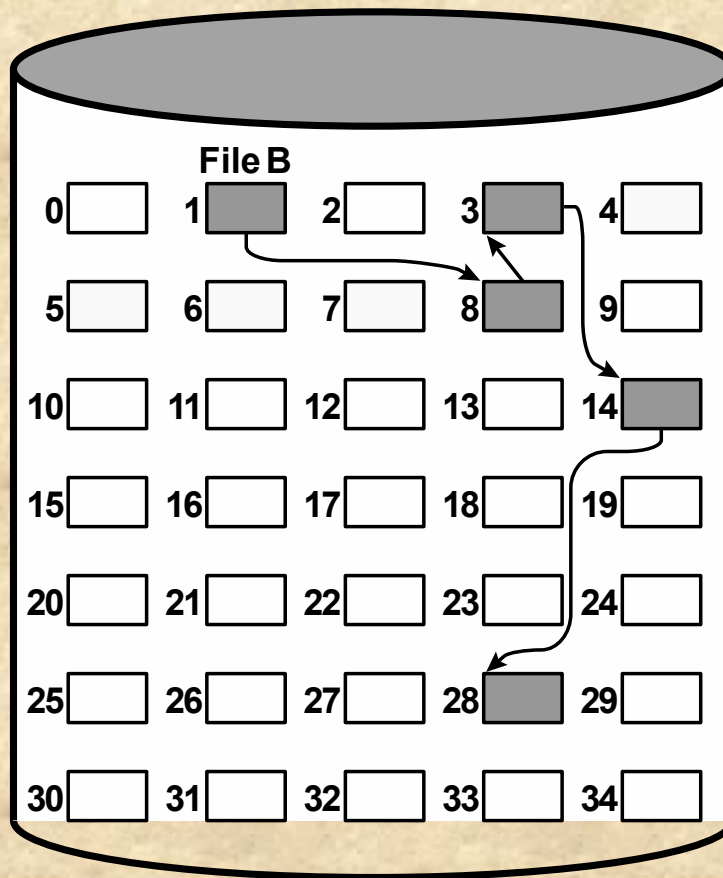
File Name	Start Block	Length
File A	2	3
File B	9	5
File C	18	8
File D	30	2
File E	26	3



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

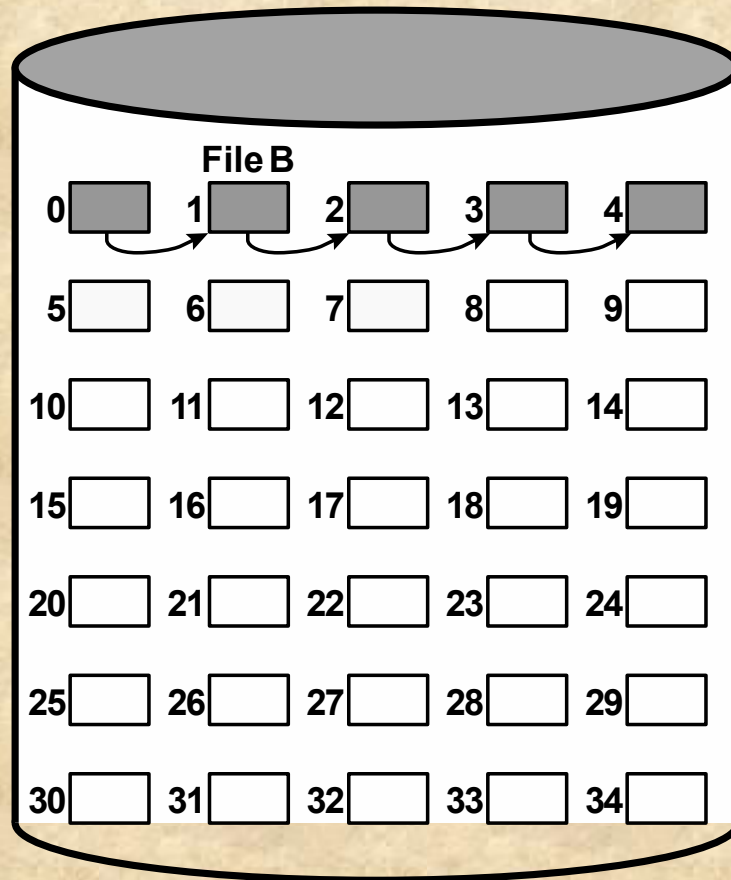
52 **Figure 12.10 Contiguous File Allocation (After Compaction)**



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
...

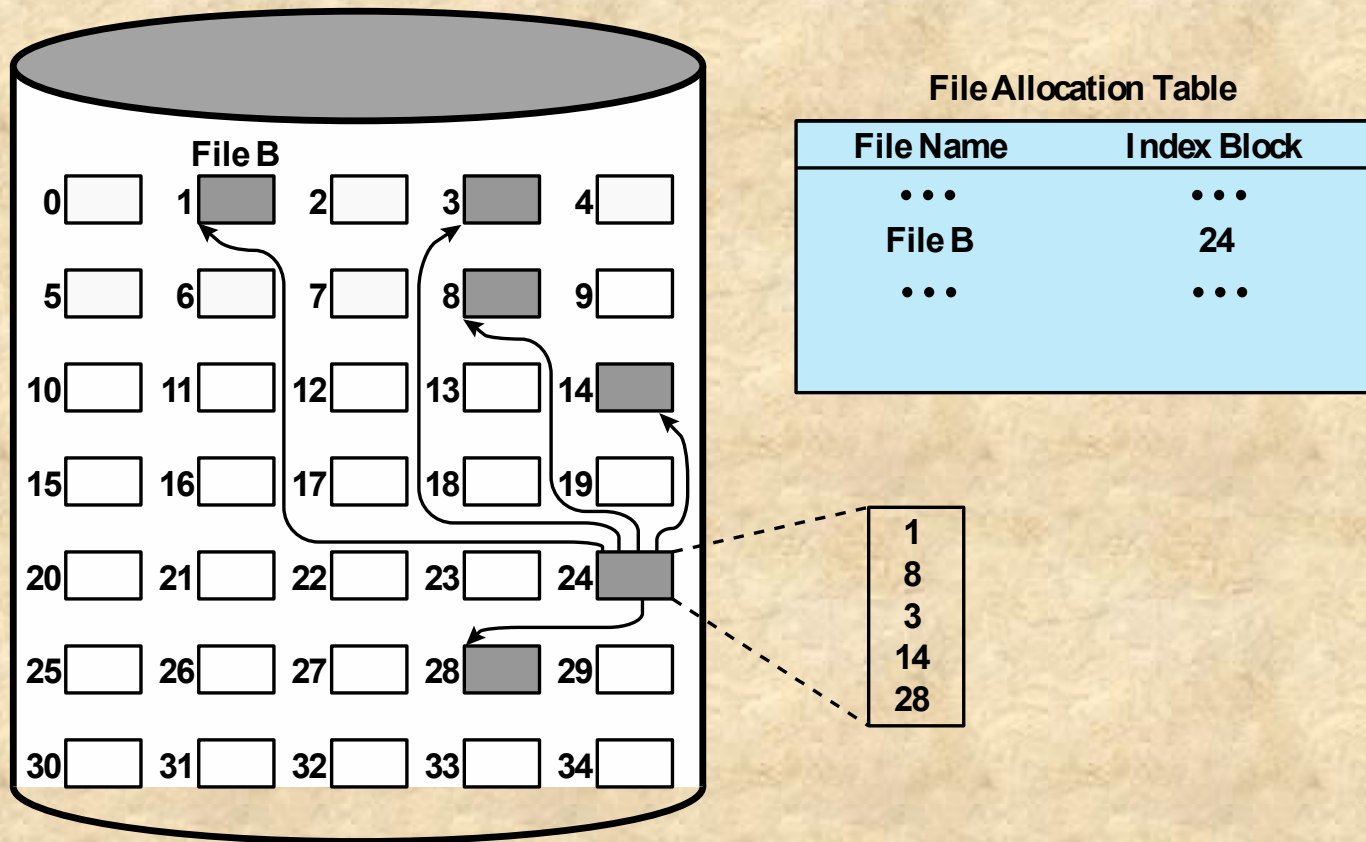
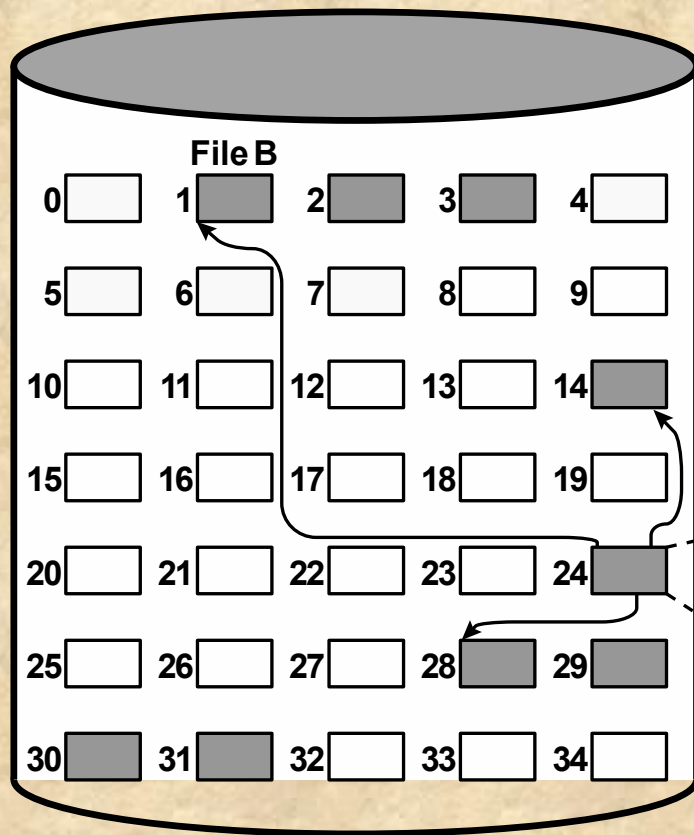


Figure 12.13 Indexed Allocation with Block Portions



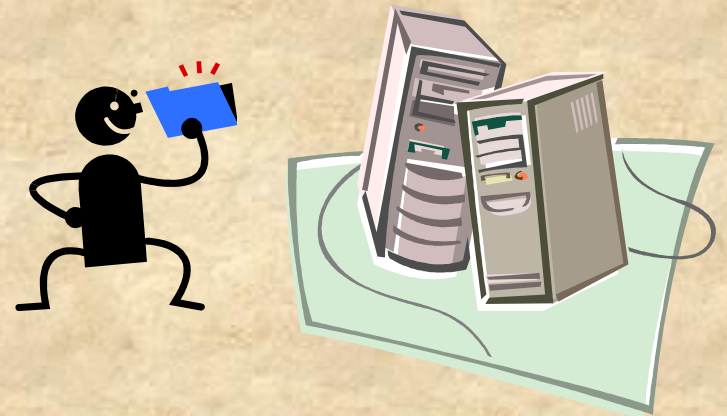
File Allocation Table	
File Name	Index Block
...	...
File B	24
...	...

Start Block	Length
1	3
28	4
14	1

56 **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

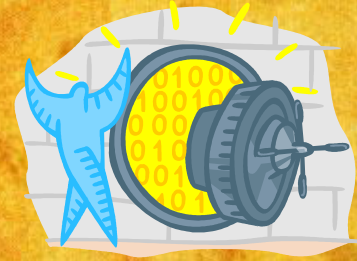
the size of the free block list is 24 or 32 times the size of the corresponding bit table and must be stored on disk

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

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 it is linked to

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



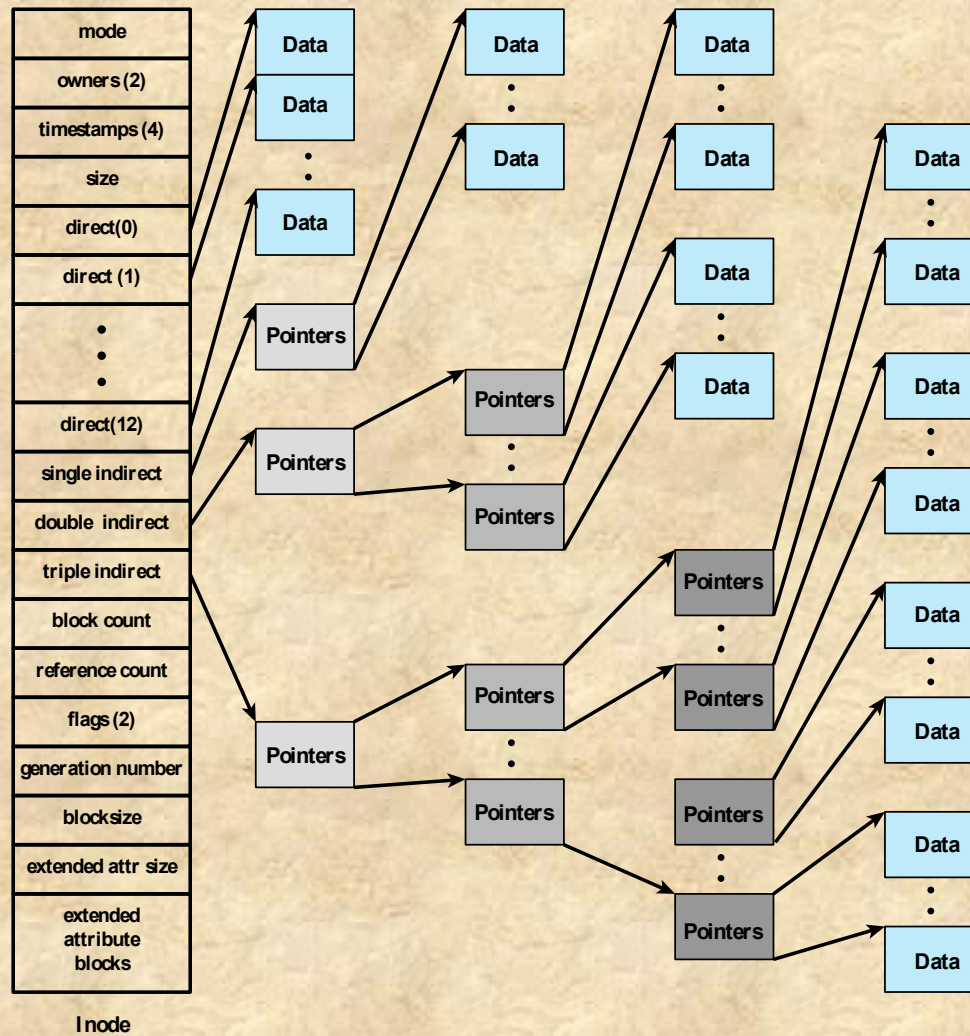


Figure 12.15 Structure of FreeBSD inode and File

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)

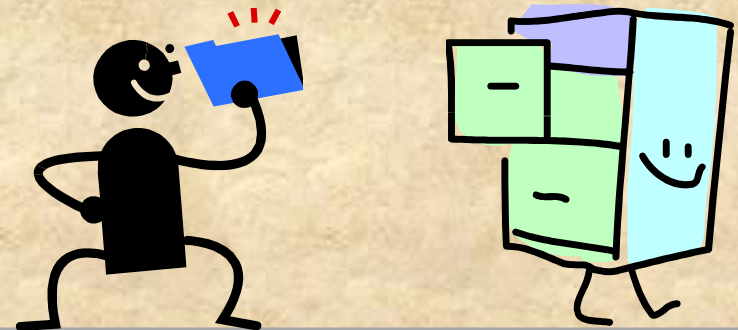
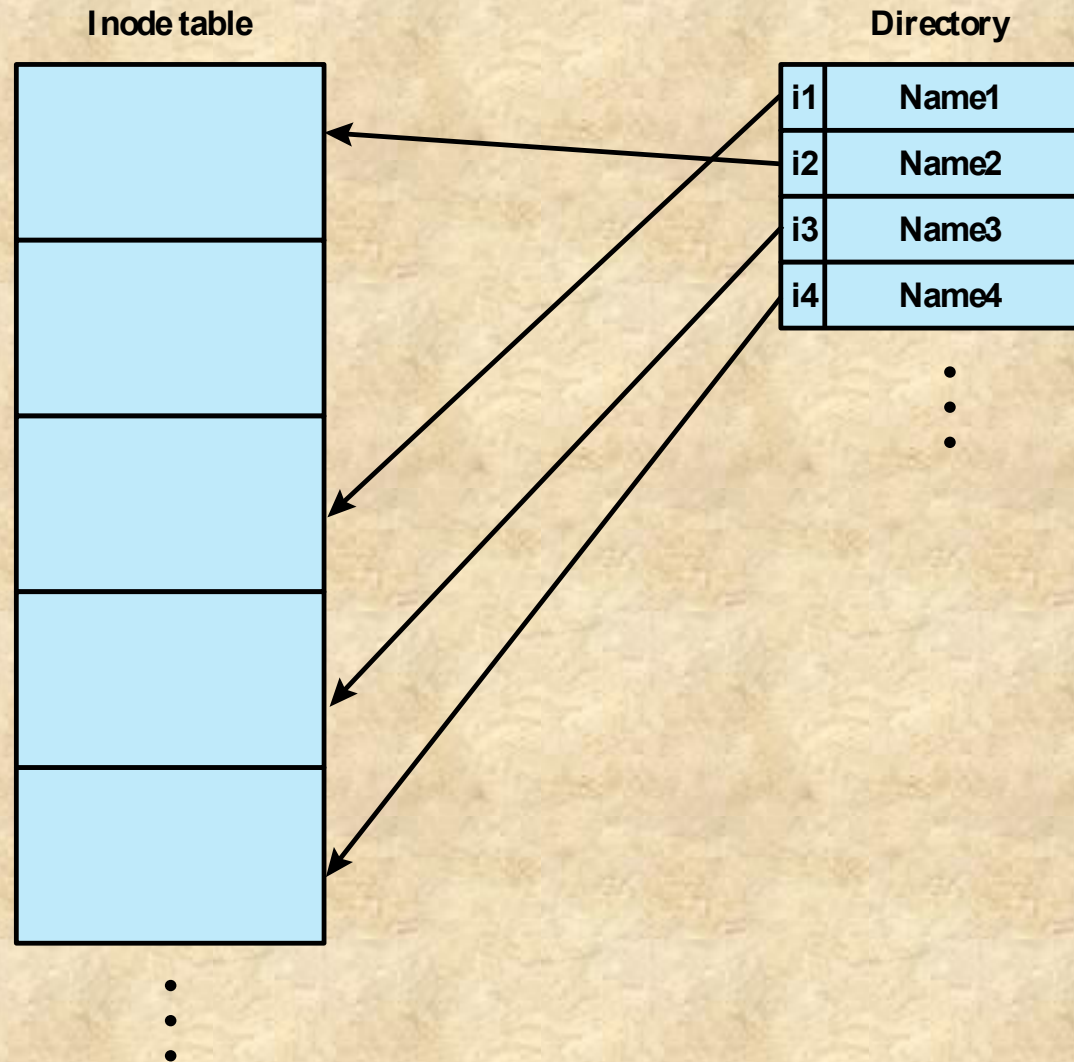


Table 12.3

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	512 512 = 256K	1G
Triple Indirect	512 256K = 128M	512G



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

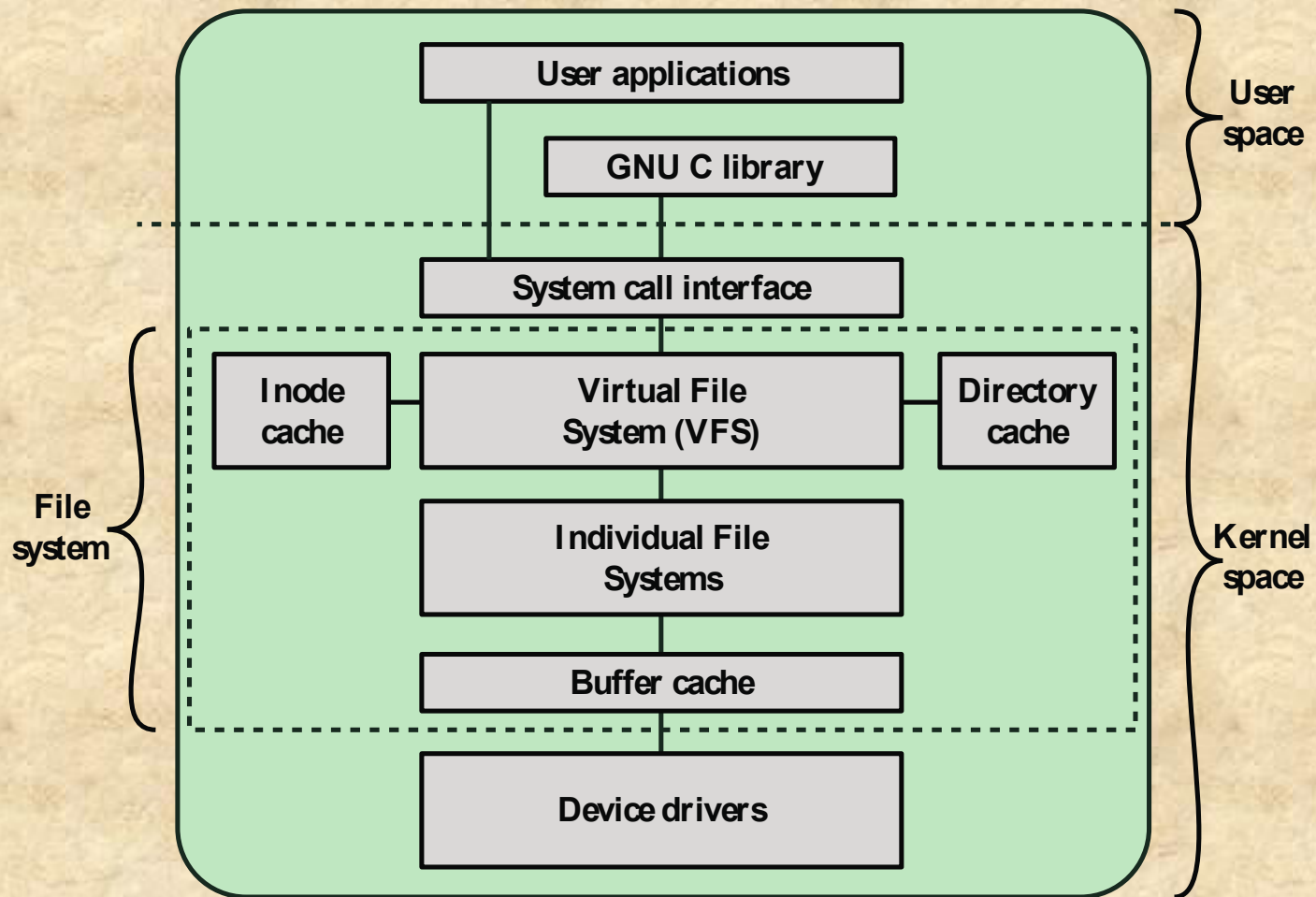


Figure 12.17 Linux Virtual File System Context

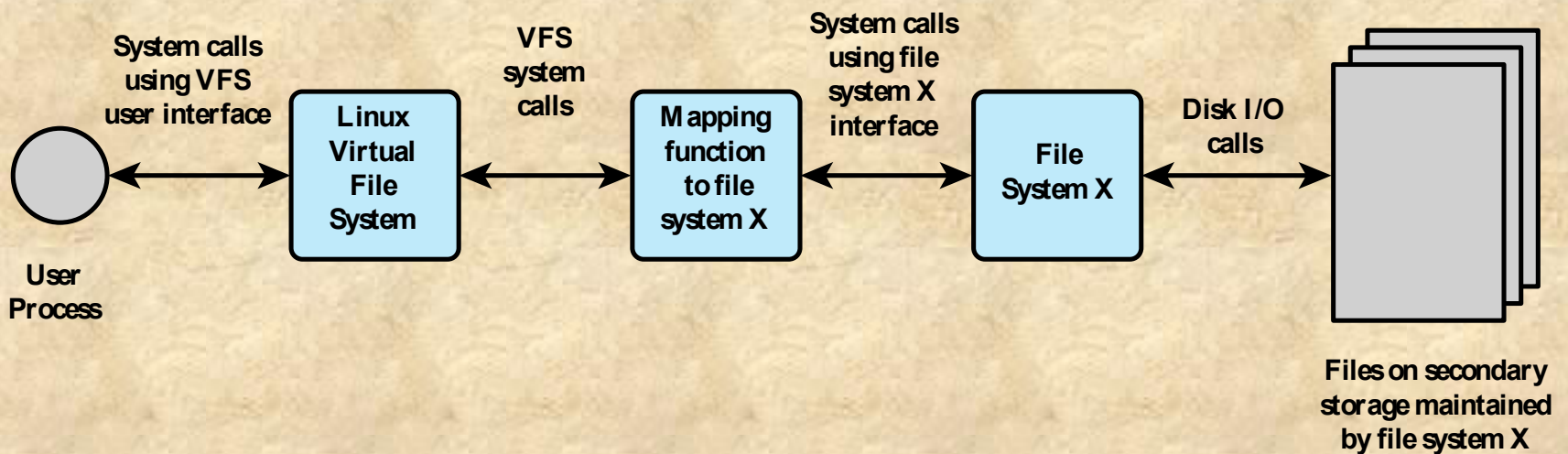


Figure 12.18 Linux Virtual File System Concept

Primary Object Types in VFS

Superblock Object

- represents a specific mounted file system

Inode Object

- represents a specific file



Dentry Object

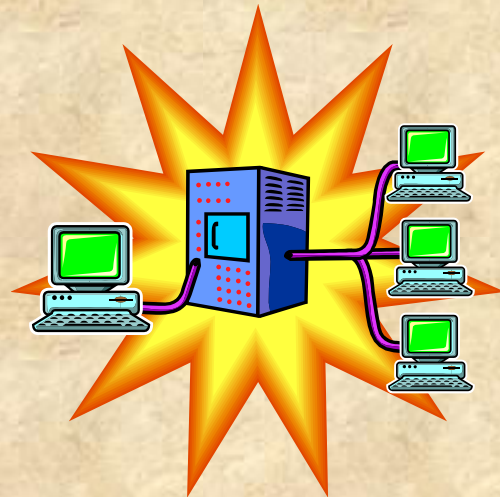
- represents a specific directory entry

File Object

- represents an open file associated with a process

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 264 bytes

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

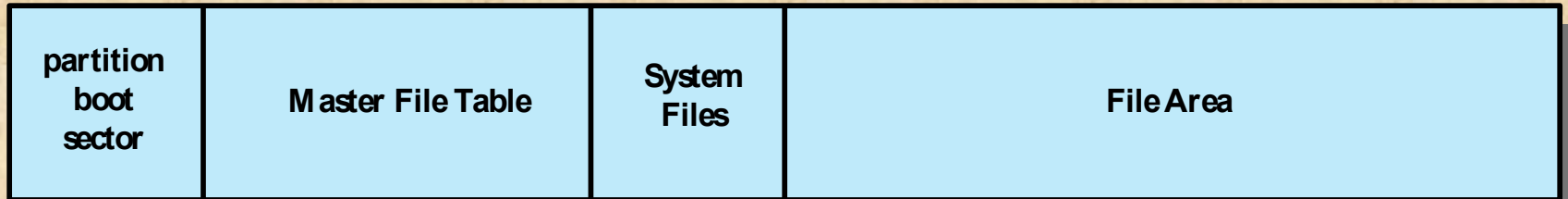


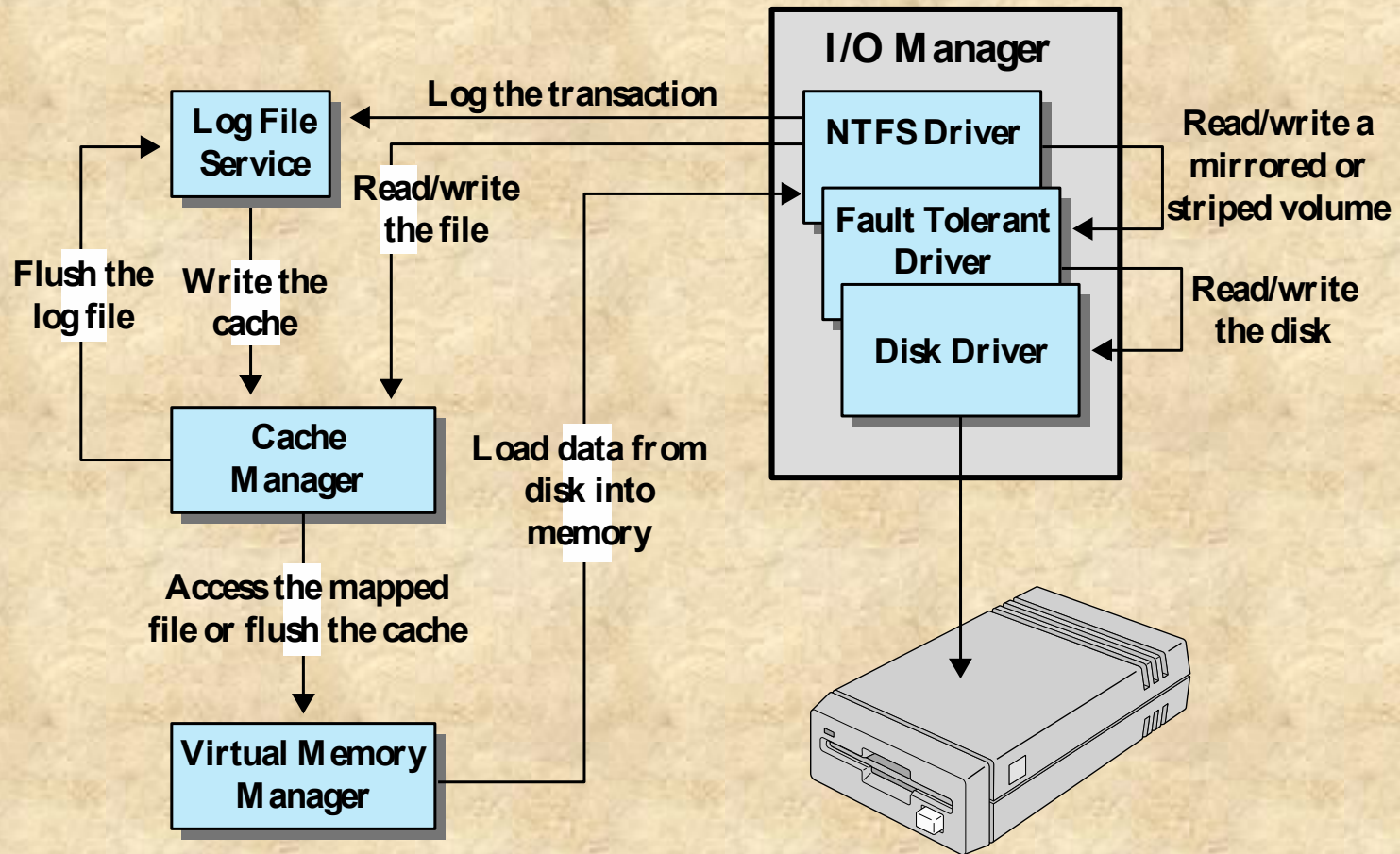
Figure 12.19 NTFS Volume Layout

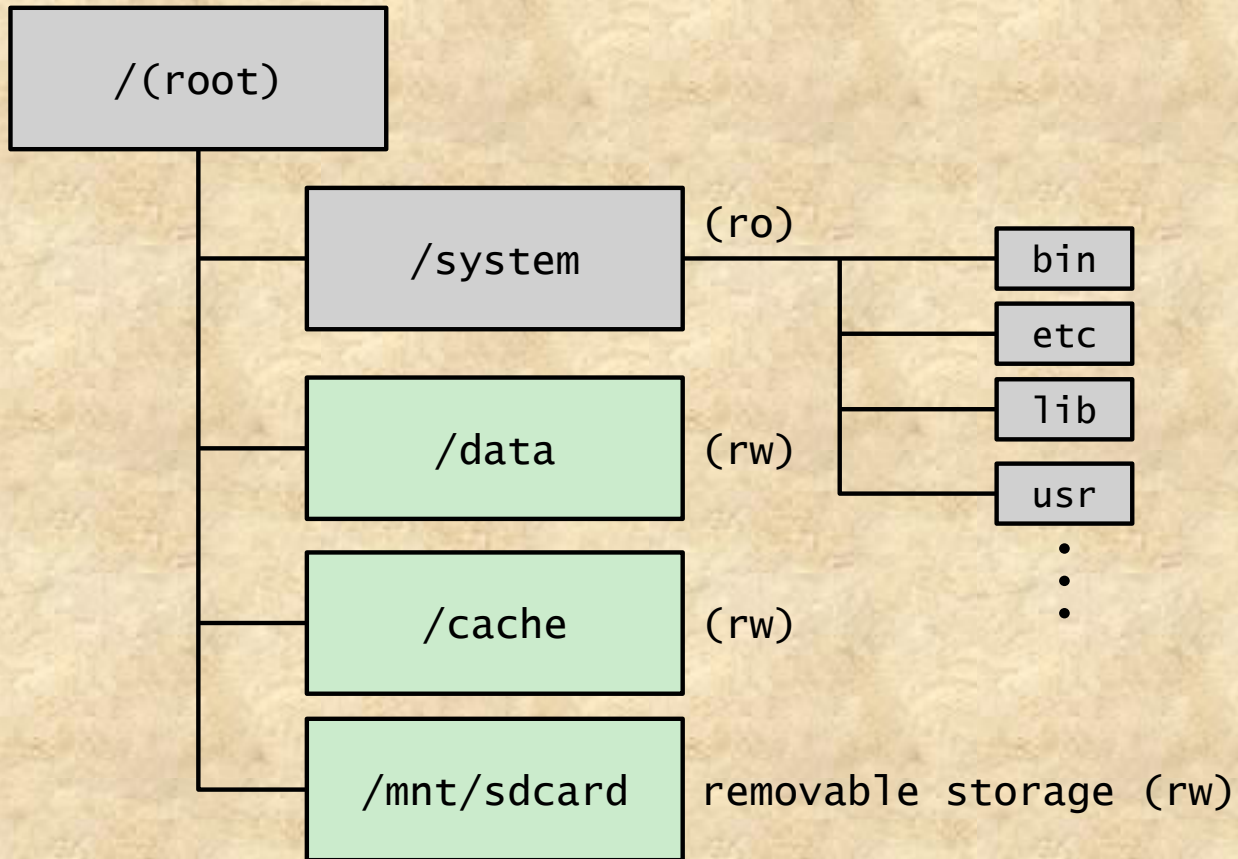
Master File Table (MFT)

- 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.





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

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 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
- 82 ■ 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