# 12. File-System Implementation

*Sungyoung Lee*

*College of Engineering*

*KyungHee University*

# Contents

- *File System Structure*
- *File System Implementation*
- *Directory Implementation*
- *Allocation Methods*
- *Free-Space Management*
- *Efficiency and Performance*
- *Recovery*
- *Log-Structured File Systems*
- *NFS*

**n** User's view on file systems:

  **ü** How files are named?

  **ü** What operations are allowed on them?

  **ü** What the directory tree looks like?

**n** Implementor's view on file systems:

  **ü** How files and directories are stored?

  **ü** How disk space is managed?

  **ü** How to make everything work efficiently and reliably?

# File-System Structure

**n** File structure

    ü Logical storage unit
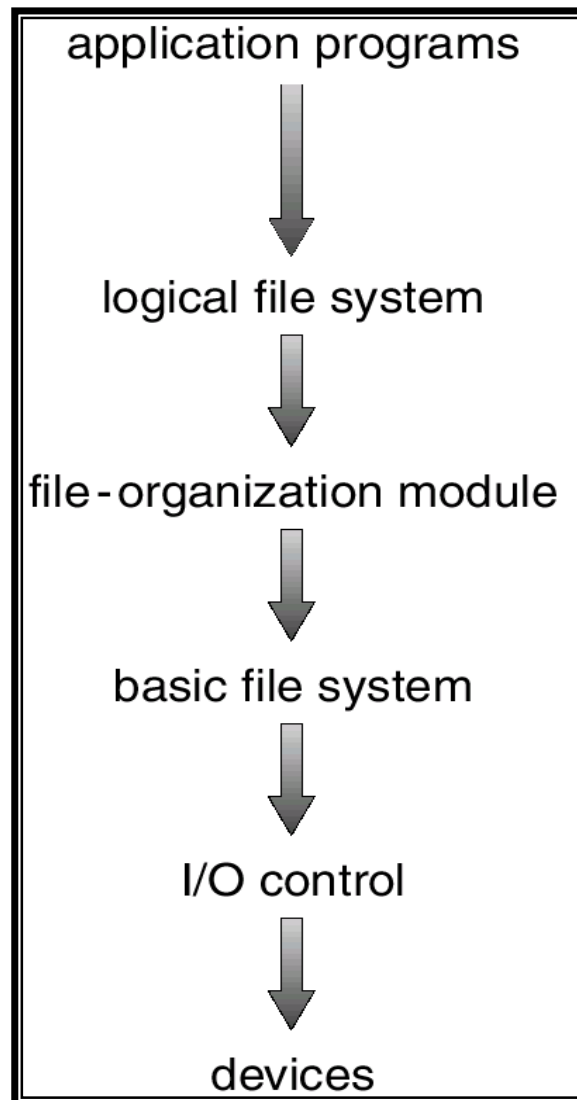
    ü Collection of related information

**n** File system resides on secondary storage (disks)
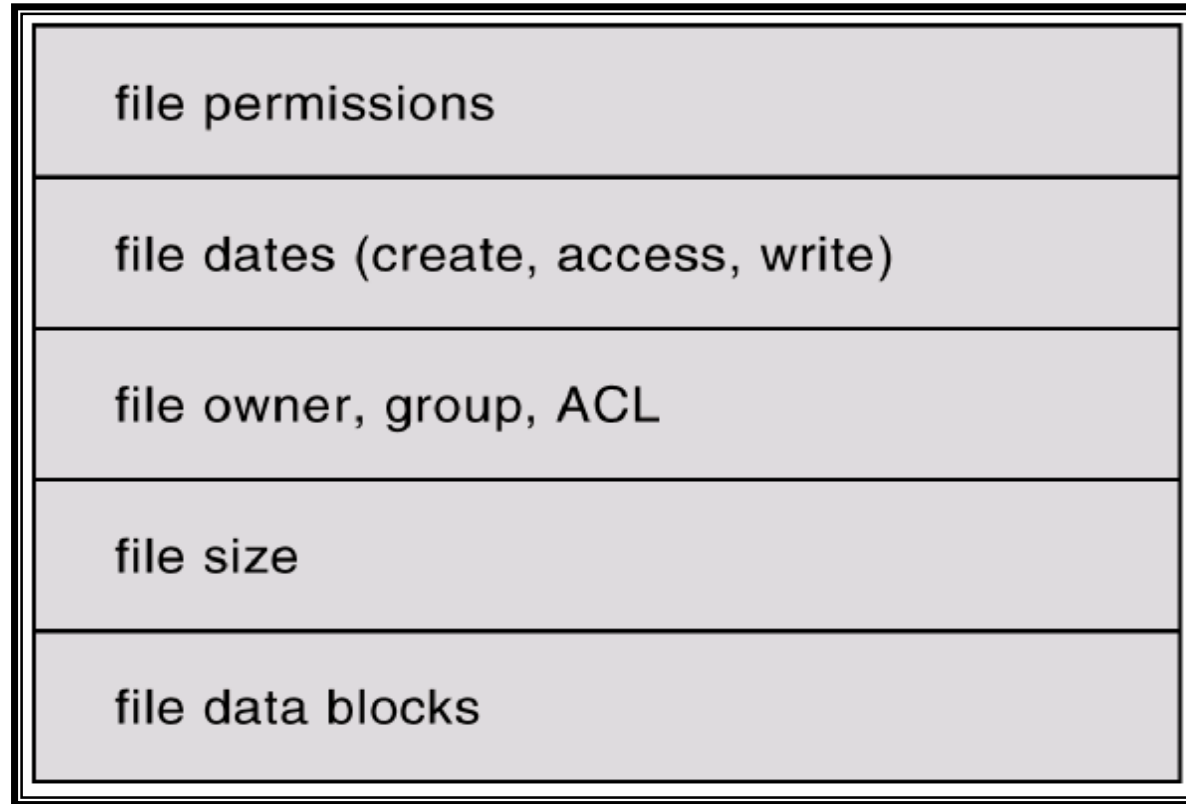
**n** File system organized into layers

**n** *File control block*

    ü storage structure consisting of information about a file

# Layered File System

# A Typical File Control Block

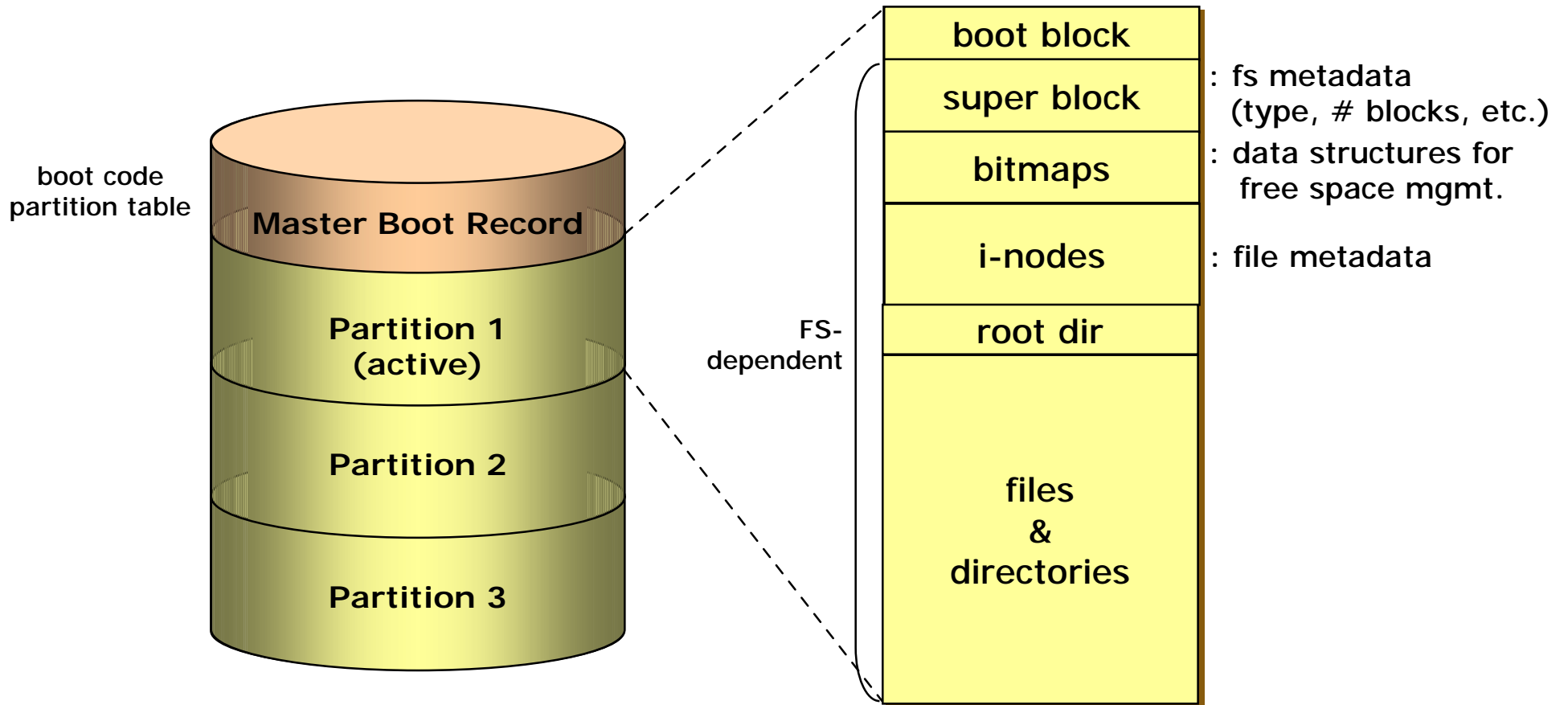| |
|---|
| file permissions |
| file dates (create, access, write) |
| file owner, group, ACL |
| file size |
| file data blocks |

**n** On-disk structure

- ü Boot control block
  - § Boot block(UFS) or Boot sector(NTFS)
- ü Partition control block
  - § Super block(UFS) or Master file table(NTFS)
- ü Directory structure
- ü File control block (FCB)
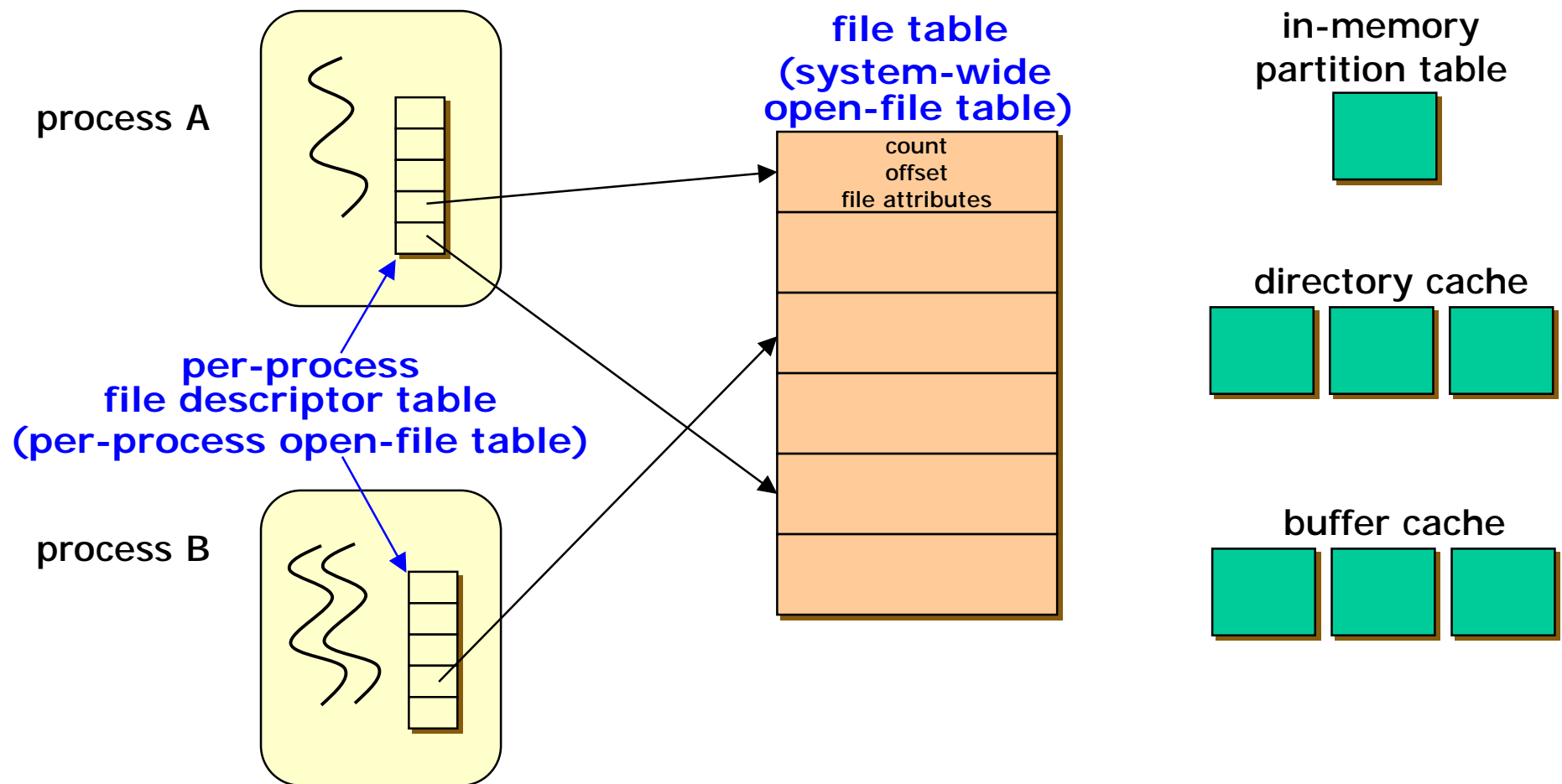  - § I-node(UFS) or In master file table(NTFS)

**n** In-memory structure

- ü In-memory partition table
- ü In-memory directory structure
- ü System-wide open file table
- ü Per-process open file table

boot code
partition table

Master Boot Record

Partition 1
(active)

Partition 2

Partition 3

FS-
dependent

boot block

super block

bitmaps

i-nodes

root dir

files
&
directories

: fs metadata
(type, # blocks, etc.)

: data structures for
free space mgmt.

: file metadata

process A

per-process
file descriptor table
(per-process open-file table)

process B

**file table
(system-wide
open-file table)**

count
offset
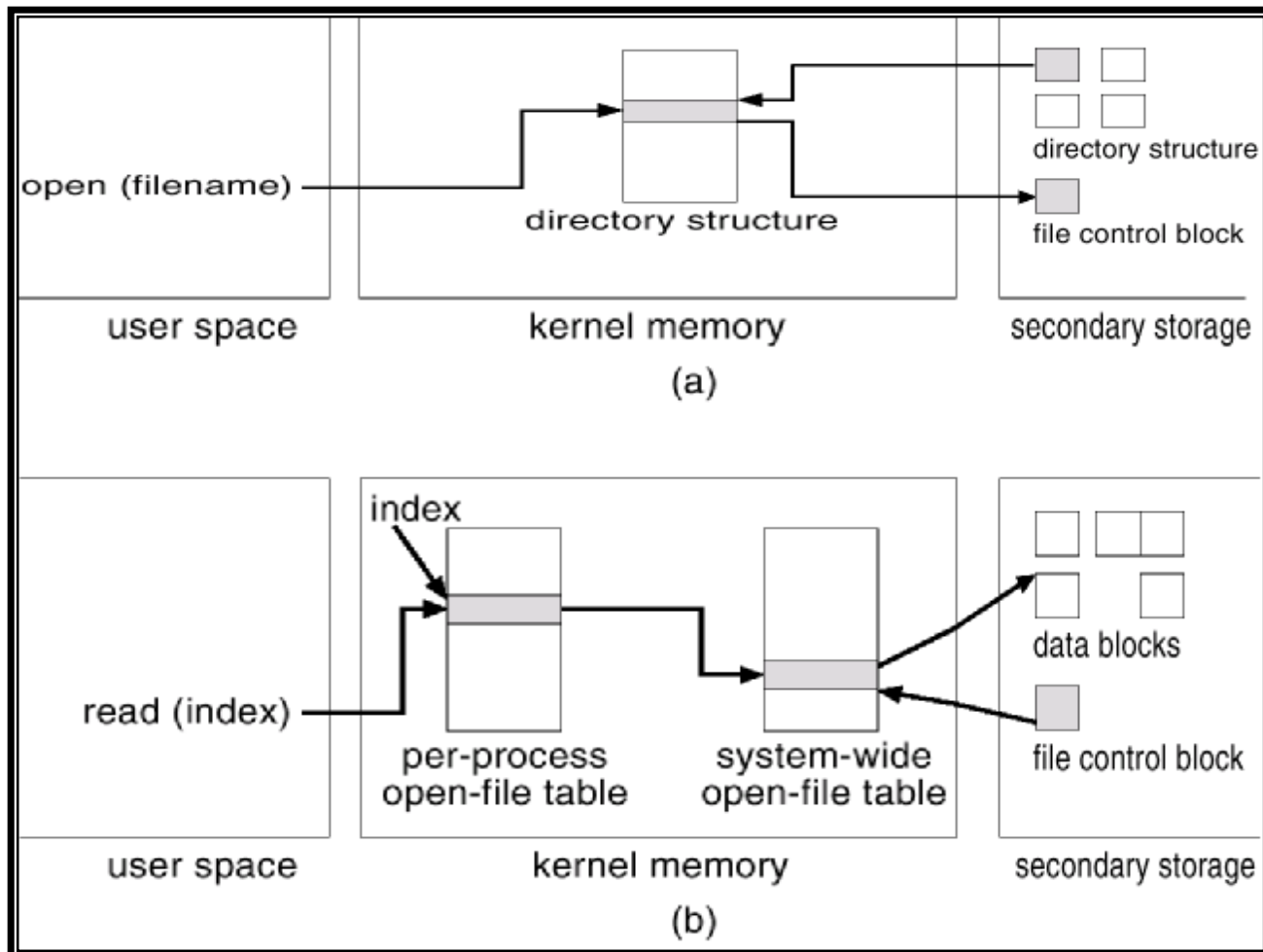file attributes

in-memory
partition table

directory cache

buffer cache

# In-Memory File System Structures

n  The following figure illustrates the necessary file system structures provided by the operating systems

n  Figure 12-3(a) refers to opening a file
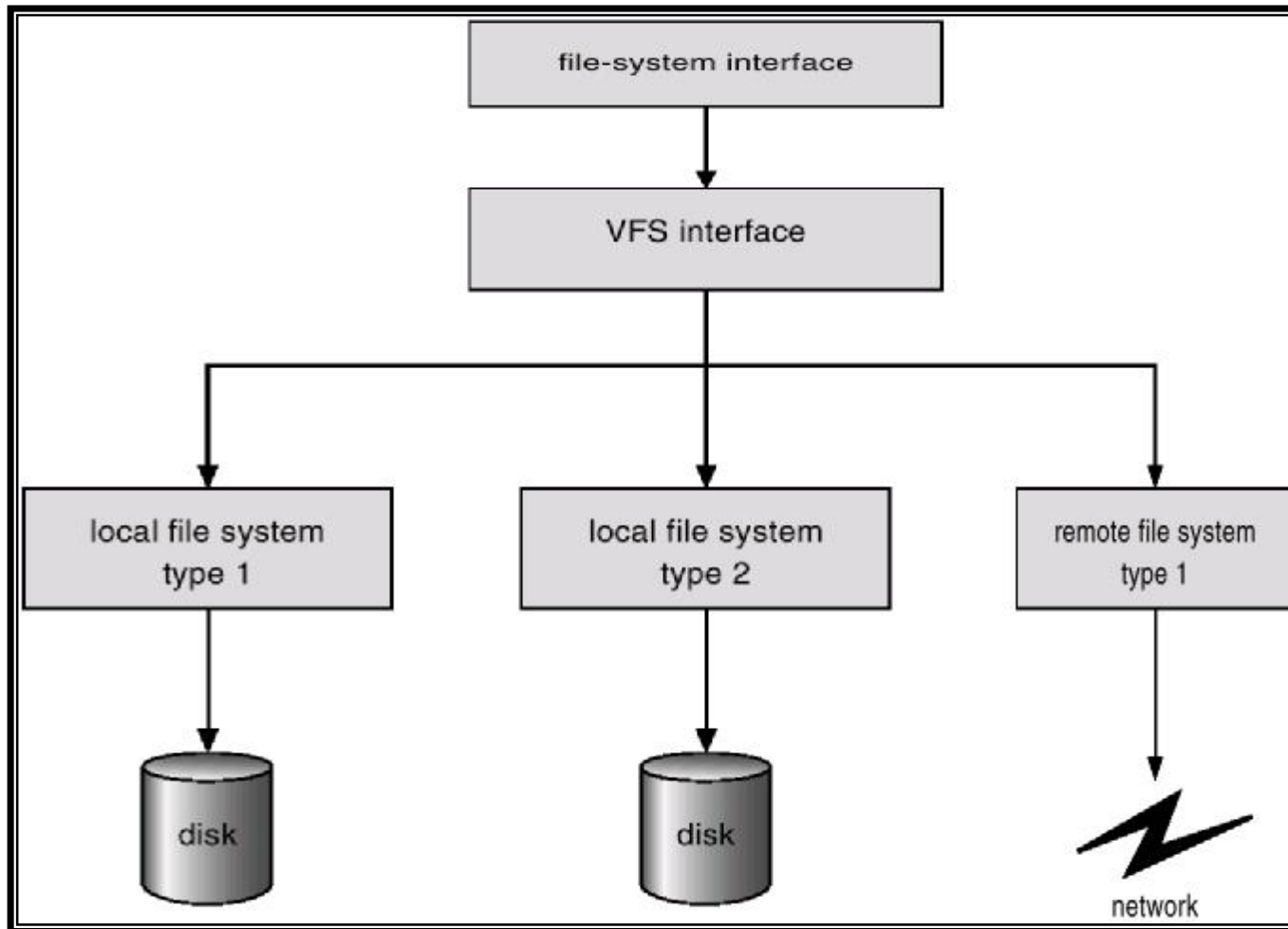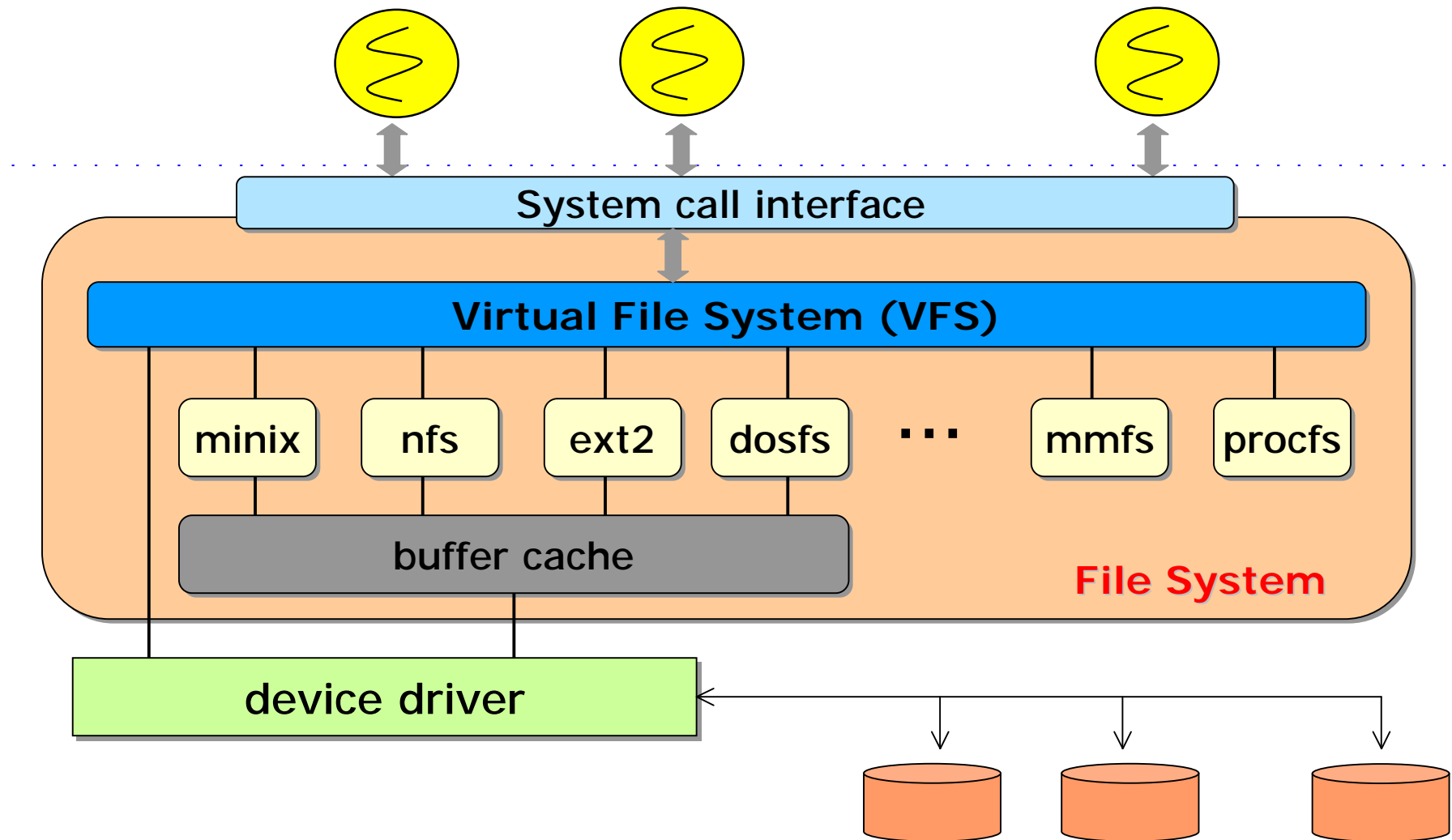
n  Figure 12-3(b) refers to reading a file

# Virtual File Systems

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

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

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

**n** Virtual File System

- **ü** Manages kernel-level file abstractions in one format for all file systems
- **ü** Receives system call requests from user-level (e.g., open, write, stat, etc.)
- **ü** Interacts with a specific file system based on mount point traversal
- **ü** Receives requests from other parts of the kernel, mostly from memory management
- **ü** Translates file descriptors to VFS data structures (such as vnode)

**n** Linux: VFS common file model

- **ü** The superblock object
    - **§** stores information concerning a mounted file system
- **ü** The inode object
    - **§** stores general information about a specific file
- **ü** The file object
    - **§** stores information about the interaction between an open file and a process
- **ü** The dentry object
    - **§** stores information about the linking of a directory entry with the corresponding file

# Directory Implementation

**n** Linear list of file names with pointer to the data blocks

    ü simple to program
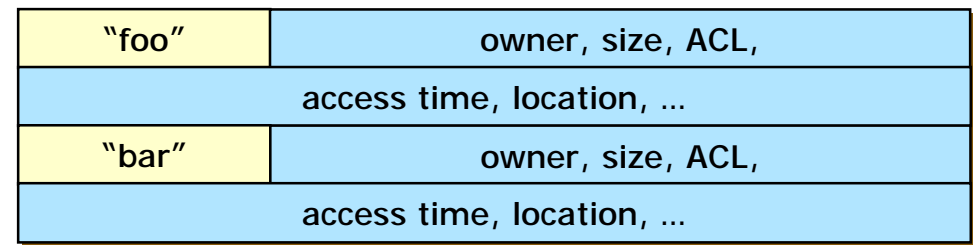
    ü time-consuming to execute

**n** Hash Table

    ü linear list with hash data structure

    ü decreases directory search time

    ü *collisions*

        § situations where two file names hash to the same location

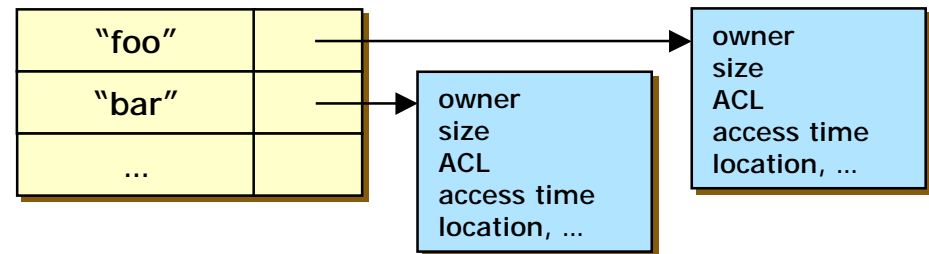    ü fixed size

**n** The location of metadata

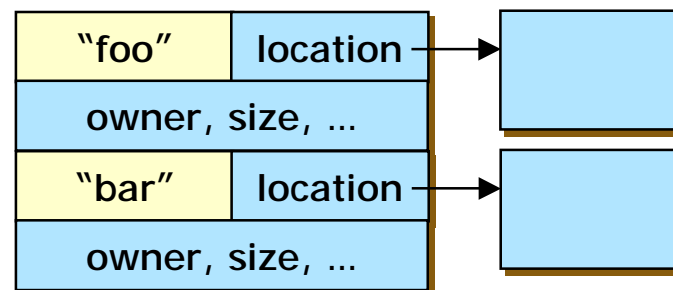   **ü** In the directory entry

| "foo" | owner, size, ACL, |
|---|---|
| access time, location, ... | |
| "bar" | owner, size, ACL, |
| access time, location, ... | |

   **ü** In the separate
data structure
(e.g., i-node)

| "foo" | | → | owner<br>size<br>ACL<br>access time<br>location, ... |
| "bar" | | → owner<br>size<br>ACL<br>access time<br>location, ... | |
| ... | | | |

   **ü** A hybrid approach

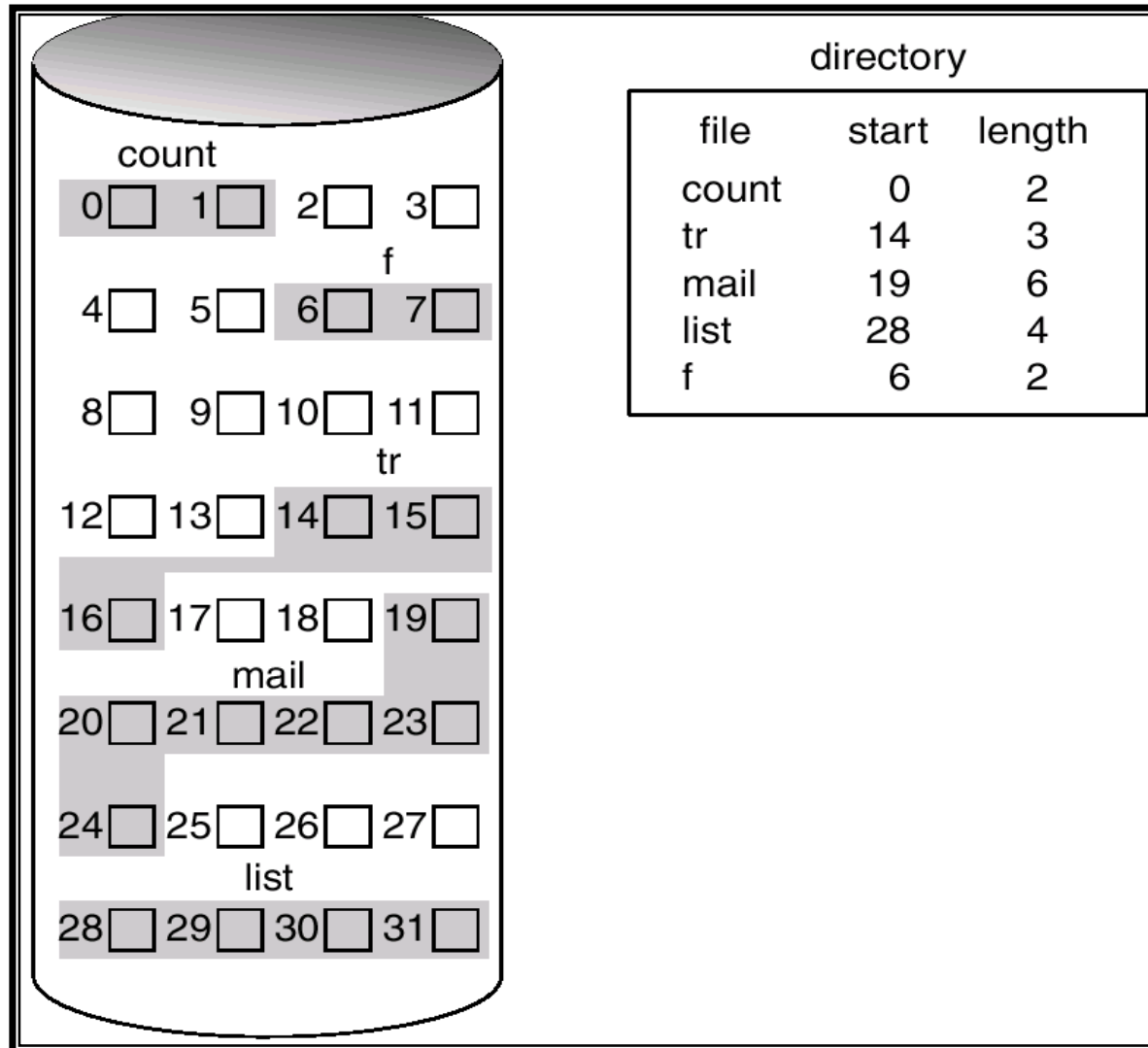| "foo" | location → | |
| owner, size, ... | | |
| "bar" | location → | |
| owner, size, ... | | |

# Allocation Methods

**n** An allocation method refers to how disk blocks are allocated for files

&#252; Contiguous allocation

&#252; Linked allocation

&#252; Indexed allocation

# Contiguous Allocation

**n** Each file occupies a set of contiguous blocks on the disk

**n** Simple
- **ü** only starting location (block #) and length (number of blocks) are required

**n** Random access

**n** Wasteful of space (dynamic storage-allocation problem)

**n** Files cannot grow

# Extent-Based Systems

**n**  Many newer file systems (I.e. Veritas File System) use a modified contiguous allocation scheme

**n**  Extent-based file systems allocate disk blocks in **extents**

**n**  An **extent** is a contiguous block of disks

    ü  Extents are allocated for file allocation

    ü  A file consists of one or more extents

**n** Advantages
- ü The number of disk seeks is minimal
- ü Directory entries can be simple:

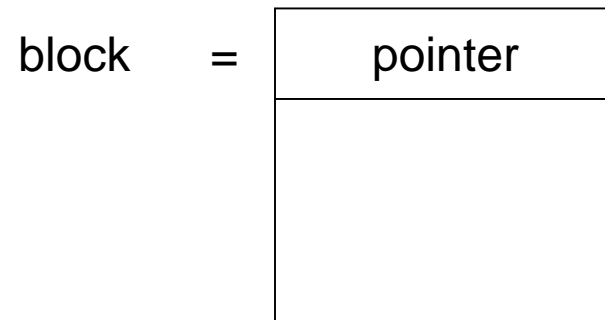  <file name, starting disk block, length, etc.>

**n** Disadvantages
- ü Requires a dynamic storage allocation: First / best fit
- ü External fragmentation: may require a compaction
- ü The file size is hard to predict and varying over time

**n** Feasible and widely used for CD-ROMS
- ü All the file sizes are known in advance
- ü Files will never change during subsequent use

# Linked Allocation

n  Each file is a linked list of disk blocks

    ü  Blocks may be scattered anywhere on the disk

block  =  | pointer |
          |         |

# Linked Allocation (Cont'd)

**n** Simple
- ü need only starting address

**n** Free-space management system
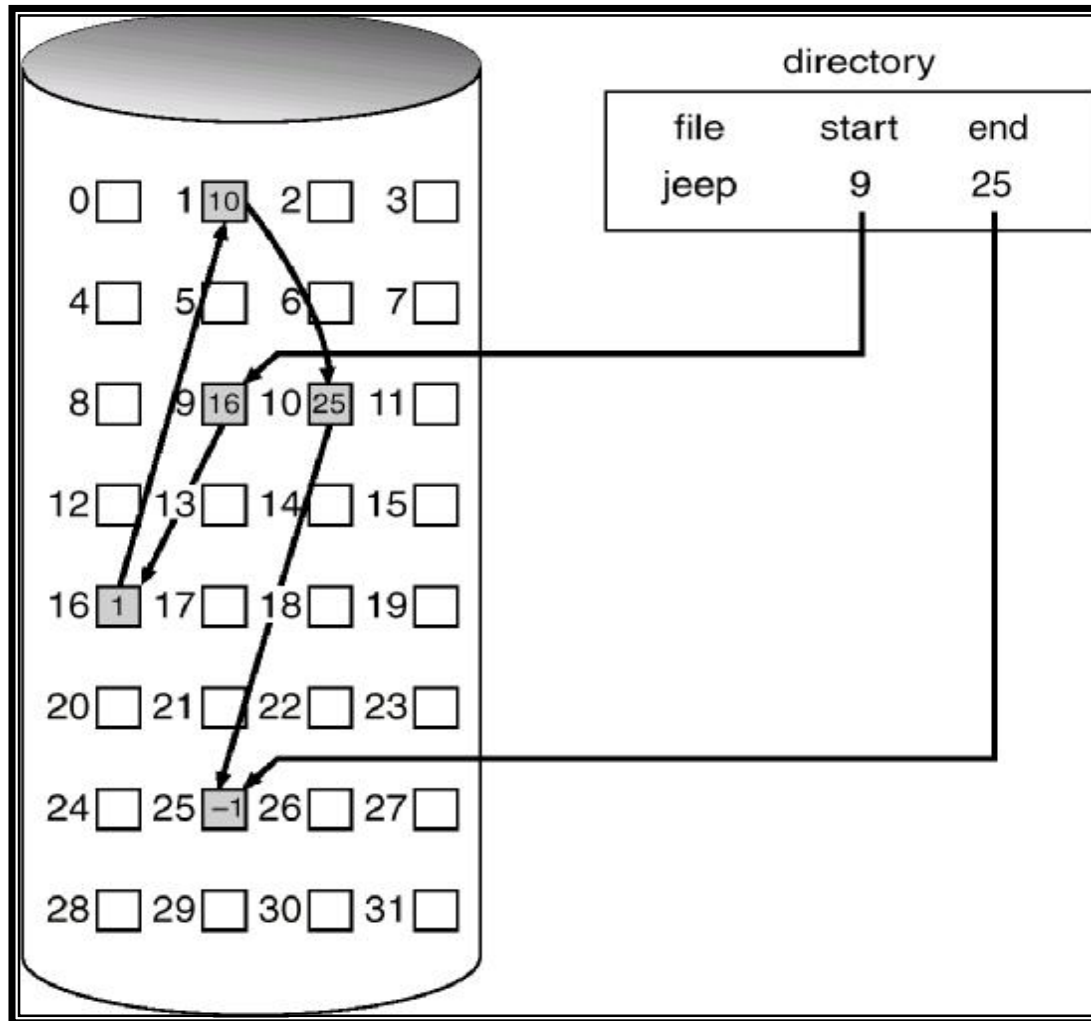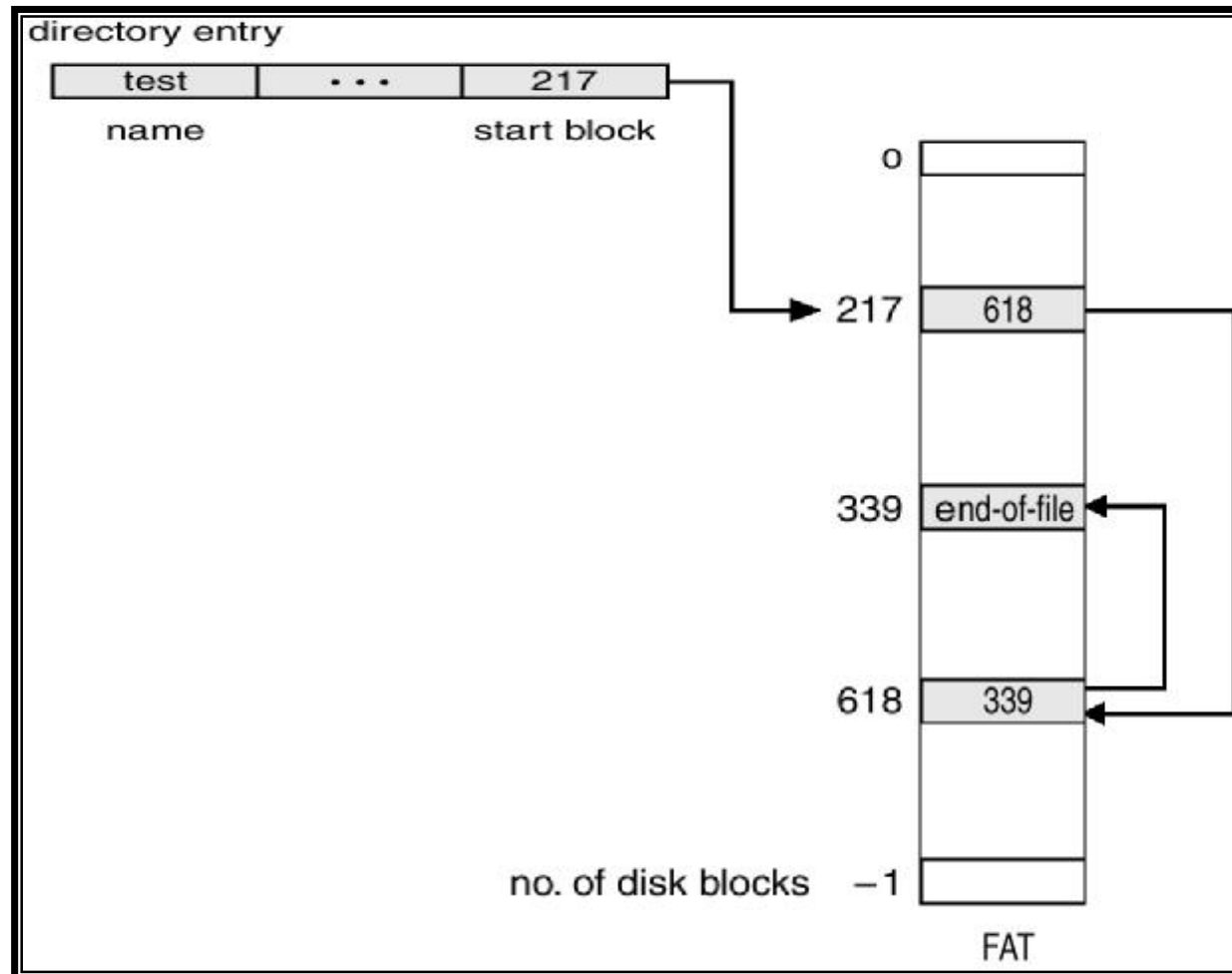- ü no waste of space

**n** No random access

**n** File-allocation table (FAT)
- ü disk-space allocation used by MS-DOS and OS/2

**n** Advantages

  ü Directory entries are simple:

  <file name, starting block, ending block, etc.>

  ü No external fragmentation

  § the disk blocks may be scattered anywhere on the disk

  ü A file can continue to grow as long as free blocks are available

**n** Disadvantages

  ü It can be used only for sequentially accessed files

  ü Space overhead for maintaining pointers to the next disk block

  ü The amount of data storage in a block is no longer a power of two because the pointer takes up a few bytes
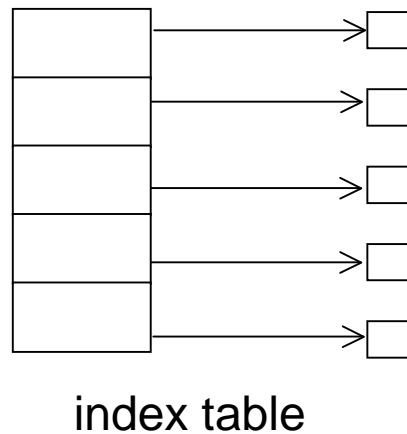
  ü Fragile: a pointer can be lost or damaged

**n** Collect blocks into multiples (clusters) and allocate the clusters to files

    ü e.g., 4 blocks / 1 cluster

**n** Advantages

    ü The logical-to-physical block mapping remains simple

    ü Improves disk throughput (fewer disk seeks)

    ü Reduced space overhead for pointers

**n** Disadvantages

    ü Internal fragmentation

# Indexed Allocation

**n** Brings all pointers together into the *index block*

**n** Logical view



index table

**n** Need index table

**n** Random access

**n** Dynamic access without external fragmentation, but have overhead of index block

M

outer-index

index table

file

# Combined Scheme:  UNIX (4K bytes per block)

**n** Advantages

ü Supports direct access, without suffering from external fragmentation

ü I-node need only be in memory when the corresponding file is open

**n** Disadvantages

ü Space overhead for indexes:

(1) Linked scheme: link several index blocks

(2) Multilevel index blocks

(3) Combined scheme: UNIX

- 12 direct blocks, single indirect block, double indirect block,

triple indirect block

# Free-Space Management

**n** Bit vector (*n* blocks)

```
 0  1   2                          n-1
┌──┬──┬──┬──┬──┬──┬─────────┬──┐
│  │  │  │  │  │  │   …     │  │
└──┴──┴──┴──┴──┴──┴─────────┴──┘
```

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

Block number calculation

(number of bits per word) *
(number of 0-value words) +
offset of first 1 bit

# Free-Space Management (Cont'd)

n Bit map requires extra space.  Example:

$$\text{block size} = 2^{12} \text{ bytes}$$

$$\text{disk size} = 2^{30} \text{ bytes (1 gigabyte)}$$

$$n = 2^{30}/2^{12} = 2^{18} \text{ bits (or 32K bytes)}$$

n Easy to get contiguous files

n Linked list (free list)
ü Cannot get contiguous space easily
ü No waste of space

n Grouping

n Counting

**n** Grouping

   ü Store the addresses of *n* free blocks in the first free block

   ü The addresses of a large number of free blocks can be found quickly

**n** Counting

   ü Keep the address of the free block and the number of free contiguous blocks

   ü The length of the list becomes shorter and the count is generally greater than 1

      § Several contiguous blocks may be allocated or freed simultaneously

# Efficiency and Performance

**n** Efficiency dependent on:
- **ü** Disk allocation and directory algorithms
- **ü** Types of data kept in file's directory entry

**n** Performance
- **ü** Disk cache
  - **§** separate section of main memory for frequently used blocks
- **ü** Free-behind and Read-ahead
  - **§** techniques to optimize sequential access
- **ü** Improve PC performance by dedicating section of memory as virtual disk, or RAM disk

# Buffer Cache

n  Applications exhibit significant locality for reading and writing files

n  Idea: cache file blocks in memory to capture locality in buffer cache (or disk cache)

  ü  Cache is system wide, used and shared by all processes
  ü  Reading from the cache makes a disk perform like memory
  ü  Even a 4MB cache can be very effective

n  Issues

  ü  The buffer cache competes with VM
  ü  Like VM, it has limited size
  ü  Need replacement algorithms again
      (References are relatively infrequent, so it is feasible to keep all the blocks in exact LRU order)

# Page Cache

n   A **page cache** caches pages rather than disk blocks using virtual memory techniques

n   Memory-mapped I/O uses a page cache

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

n   This leads to the following figure

# Unified Buffer Cache

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

**n** Synchronous writes are very slow

**n** Asynchronous writes (or write-behind, write-back)
- ü Maintain a queue of uncommitted blocks
- ü Periodically flush the queue to disk
- ü Unreliable: metadata requires synchronous writes (with small files, most writes are to metadata)

- **n** File system predicts that the process will request next block
    - ü File system goes ahead and requests it from the disk
    - ü This can happen while the process is computing on previous block, overlapping I/O with execution
    - ü When the process requests block, it will be in cache

- **n** Compliments the disk cache, which also is doing read ahead

- **n** Very effective for sequentially accessed files

- **n** File systems try to prevent blocks from being scattered across the disk during allocation or by restructuring periodically

- **n** Cf) Free-behind

# Block Size Performance vs. Efficiency

**n** Block size

    **ü** Disk block size vs. file system block size

    **ü** The median file size in UNIX is about 1KB



(All files are 2KB)

# Recovery

**n** Consistency checking
  - ü compares data in directory structure with data blocks on disk, and tries to fix inconsistencies

**n** Use system programs to *back up* data from disk to another storage device (floppy disk, magnetic tape)

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

**n** **File system consistency**

ü File system can be left in an inconsistent state if cached blocks are not written out due to the system crash

ü It is especially critical if some of those blocks are i-node blocks, directory blocks, or blocks containing the free list

ü Most systems have a utility program that checks file system consistency

§ Windows: scandisk

§ UNIX: fsck

# Log Structured File Systems

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

n All transactions are written to a **log**
  ü A transaction is considered **committed** once it is written to the log
  ü However, the file system may not yet be updated

n The transactions in the log are asynchronously written to the file system
  ü When the file system is modified, the transaction is removed from the log

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

n **Journaling file systems**

   ü Fsck'ing takes a long time, which makes the file system restart slow in the event of system crash

   ü Record a log, or journal, of changes made to files and directories to a separate location (preferably a separate disk)

   ü If a crash occurs, the journal can be used to undo any partially completed tasks that would leave the file system in an inconsistent state

   ü IBM JFS for AIX, Linux

      Veritas VxFS for Solaris, HP-UX, Unixware, etc.

      SGI XFS for IRIX, Linux

      Reiserfs, ext3 for Linux

      NTFS for Windows

# The Sun Network File System (NFS)

n  An implementation and a specification of a software system for accessing remote files across LANs (or WANs)

n  The implementation is part of the Solaris and SunOS operating systems running on Sun workstations using an unreliable datagram protocol (UDP/IP protocol and Ethernet)

# NFS (Cont'd)

**n** Interconnected workstations viewed as a set of independent machines with independent file systems, which allows sharing among these file systems in a transparent manner

- ü A remote directory is mounted over a local file system directory. The mounted directory looks like an integral subtree of the local file system, replacing the subtree descending from the local directory

- ü Specification of the remote directory for the mount operation is nontransparent; the host name of the remote directory has to be provided. Files in the remote directory can then be accessed in a transparent manner

- ü Subject to access-rights accreditation, potentially any file system (or directory within a file system), can be mounted remotely on top of any local directory

# NFS (Cont'd)

n NFS is designed to operate in a heterogeneous environment of different machines, operating systems, and network architectures

    ü The NFS specifications independent of these media

n This independence is achieved through the use of RPC primitives built on top of an External Data Representation (XDR) protocol used between two implementation-independent interfaces

n The NFS specification distinguishes between the services provided by a mount mechanism and the actual remote-file-access services

# Three Independent File Systems

Mounts         Cascading mounts

# NFS Mount Protocol

- **Establishes initial logical connection between server and client**
- **Mount operation includes name of remote directory to be mounted and name of server machine storing it**
    - Mount request is mapped to corresponding RPC and forwarded to mount server running on server machine
    - *Export list* – specifies local file systems that server exports for mounting, along with names of machines that are permitted to mount them
- **Following a mount request that conforms to its export list, the server returns a *file handle* — a key for further accesses**
- **File handle**
    - A file-system identifier, and an inode number to identify the mounted directory within the exported file system
- **The mount operation changes only the user's view and does not affect the server side**

# NFS Protocol

**n** Provides a set of remote procedure calls for remote file operations

**n** The procedures support the following operations:

    ü searching for a file within a directory

    ü reading a set of directory entries

    ü manipulating links and directories

    ü accessing file attributes

    ü reading and writing files

**n** NFS servers are *stateless*

    ü Each request has to provide a full set of arguments

**n** Modified data must be committed to the server's disk before results are returned to the client (lose advantages of caching)

**n** The NFS protocol does not provide concurrency-control mechanisms

# Three Major Layers of NFS Architecture

**n** UNIX file-system interface

   **ü** Based on the **open, read, write**, and **close** calls, and file descriptors

**n** *Virtual File System* (VFS) layer

   **ü** Distinguishes local files from remote ones, and local files are further distinguished according to their file-system types

   **ü** The VFS activates file-system-specific operations to handle local requests according to their file-system types

   **ü** Calls the NFS protocol procedures for remote requests

**n** NFS service layer

   **ü** Bottom layer of the architecture

   **ü** Implements the NFS protocol

# NFS Path-Name Translation

n   Performed by breaking the path into component names and performing a separate NFS lookup call for every pair of component name and directory vnode

n   To make lookup faster, a directory name lookup cache on the client's side holds the vnodes for remote directory names

# NFS Remote Operations

- **n** Nearly one-to-one correspondence between regular UNIX system calls and the NFS protocol RPCs (except opening and closing files)
- **n** NFS adheres to the remote-service paradigm, but employs buffering and caching techniques for the sake of performance
- **n** File-blocks cache
    - ü When a file is opened, the kernel checks with the remote server whether to fetch or revalidate the cached attributes
    - ü Cached file blocks are used only if the corresponding cached attributes are up to date
- **n** File-attribute cache
    - ü The attribute cache is updated whenever new attributes arrive from the server
- **n** Clients do not free delayed-write blocks until the server confirms that the data have been written to disk

# Appendix) Implementation Examples

**n** Fast file system (FFS)

- ü The original Unix file system (70's) was very simple and straightforwardly implemented:
  - § Easy to implement and understand
  - § But very poor utilization of disk bandwidth (lots of seeking)
- ü BSD Unix folks redesigned file system called FFS
  - § McKusick, Joy, Fabry, and Leffler (mid 80's)
  - § Now it is the file system from which all other UNIX file systems have been compared
- ü The basic idea is aware of disk structure
  - § Place related things on nearby cylinders to reduce seeks
  - § Improved disk utilization, decreased response time

**n** Data and i-node placement

ü Original Unix FS had two major problems:

(1) Data blocks are allocated randomly in aging file systems

§ Blocks for the same file allocated sequentially when FS is new

§ As FS "ages" and fills, need to allocate blocks freed up when other files are deleted

§ Problem: Deleted files essentially randomly placed

§ So, blocks for new files become scattered across the disk

(2) i-nodes are allocated far from blocks

§ All i-nodes at the beginning of disk, far from data

§ Traversing file name paths, manipulating files and directories require going back and forth from i-nodes to data blocks

ü Both of these problems generate many long seeks!

**n** Cylinder groups

- ü BSD FFS addressed these problems using the notion of a cylinder group
- ü Disk partitioned into groups of cylinders
- ü Data blocks from a file all placed in the same cylinder group
- ü Files in same directory placed in the same cylinder group
- ü I-nodes for files allocated in the same cylinder group as file's data blocks

**n** Free space reserve

  ü If the number of free blocks falls to zero, the file system throughput tends to be cut in half, because of the inability to localize blocks in a file

  ü A parameter, called free space reserve, gives the minimum acceptable percentage of file system blocks that should be free

  ü If the number of free blocks drops below this level, only the system administrator can continue to allocate blocks

  ü Normally 10%; this is why df may report > 100%

**n** Fragments

  **ü** Small blocks (1KB) caused two problems:

  - **§** Low bandwidth utilization

  - **§** Small max file size (function of block size)

  **ü** FFS fixes by using a larger block (4KB)

  - **§** Allows for very large files
    (1MB only uses 2 level indirect)

  - **§** But introduces internal fragmentation: there are many small files (i.e., < 4KB)

  **ü** FFS introduces "fragments" to fix internal fragmentation

  - **§** Allows the division of a block into one or more fragments (1K pieces of a block)

**n** Media failures

ü Replicate master block (superblock)

**n** File system parameterization

ü Parameterize according to disk and CPU characteristics

§ Maximum blocks per file in a cylinder group

§ Minimum percentage of free space

§ Sectors per track

§ Rotational delay between contiguous blocks

§ Tracks per cylinder, etc.

ü Skip according to rotational rate and CPU latency

**n** History

- ü Evolved from Minix file system
    - § Block addresses are stored in 16bit integers – maximal file system size is restricted to 64MB
    - § Directories contain fixed-size entries and the maximal file name was 14 characters
- ü Virtual File System (VFS) is added
- ü Extended Filesystem (Ext FS), 1992
    - § Added to Linux 0.96c
    - § Maximum file system size was 2GB, and the maximal file name size was 255 characters
- ü Second Extended File-system (Ext2 FS), 1994

# Linux Ext2 File System (Cont'd)

**n** Ext2 Features

- ü Configurable block sizes (from 1KB to 4KB)
    - § depending on the expected average file size
- ü Configurable number of i-nodes
    - § depending on the expected number of files
- ü Partitions disk blocks into groups
    - § lower average disk seek time
- ü Preallocates disk data blocks to regular files
    - § reduces file fragmentation
- ü Fast symbolic links
    - § If the pathname of the symbolic link has 60 bytes or less, it is stored in the i-node
- ü Automatic consistency check at boot time

**n** Disk layout

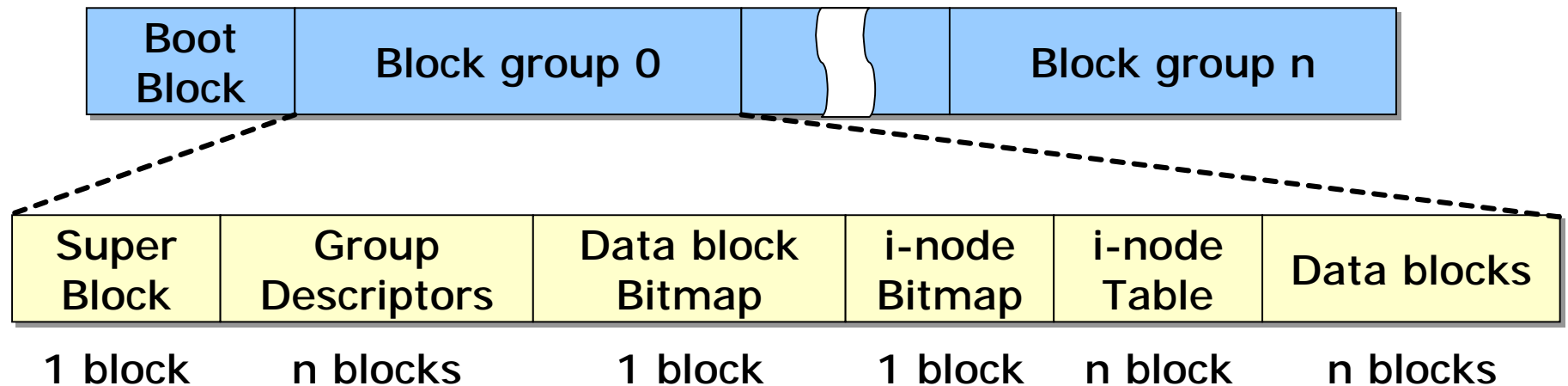    ü Boot block

        § reserved for the partition boot sector

    ü Block group

        § Similar to the cylinder group in FFS

        § All the block groups have the same size and are stored sequentially

| Boot Block | Block group 0 | { | Block group n |
|------------|---------------|---|---------------|

| Super Block | Group Descriptors | Data block Bitmap | i-node Bitmap | i-node Table | Data blocks |
|-------------|-------------------|-------------------|--------------|--------------|-------------|
| 1 block | n blocks | 1 block | 1 block | n block | n blocks |

**n Block group**

- ü Superblock: stores file system metadata
    - § Total number of i-nodes
    - § File system size in blocks
    - § Free blocks / i-nodes counter
    - § Number of blocks / i-nodes per group
    - § Block size, etc.
- ü Group descriptor
    - § Number of free blocks / i-nodes / directories in the group
    - § Block number of block / i-node bitmap, etc.
- ü Both the superblock and the group descriptors are duplicated in each block group
    - § Only those in block group 0 are used by the kernel
    - § fsck copies them into all other block groups
    - § When data corruption occurs, fsck uses old copies to bring the file system back to a consistent state

**n** Block group size

- ü The block bitmap must be stored in a single block
    - § In each block group, there can be at most 8xb blocks, where b is the block size in bytes
- ü The smaller the block size, the larger the number of block groups
- ü Example: 8GB Ext2 partition with 4KB block size
    - § Each 4KB block bitmap describes 32K data blocks
      = 32K * 4KB = 128MB
    - § At most 64 block groups are needed

**n** Directory structure

| | inode | record length | name length | file type | name |
|---|---|---|---|---|---|
| 0 | 21 | 12 | 1 | 2 | . \0 \0 \0 |
| 12 | 22 | 12 | 2 | 2 | . . \0 \0 |
| 24 | 53 | 16 | 5 | 2 | h o m e 1 \0 \0 \0 |
| 40 | 67 | 28 | 3 | 2 | u s r \0 |
| 62 | 0 | 16 | 7 | 1 | o l d f i l e \0 |
| 68 | 34 | 12 | 3 | 2 | b i n \0 |

*<deleted file>*