# Operating System: Chap11 File System Implementation

National Tsing-Hua University

2016, Fall Semester
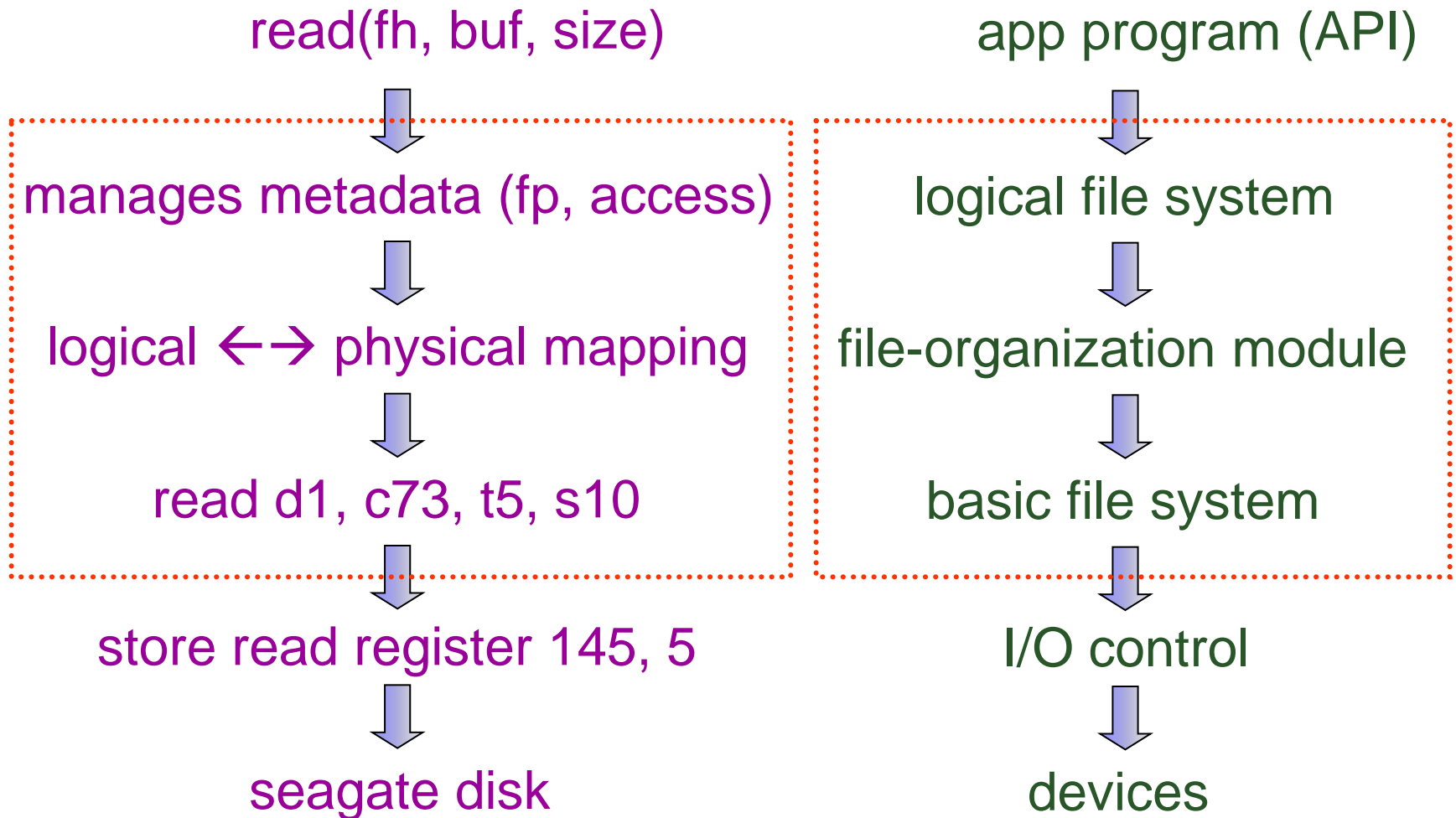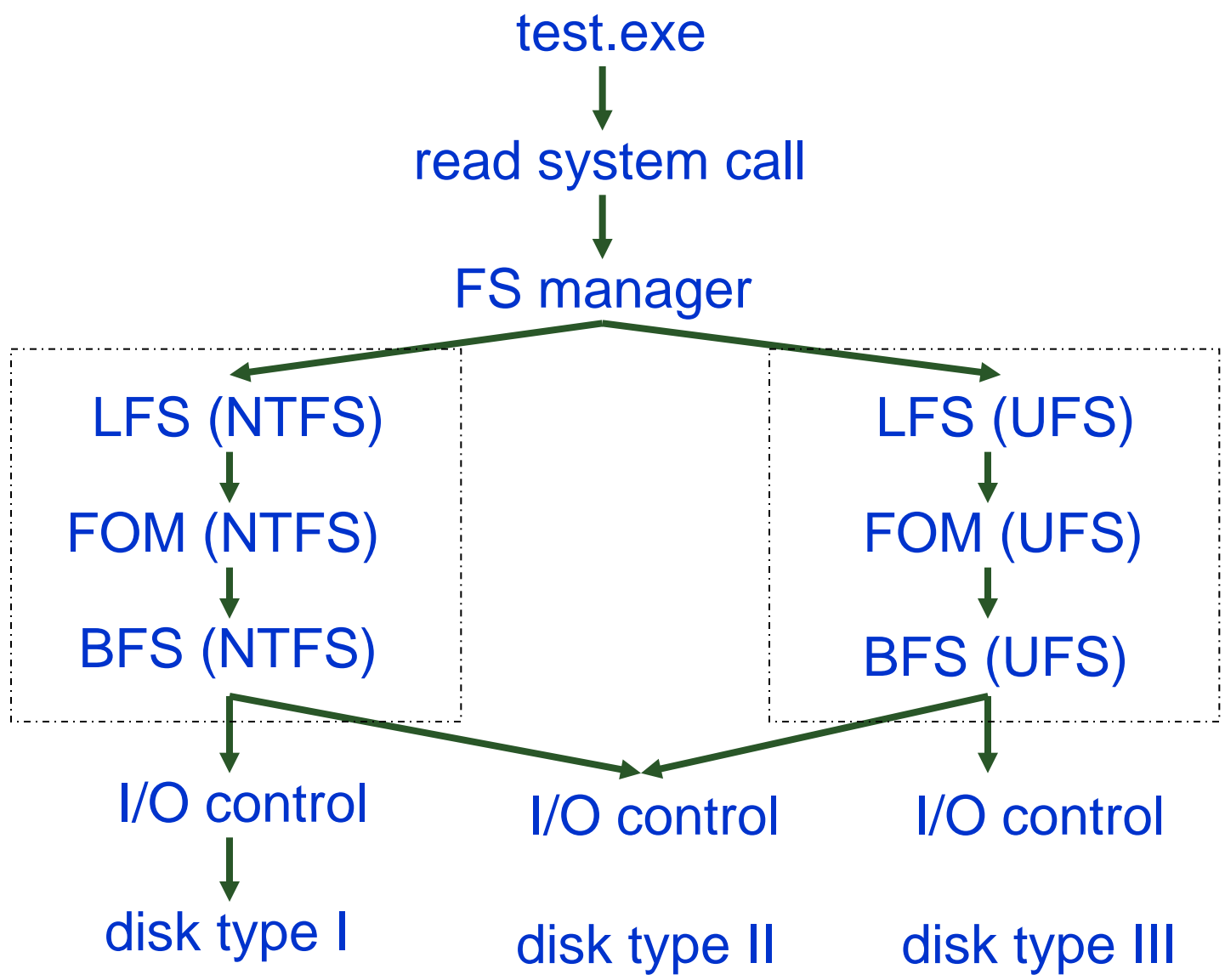
# Overview

- File-System Structure

- File System Implementation

- Disk Allocation Methods

- Free-Space Management

# File-System Structure

- I/O transfers between memory and disk are performed in units of **blocks**
  - one block is one or more **sectors**
  - one sector is usually 512 bytes
- One OS can support more than 1 FS types
  - NTFS, FAT32
- Two design problems in FS
  - interface to user programs
  - interface to physical storage (disk)

# Layered File System

read(fh, buf, size)                    app program (API)

manages metadata (fp, access)          logical file system

logical ←→ physical mapping            file-organization module

read d1, c73, t5, s10                  basic file system

store read register 145, 5             I/O control

seagate disk                           devices

test.exe

read system call

FS manager

LFS (NTFS)

FOM (NTFS)

BFS (NTFS)

LFS (UFS)

FOM (UFS)

BFS (UFS)

I/O control

I/O control

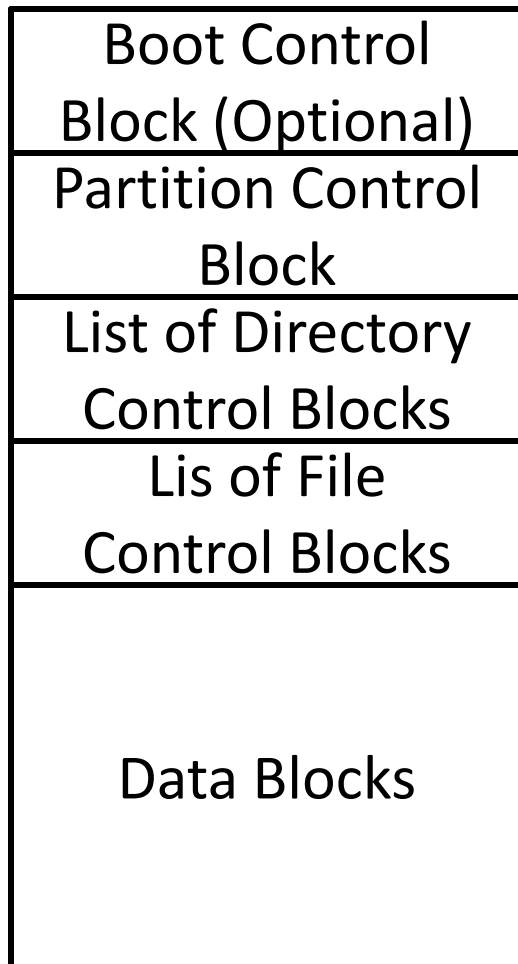I/O control

disk type I

disk type II

disk type III
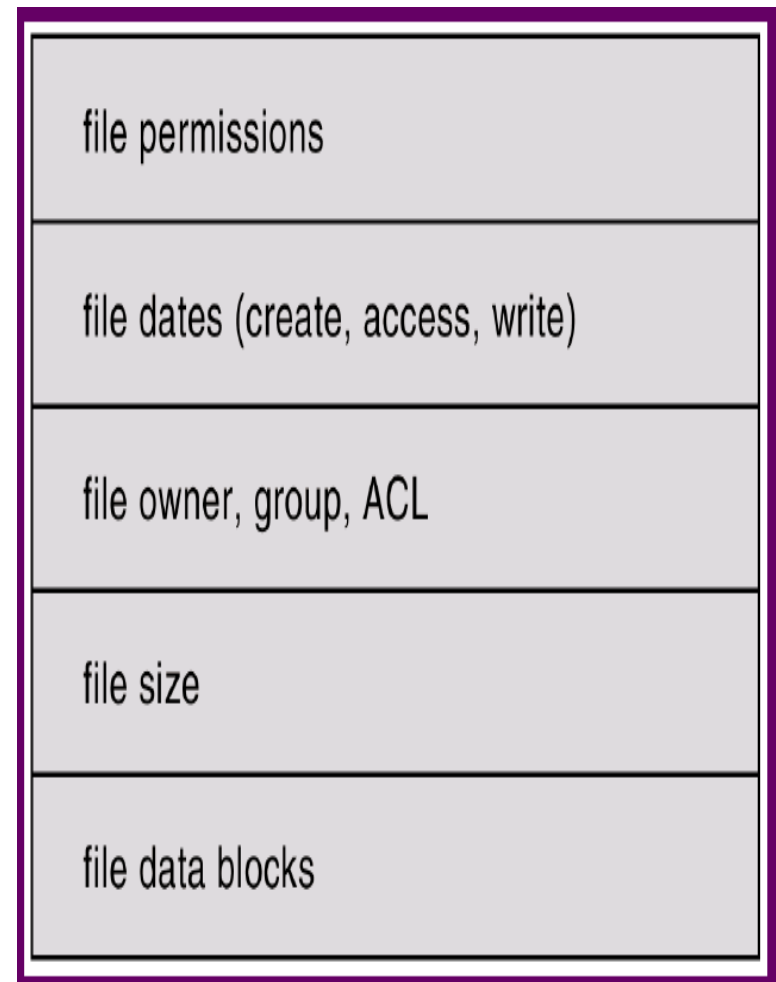
# File-System Implementation

# On-Disk Structure

- **Boot control block** (per partition): information needed to boot an OS from that partition
  - typical the first block of the partition (empty means no OS)
  - UFS (Unix File Sys.): **boot block**,  NTFS: partition boot sector
- **Partition control block** (per partition): partition details
  - details: # of blocks, block size, free-block-list, free FCB pointers, etc
  - UFS: **superblock**, NTFS: Master File Table
- **File control block** (per file): details regarding a file
  - details: permissions, size, location of data blocks
  - UFS: **inode**, NTFS: stored in MFT (relational database)
- **Directory structure** (per file system): organize files
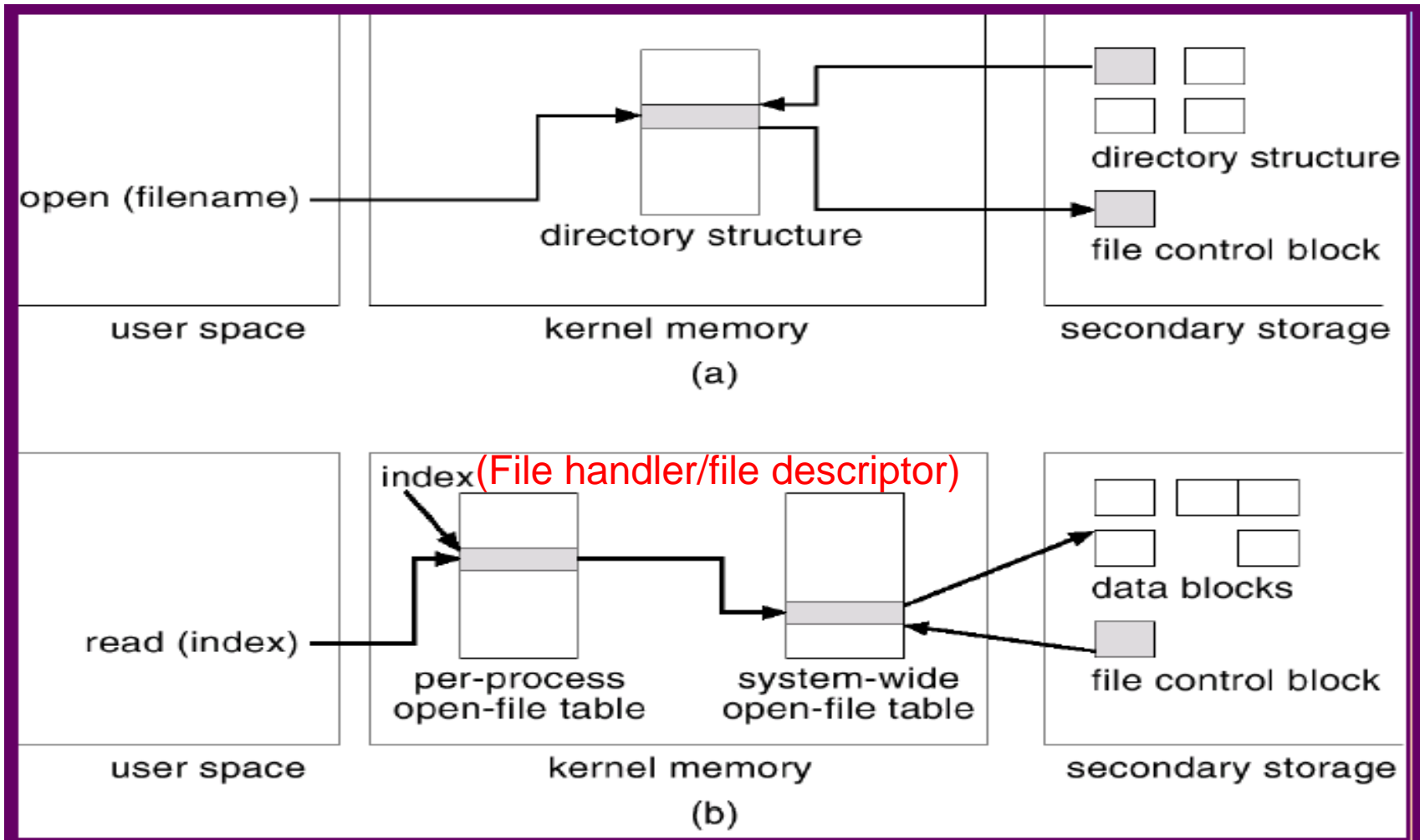
# On-Disk Structure

## Partition

| |
|---|
| Boot Control Block (Optional) |
| Partition Control Block |
| List of Directory Control Blocks |
| Lis of File Control Blocks |
| Data Blocks |

## File Control Block (FCB)

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

# In-Memory Structure

- **in-memory partition table**: information about each **mounted partition**

- **in-memory directory structure**: information of **recently accessed directories**

- **system-wide open-file table**: contain a copy of each **opened file's FCB**

- **per-process open-file table**: **pointer (file handler/descriptor)** to the corresponding entry in the above table
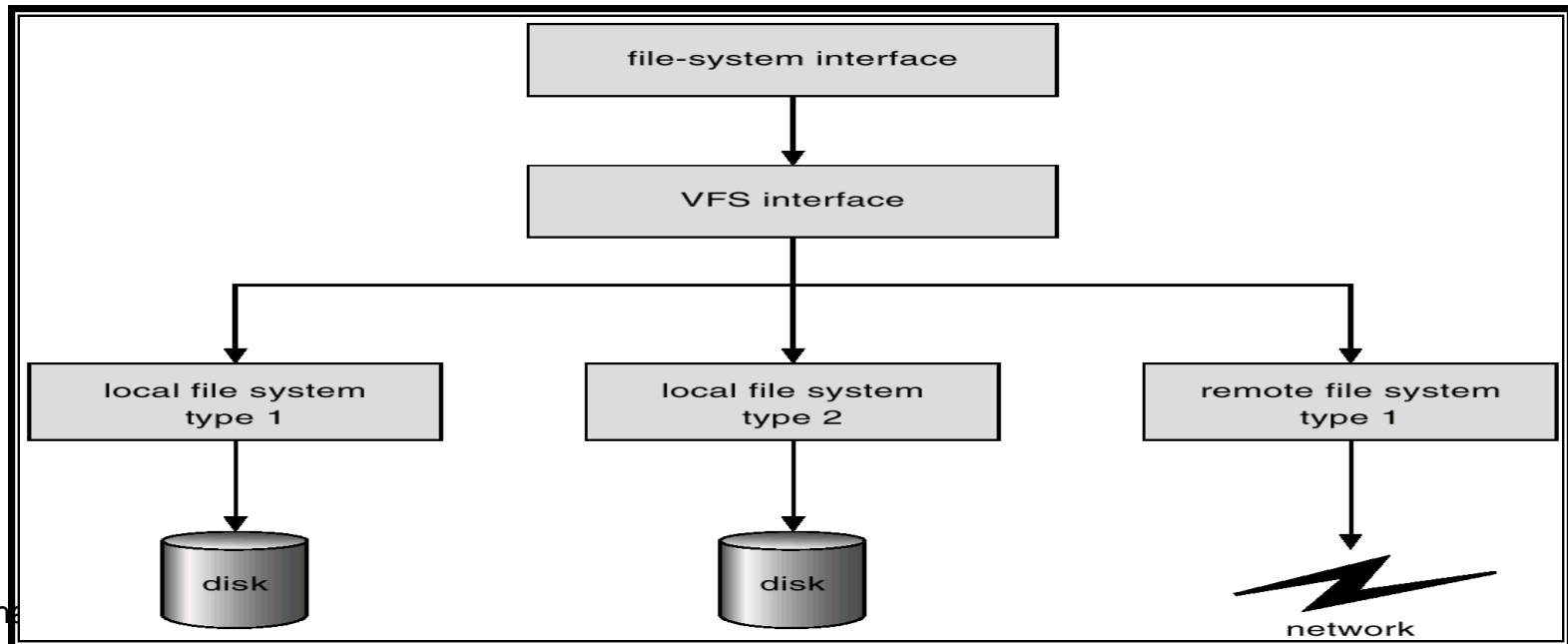
# File-Open & File-Read

# File Creation Procedure

1. OS allocates a new **FCB**

2. Update **directory structure**

    1. OS reads in the corresponding directory structure into memory

    2. Updates the dir structure with the new file name and the FCB

    3. (After file being closed), OS writes back the directory structure back to disk

3. The file appears in user's dir command

# Virtual File System

- VFS provides an object-oriented way of implementing file systems
- VFS allows the same **system call interface** to be used for different types of FS
- VFS calls the appropriate FS routines based on the partition info

# Virtual File System

- Four main object types defined by Linux VFS:
  - inode ➜ an individual file
  - file object ➜ an open file
  - superblock object ➜ an entire file system
  - dentry object ➜ an individual directory entry
- VFS defines a set of operations that must be implemented (e.g. for file object)
  - int open(…) ➜ open a file
  - ssize_t read() ➜ read from a file

# Directory Implementation

- Linear lists
  - list of file names with pointers to data blocks
  - easy to program but poor performance
    - insertion, deletion, searching

- Hash table – linear list w/ hash data structure
  - constant time for searching
  - linked list for collisions on a hash entry
  - hash table usually has fixed # of entries

# Review Slides ( I )

- Transfer unit between memory and disk?
- App ➔ LFS ➔ FOM ➔ BFS ➔I/O Control ➔Devices
- On-disk structure
  - Boot control block, Partition control block
  - File control block, Directory structure
- In-memory
  - Partition table, Directory structure
  - System-wide open-file table
  - Per-process open-file table
- Steps to open file, read/write file and create file?
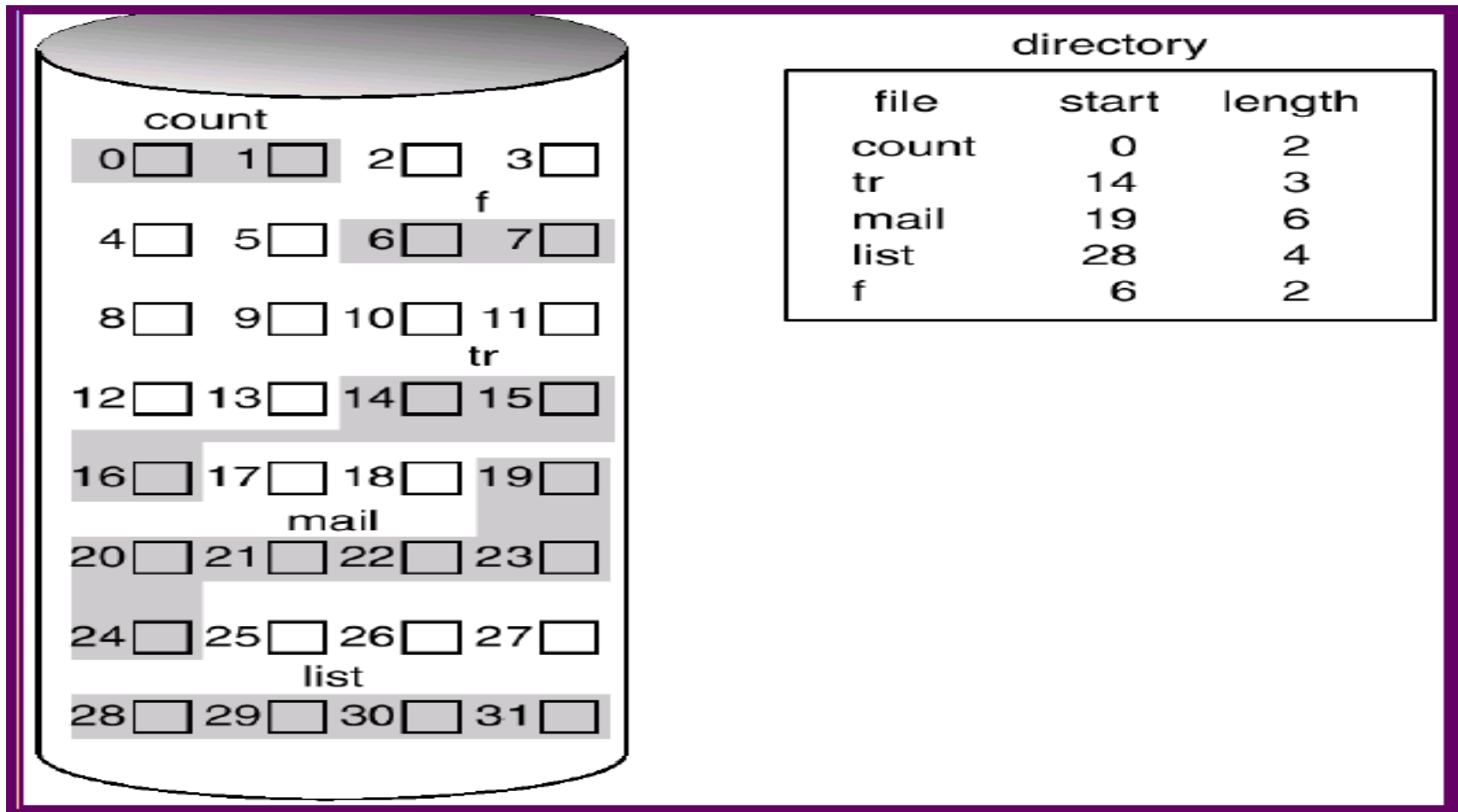- Purpose of VFS?

# Allocation Methods

# Outline

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

- Allocation strategy:
  - ➢ **Contiguous allocation**
  - ➢ **Linked allocation**
  - ➢ **Indexed allocation**

# Contiguous Allocation



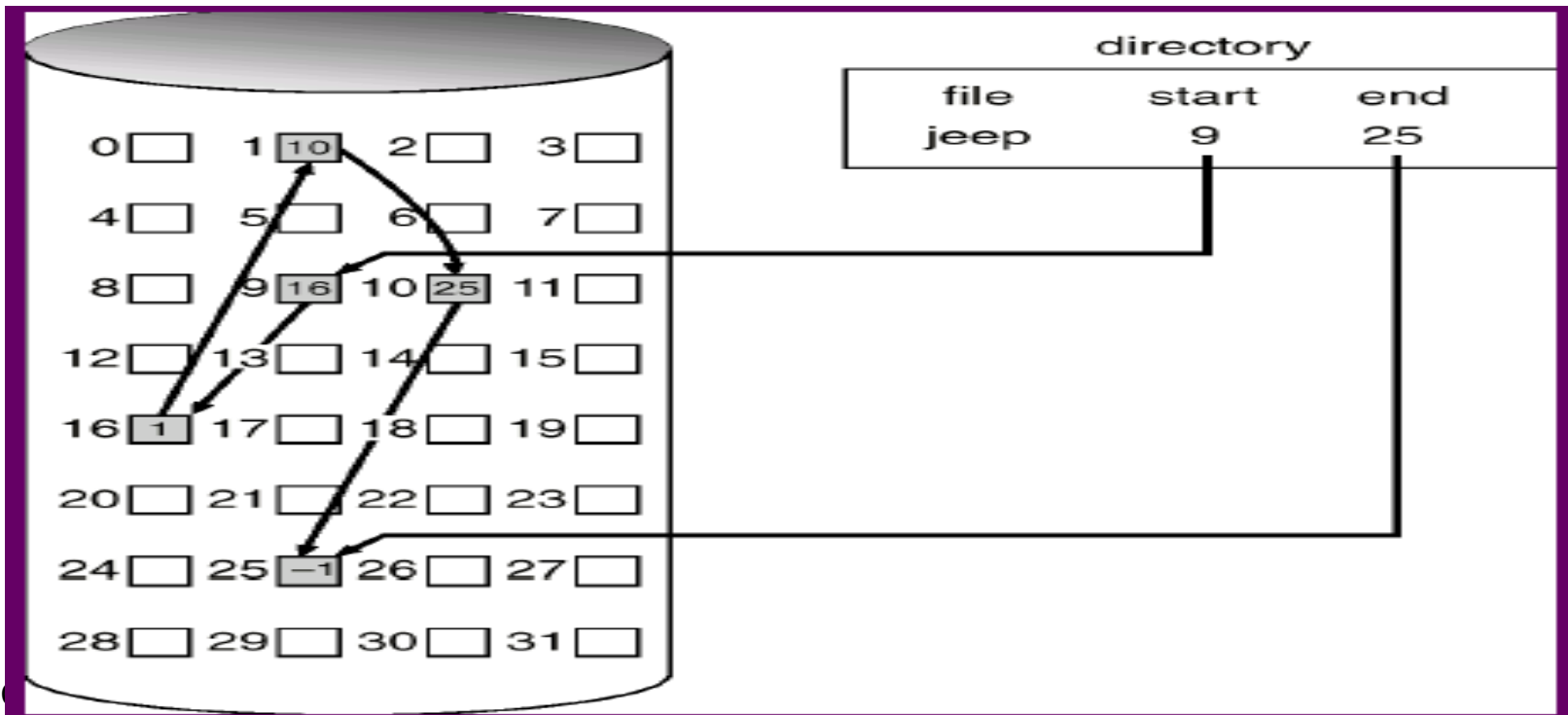| file | start | length |
|------|-------|--------|
| count | 0 | 2 |
| tr | 14 | 3 |
| mail | 19 | 6 |
| list | 28 | 4 |
| f | 6 | 2 |

# Contiguous Allocation

- Each file occupies a set of contiguous blocks
  - \# of disk seeks is minimized
  - The dir entry for each file = (starting #, size)
- Both sequential & random access can be implemented efficiently
- Problems
  - External fragmentation ➜ compaction
  - File cannot grow ➜ extend-based FS

# Extent-Based File System

- Many newer file system use a modified contiguous allocation scheme

- Extent-based file systems allocate disk blocks in extents

- An extent is a **contiguous blocks** of disks
  - A file contains one or more extents
  - An extent: (starting block #, length, **pointer to next extent**)
  - ☹ Random access become more costly
  - ☹ Both internal & external fragmentation are possible

# Linked Allocation

■ Each file is a linked list of blocks
  ➢ Each block contains a pointer to the next block
  ➔ data portion: block size – pointer size
■ File read: following through the list

# Linked Allocation

- **Problems**
  - Only good for **sequential-access** files
    - Random access requires traversing through the link list
    - Each access to a link list is a disk I/O (because link pointer is stored inside the data block)
  - Space required for pointer (4 / 512 = 0.78%)
    - solution: unit = cluster of blocks
      - → internal fragmentation
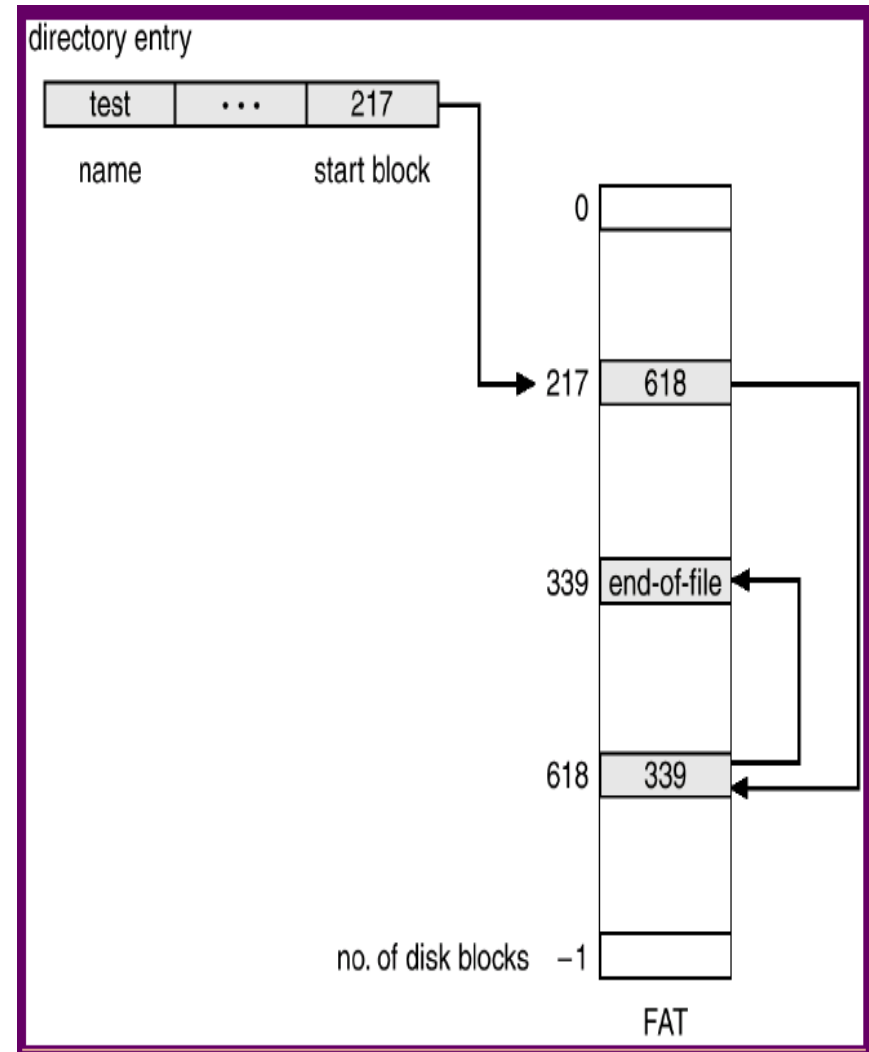  - Reliability
    - One missing link breaks the whole file

# FAT (File Allocation Table) file system

- **FAT32**
  - ➢ Used in MS/DOS & OS/2
  - ➢ Store all links in a table
  - ➢ 32 bits per table entry
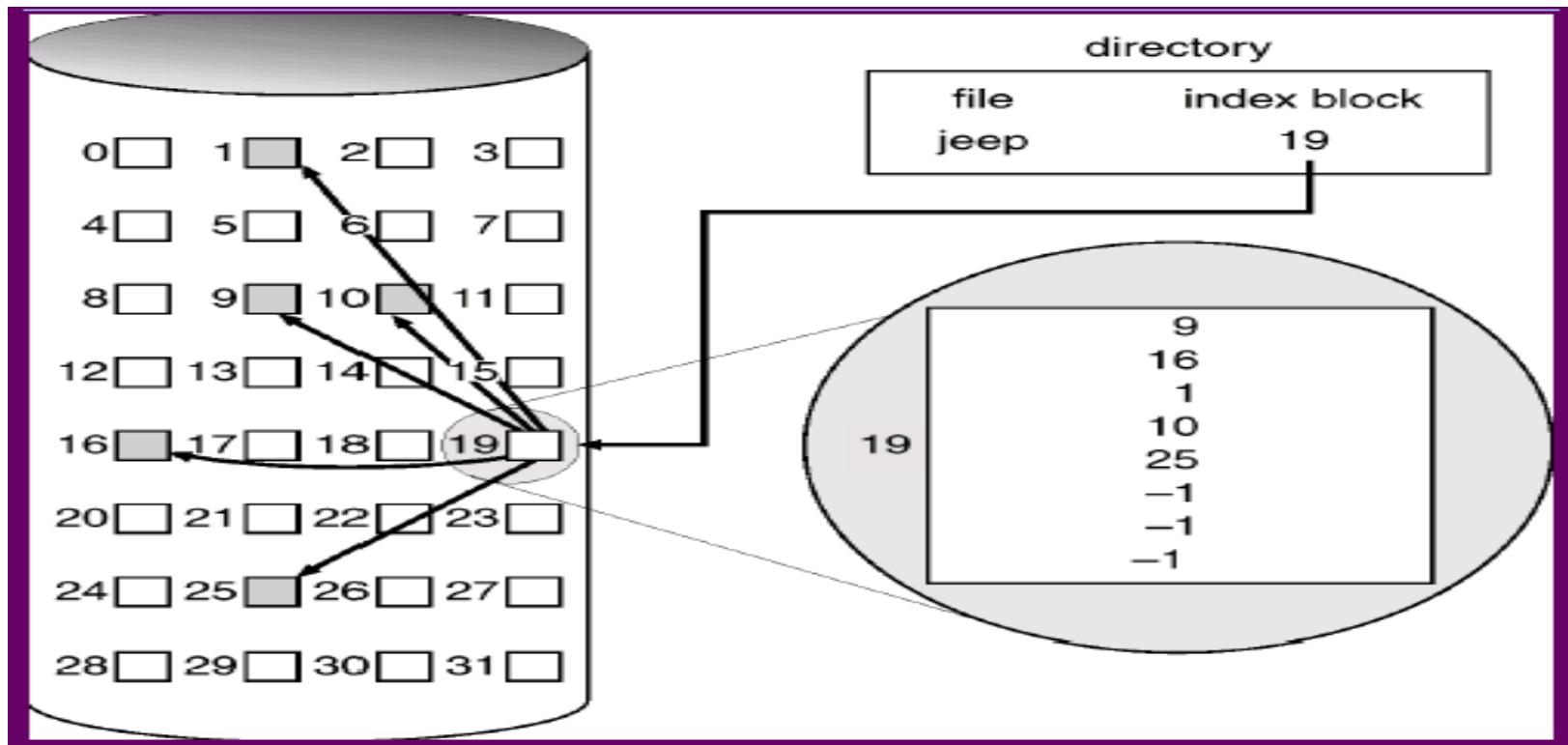  - ➢ located in a section of disk at the beginning of each partition
- **FAT(table) is often cached in memory**
  - ➢ Random access is improved
  - ➢ Disk head find the location of any block by reading FAT



directory entry

| test | . . . | 217 |
| --- | --- | --- |

name            start block

0

217    618

339    end-of-file

618    339

no. of disk blocks   −1

FAT

# Indexed Allocation Example

- The directory contains the address of the file index block
- Each file has its own index block
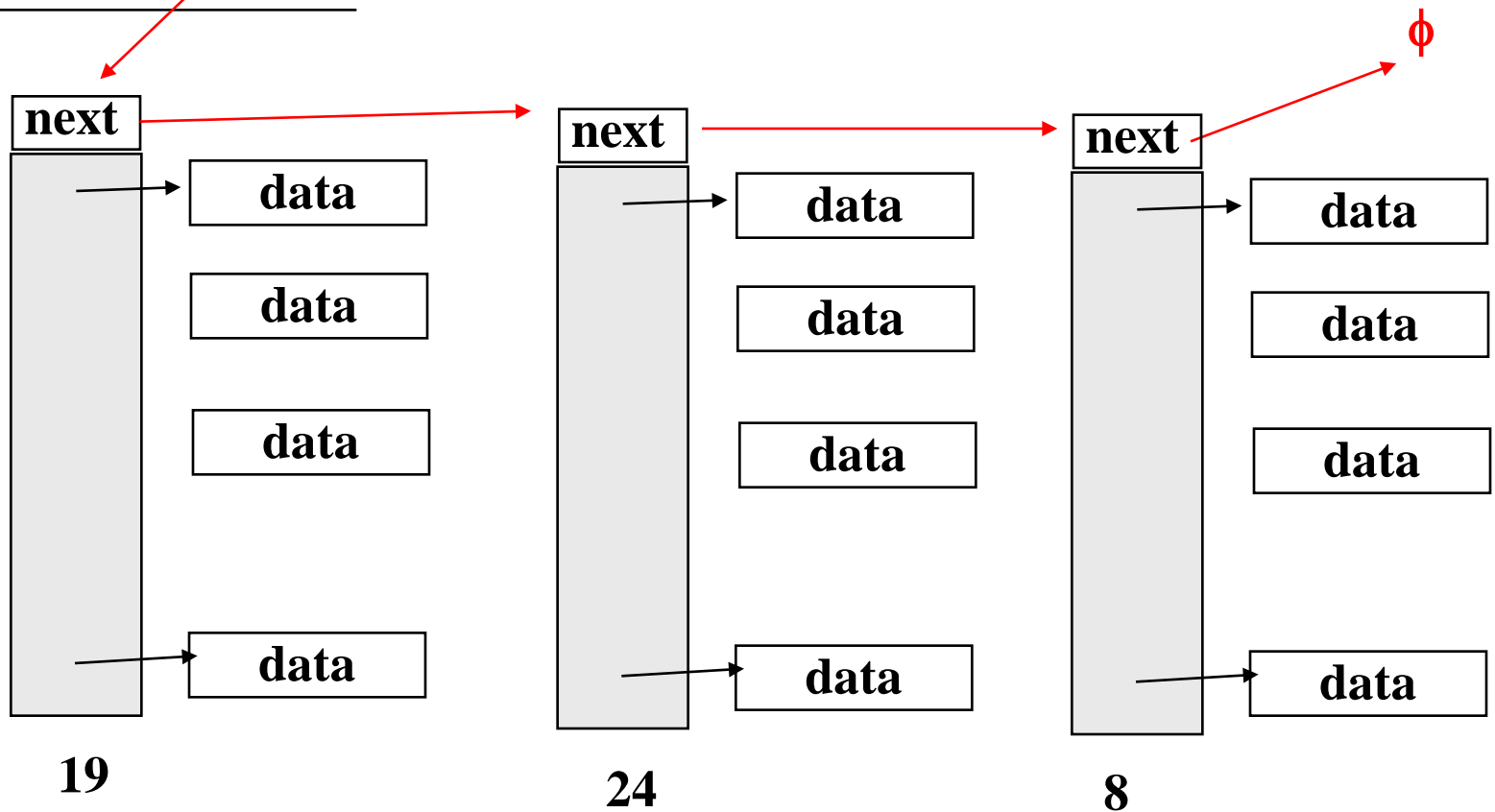- Index block stores block # for file data

# Indexed Allocation

- Bring all the pointers together into one location: the index block (one for each file)

☺:      1. Implement **direct and random access efficiently**

         2. No external fragmentation
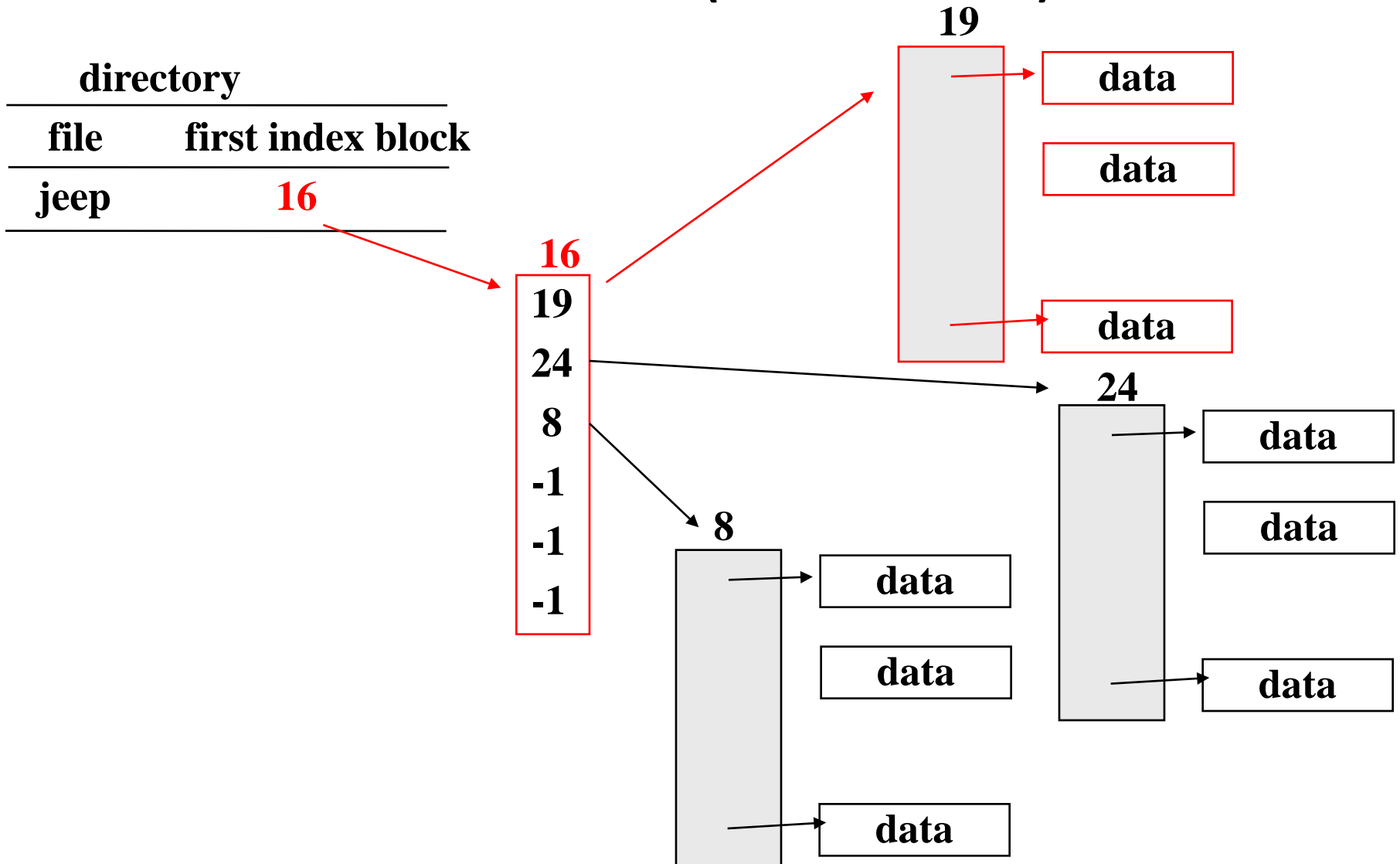
         2. Easy to create a file (no allocation problem)

☹:      1. Space for index blocks

         2. How large the index block should be ?

- **linked scheme**
- **multilevel index**
- **combined scheme** (inode in BSD UNIX)

# Linked Indexed Scheme

**directory**

| file | first index block |
|------|-------------------|
| jeep | **19** |



19          24          8

# Multilevel Scheme (two-level)

**directory**

| file | first index block |
|------|-------------------|
| jeep | **16** |

**16**
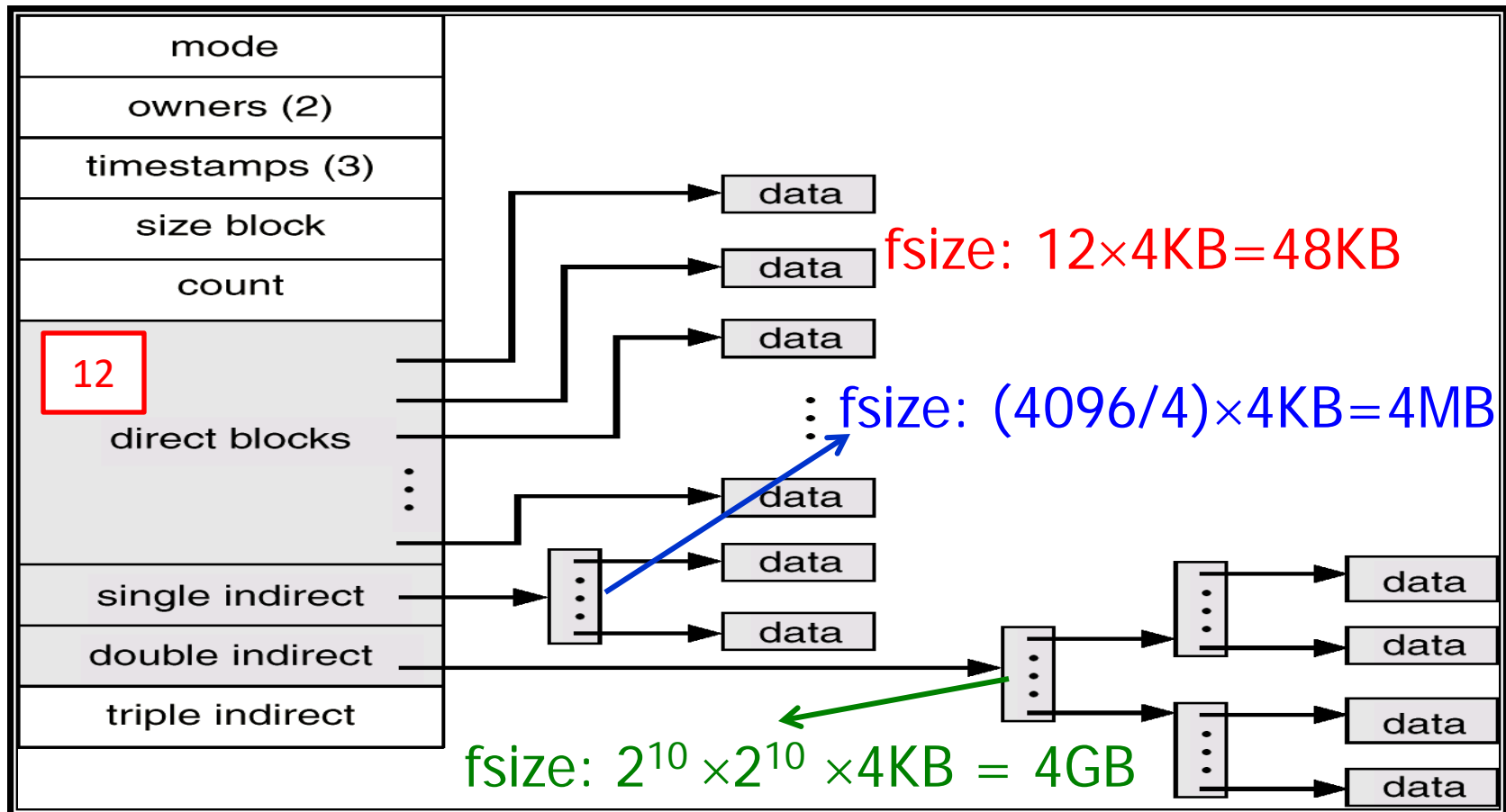19
24
8
-1
-1
-1

**19**
data
data
data

**24**
data
data
data

**8**
data
data
data

# Combined Scheme: UNIX inode

- File pointer:4B (32bits)➜reach only 4GB ($2^{32}$) files
- Let each data/index block be 4KB



fsize: $12 \times 4KB = 48KB$

fsize: $(4096/4) \times 4KB = 4MB$

fsize: $2^{10} \times 2^{10} \times 4KB = 4GB$

# Free-Space Management

# Free Space

- Free-space list: records all free disk blocks
- Scheme
  - Bit vector
  - Linked list (same as linked allocation)
  - Grouping (same as linked index allocation)
  - Counting (same as contiguous allocation)
- File systems usually manage free space in the same way as a file
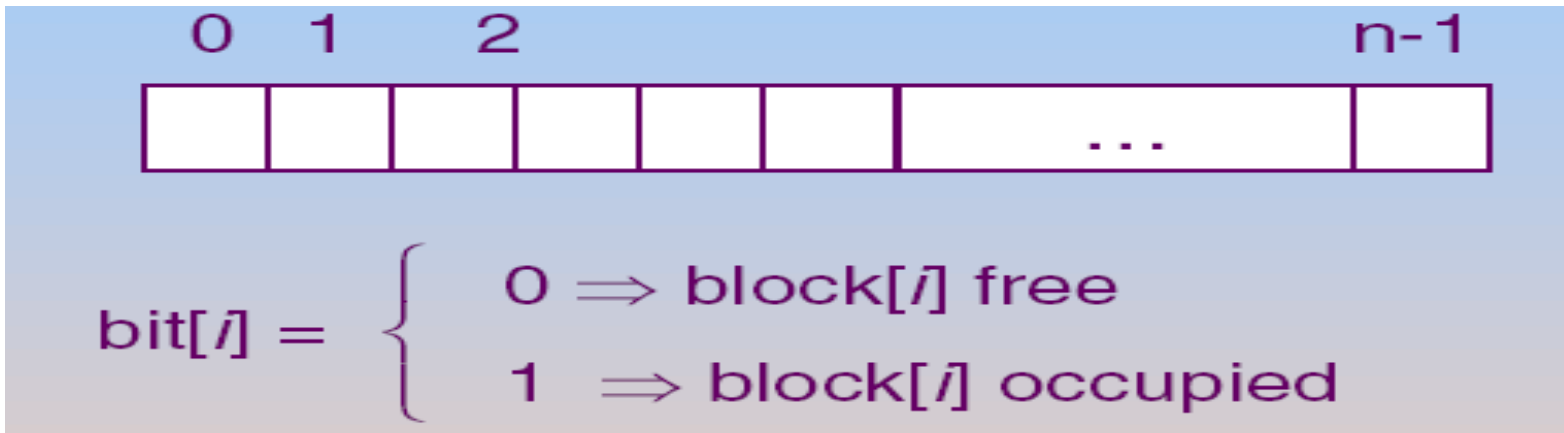
# Bit vector

- Bit Vector (bitmap): one bit for each block
  - e.g. 001111001111111110011001100000.....
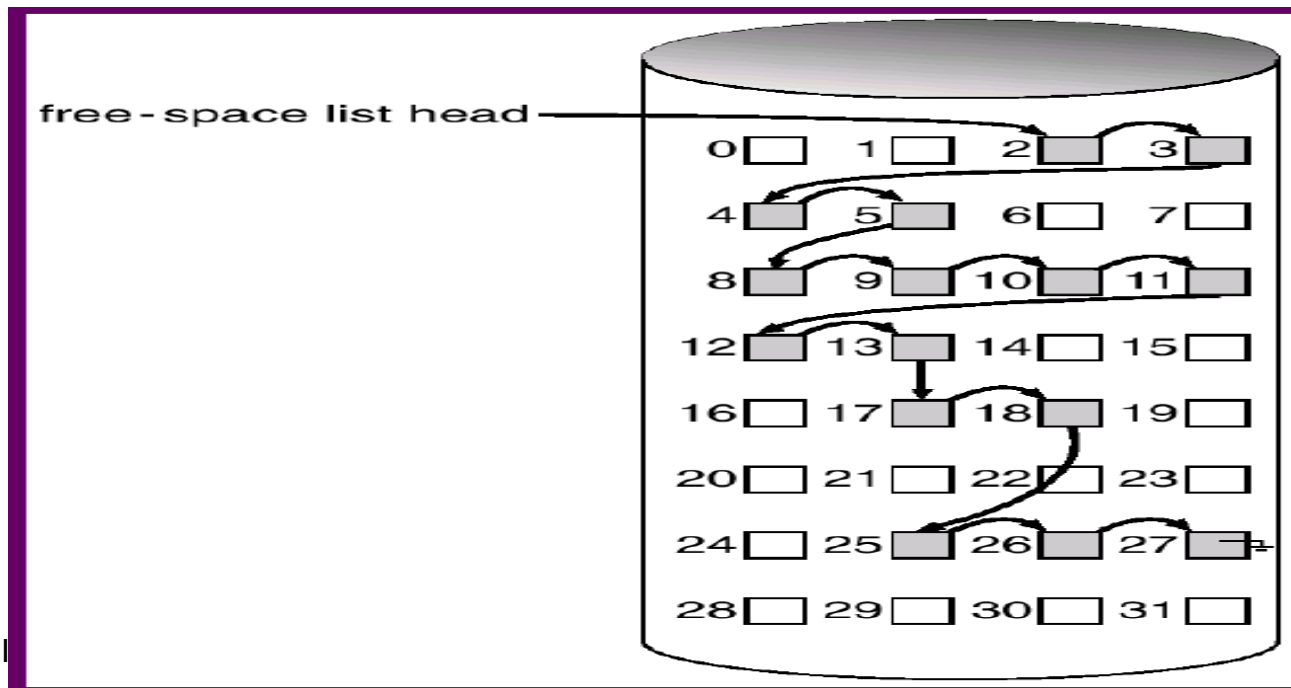- ☺: simplicity, efficient

  (HW support bit-manipulation instruction)
- ☹:bitmap must be cached for good performance
  - A 1-TB(4KB block) disk needs 32MB bitmap



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

# Linked List

- Link together all free blocks (same as linked allocation)
- Keep the first free block pointer in a special location on the disk and caching in memory
- Traversing list could be inefficient
  - No need for traversing; Put all link-pointers in a table(FAT)

# Grouping & Counting

- **Grouping (Same as linked-index allocation)**
  - store address of n free blocks in the 1st block
  - the first (n-1) pointers are free blocks
  - the last pointers is another grouping block
- **Counting (Same as contiguous allocation)**
  - keep the address of the first free block and # of contiguous free blocks

# Review Slides ( II )

- Allocation:
    - Contiguous file allocation? Extent-based file system?
    - Linked allocation?
    - Indexed allocation?
        - Linked scheme
        - multilevel index allocation
        - Combine scheme
- Free space:
    - Bit vector, linked list, counting, grouping

# Reading Material & HW

- Chap 11
- Problems:
  - 11.1, 11.2, 11.3, 11.4, 11.7, 11.8