



國立陽明交通大學

NATIONAL YANG MING CHIAO TUNG UNIVERSITY

Institute of Artificial Intelligence Innovation

Department of Computer Science

Operating System

Lecture 10: File System Implementation

Shuo-Han Chen 陳碩漢

shch@nycu.edu.tw

Wed. 10:10 - 12:00 EC115 +

Fri. 11:10 – 12:00 Online



Course Schedule

W	Date	Lecture	Online	Homework
1	Sept. 4	Lec00: Course Overview & Historical Prospective		
2	Sept. 11	Lec01: Introduction	V	
3	Sept. 18	Lec02: OS Structure	V	HW01 Due 10/5
4	Sept. 25	Lec03: Processes Concept	X	
5	Oct. 2	Typhoon – No class	V	
6	Oct. 9	Lec07: Memory Management	V	
7	Oct. 16	Lec08: Virtual Memory Management	V	HW02 Due 11/2
8	Oct. 23	Lec04: Multithreaded Programming	V	
9	Oct. 30	Midterm Exam		
10	Nov. 6	Lec05: Process Scheduling	V	Let's take a breath
11	Nov. 13	Lec06: Process Synchronization & Deadlocks	X	HW03
12	Nov. 20	School Event – No class	V	
13	Nov. 27	Lec09: File System Interface	V	
14	Dec. 4	Lec10: File System Implementation	V	HW04
15	Dec. 11	Lec11: Mass Storage System & Lec12: IO Systems	V	
16	Dec. 18	School Final Exam		

Overview

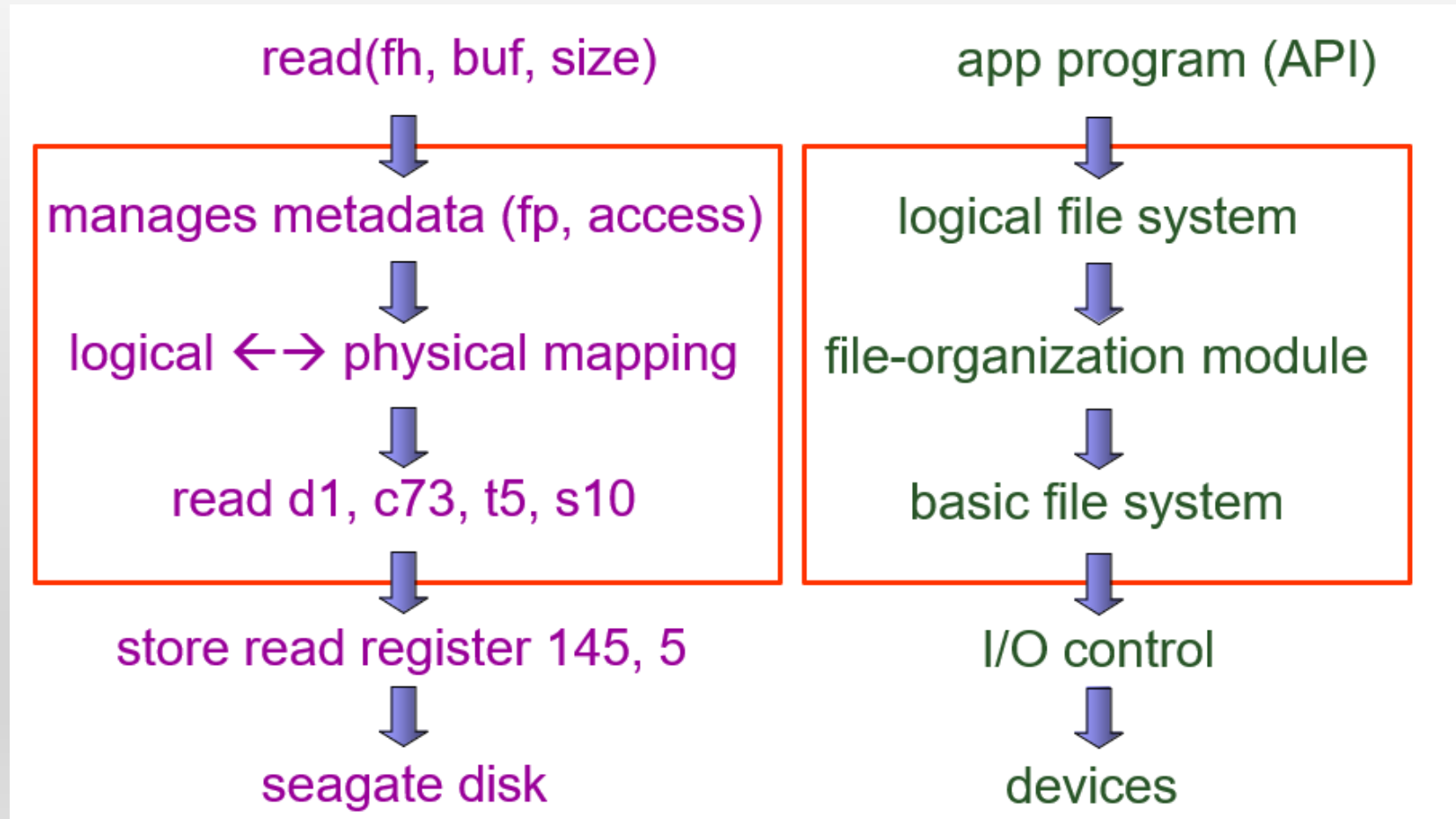
- File-System Structure
- File System Implementation
- Disk Allocation Methods
- Free-Space Management

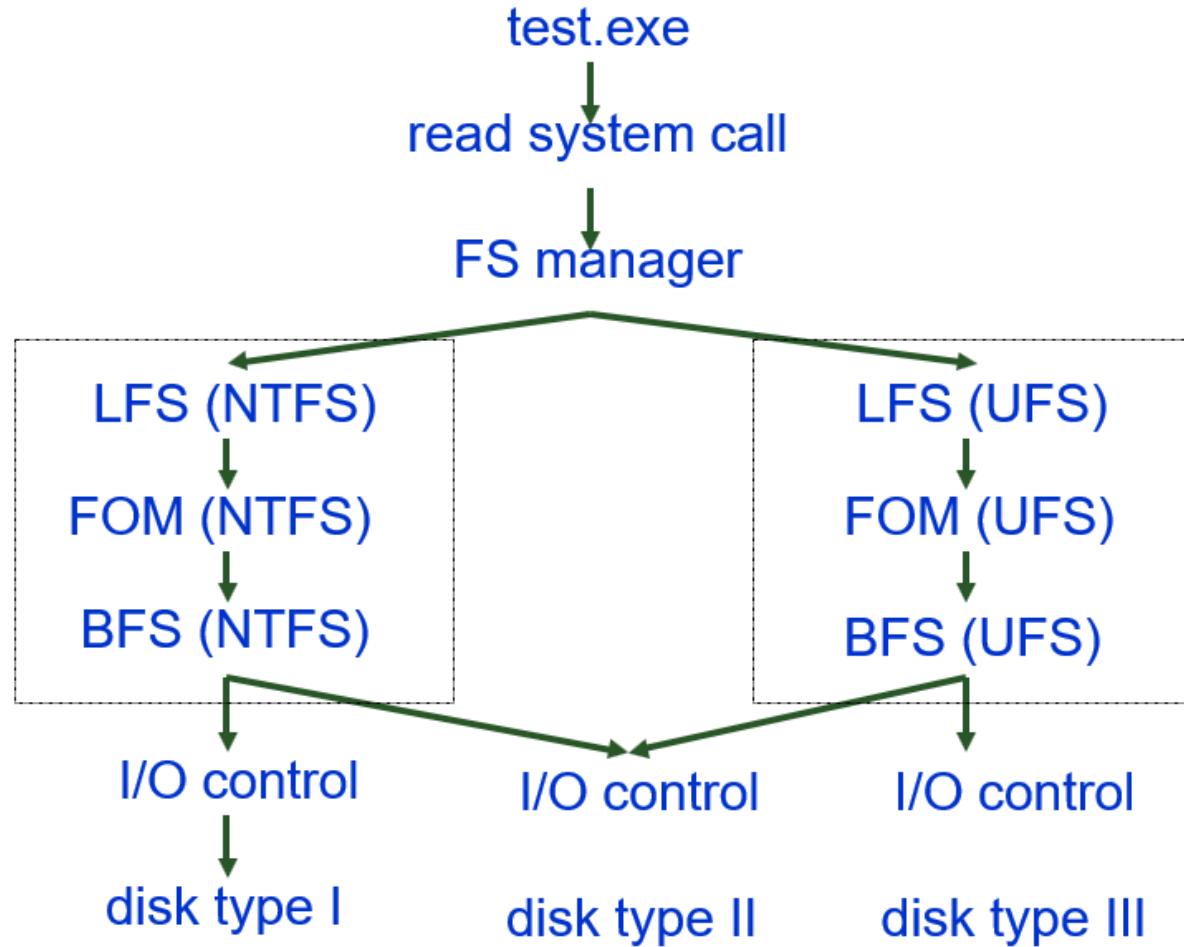


File-System Structure

- 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





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
List of File Control Blocks
Data Blocks

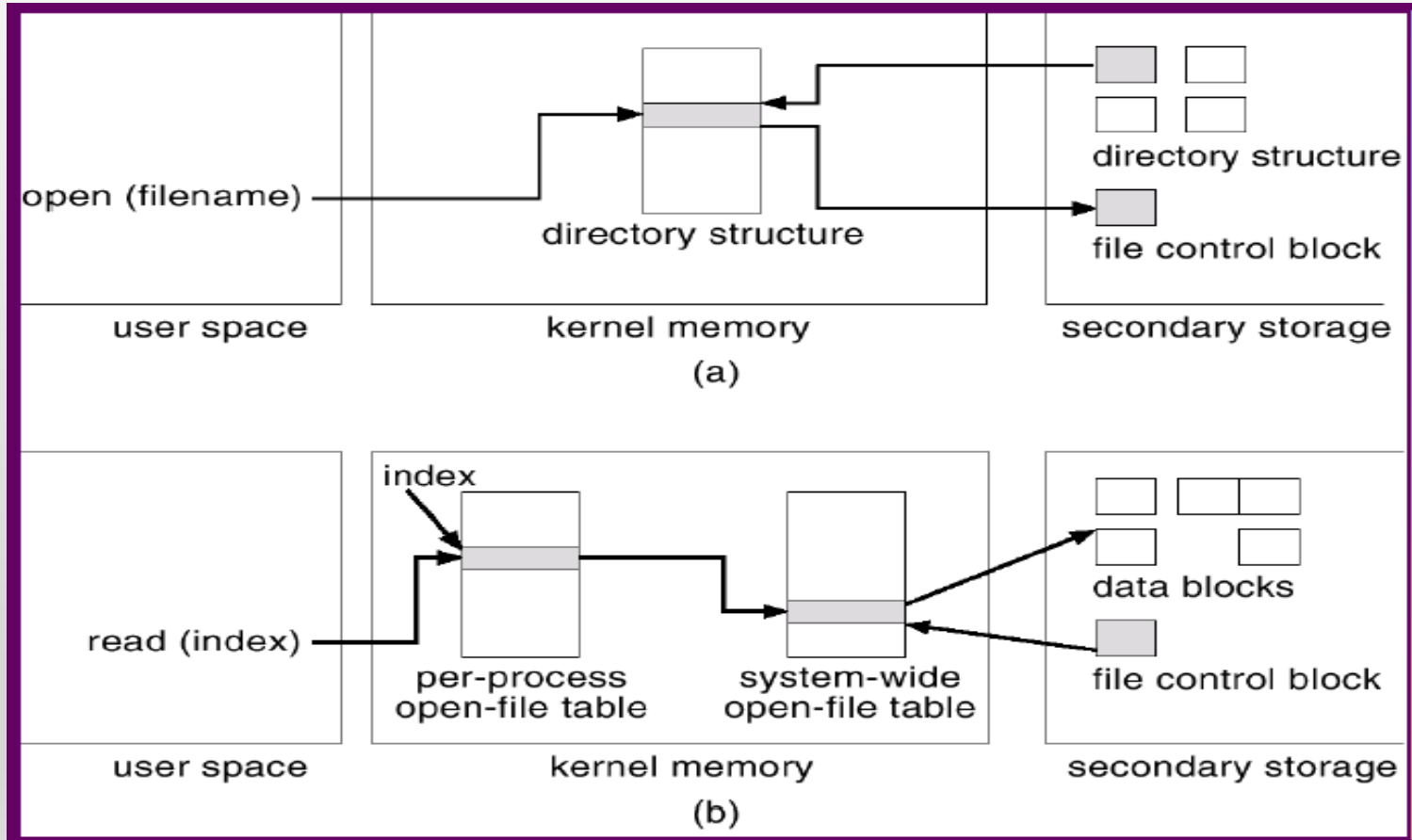
File Control Block (FCB)

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

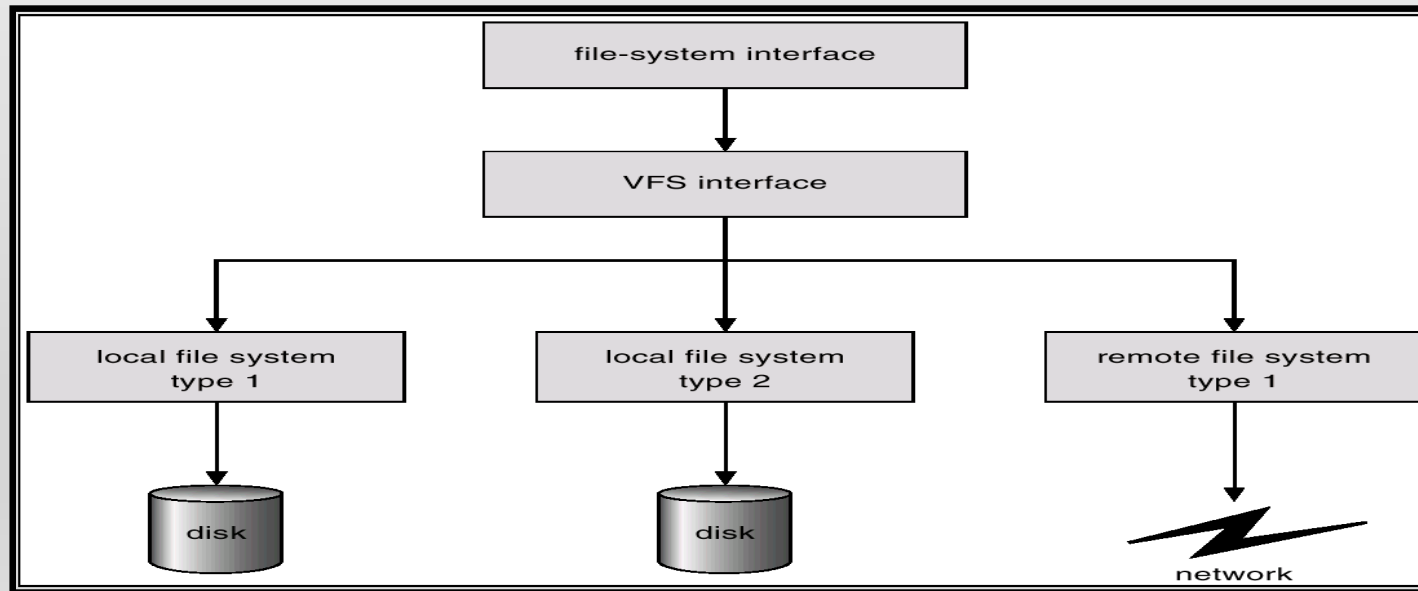


File Creation Procedure

- OS allocates a new FCB
- 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
- 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)

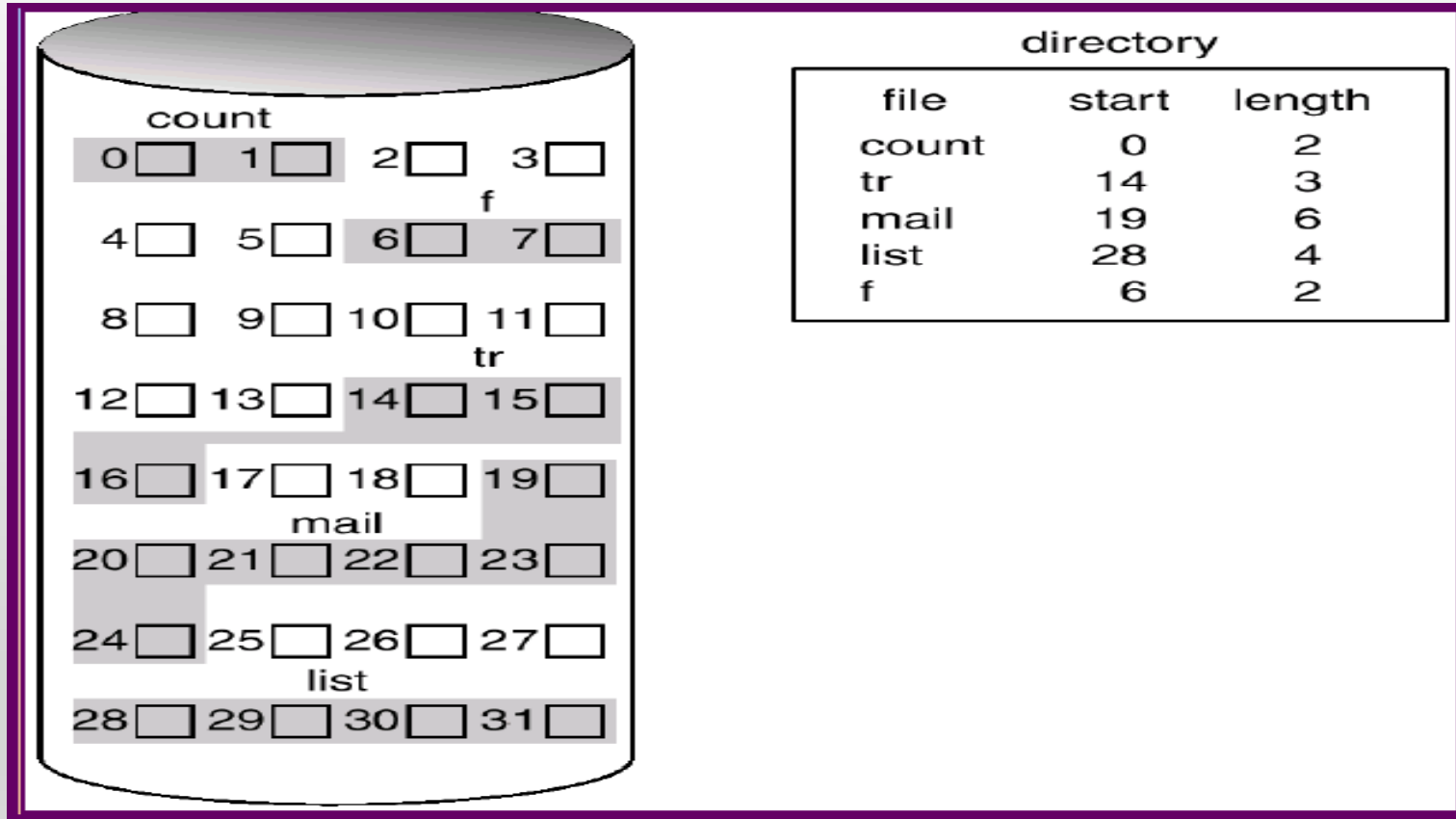
- Transfer unit between memory and disk?
- App -> LFS -> FOM -> BFS -> I/O Control -> Devices
- On-disk structure
 - Boot control block, Partition control block
 - **File control block**, Directory structure
- In-memory
 - Partition table, Directory structure
 - System-wide open-file table
 - Per-process 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



Contiguous Allocation

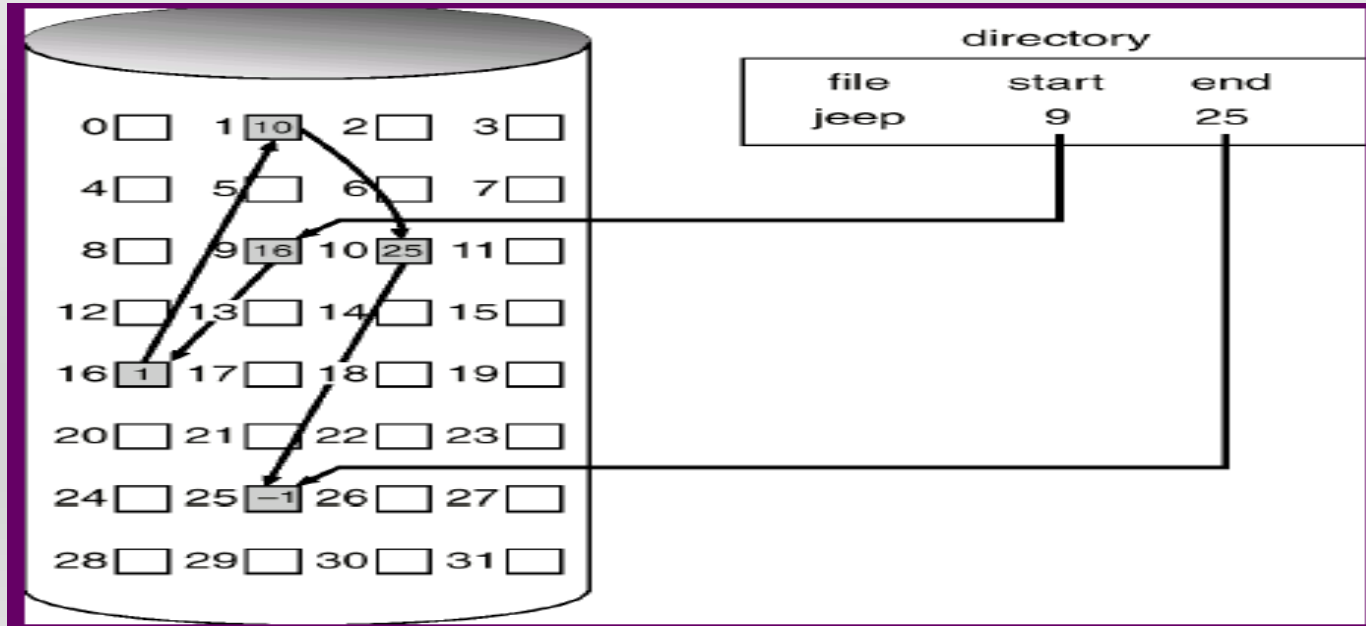
- Each file occupies a set of contiguous blocks
 - # of disk seeks is minimized
 - The dir entry for each file = (starting #, size)
- Both sequential & random access can be implemented efficiently
- Problems
 - External fragmentation -> compaction
 - File cannot grow -> extend-based FS

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

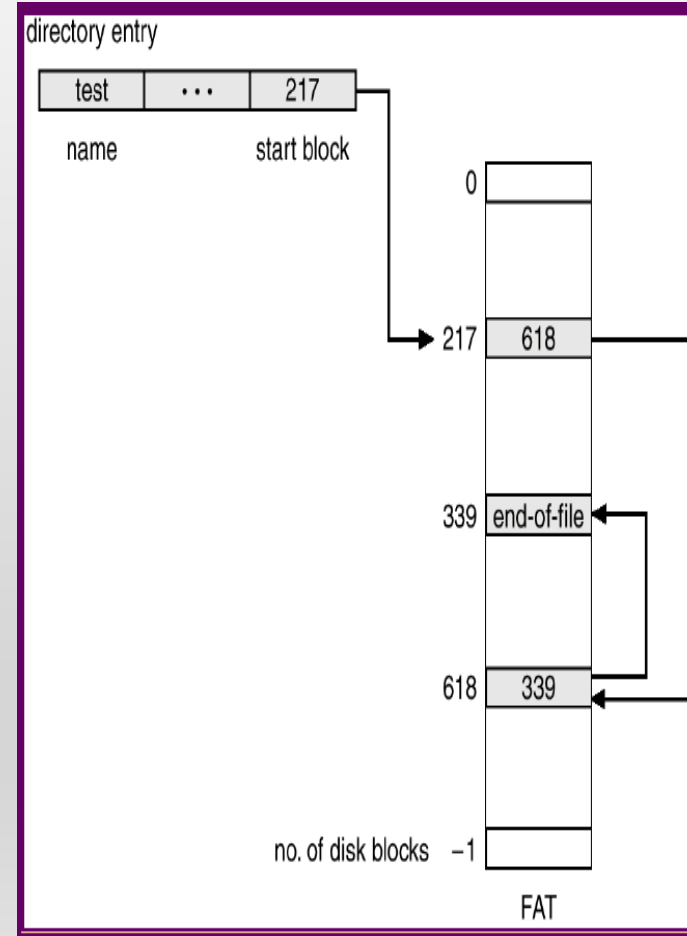


Linked Allocation

- 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

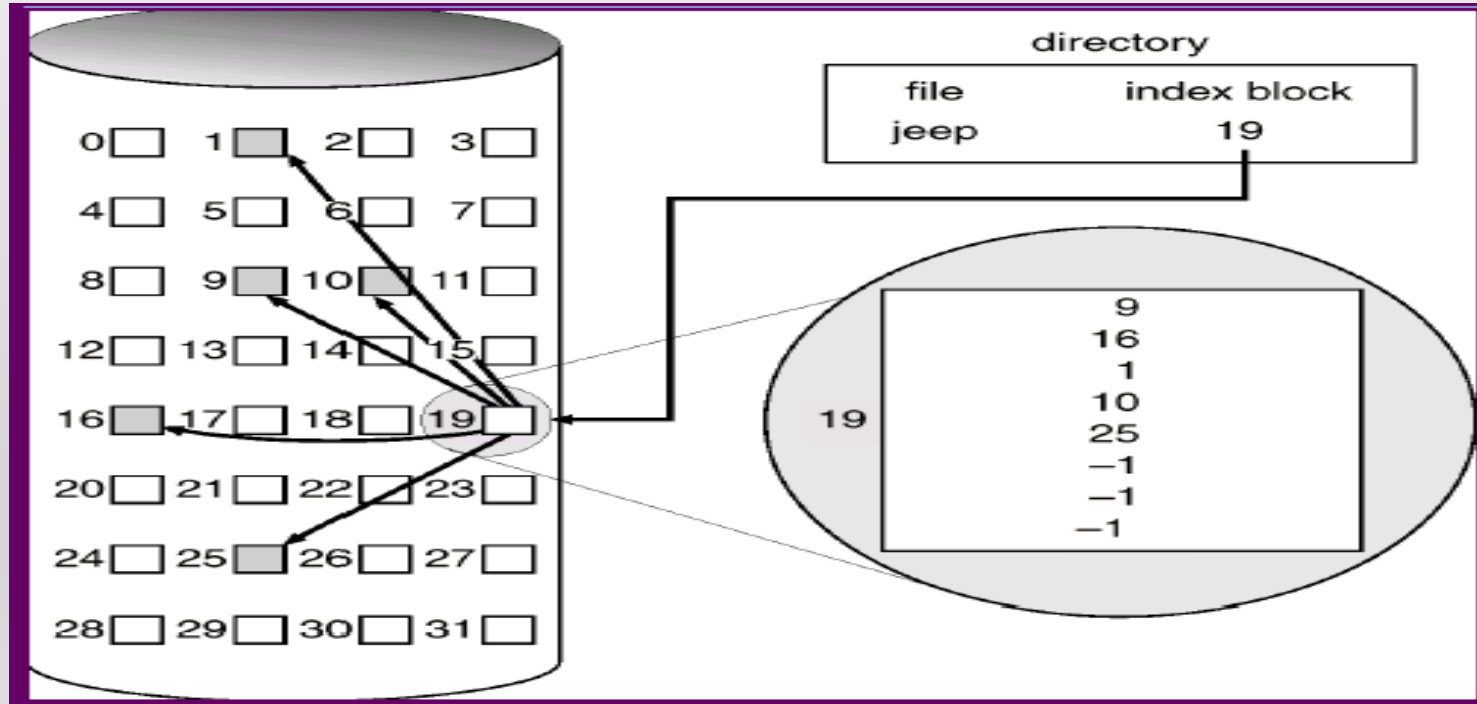
FAT (File Allocation Table) file system

- 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



Indexed Allocation Example

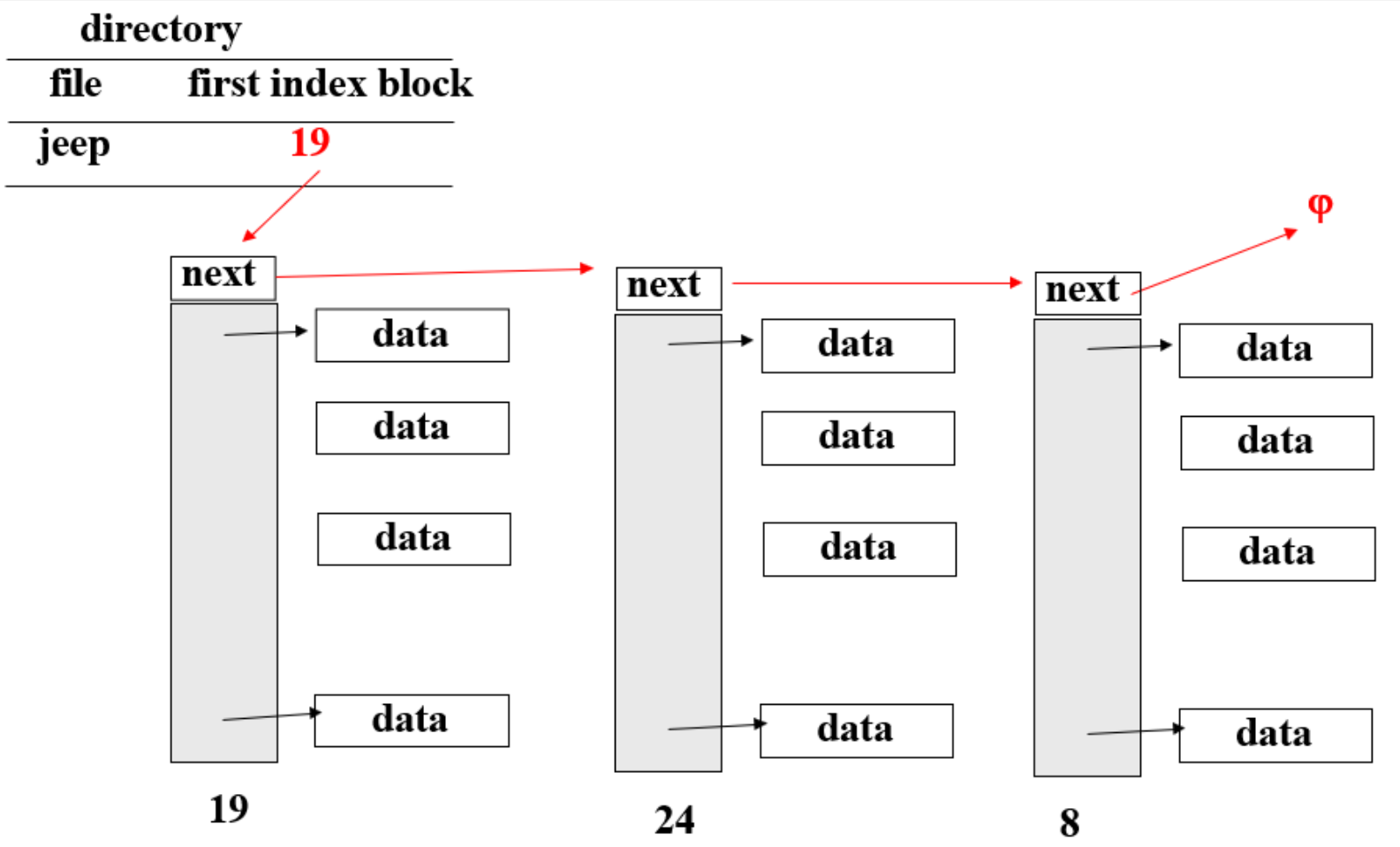
- The directory contains the address of the file index block
- Each file has its own index block
- Index block stores **block #** for file data



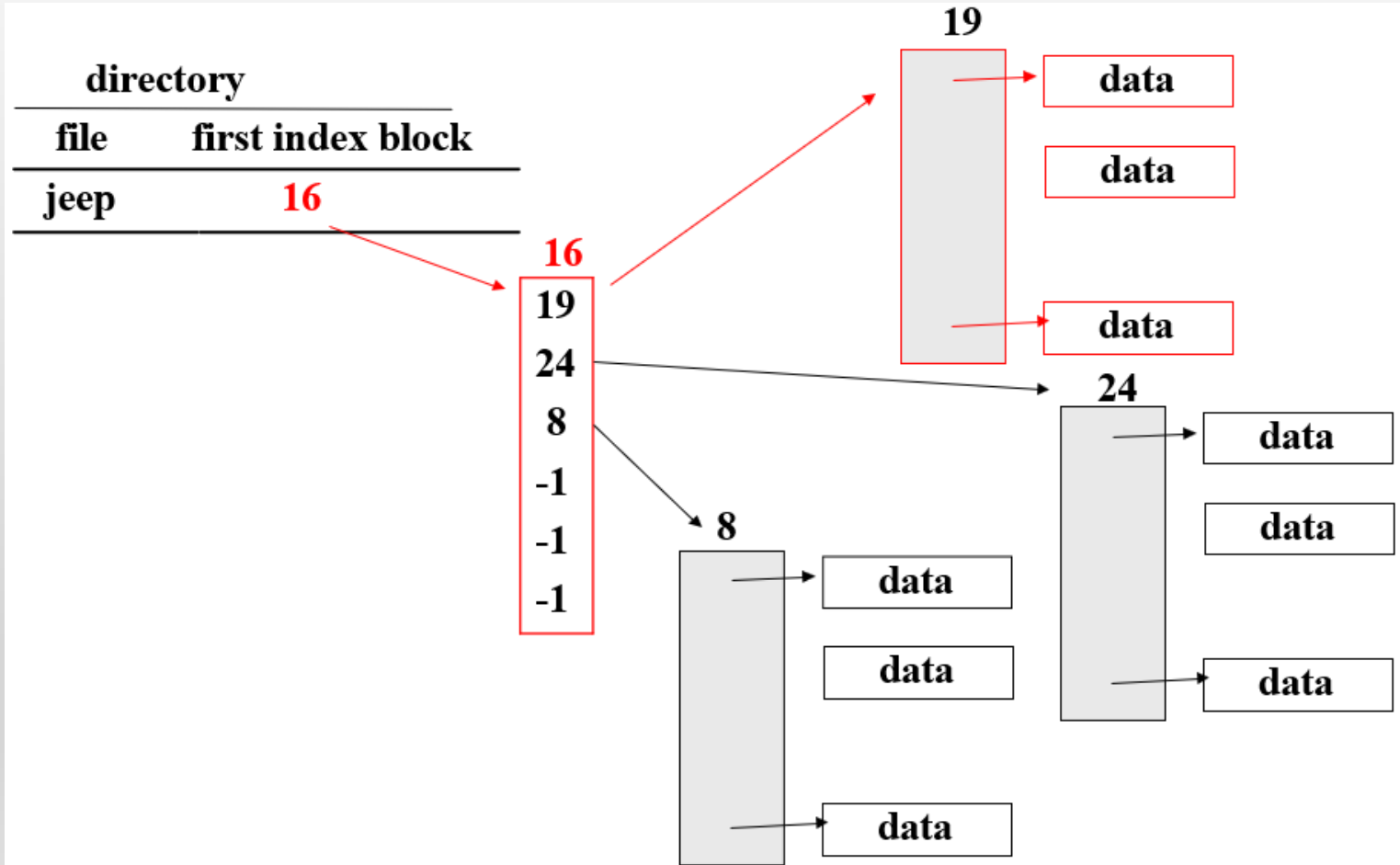
Indexed Allocation

- Bring all the pointers together into one location: the **index block** (one for each file)
- Good:
 1. Implement direct and random access efficiently
 2. No external fragmentation
 3. Easy to create a file (no allocation problem)
- Bad:
 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

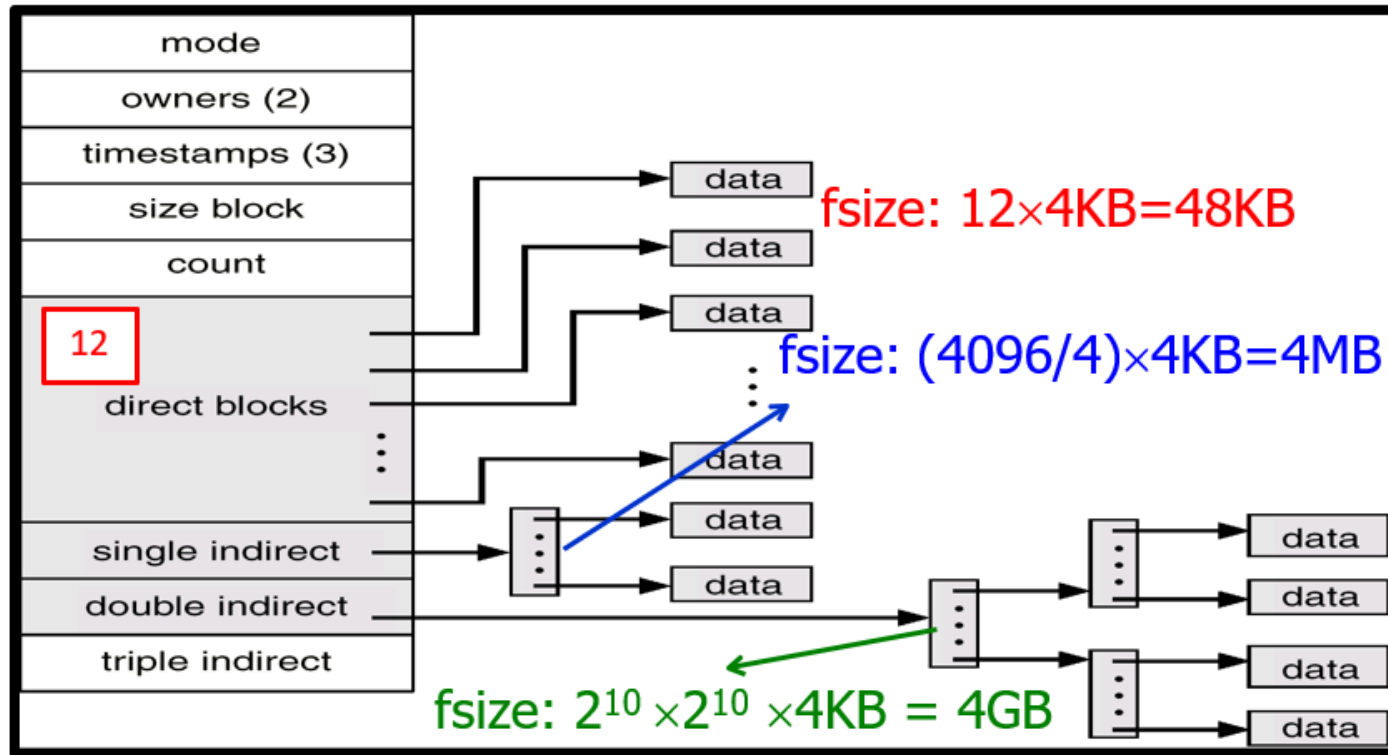


Multilevel Scheme (two-level)



Combined Scheme: UNIX inode

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



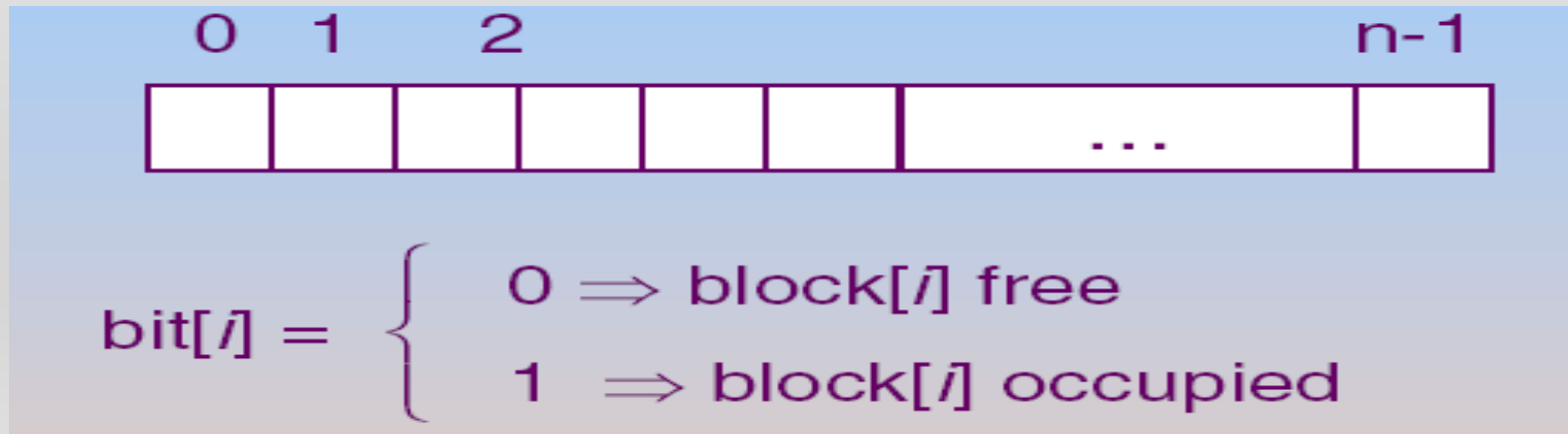
Free-Space Management

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

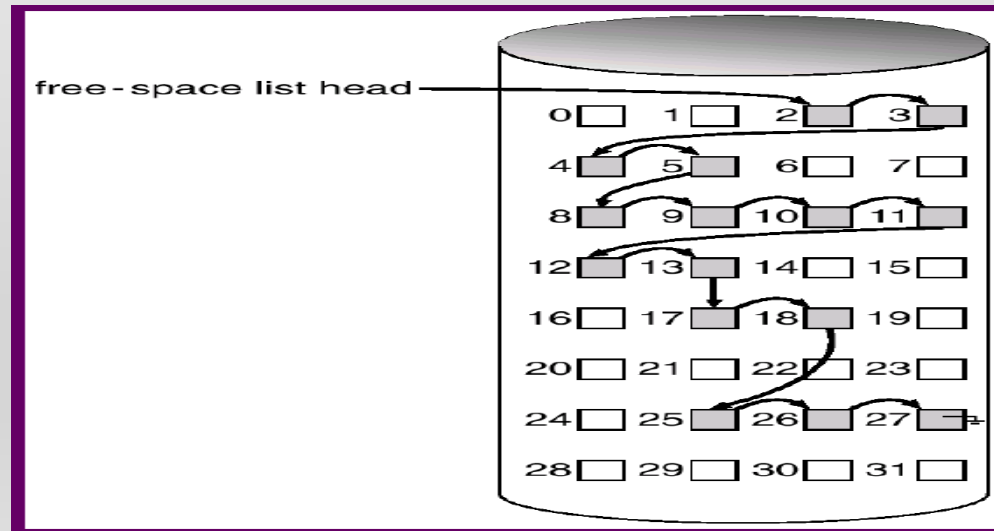
Bit vector

- **Bit Vector (bitmap)**: one bit for each block
 - e.g. 00111100111111111001110011000000.....
- Good: simplicity, efficient
 - (HW support bit-manipulation instruction)
- Bad: bitmap must be cached for good performance
 - A 1-TB(4KB block) disk needs 32MB bitmap



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 pointer 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

Q & A

Thank you for your attention