

# Operating system

## Part X: File System (Interface & Implementation)

By KONG LingBo ( 孔令波 )

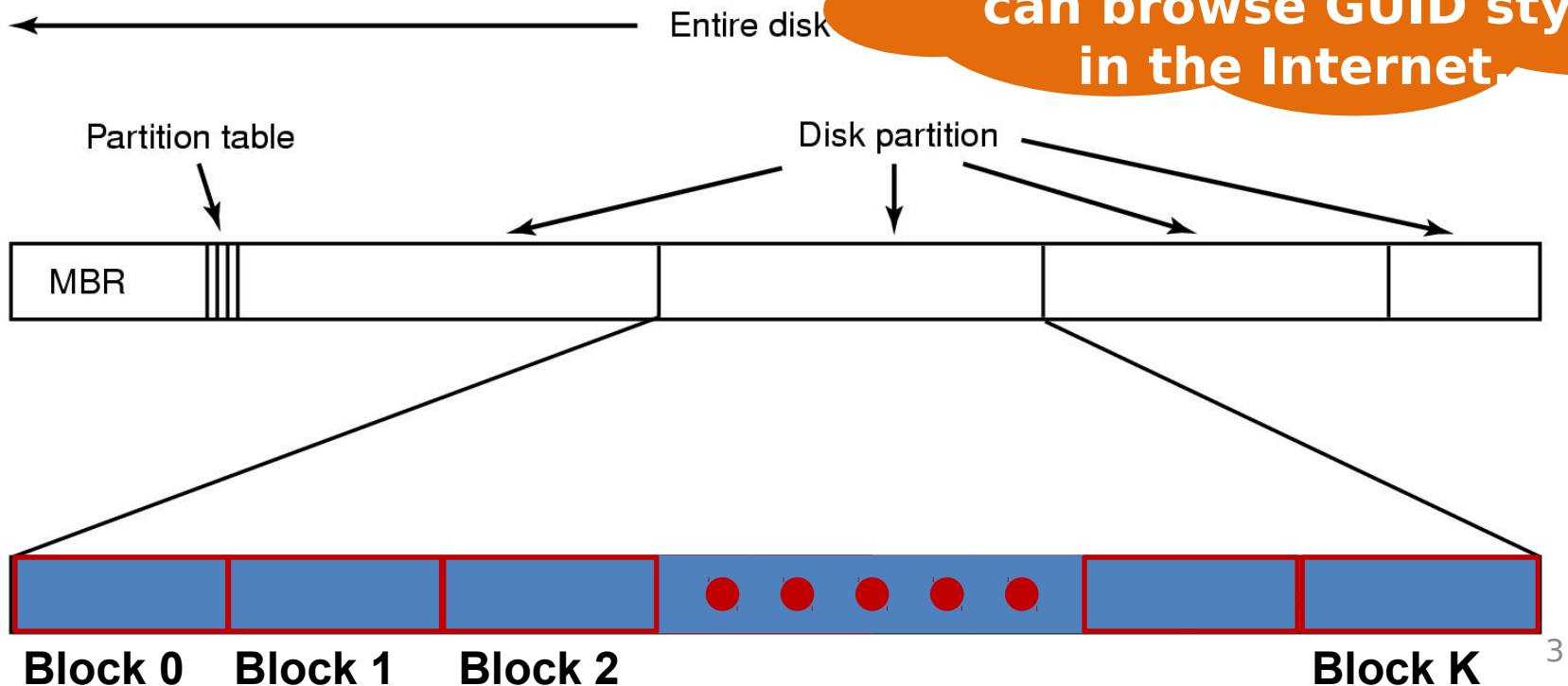
# Goals

- Know the basic concepts related to file system
  - File, directory,
- Know the implementation techniques
  - File organization
  - Directory implementation
- Samples of File systems
  - NFS, NTFS ...

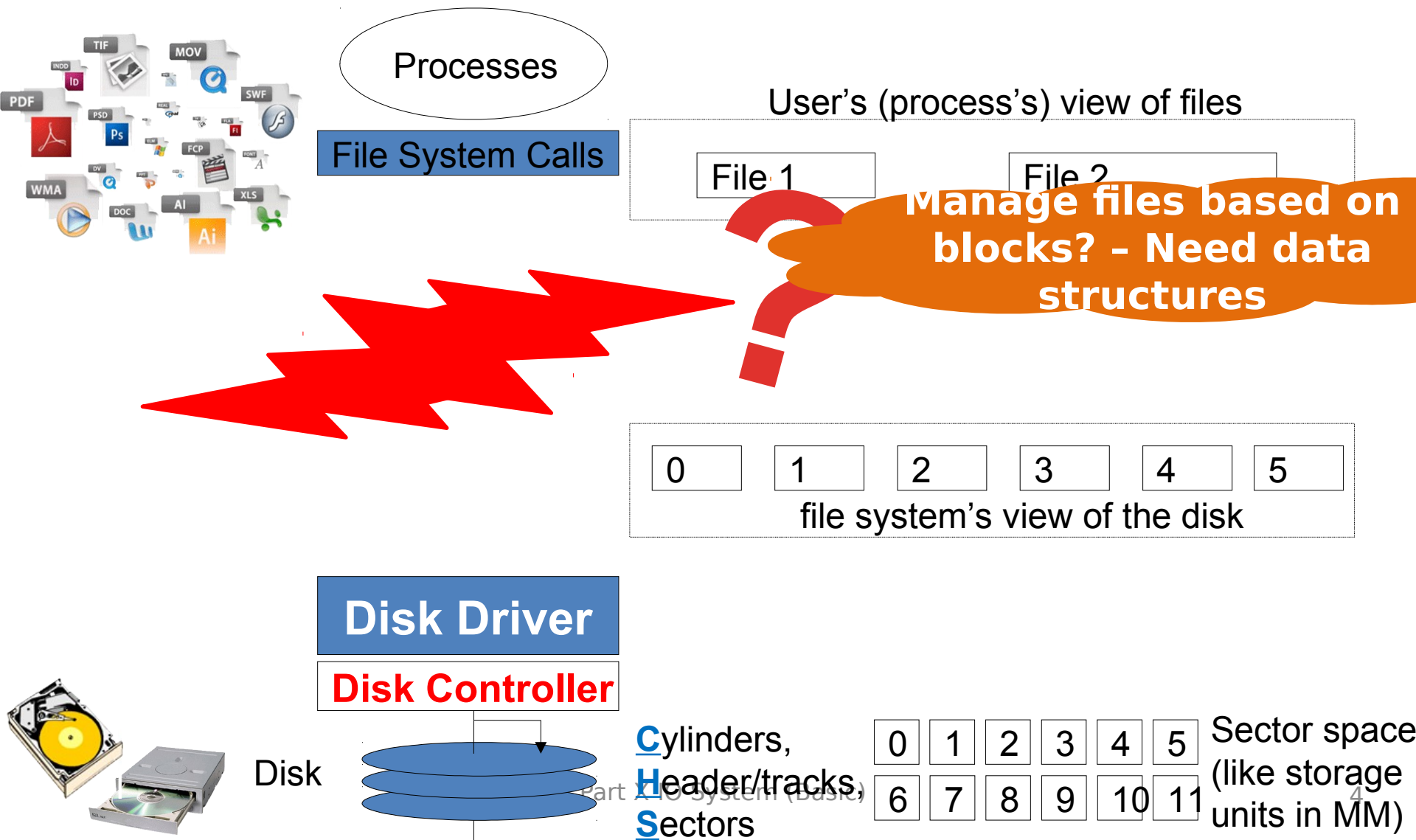
# We've learned - Disk Space Organization

- Disk can be **PARTITIONED**
  - Each partition can have a different OS and/or different file system
  - One partition can be swap space
- Each partition has

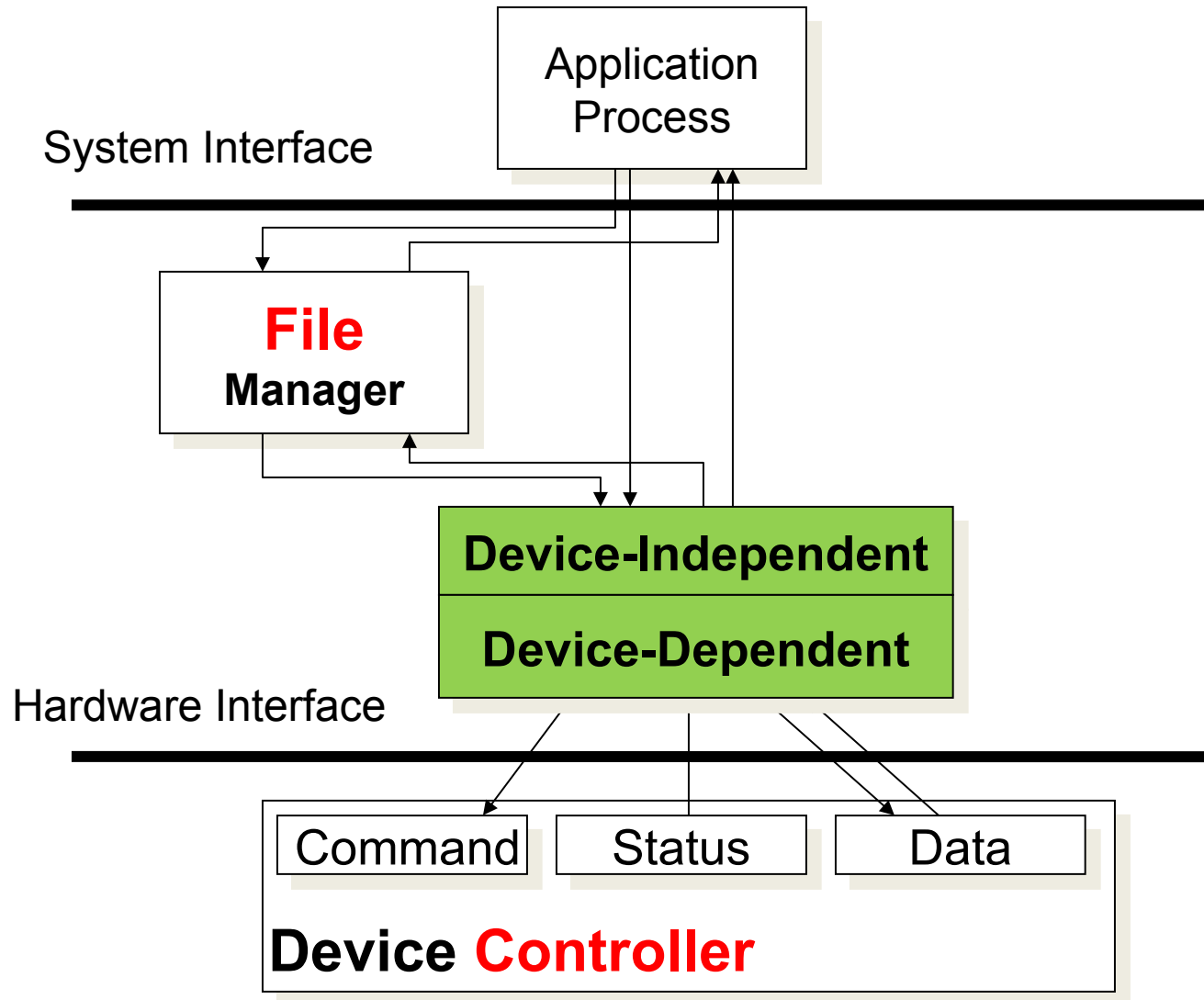
**MBR style is the traditional way to partition the disk. You can browse GUID style in the Internet.**



# We've known - linear addressed block space for files



# File & Device Management



# Problems?

- Like Mapping 1, we need following information – data structures & Related Operations
  - Organize blocks into semantic regions
    - Boot block, superblock, directory
  - Free space
    - You need know where available blocks are.
  - Map a file (collection of bits) into blocks
    - Needed blocks, and how to indicate them?
  - How to organize so many files?
    - Directory

- Mapping 2: File to HDisk (Linear address space to block Space)
  - Organize blocks into semantic regions
    - Boot block, superblock, directory
  - Free space
    - You need know where available blocks are.
  - Map a file (collection of bits) into blocks
    - Needed blocks, and how to indicate them?
  - How to organize so many files?
    - Directory

# Boot block (or sector)

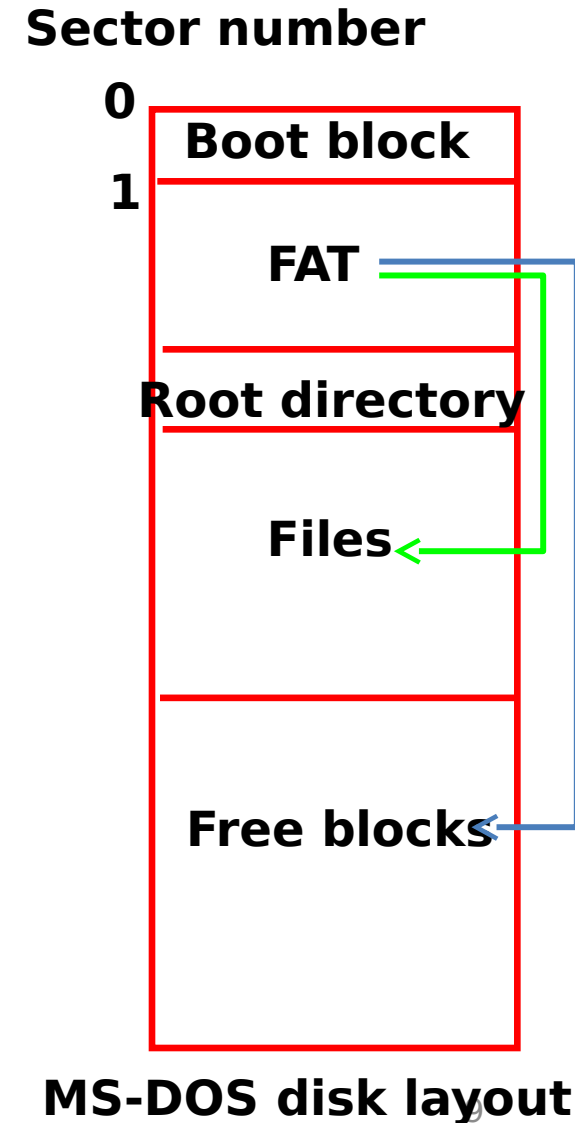
Some old system uses only one sector

- Like the MBR for the hard disk, **each partition usually uses the 1<sup>st</sup> block for special usage, called “boot block”**
  - to record some critical information of how those blocks are used, such as regions for
    - Free space, files, directory, ...
- If the partition contains OS, the boot block also records the information to locate the “bootstrap loader”



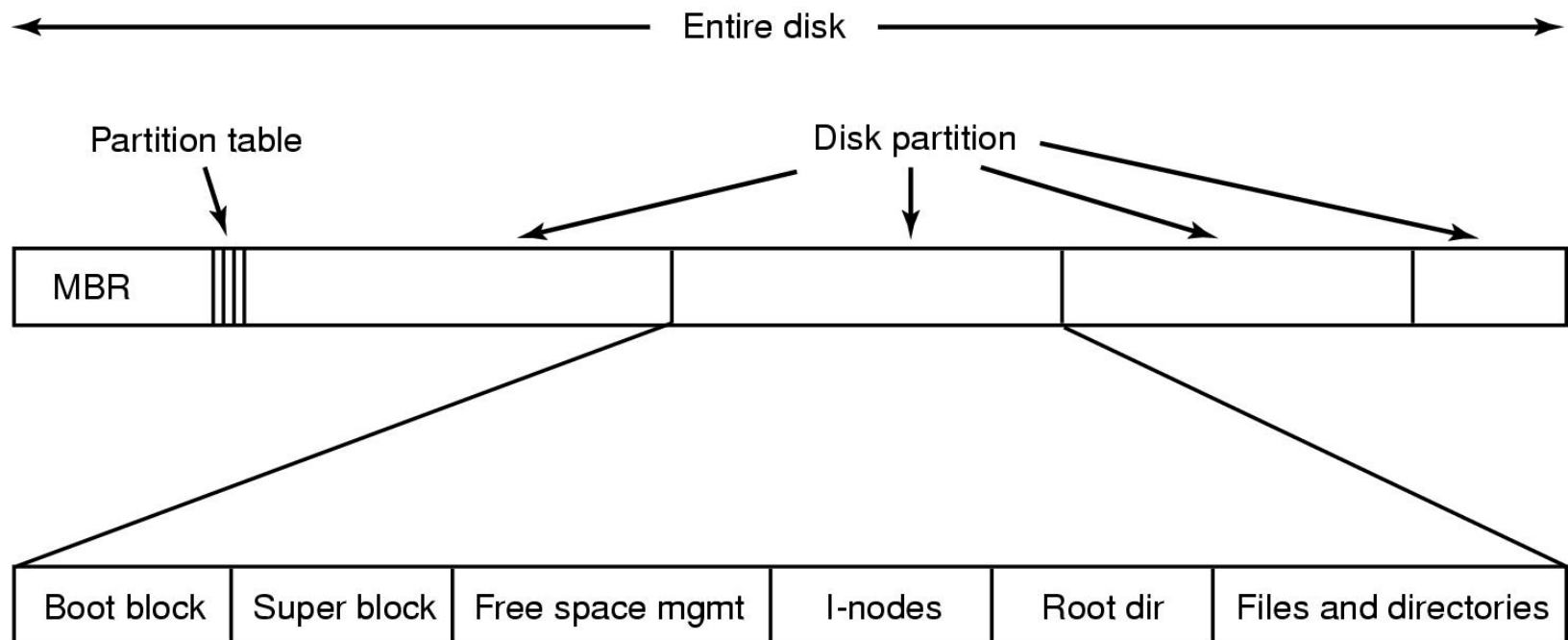
# We have many partition layouts

- MS-DOS layout
  - System disk contains boot block in first block of each partition
    - Boot block has bootstrap executable
  - Bootstrap loader copies bootstrap program into memory
    - Bootstrap program initializes all registers, finds OS kernel on disk, loads OS, and jumps to the initial address of OS



# Cont'

- UNIX like



A possible file system layout *with a **Unix** partition*

- Mapping 2: File to HDisk (Linear address space to block Space)
  - Organize blocks into semantic regions
    - Boot block, superblock, directory
  - Free space
    - You need know where available blocks are.
  - Map a file (collection of bits) into blocks
    - Needed blocks, and how to indicate them?
  - How to organize so many files?
    - Directory

# Free Space Management

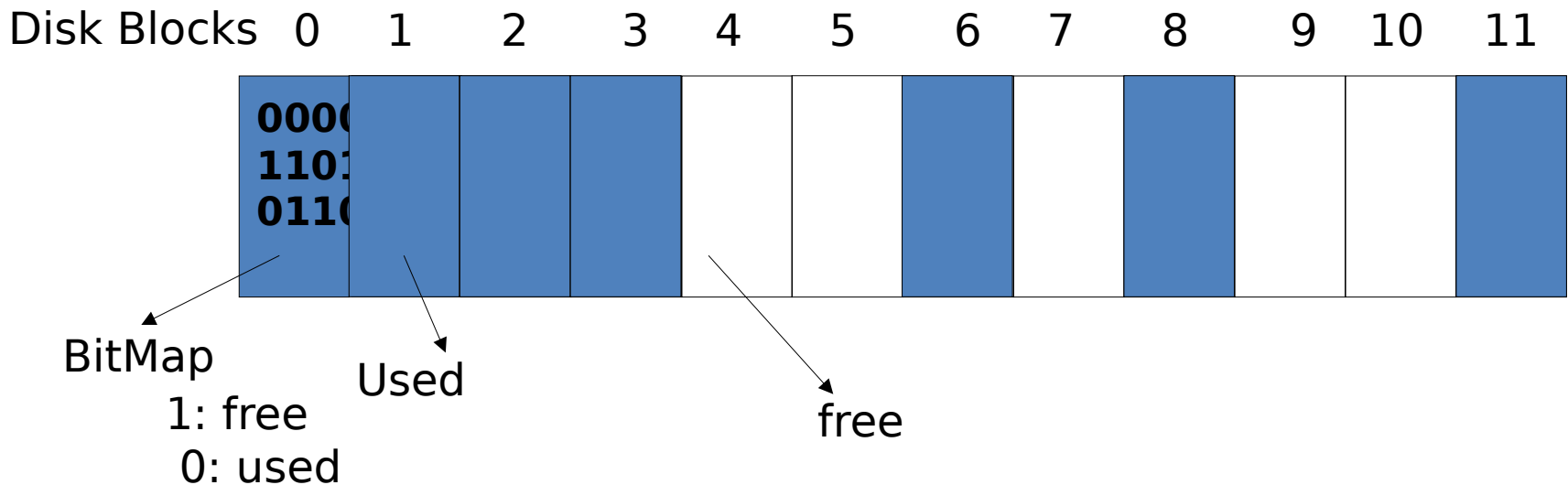
- How can we keep track of **free blocks** of the disk?
  - Which blocks are free?
- We need this info when we want to allocate a new block to a file.
  - Allocate a block that is free.
- There are several methods to keep track of free blocks:
  - Bit vector (bitmap) method
  - Linked list method
  - Grouping
  - Counting

# Free-Space Management:

## Bit Vector (Bit map: 位图)

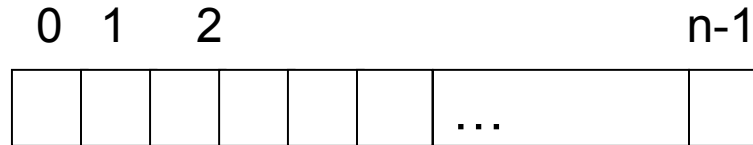
- We have a bit vector (bitmap) where we have **one bit per block** indicating if the block is used or free.
- If the block is free the corresponding bit can be 1 else it can be 0 (or vice versa).

Example:



# Cont'

- Bit vector ( $n$  blocks in disk)



$\text{bit}[i] =$ 
  
 0  $\rightarrow$  block $[i]$  used (or vice versa)
   
 1  $\rightarrow$  block $[i]$  free

Finding a free block (i.e. its number)

Start searching from the beginning of the bitmap:  
Search for the first 1

First Free Block Number =

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

```

0000000000000000
0000000000000000
0000000000000000
0000000010000000
0000000000000000
0000000000011000
0001100010000000
0000000011110000
  
```

$$3 \times 16 + 8 = \mathbf{56}$$

# Simple check

- If the block size is 4KB, and the hard disk size (or a partition) is 500 GB, how many blocks are used to store the bitmap?
  - $1 \text{ GB} = 2^{30} \text{ B}$
- You'll meet this kind of questions all the time
  - How many blocks are used to store the needed information?
  - Do you know in HDD, GB means  $10^6 \text{ B}$ ?

- Mapping 2: File to HDisk (Linear address to block Space)
  - Organize blocks into semantic regions
    - Boot block, superblock, directory
  - Free space
    - You need know where available blocks are.
  - Map a file (collection of bits) into blocks
    - Needed blocks, and how to indicate them?
  - How to organize so many files?
    - Directory



# File System Implementation:

## File Space Allocation

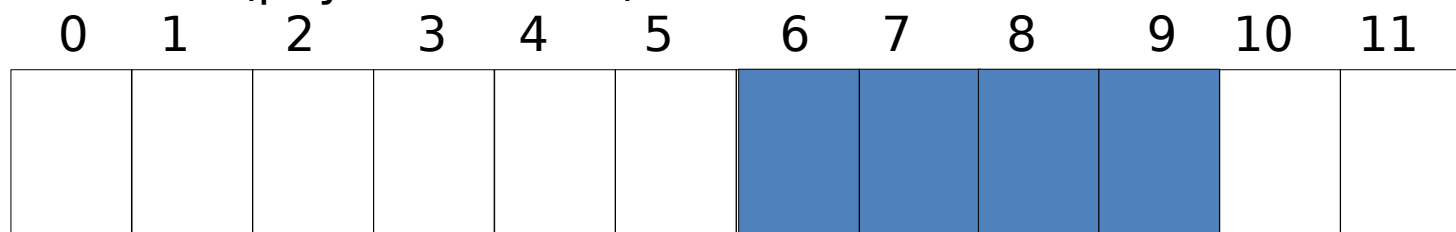
- Goals
  - Fast sequential access
  - Fast random access
  - Ability to dynamically grow
  - Minimum fragmentation
- Standard schemes
  - Contiguous allocation (fixed)
  - Linked list allocation
  - Linked list with file allocation table (FAT)
  - Linked list with Indexing (I-nodes)

# Contiguous Allocation

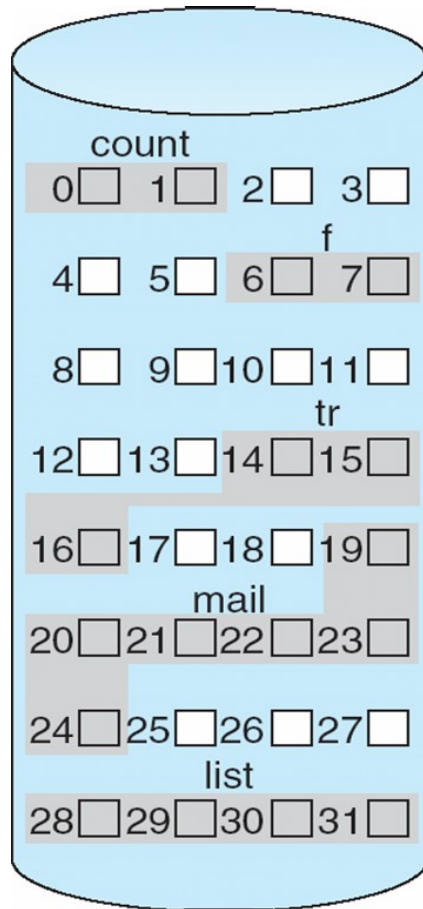
- Each file occupies a set of contiguous blocks on the disk
- + Simple – only starting location (block #) and length (number of blocks) are required to find out the *disk data blocks of file*
- + Random access is fast
- - Wasteful of space (dynamic storage-allocation problem)
- - File **file data**

Start address = 6  
Number of blocks = 4

disk blocks (physical blocks)



# Contiguous Allocation of Disk Space

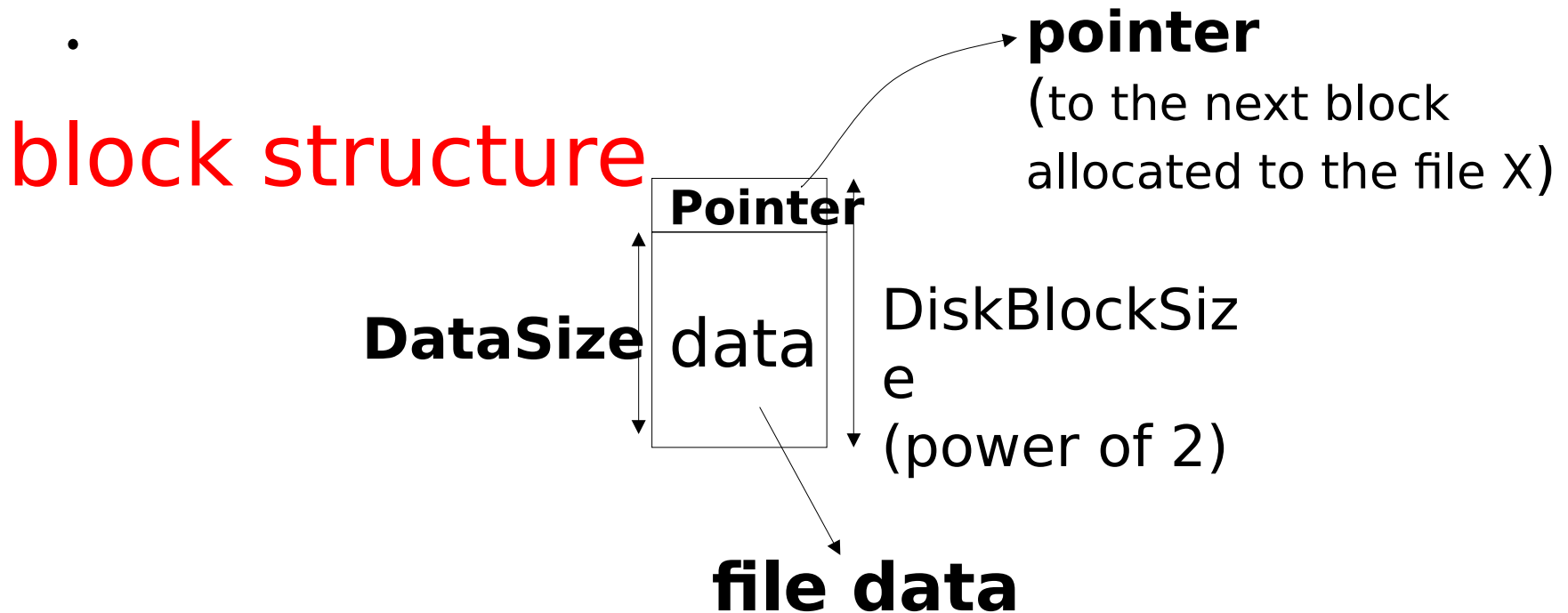


directory

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

# Linked Allocation

- Each file is a **linked list of disk blocks**: blocks may be scattered anywhere on the disk



data size in a disk block is no longer a power

# Linked Allocation (cont' )

File X

File starts at disk block 5

pointer

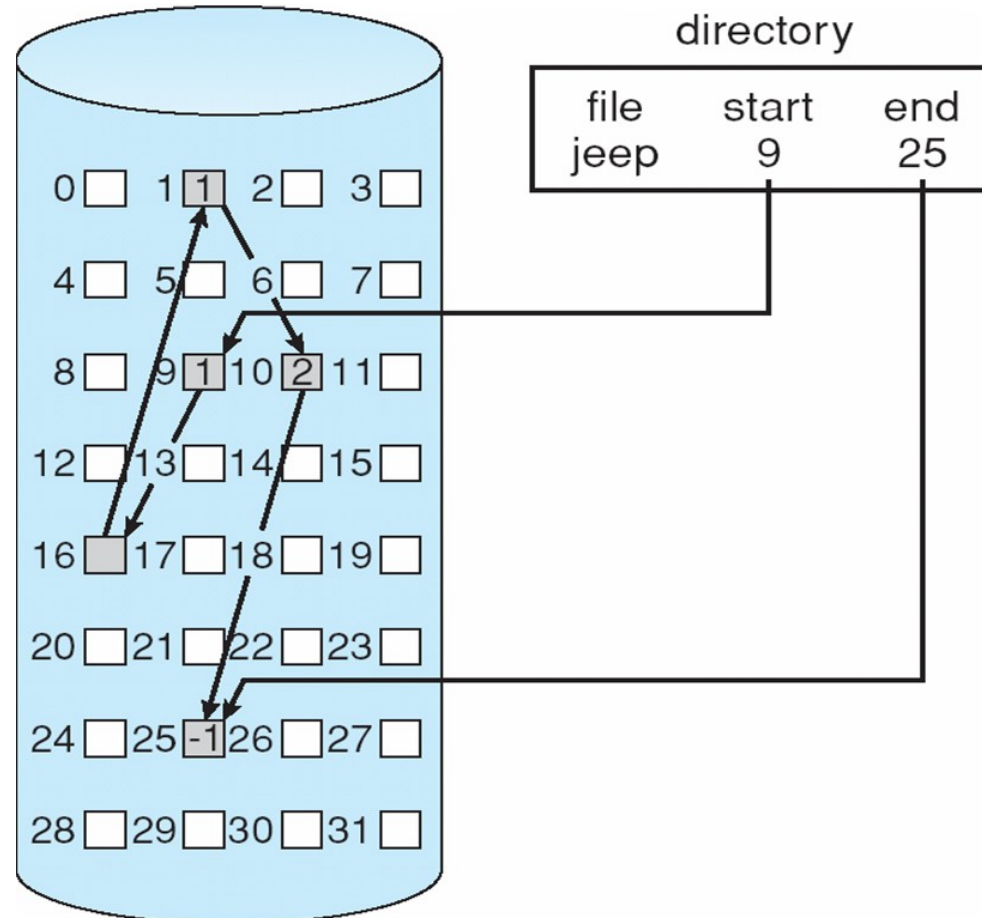
**disk blocks (physical blocks)**

0 1 2 3 4 5 6 7 8 9 10 11

			8		3			10			

data

# Linked Allocation



# Linked Allocation (cont.)

- Slow - defies principle of locality
  - Need to read through many records sequentially to find the record
- Not very reliable
  - System crashes can corrupt pointers
- Important variation on linked allocation method
  - File-allocation table (FAT) - disk-space allocation used by MS-DOS and OS/2.

Just as those shortcomings of linked list you learned in DSA

# File Allocation Table

- The File Allocation Table (**FAT**) is **a variation to the linked allocation** method used to support direct access
  - disk-space allocation used by MS-DOS and OS/2.
- The FAT file system is a simple file system originally designed for small disks and simple folder structures.
  - named for its method of organization, the **file allocation table**, which resides at the beginning of the volume.
- To protect the volume, two copies of the table are kept, in case one becomes damaged.
  - In addition, the file allocation tables and the root folder must be stored in a fixed location so that the files needed to start the system can be correctly located

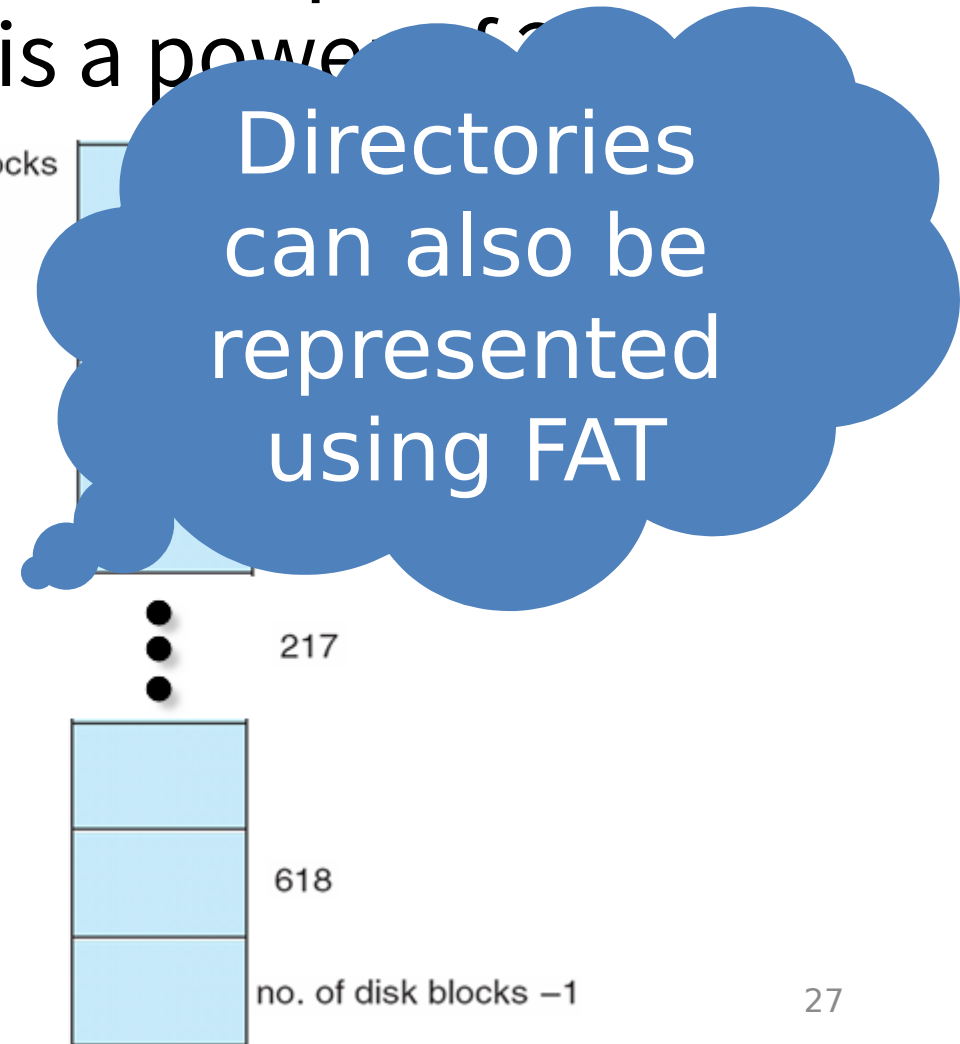
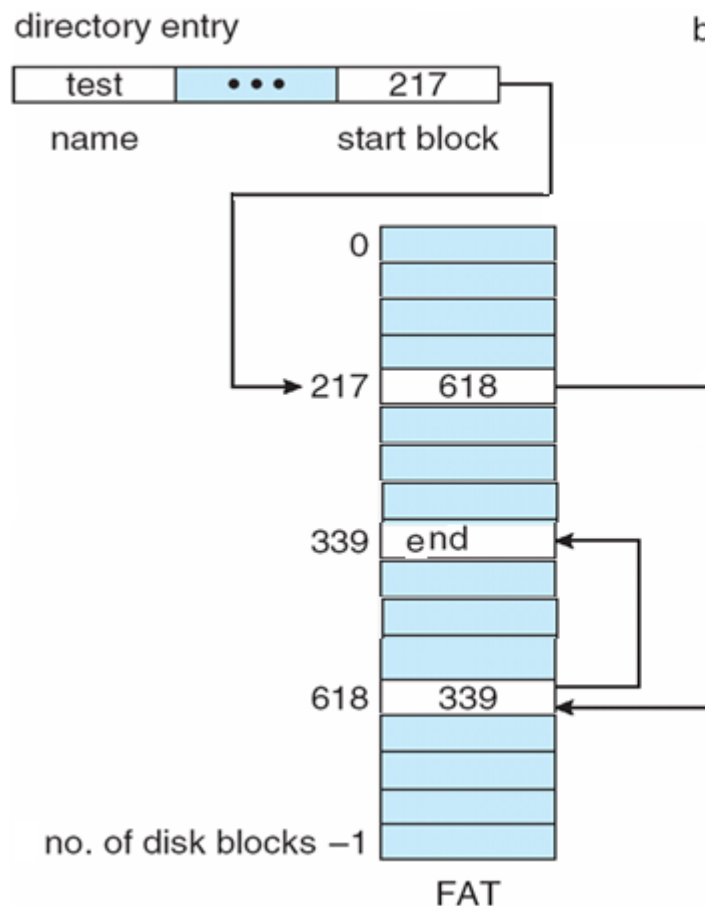


# File Allocation Table

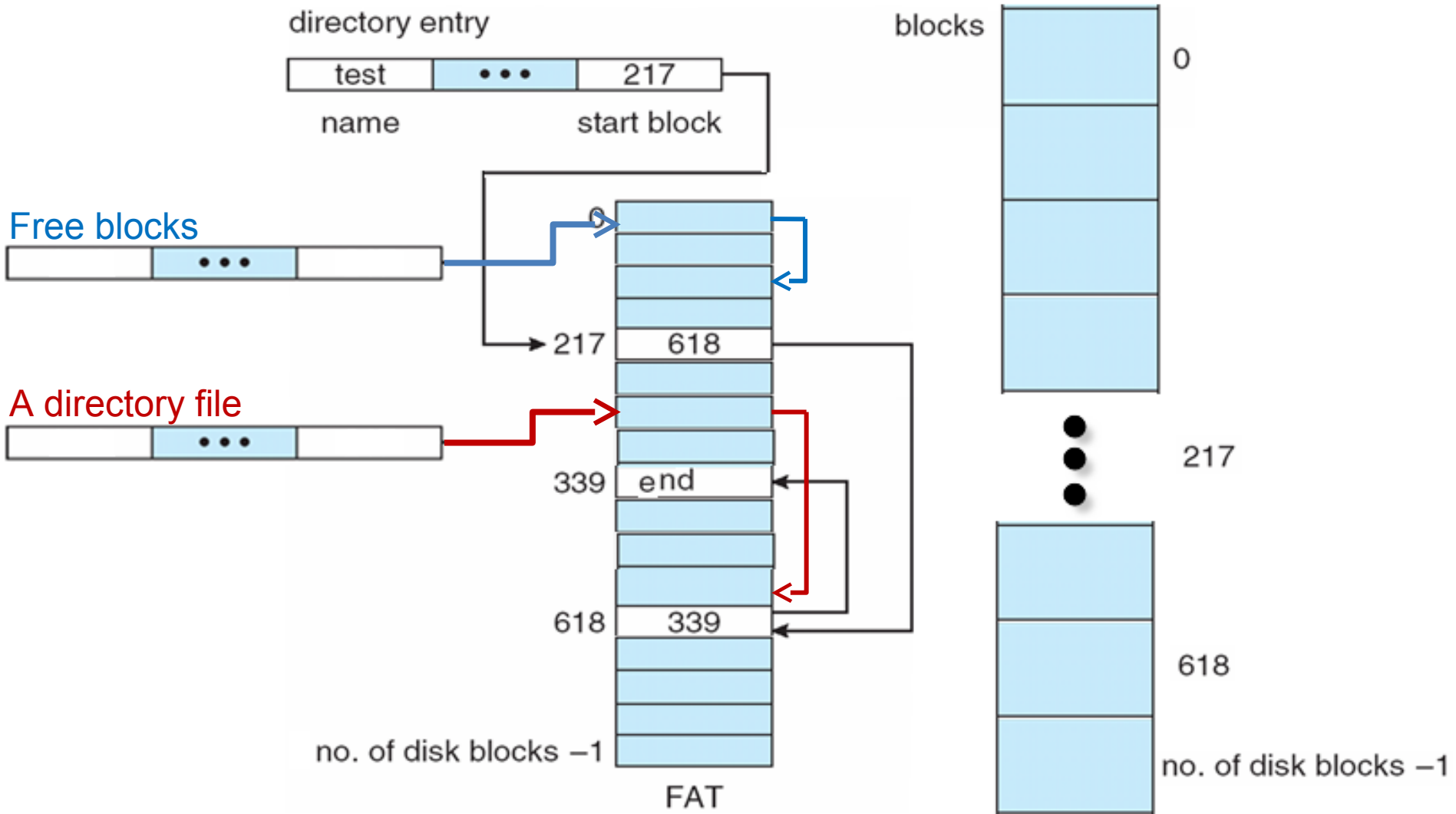
- The File Allocation Table (FAT) has many versions: FAT12, 16, 32 ...
  - **FAT12** is only seen on floppy disks and very small storage media, while **FAT16** is the older version of FAT from the Windows 95 days, and **FAT32** is newer, from the Windows 98 days.
- NT, 2000, XP, Vista, and Windows 7 can use all the FAT file systems, plus the **NTFS** (New Technology File System)

# Cont'

- **Pointers are kept in a table (FAT)**
- Data Block does not hold a pointer; hence data size in a block is a power of 2



# Cont'



# Cont'

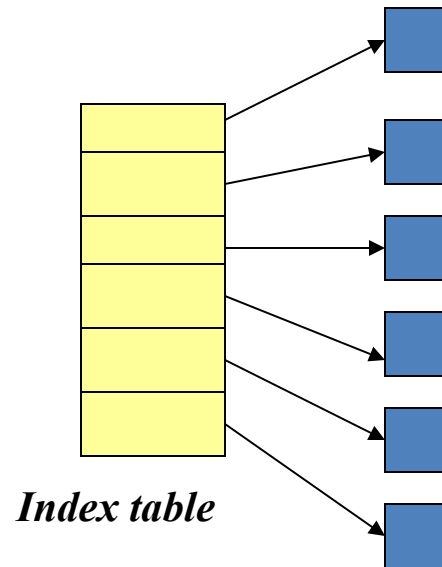
- Of course the FAT is also stored in blocks
- So the size limit of a file is determined by
  - number of FAT entries
  - Size of block
- Given:
  - Block size = 4096 bytes
  - Length of a entry in FAT = 8 bytes
  - ② one block can contain  $4096/8 = 512$  entries
  - So the size of a file by using only one block could be:  $512 * 4KB = 2048KB = 2GB$

# Questions

- A hard disk has 40G, its each block size is 1 K , and each table entry of FAT needs 20 bits , then Its FAT (File Allocation Table) need ( ) memory space
  - A ) 100M                      B ) 120M
  - C ) 140M                      D ) 160M

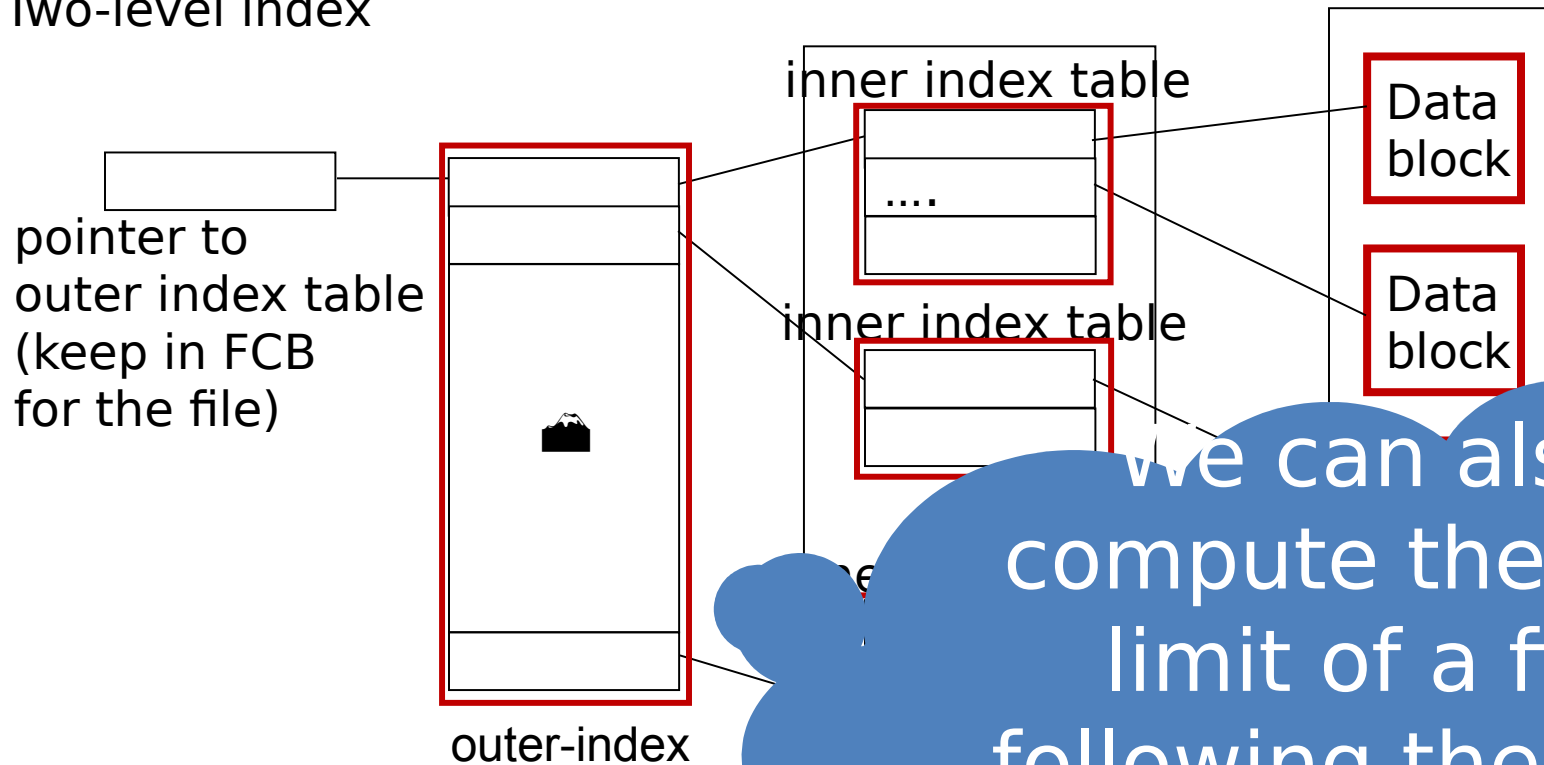
# Indexed Allocation

- Brings all pointers together into the index block.
- Logical view



# Indexed Allocation – Mapping (Cont.)

Two-level index



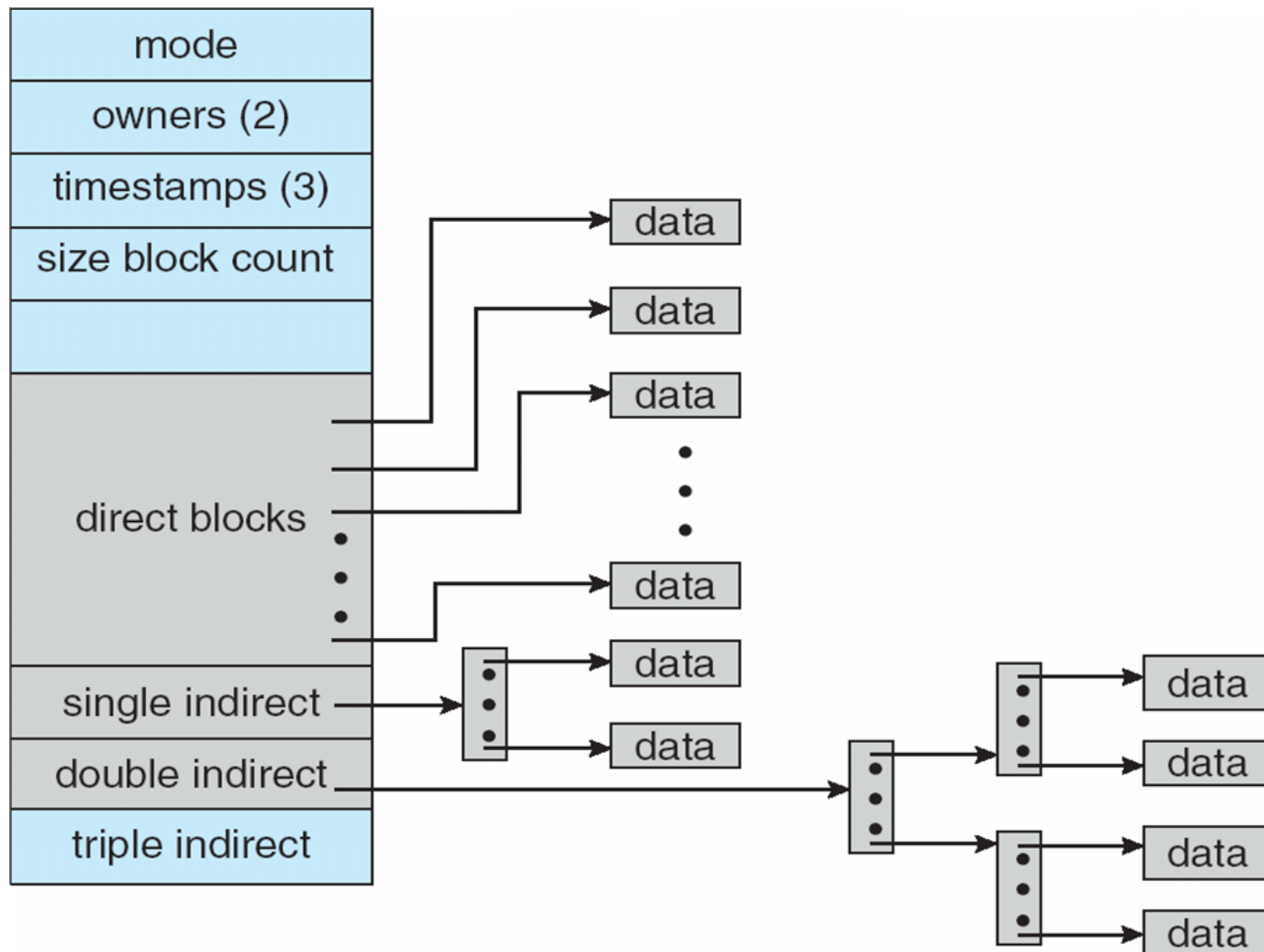
we can also compute the size limit of a file following the idea introduced in FAT. Attention to the levels

# Indexed Allocation - Mapping

- Mapping from logical to physical in a file of unbounded length.
- Linked scheme -
  - Link blocks of index tables (no limit on size)
- Multilevel Index
  - E.g. Two Level Index - first level index block points to a set of second level index blocks, which in turn point to file blocks.
  - Increase number of levels based on maximum file size desired.
  - Maximum size of file is bounded.

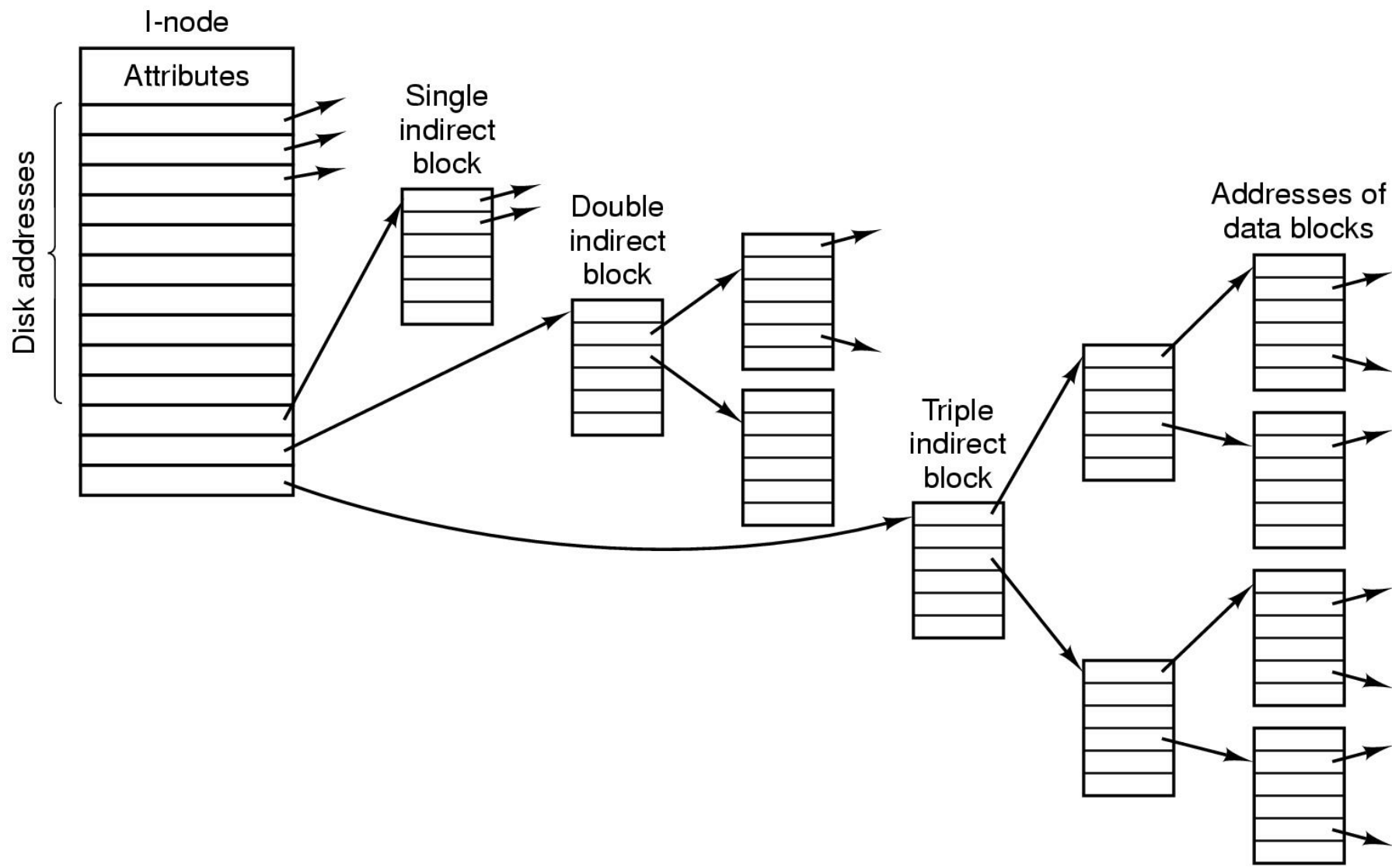


# Combined Scheme: UNIX (4K bytes per block)

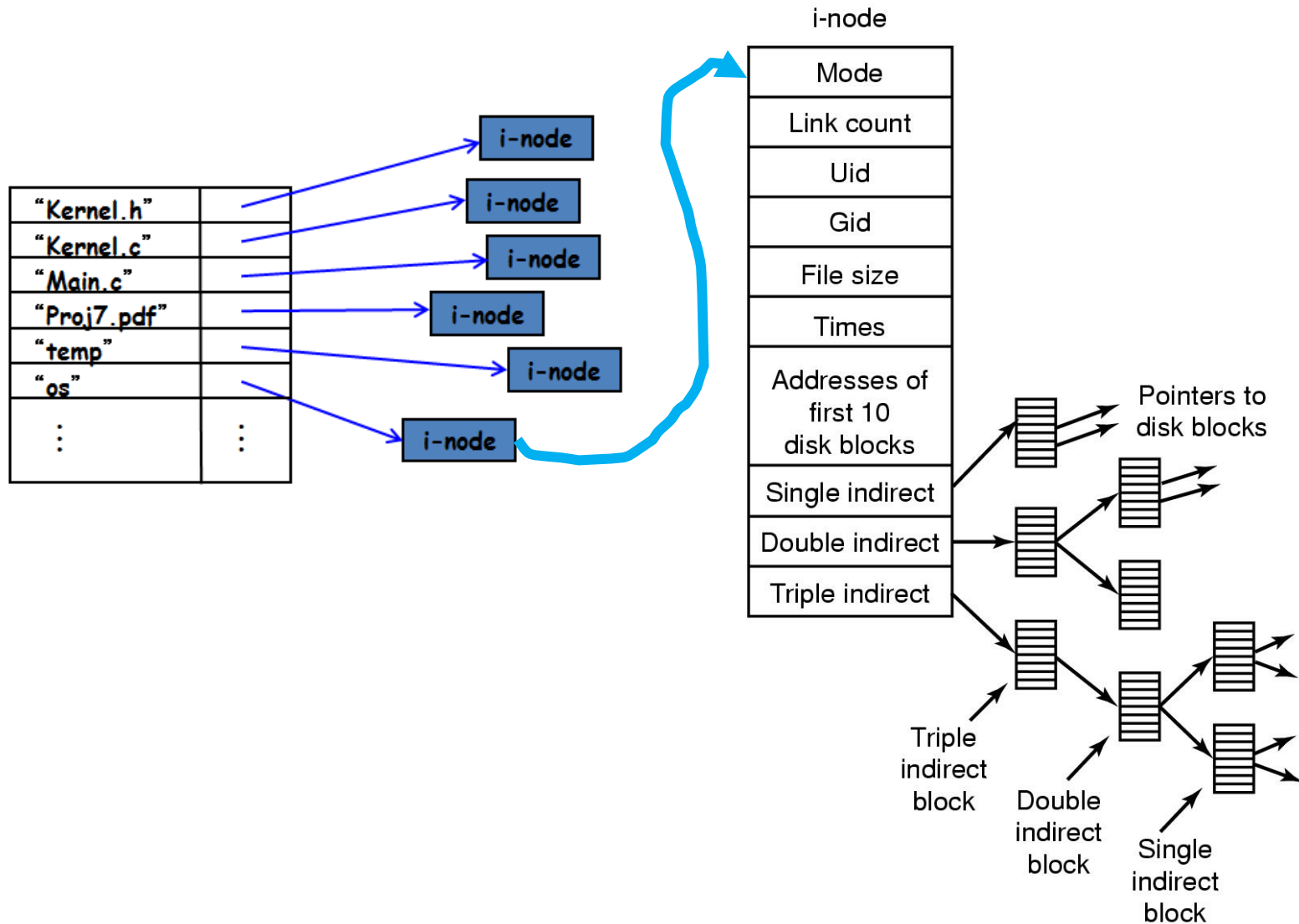


FCB - in Unix this is called **i-Node**

# Unix i-node



# The UNIX I-node



- Assume that the index-node of a file has 7 pointers, among which 4 are for data blocks, 2 pointers for “indirect block”, and 1 pointer for “doubly indirect block”.
- The size of both data block and index block are 256KB. The length of a pointer is 4 Bytes.
- What is the max file size this file system can support?

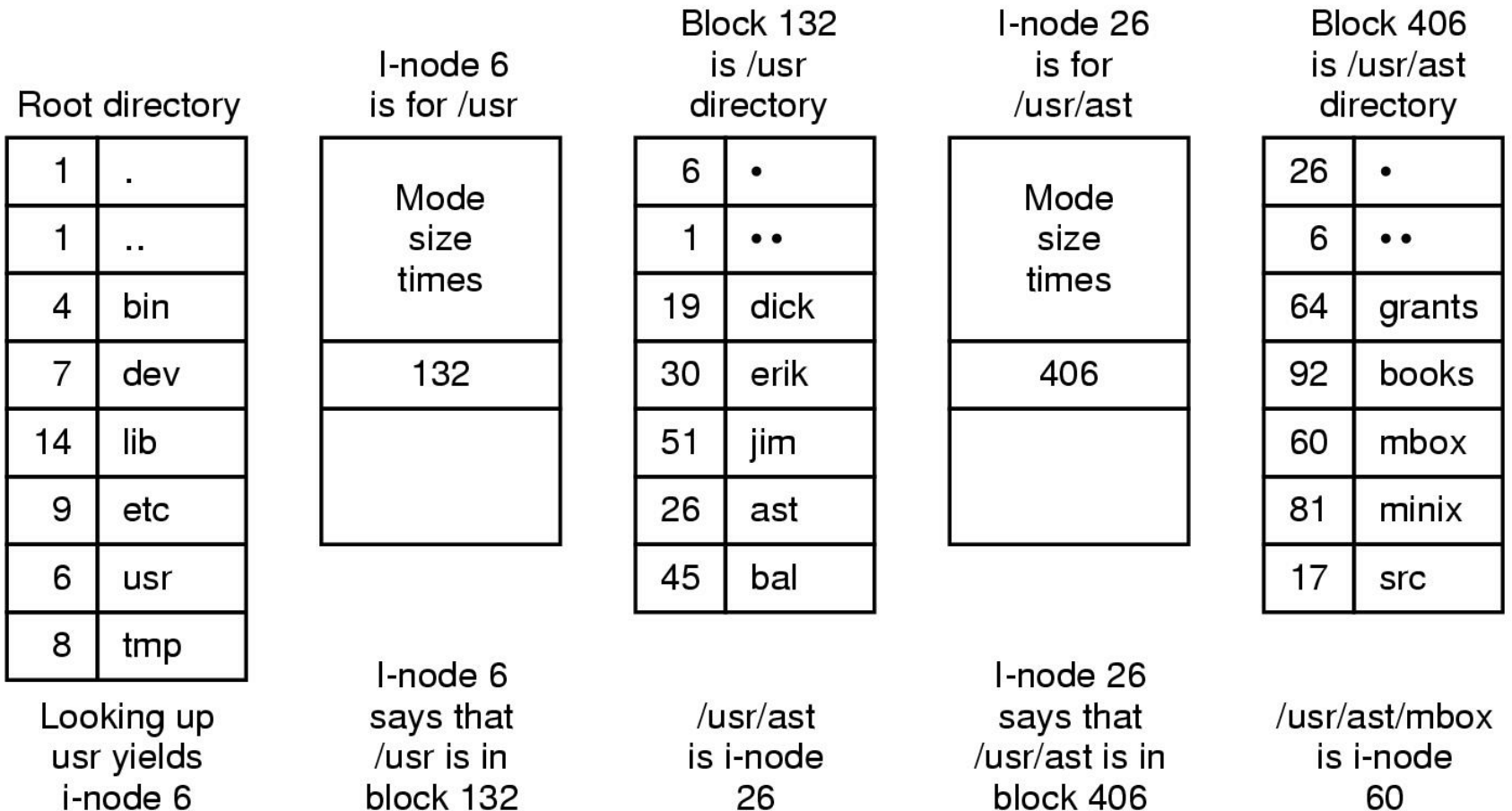
# Entry Lookup

- ❑ In Unix systems
  - The **superblock** (among other things) has the location of the i-node which represents the root directory
- ❑ Once the root directory is located a search through the directory tree finds the desired directory entry
- ❑ The directory entry provides the information needed to find the disk blocks for the requested file

# Entry Lookup

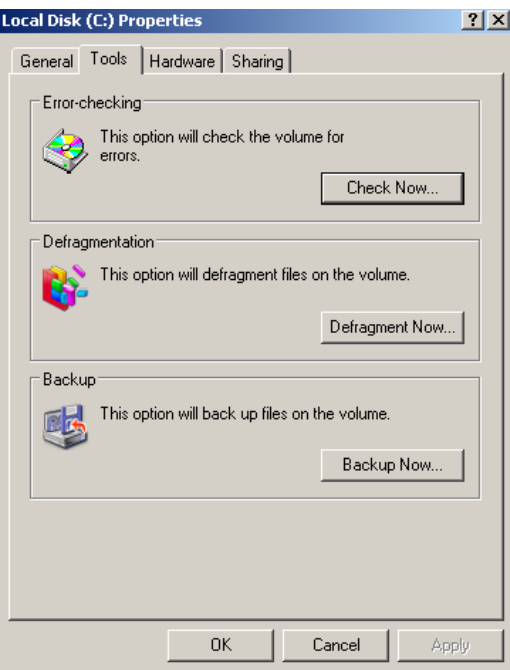
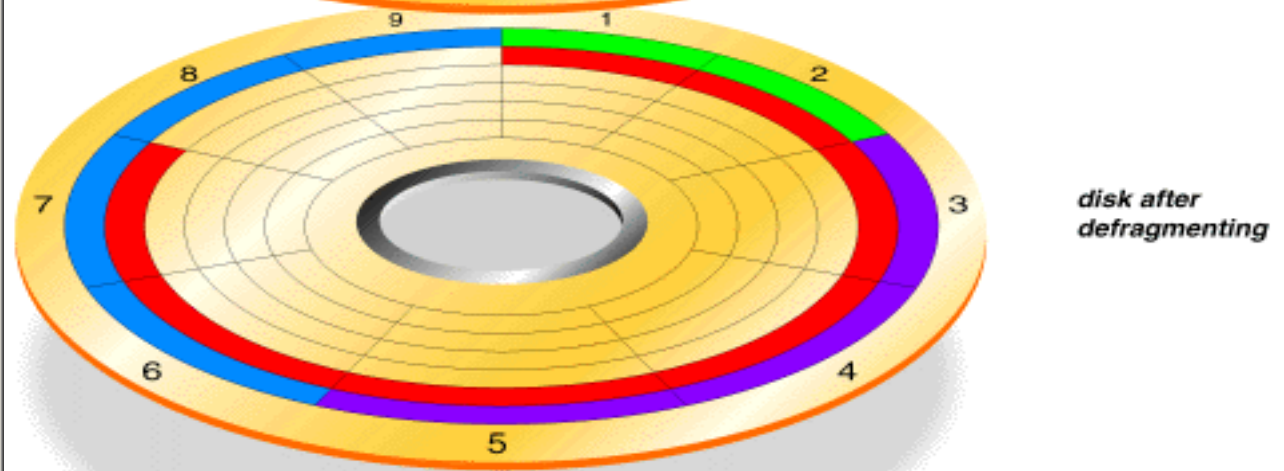
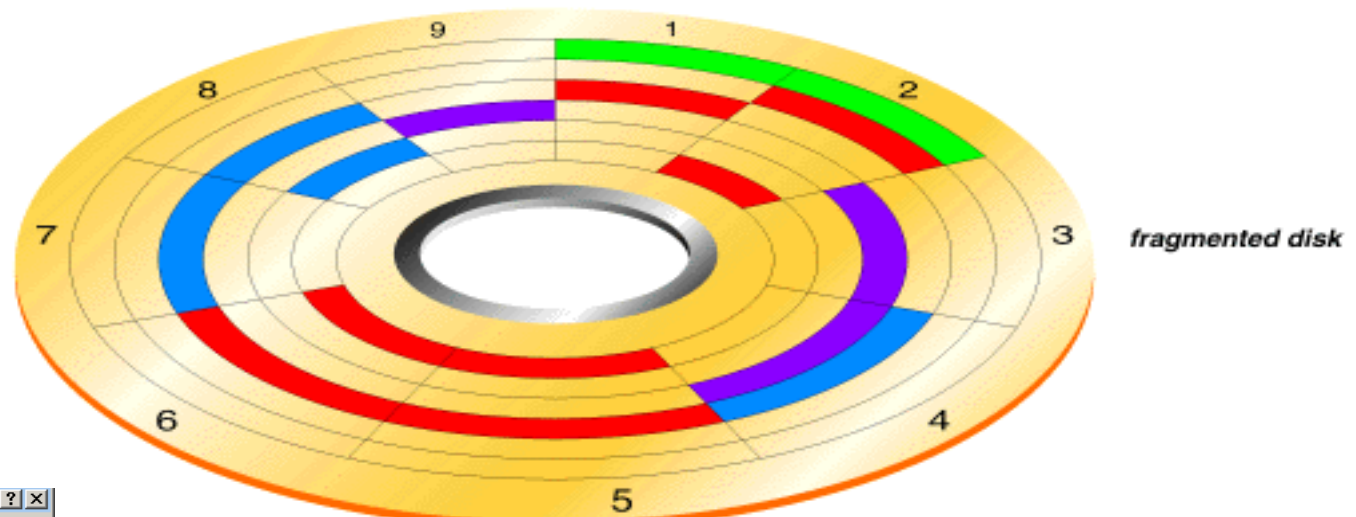
- ❑ This discussion focuses on Unix-related file systems
- ❑ When a file is opened, the file system must take the file name supplied and locate its disk blocks
- ❑ Let's see how this is done for the path name */usr/ast/mbox*

# Looking up for an entry



The steps in looking up `/usr/ast/mbox`

# Fragmentation and Defragmenting for disks

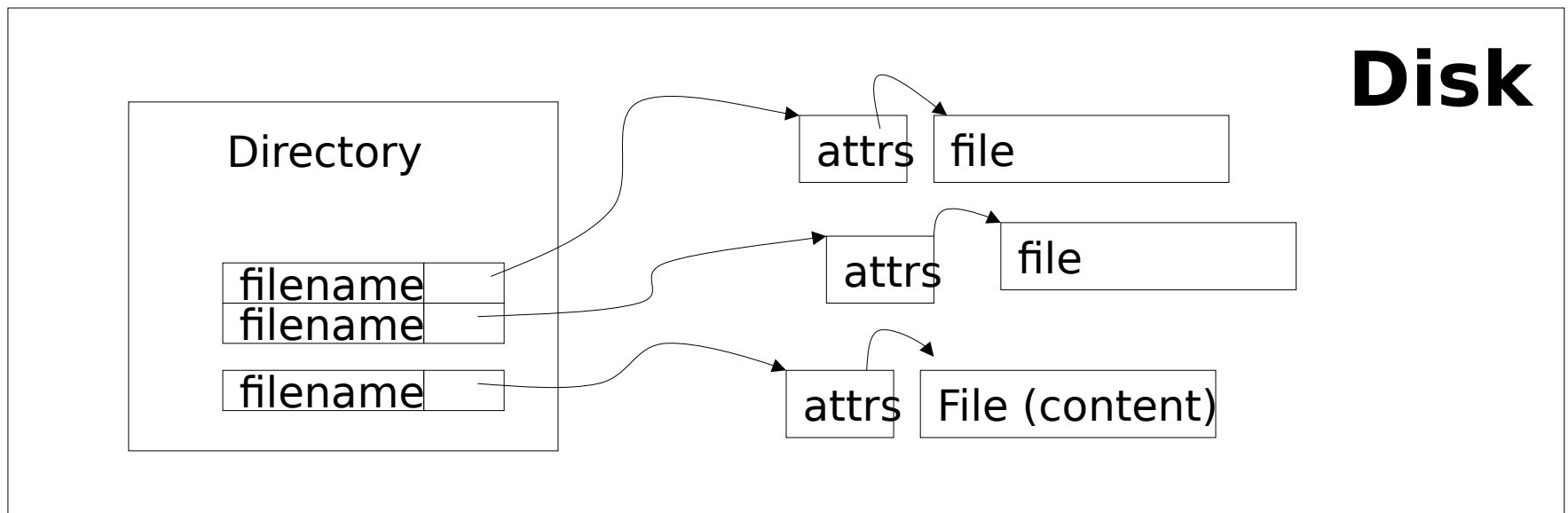




- Mapping 2: File to HDisk (Linear address space to block Space)
  - Organize blocks into semantic regions
    - Boot block, superblock, directory
  - Free space
    - You need know where available blocks are.
  - Map a file (collection of bits) into blocks
    - Needed blocks, and how to indicate them?
  - How to organize so many files?
    - Directory

# Files and Directories

- There are two basic things that are stored on disk as part of the area controlled by the file system
  - **files** (store content)
  - **directory** information (can be a tree): keeps info about files, their attributes or locations

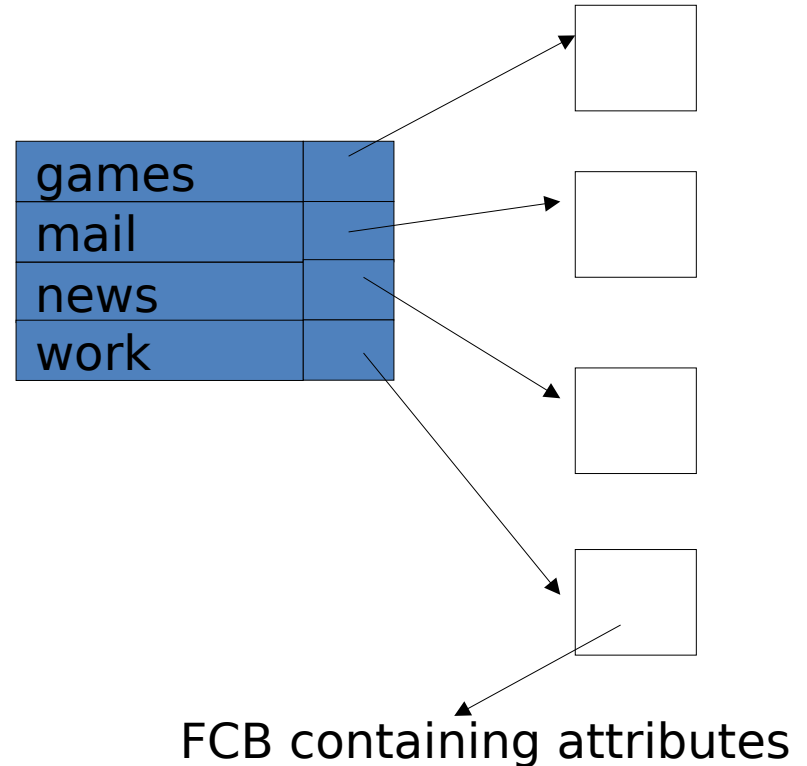


# Directory Implementation: directory entries

games	attributes
mail	attributes
news	attributes
work	attributes

a directory with fixed sized entries

attributes include location  
info for data blocks  
of the file

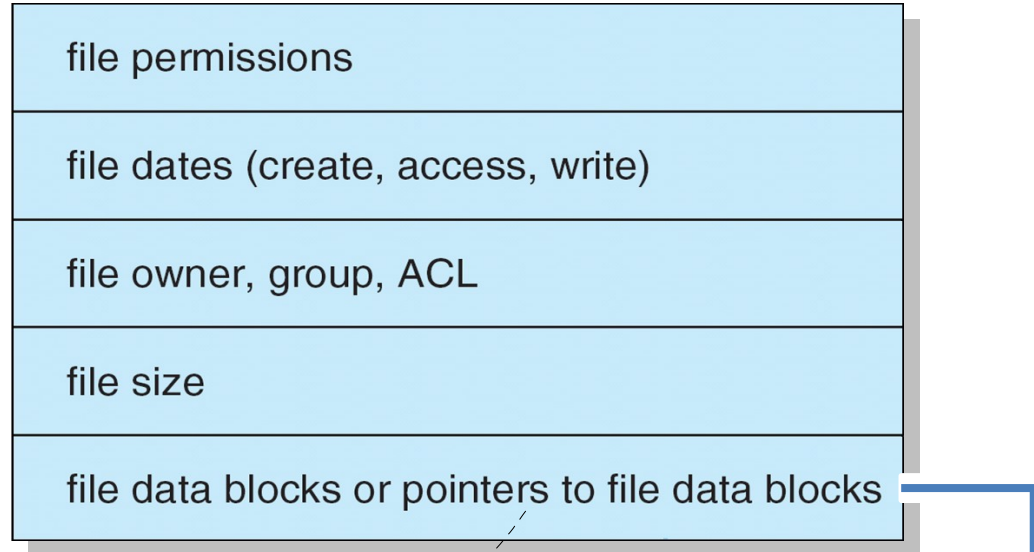


Using fixed sized names

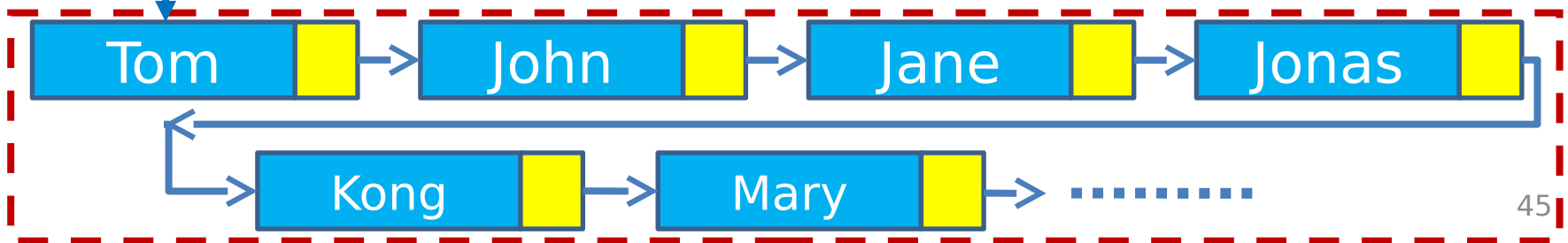
# A Typical File Control Block

Filename=X | info about locating the FCB | Directory entry

File Control Block of a file with filename X



File Data Blocks of X



# Cont'

- Linear list of file names with pointers to the data blocks
  - simple to program
  - time-consuming to execute - linear search to find entry.
  - Sorted list helps - allows binary search and decreases search time.
- Hash Table - linear list with hash data structure
  - decreases directory search time
  - collisions - situations where two file names hash to the same location.
  - Each hash entry can be a linked list - resolve collisions by adding new entry to linked list.

# Of course, you should provide File Operations

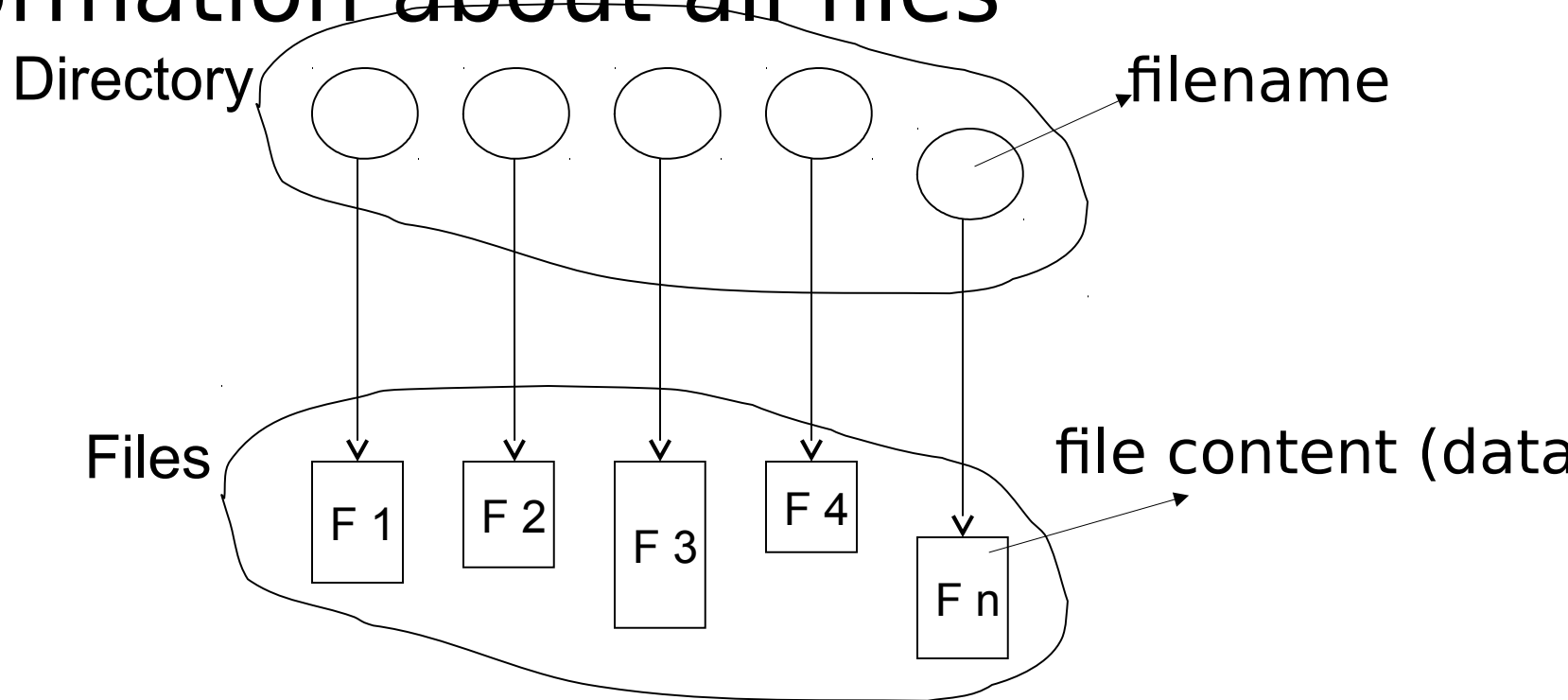
- File is an **abstract data type**
- **Common Operations that are supported by the Operating System:**
  - Create
  - Write
  - Read
  - Reposition within file
  - Delete
  - Truncate
- **Open( $F_i$ )** – search the directory structure on disk for entry  $F_i$ , and move the content of entry to memory
- **Close ( $F_i$ )** – move the content of entry  $F_i$  in memory to directory structure on disk

## Directory Implementation

- A directory entry provides the info needed to find the disk data blocks of a file
  - disk address of first block and size
  - address of first block
  - number of associated i-node
- File attributes can be stored in the directory entry (Windows) or in its i-node (Unix)
- File name and support of variable length and long file names (255 chars)

# Directory Structure

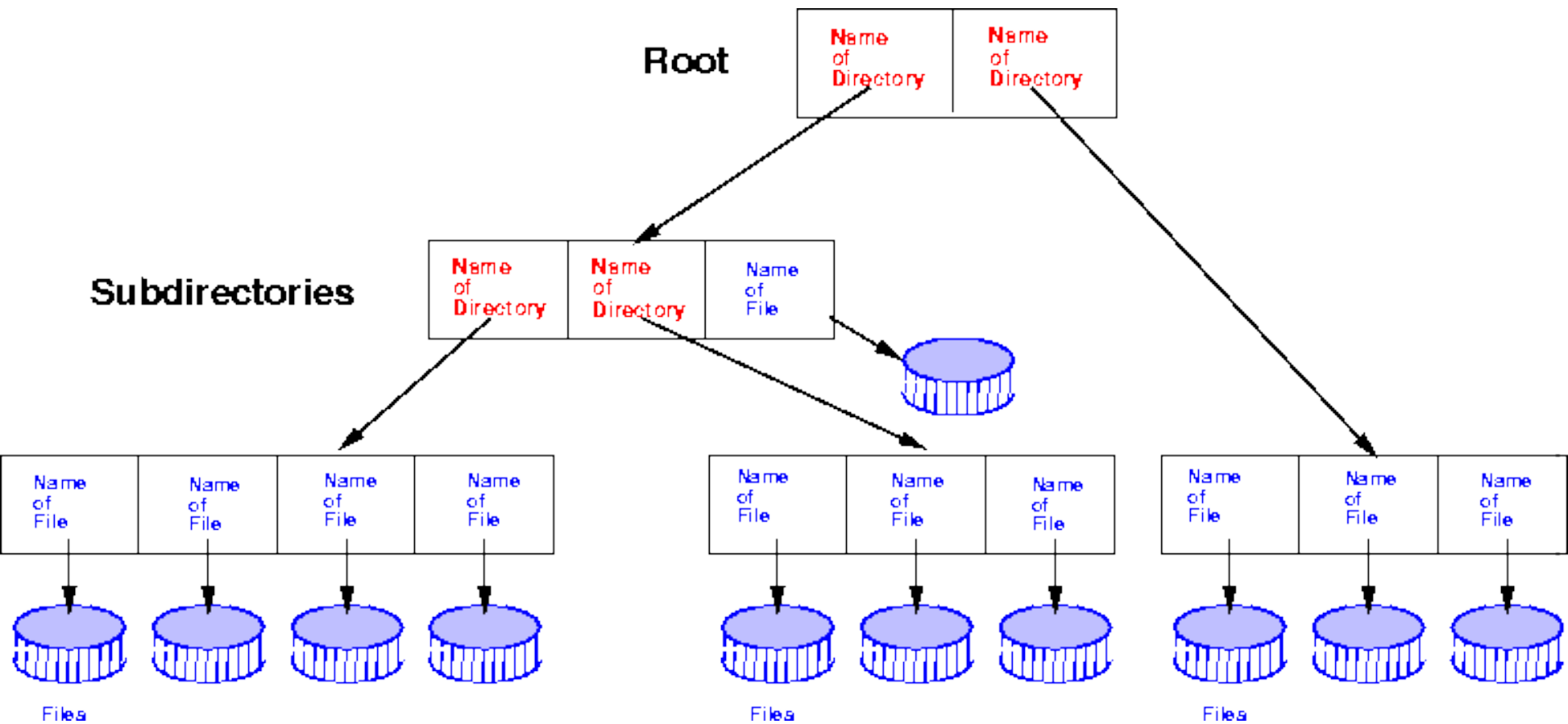
- A collection of nodes containing information about all files



- Both the directory structure and the files reside on disk
- Backups of these two structures could be kept on tapes



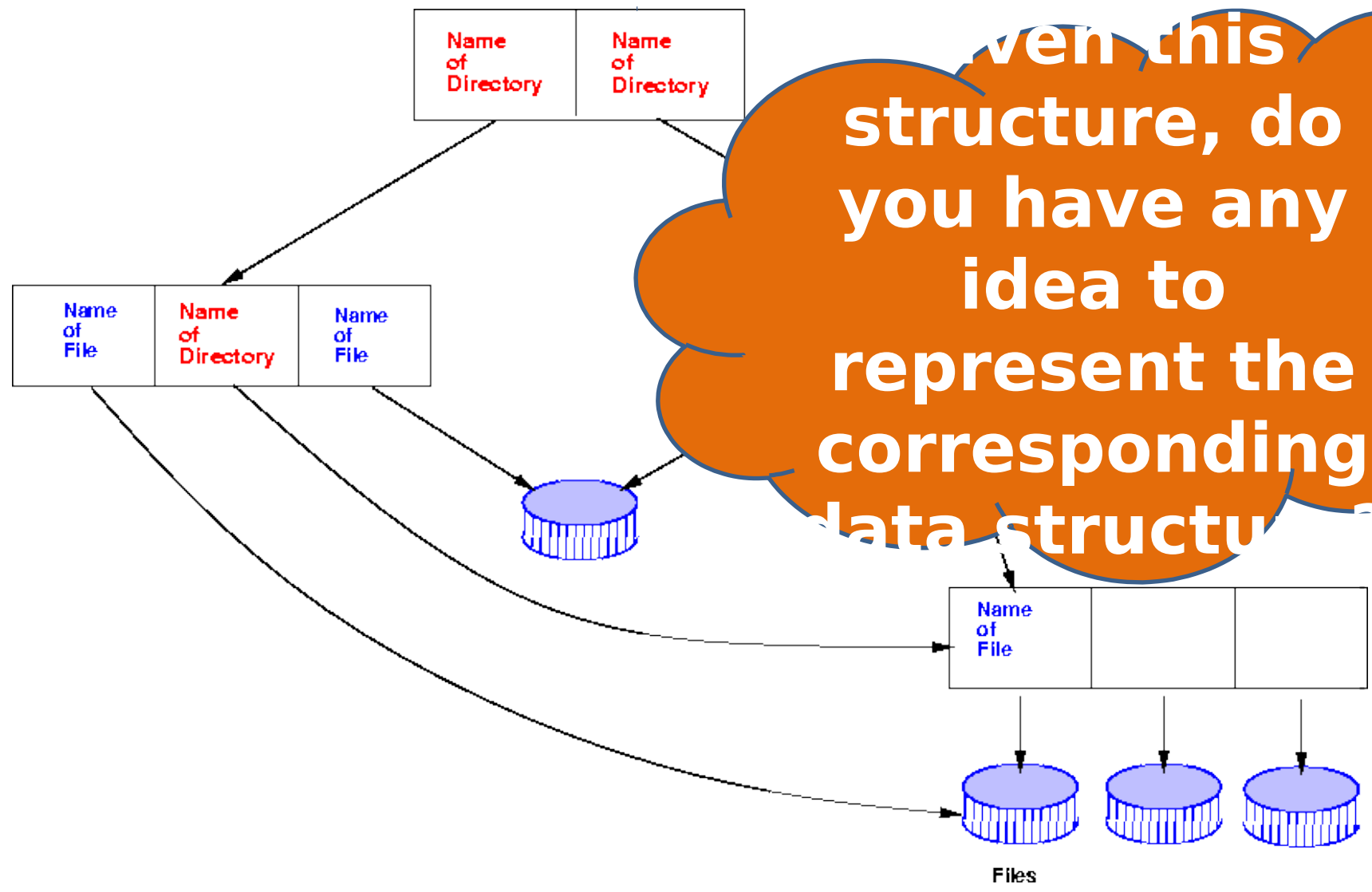
Directory Structure:  
Tree structured Directories



# Cont'

- Absolute or relative path name
  - Absolute from root
  - Relative paths from current working directory pointer.
- Creating a new file is done in current directory
- Creating a new subdirectory is done in current directory, e.g. `mkdir <dir-name>`
- Delete a file , e.g. `rm file-name`
- Deletion of directory
  - Option 1 : Only delete if directory is empty
  - Option 2: delete all files and subdirectories under directory

Directory Structure:  
Acyclic Graph Directories [无环图目录]



even this structure, do you have any idea to represent the corresponding data structure?

# Cont'

- Acyclic graphs allow sharing
- Implementation by links
  - Links are pointers to other files or subdirectories
  - Symbolic links or relative path name
    - Directory entry is marked as a link and name of real file/directory is given. Need to resolve link to locate file.
- Implementation by shared files
  - Duplicate information in sharing directories
  - Original and copy indistinguishable.
  - Need to maintain consistency if one of them is modified.

**In summary**

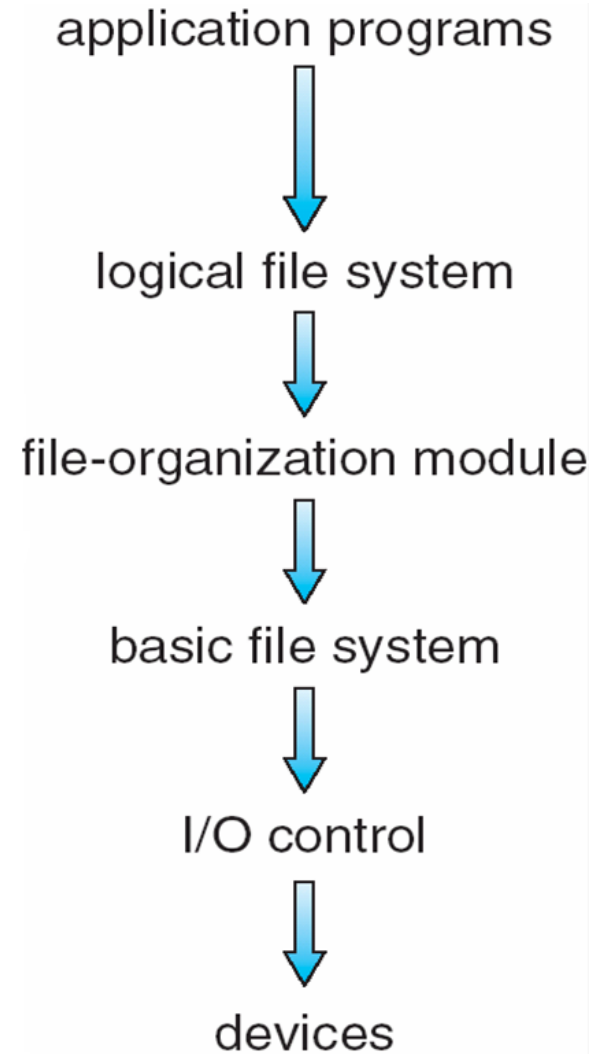


# File system [ 文件系统 ]

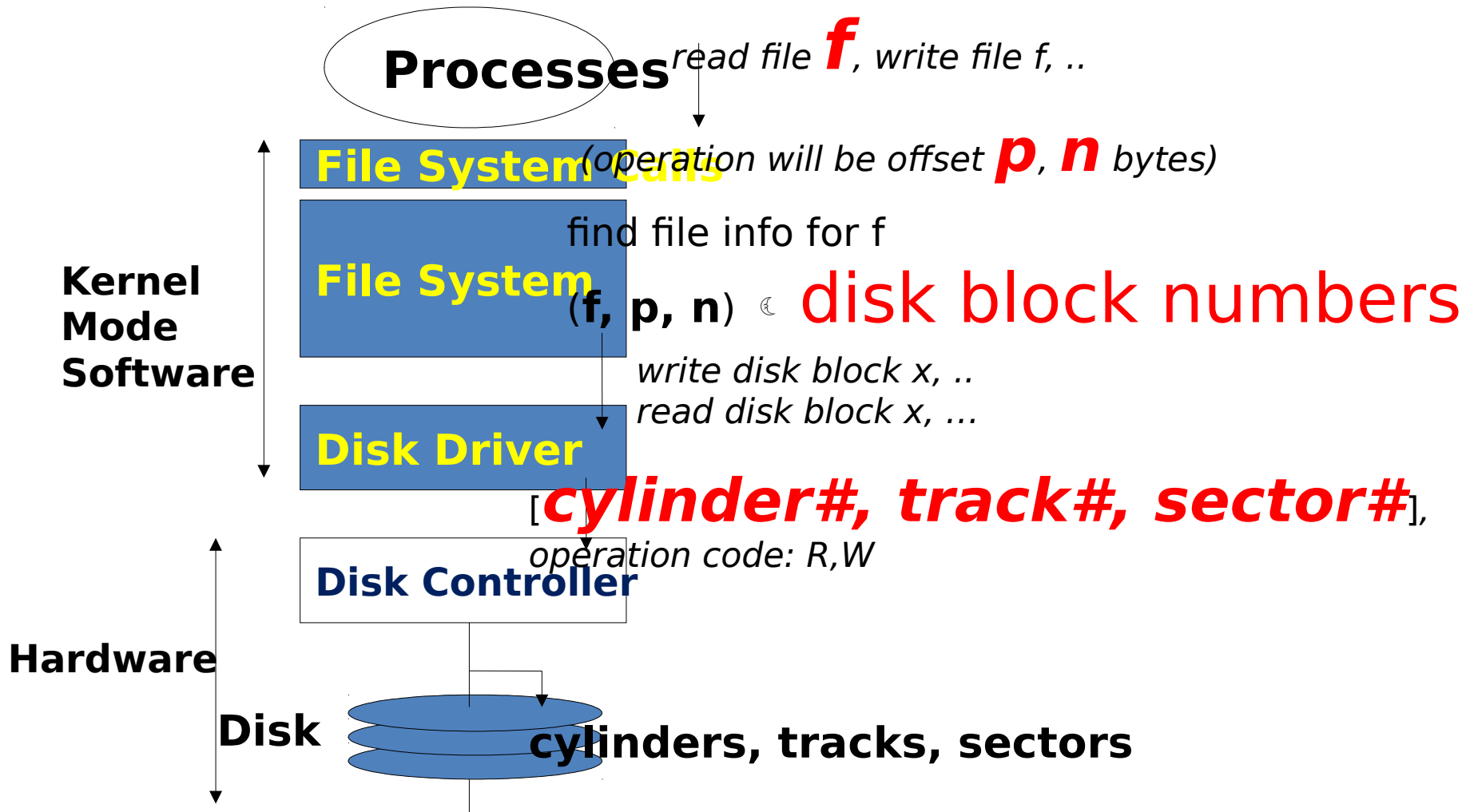
- The collection of algorithms and data structures which perform the translation from **logical file operations** (system calls) to **actual physical storage of information**
  - Provide storage of data and manipulation
  - Guarantee consistency of data and minimize errors
  - Optimize performance (system and user)
  - Eliminate data loss (data destruction)
  - Support variety of I/O devices
  - Provide a standard user interface
  - Support multiple users

# File System

- The file system should provide an **efficient** implementation of the interface
  - storing, locating, retrieving data
- The problem: define **data structures** and **algorithms** to map the **logical File System** onto the disk
  - some data structures live on disk
  - some data structures live (temporarily) in memory
- Typical layer organization:
  - Good for modularity and code reuse
  - Bad for overhead



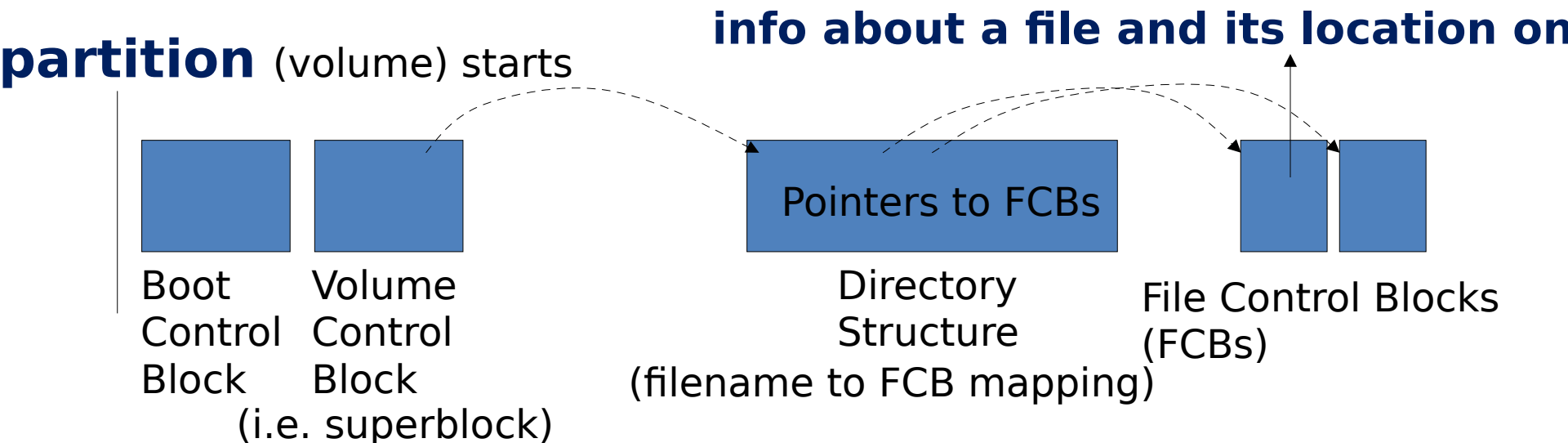
# Layered Software

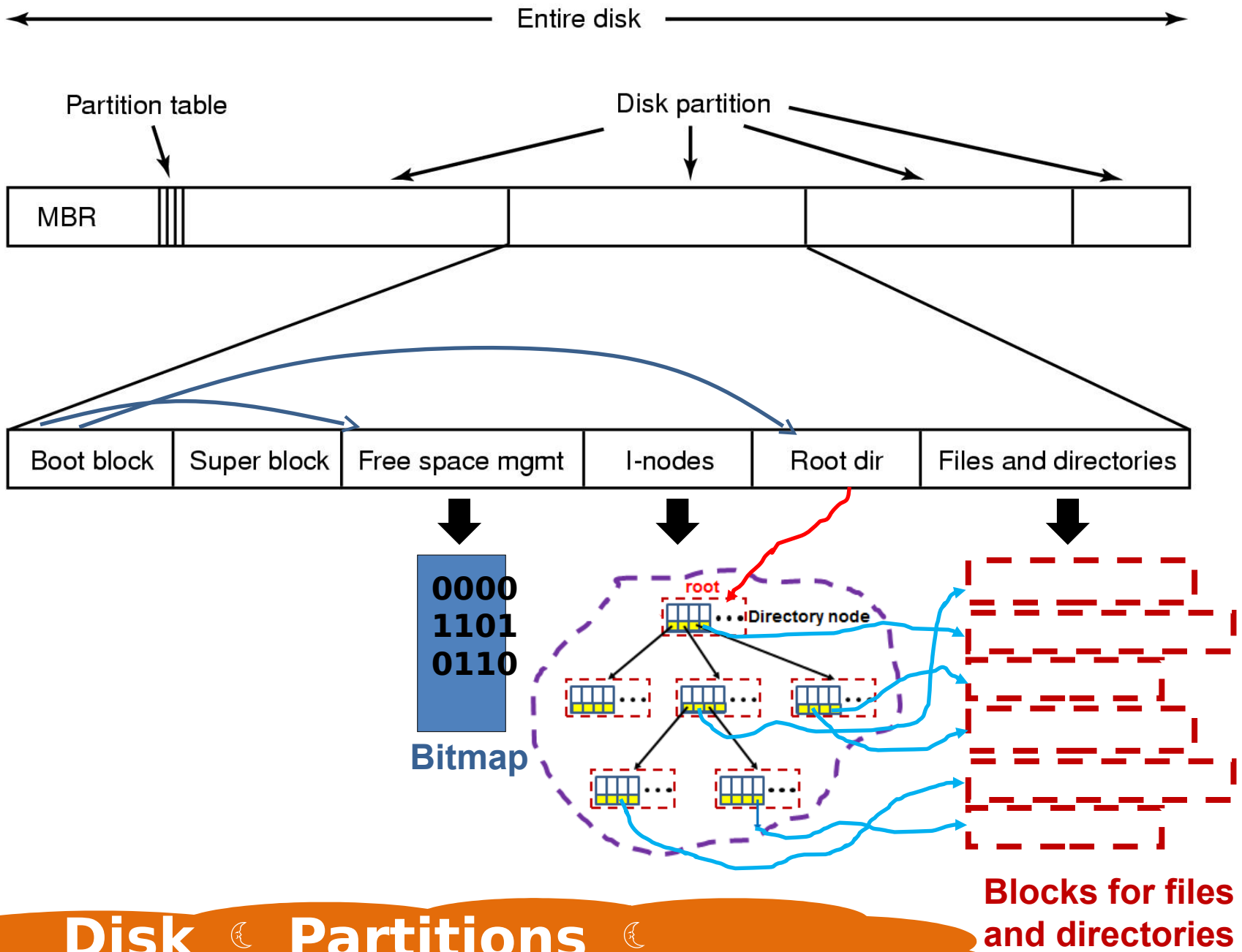




# File System Implementation

- **Major On-disk Structures (information):**
  - **Boot control block** contains info needed by system to boot OS from that volume
  - **Volume control block** contains volume details
  - Directory structure organizes the files
  - Per-file **File Control Block (FCB)** contains many details about the file





**Disk** ☾ **Partitions** ☾  
**Blocks** ☾ **FS**

# A Directory

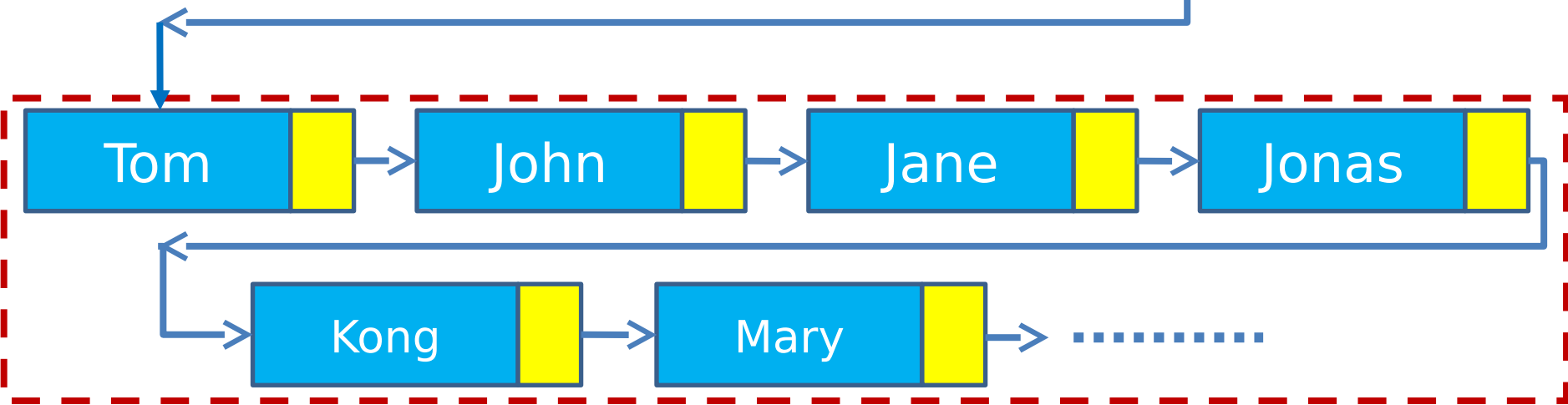
File space

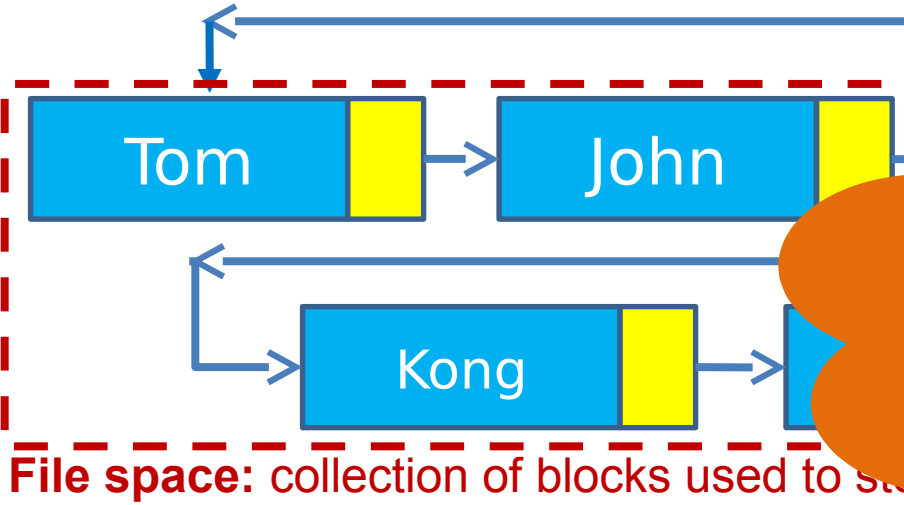
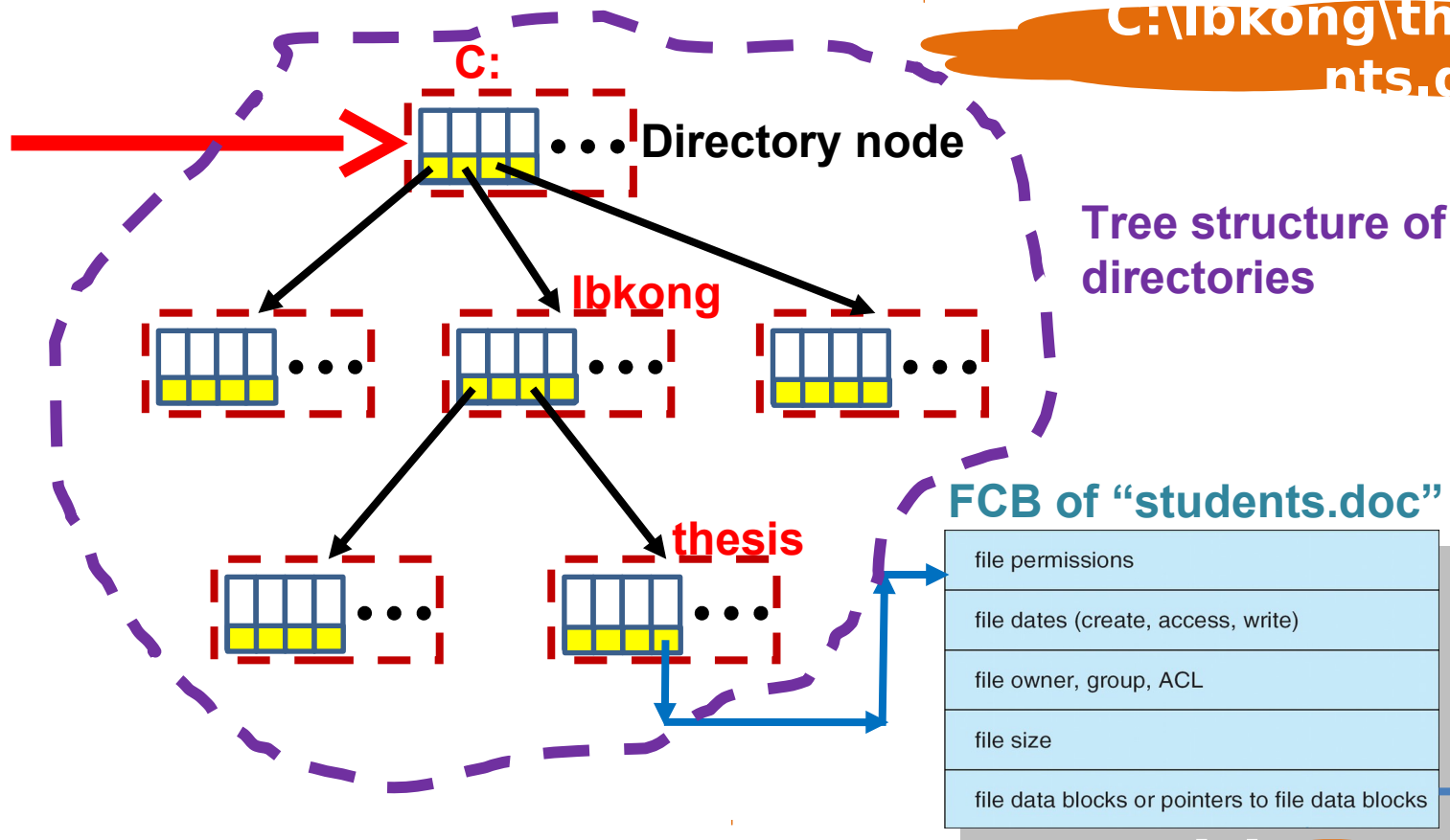
Header  
(pointer)  
column

game	
mail	
news	
work	

file permissions
file dates (create, access, write)
file owner, group, ACL
file size
file data blocks or pointers to file data blocks

**File  
Control  
Block of  
a file  
with  
filenam  
e X**





With your DSA experience, you should consider the related operations on this data structure. Could you imagine how to locate "students.doc"? And other file operation??

# Compared with MM

To run our program, the needed pages should be fetched into MM  
(PCB maintains the reference for our program file)

(OS maintains the **mapping info** between program pages and blocks)



Addressed **frame/page** space  
(usually 4KB)



Addressed **Storage unit** space  
(1 byte)

Our program exists first as file  
which may contain many **pages**

(filename ↪ FCB ↪ Efficient data structure (like Hash, i-node,...) ↪ **Blocks**)



Addressed **block** space  
(usually 4KB)



Addressed **Sector** space  
(512 bytes)



# Suggest for deeper study

- You can try to implement a simulation program to support file + directory
  - Data structure + Operators
    - New file, Delete, Rename
- Continue to consider concurrent access of a file by many users.
  - What should we consider ? – Data structure + Operators
- Have you ever considered how to store an array of structs into a file? – Don't forget to reconstruct their structure when reading into Memory

**helpful for your understanding about the implementation of DBMS**