# File Systems

Kartik Gopalan

Chapter 4

From Tanenbaum's Modern Operating System

# What is a File System?

- File system is the OS component that organizes data on the raw storage device.

- Data, by itself, is just a meaningless sequence of bits and bytes.

- Metadata is the information that describes the data.
  - Also called attributes.
  - Without meta-data, data itself will be incomprehensible.

- A File System defines
  - Format of the data objects.
  - The format/meaning of meta-data associated with each data object. E.g. File name, permissions, and size of a file.
  - Location of the individual data blocks of for each data object.
  - A framework to manage free space on the raw storage.

# Metadata — Examples

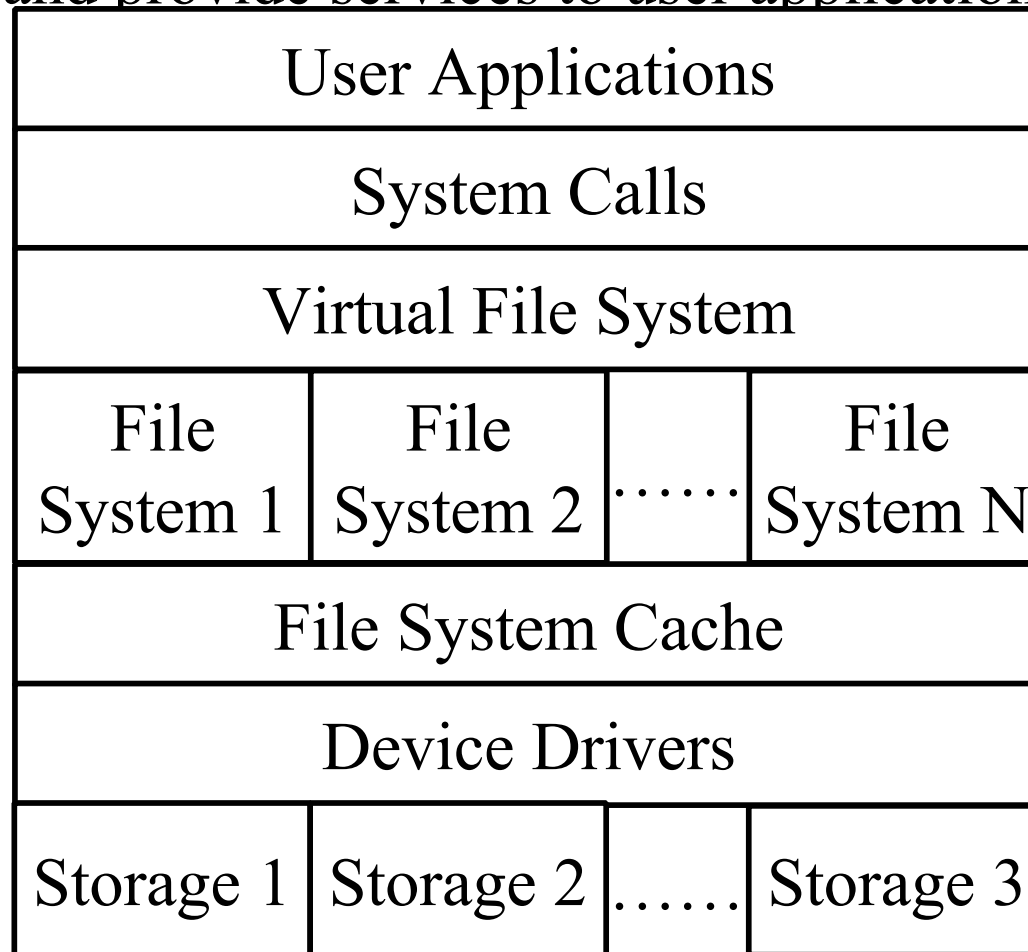| Attribute | Meaning |
|---|---|
| Protection | Who can access the file and in what way |
| Password | Password needed to access the file |
| Creator | ID of the person who created the file |
| Owner | Current owner |
| Read-only flag | 0 for read/write; 1 for read only |
| Hidden flag | 0 for normal; 1 for do not display in listings |
| System flag | 0 for normal files; 1 for system file |
| Archive flag | 0 for has been backed up; 1 for needs to be backed up |
| ASCII/binary flag | 0 for ASCII file; 1 for binary file |
| Random access flag | 0 for sequential access only; 1 for random access |
| Temporary flag | 0 for normal; 1 for delete file on process exit |
| Lock flags | 0 for unlocked; nonzero for locked |
| Record length | Number of bytes in a record |
| Key position | Offset of the key within each record |
| Key length | Number of bytes in the key field |
| Creation time | Date and time the file was created |
| Time of last access | Date and time the file was last accessed |
| Time of last change | Date and time the file has last changed |
| Current size | Number of bytes in the file |
| Maximum size | Number of bytes the file may grow to |

# File Types

- Regular files
  - Contains actual data

- Directories
  - Files that help locate other files

- Character Special Files
  - Typically for byte-oriented I/O operations

- Block special files
  - For block-oriented I/O operations
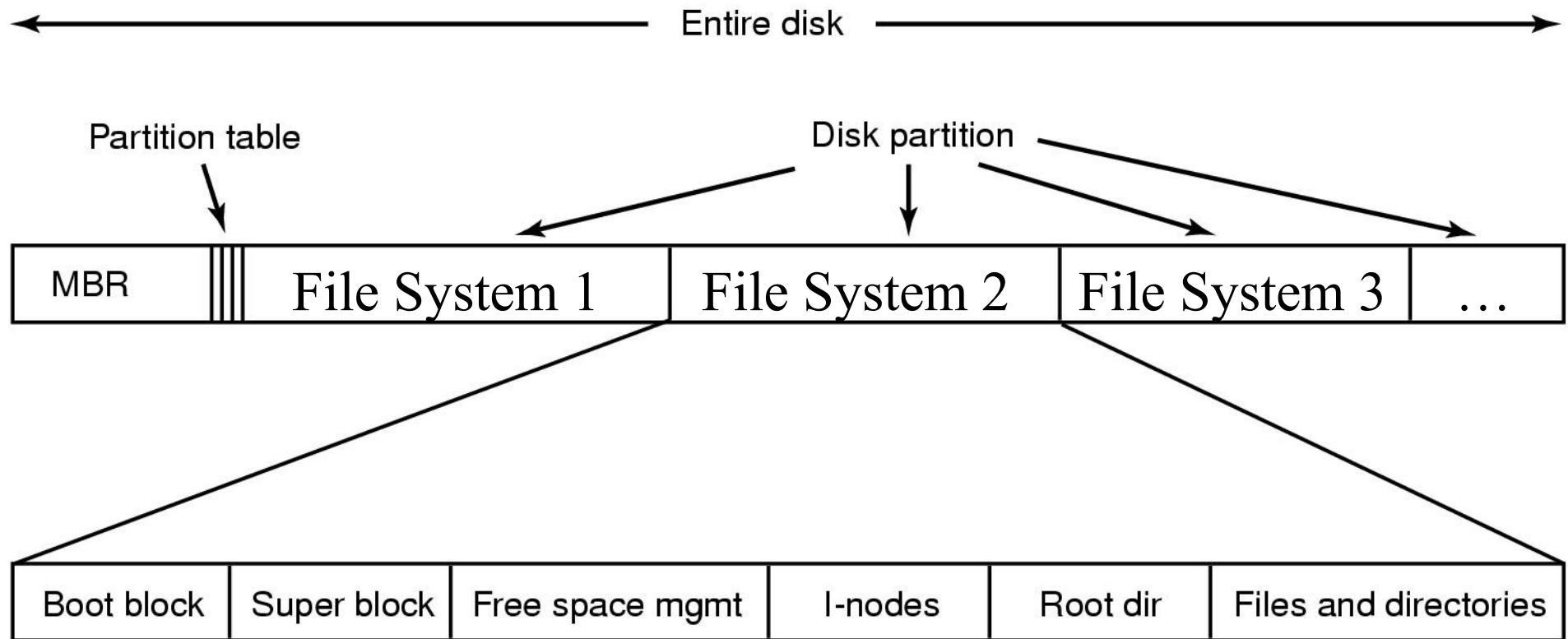
# Basic File System Operations

- Create a file
- Open an existing file
- Write to a file
- Read from a file
- Seek to somewhere in a file
- Close an open file
- Delete a file
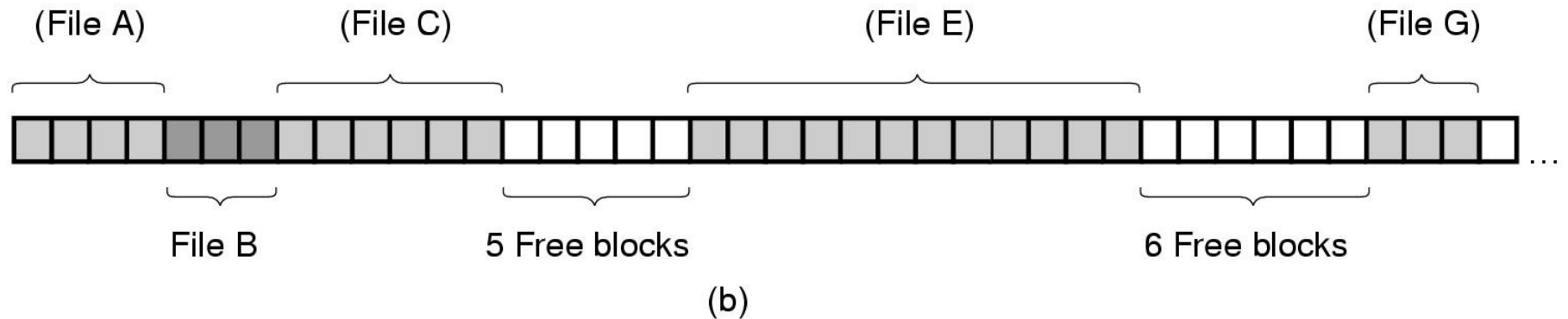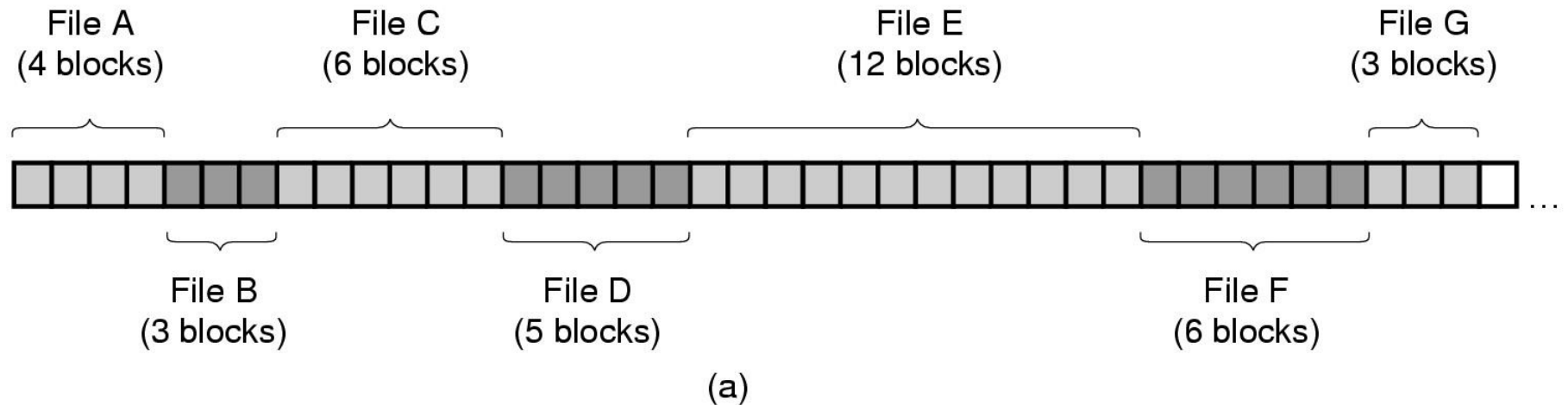
# Virtual File System (VFS)

- VFS provides
  1. A common system call interface to user applications to access different file systems implemented in the OS.
  2. A common interface to file systems to "plug into" the operating system and provide services to user applications.

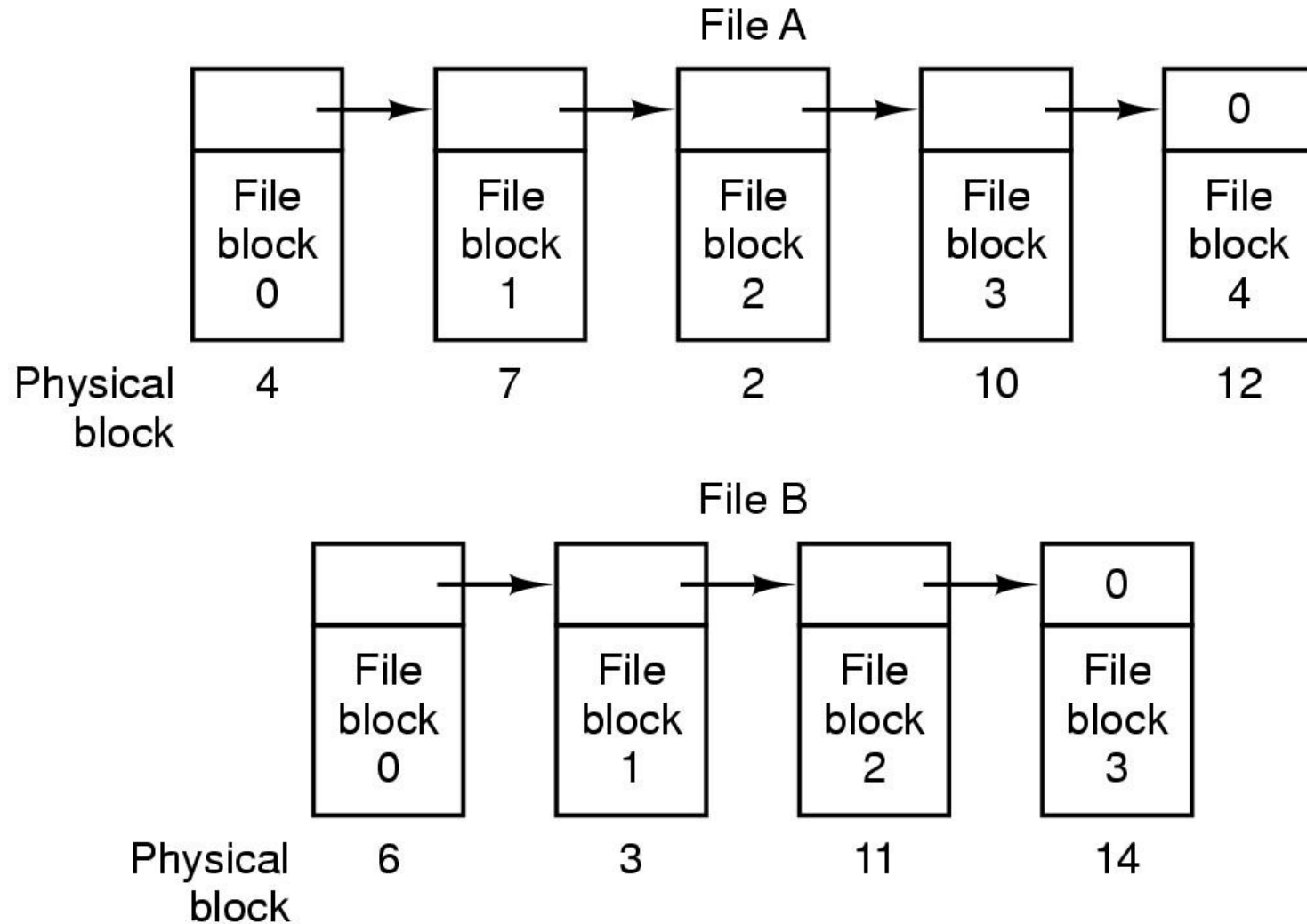| User Applications | | | |
|:---:|:---:|:---:|:---:|
| System Calls | | | |
| Virtual File System | | | |
| File System 1 | File System 2 | ······ | File System N |
| File System Cache | | | |
| Device Drivers | | | |
| Storage 1 | Storage 2 | ······ | Storage 3 |

# Partitions and File-system Layout

# Organizing Files on Disk: Contiguous Allocation

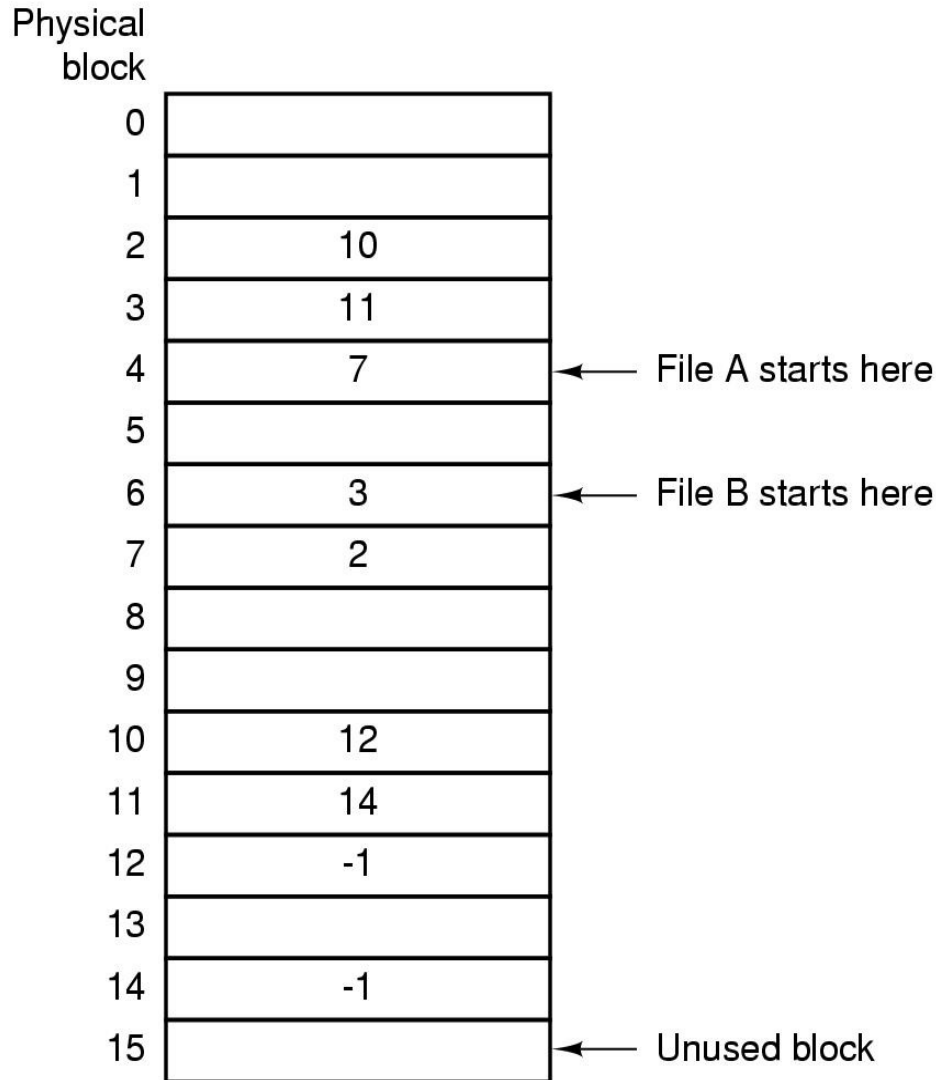

(a)

(a) Contiguous allocation of disk space for 7 files

(b) State of the disk after files D and E have been removed

# Organizing Files on Disk: Singly Linked List of Blocks

File A

| | | | | |
|---|---|---|---|---|
| | | | | 0 |
| File block 0 | File block 1 | File block 2 | File block 3 | File block 4 |

Physical block    4        7        2        10        12

File B

| | | | |
|---|---|---|---|
| | | | 0 |
| File block 0 | File block 1 | File block 2 | File block 3 |

Physical block    6        3        11        14

- Advantage: Logically contiguous blocks can be discontiguous on disk
- Disadvantage: Random seeks are expensive. Requires traversal from the start.
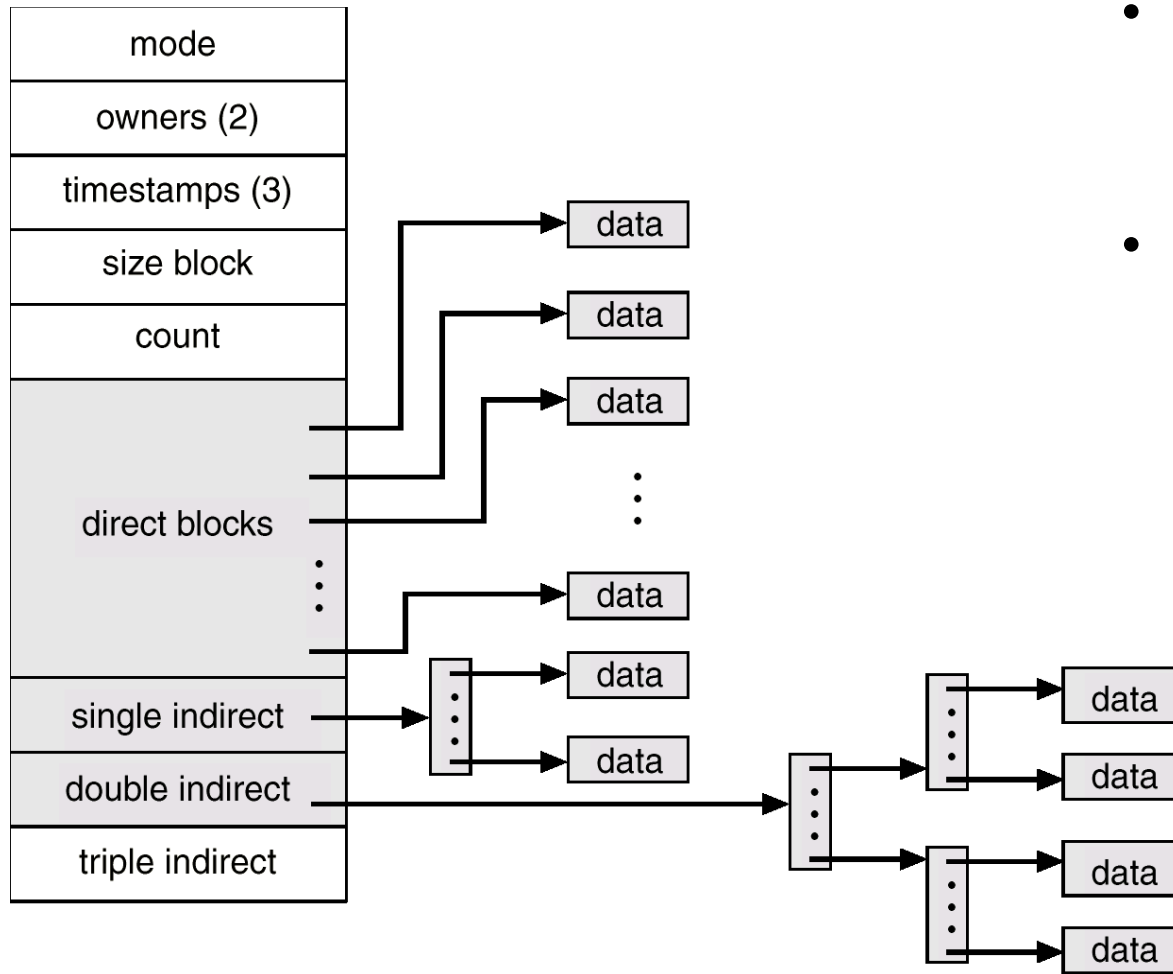
# Organizing Files on Disk: File Allocation Table

Physical
block

| Block | Value | |
|-------|-------|---|
| 0 | | |
| 1 | | |
| 2 | 10 | |
| 3 | 11 | |
| 4 | 7 | ← File A starts here |
| 5 | | |
| 6 | 3 | ← File B starts here |
| 7 | 2 | |
| 8 | | |
| 9 | | |
| 10 | 12 | |
| 11 | 14 | |
| 12 | -1 | |
| 13 | | |
| 14 | -1 | |
| 15 | | ← Unused block |

- Linked list allocation using a file allocation table in RAM

- Doesn't need a separate "next" pointer within each block

- But random seeks are still expensive
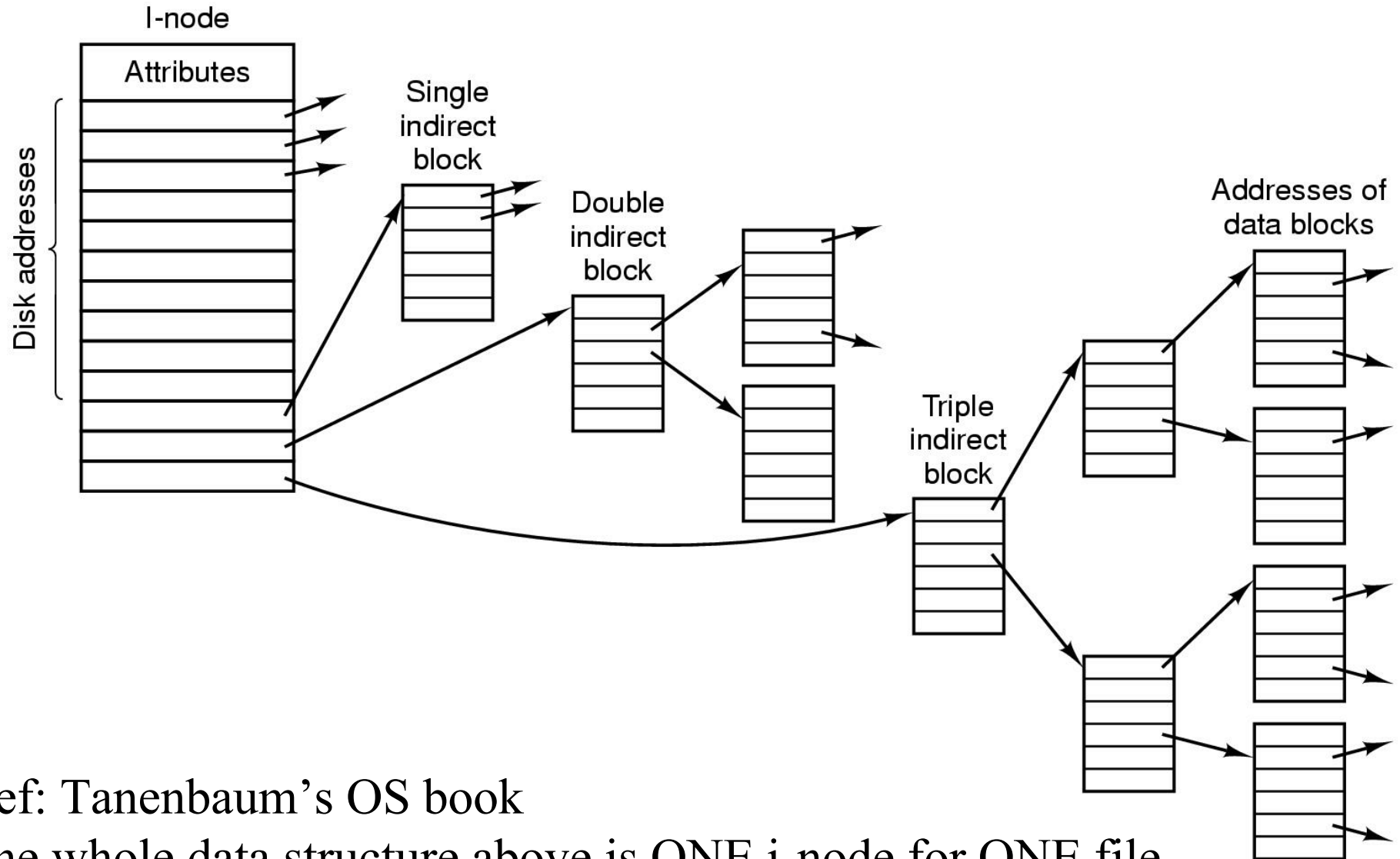
# i-nodes (index nodes)

- Each file is described by an i-node
- i-node stores the metadata for the file
- Metadata = file info + location of data blocks
- i-nodes are stored on the disk

- An i-node is to a file what a page-table is to a virtual address space

  - Page table maps
    - virtual page number in a virtual address space ——> physical page frame number in DRAM
  - i-node maps
    - logical block number in a file ——> physical block location on disk
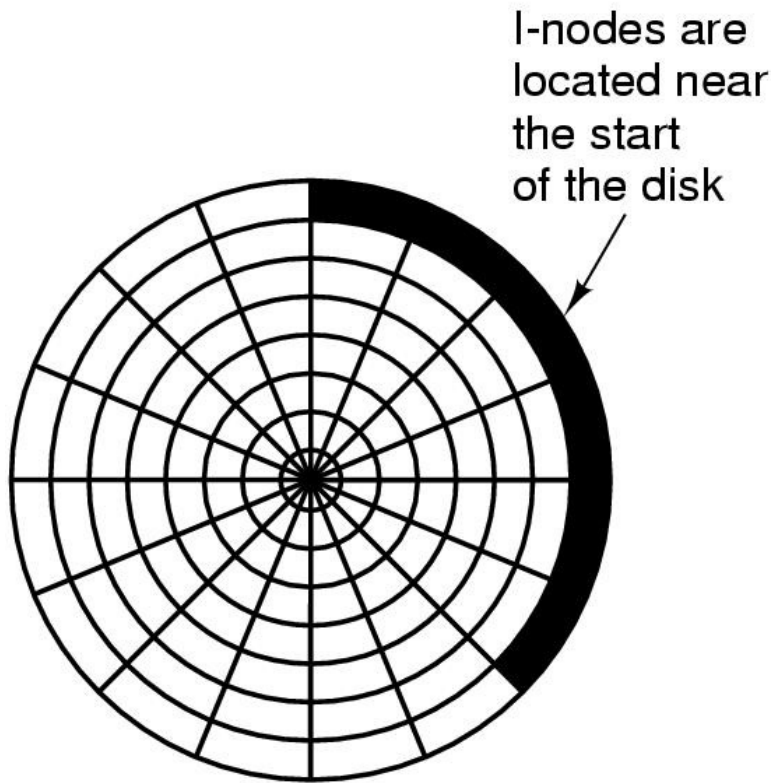
# Unix i-node (index node)



- Small files can be accessed quickly

- If each block is 4KB
  - First 48KB of file reachable from 12 direct blocks
  - Next 4MB available from single-indirect blocks
  - Next 4GB available from double-indirect blocks
  - Next 4TB available through the triple-indirect blocks

- Any block can be found with at most 3 disk accesses
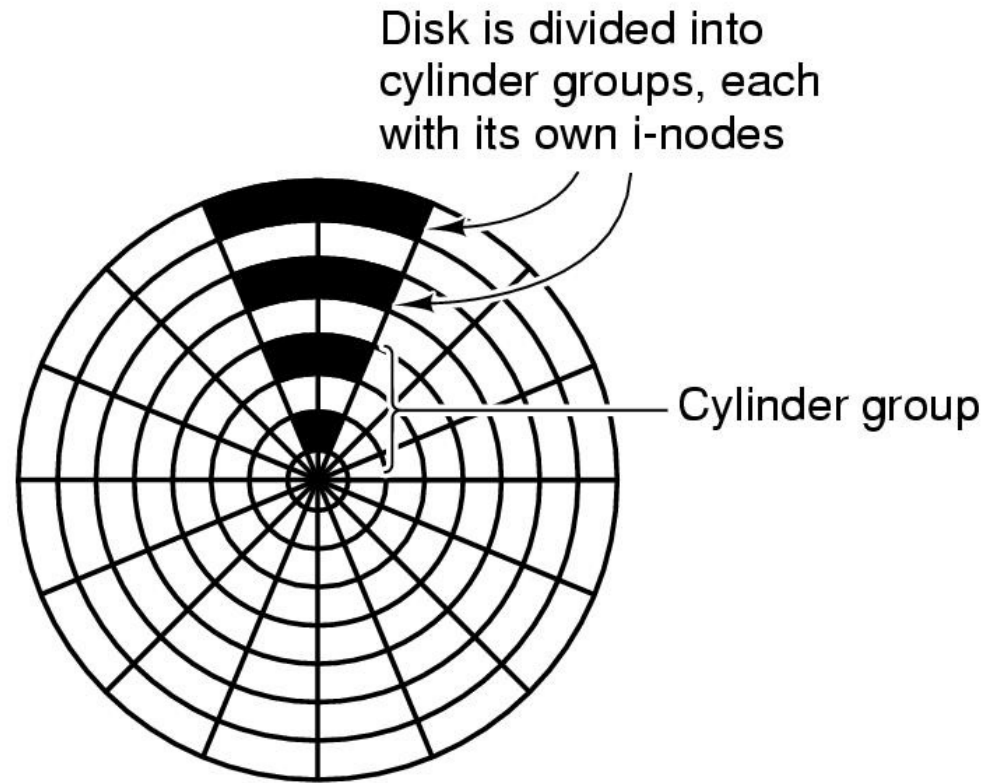
# Another view of a UNIX i-node



- Ref: Tanenbaum's OS book
- The whole data structure above is ONE i-node for ONE file
- i-node grows dynamically as the file grows
- Just like page-tables, i-node tracks a giant array broken up into many pieces
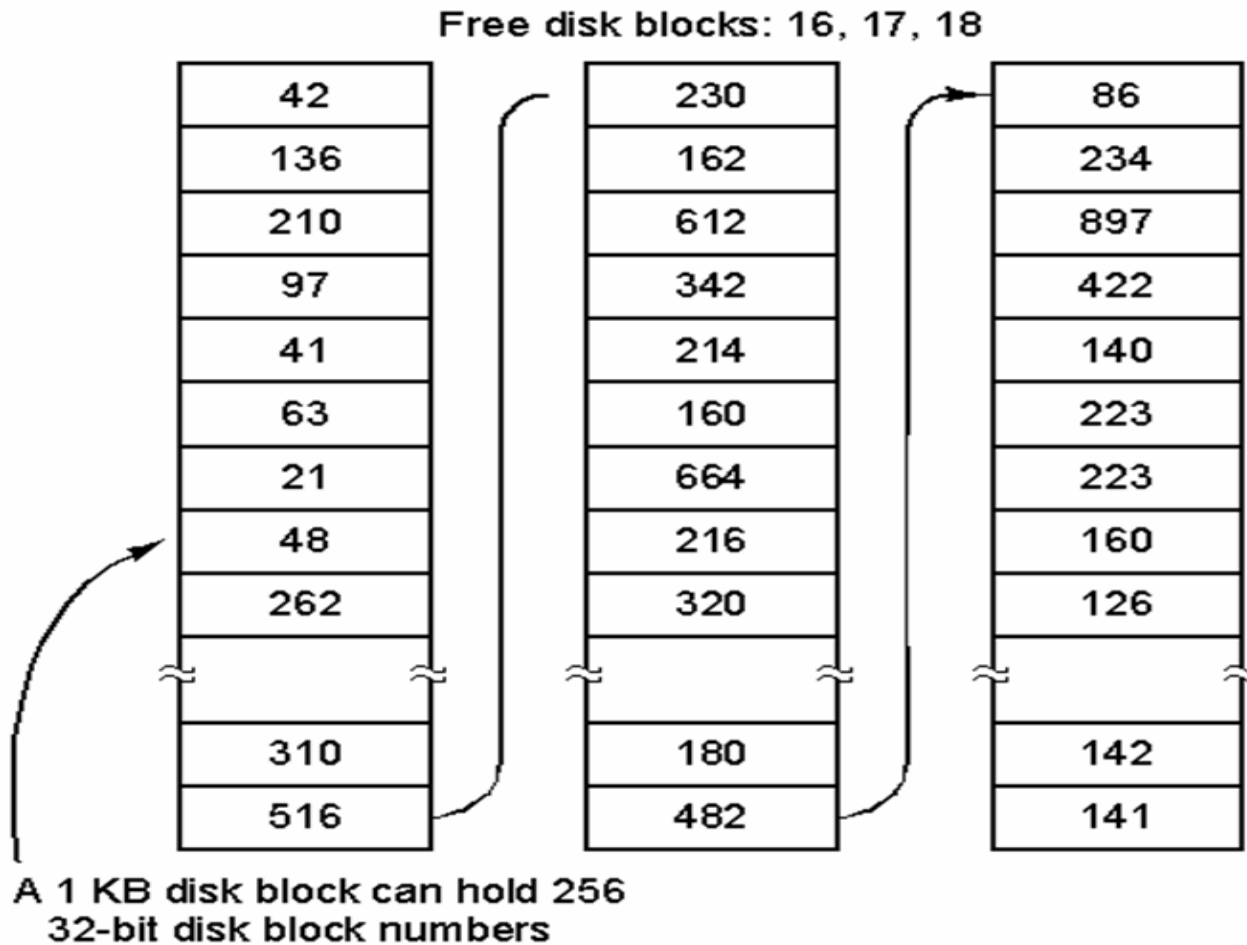
# Performance impact of i-node placement



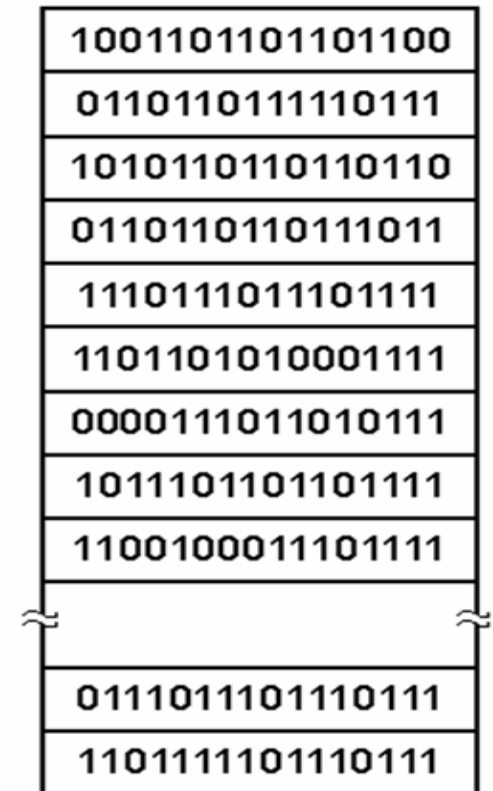- I-nodes placed at the start of the disk
- Longer average seek times

- Each track on disk has its own blocks and i-nodes
- Shorter average seek times

# Managing Free Space on Disk

Free disk blocks: 16, 17, 18

| | | |
|---|---|---|
| 42 | 230 | 86 |
| 136 | 162 | 234 |
| 210 | 612 | 897 |
| 97 | 342 | 422 |
| 41 | 214 | 140 |
| 63 | 160 | 223 |
| 21 | 664 | 223 |
| 48 | 216 | 160 |
| 262 | 320 | 126 |
| 310 | 180 | 142 |
| 516 | 482 | 141 |

A 1 KB disk block can hold 256
32-bit disk block numbers

(a) Linked list

```
1001101101101100
0110110111110111
1010110110110110
0110110110111011
1110111011101111
1101101010001111
0000111011010111
1011101101101111
1100100011101111
0111011101110111
1101111101110111
```
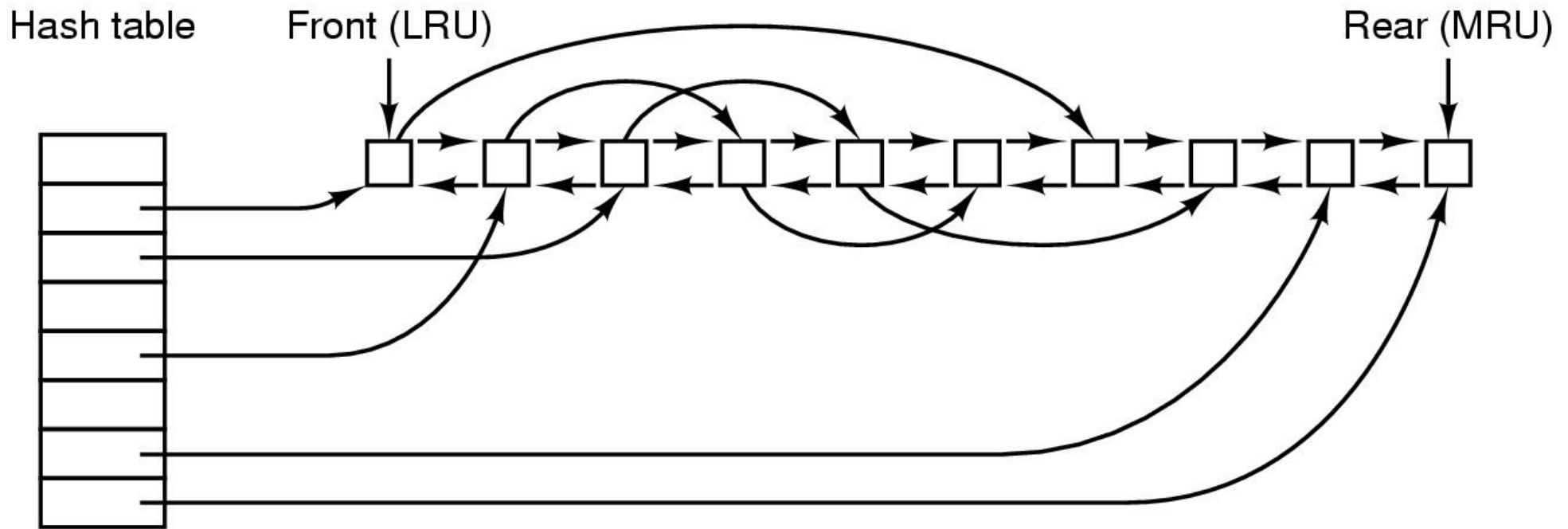
A bit map

(b)

# File System Cache

- Small area in main memory that stores frequently accessed data blocks in the file system.

- Before accessing the disk, look in the FS cache.

- If data block is in FS cache, no need to go to disk.

- Periodically, purge the cache of infrequently used data blocks.

- Claim: If the cache works well, then most I/O accesses to the physical disk will be writes. Why?

# Data Structure for File-System Cache



- Cache Lookup:
  - Pages in cache are looked up via a hash table for fast access.
- Cache Eviction:
  - Another doubly-linked list maintained to identify least-recently used (LRU) victim pages that are periodically purged.
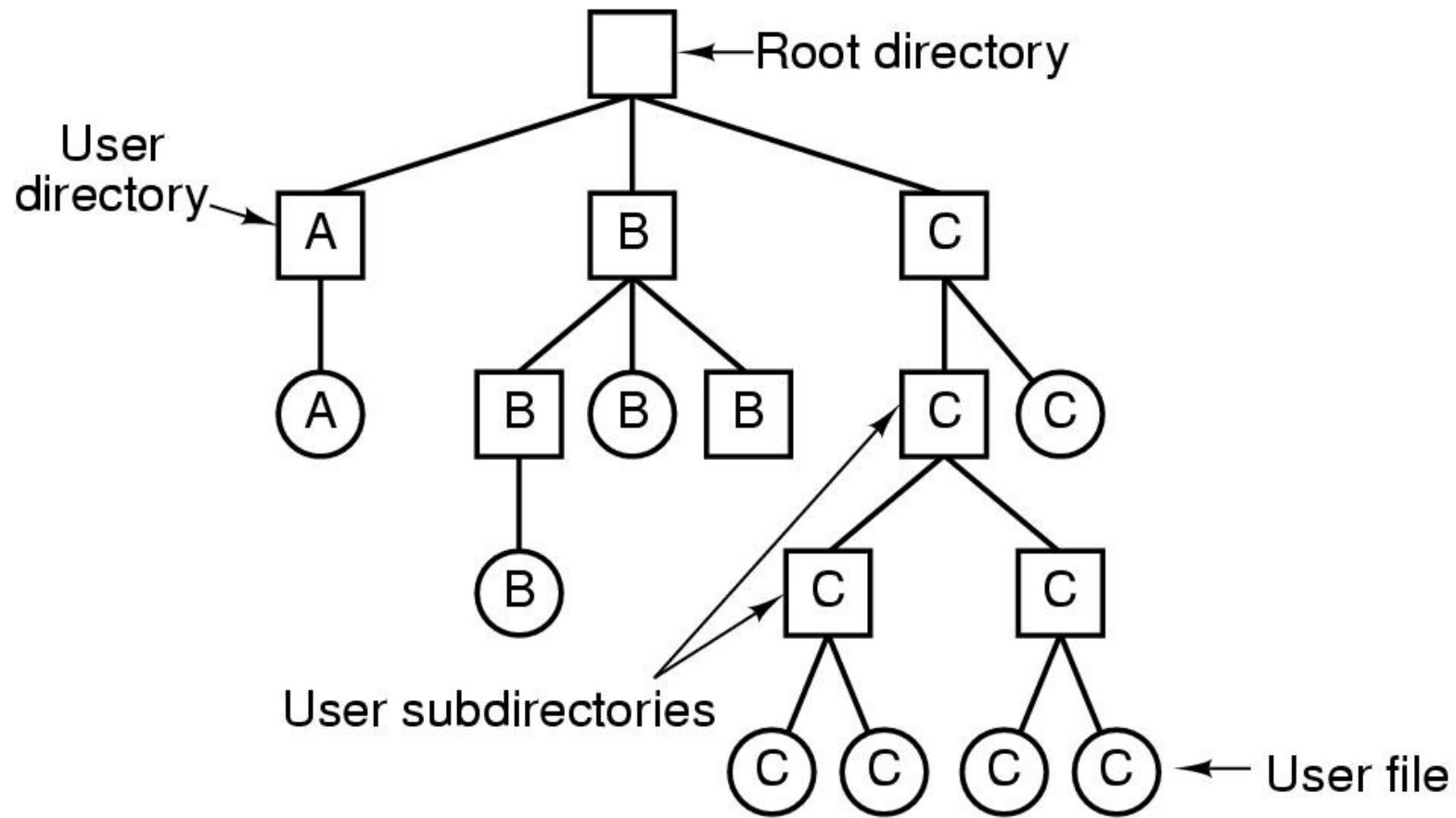  - Is the victim page dirty? Then write to disk. Else discard.

# Virtual memory page cache and FS cache

- Often they are managed in a unified manner

- Meaning: one common page-cache is used for managing pages for both virtual memory and file system.

- For example, Linux maintains one common set of data structures to keep track of active and inactive (LRU) pages.
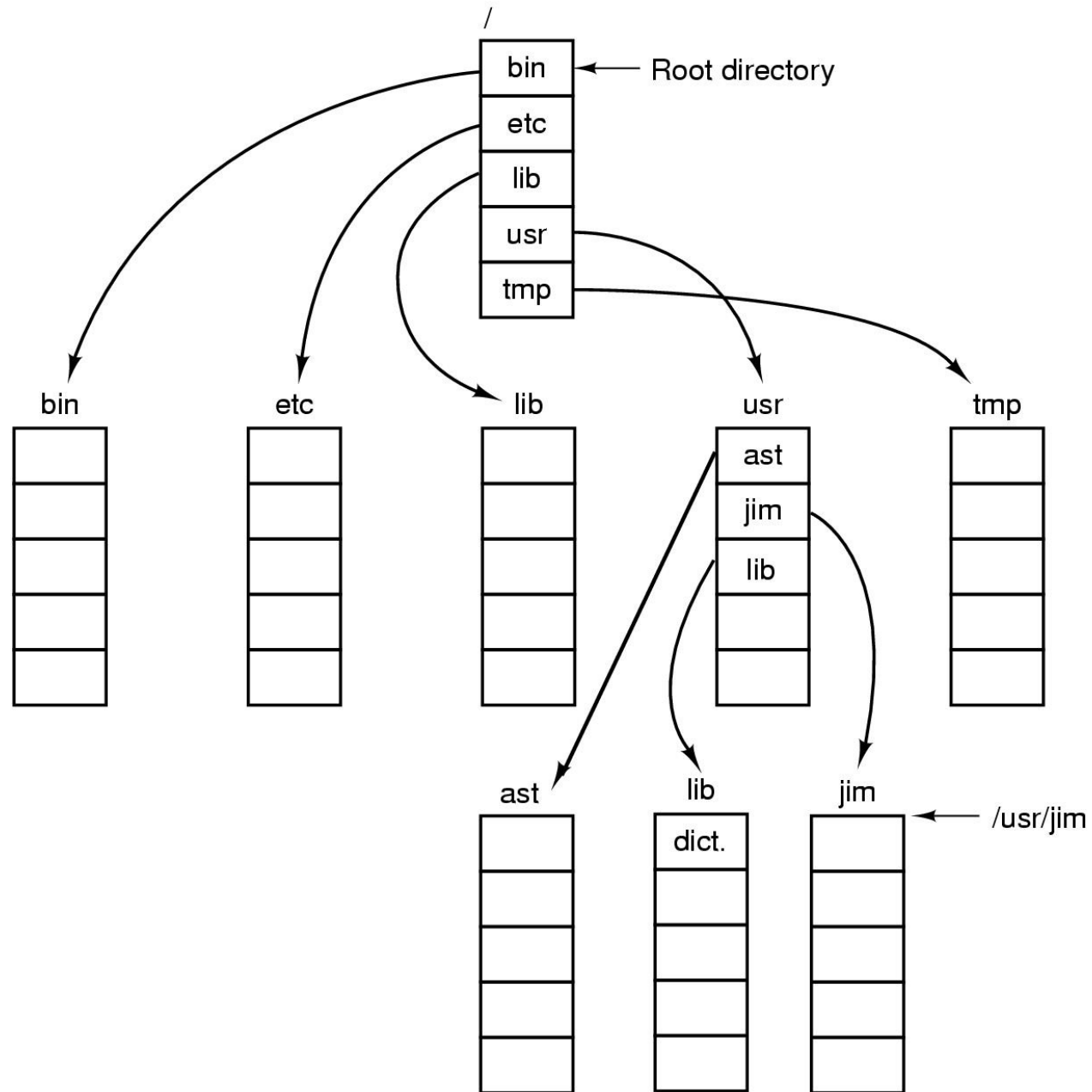
# Log-Structured File Systems

- With CPUs faster, memory larger
  - disk caches are also getting larger
  - increasing number of read requests come from file system cache
  - Thus, most disk accesses will be writes

- LFS treats the entire disk as a log
  - all writes are initially buffered in memory
  - periodically commit the writes to the end of the disk log
  - When file is opened, locate i-node, then find blocks
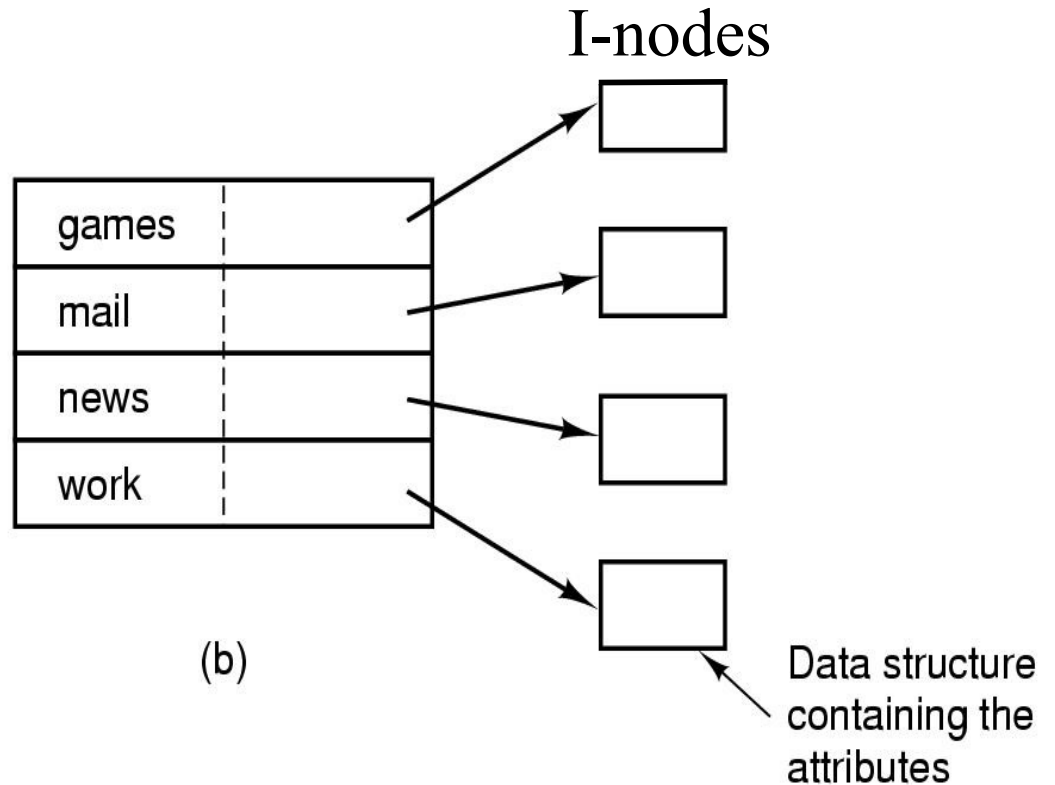
# Hierarchical Directory Systems

# Path Names

# Implementing Directories


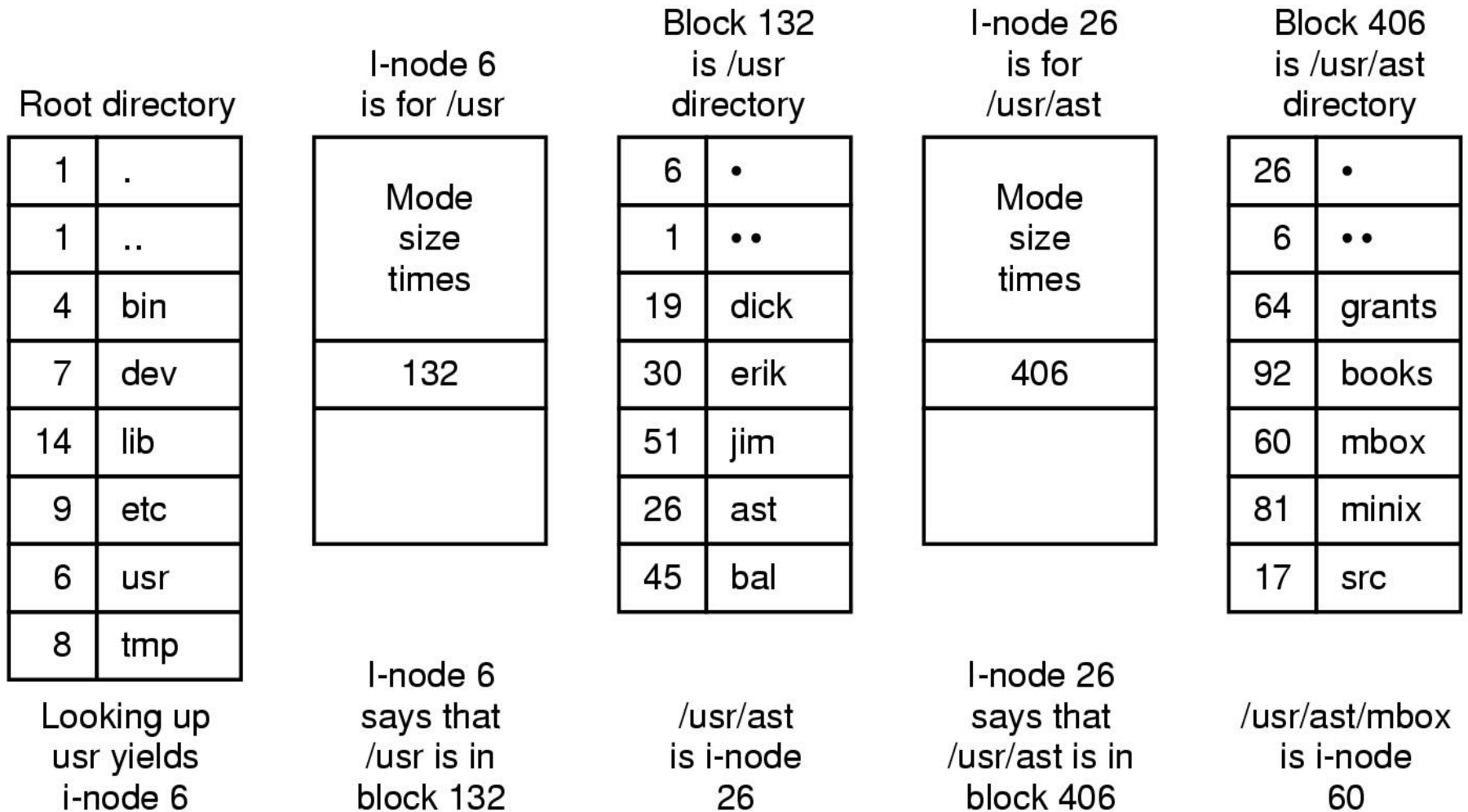
(a)

(b)

I-nodes

Data structure containing the attributes

(a) A simple directory with fixed size entries. Disk addresses and attributes are stored in directory entry

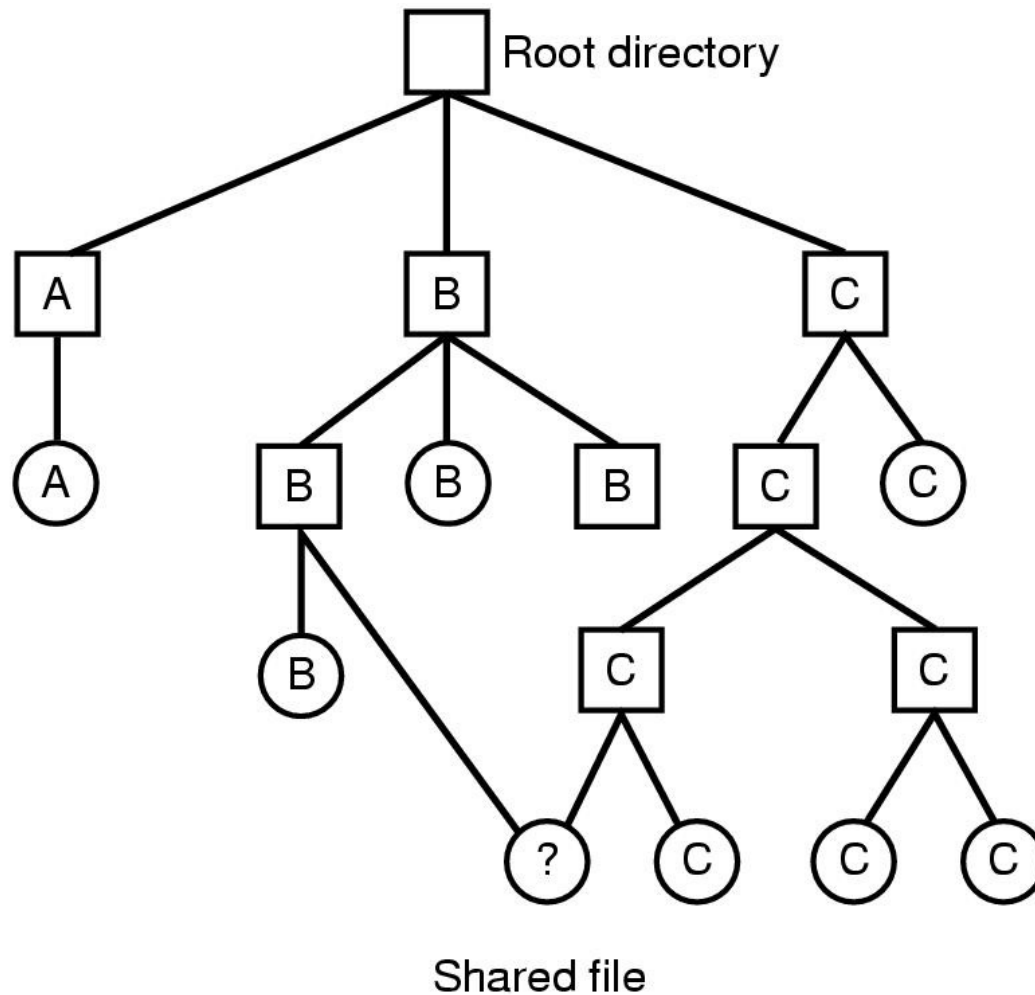(b) Directory in which each entry just refers to an i-node

# Path lookup in a typical Unix File System

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

| Root directory | | I-node 6<br>is for /usr | Block 132<br>is /usr<br>directory | | I-node 26<br>is for<br>/usr/ast | Block 406<br>is /usr/ast<br>directory | |
|---|---|---|---|---|---|---|---|
| 1 | . | | 6 | • | | 26 | • |
| 1 | .. | Mode<br>size<br>times | 1 | •• | Mode<br>size<br>times | 6 | •• |
| 4 | bin | | 19 | dick | | 64 | grants |
| 7 | dev | 132 | 30 | erik | 406 | 92 | books |
| 14 | lib | | 51 | jim | | 60 | mbox |
| 9 | etc | | 26 | ast | | 81 | minix |
| 6 | usr | | 45 | bal | | 17 | src |
| 8 | tmp | | | | | | |

Looking up
usr yields
i-node 6

I-node 6
says that
/usr is in
block 132

/usr/ast
is i-node
26

I-node 26
says that
/usr/ast is in
block 406

/usr/ast/mbox
is i-node
60

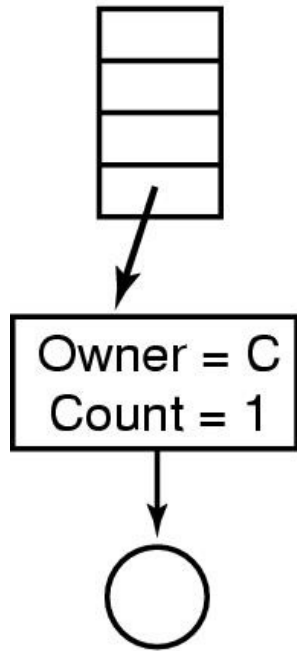Next, lookup i-node 60 to locate the data blocks of the file /usr/ast/mbox

# Shared Files — Hard Links



A file shared between two directories
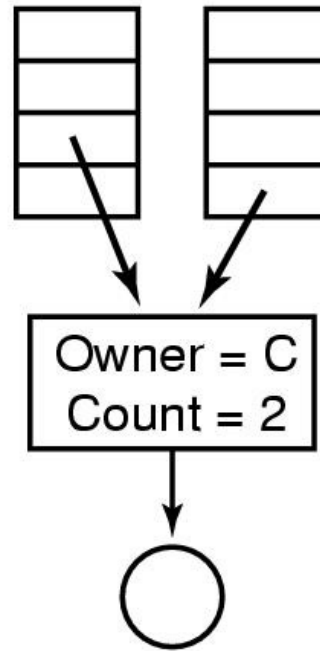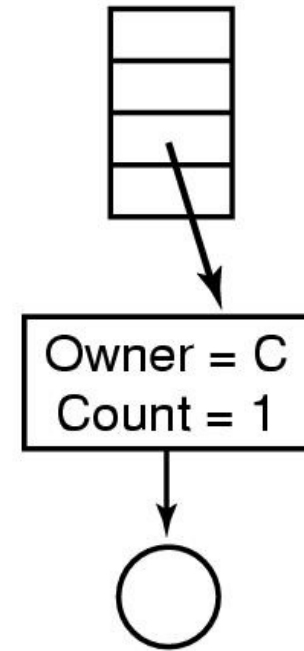
# Shared Files – Hard Links



C's directory

Owner = C
Count = 1

(a)

Situation prior to linking

B's directory   C's directory

Owner = C
Count = 2

(b)

After the hard link is created

B's directory

Owner = C
Count = 1

(c)

After the original owner removes the file

File still exists, but with count decremented

# Quiz on Inodes

- All blocks in a disk are of size 4KB (4096 bytes). A disk block can store either data or metadata (but not both).

- Each block address, i.e. a block's location on the disk, is 8-bytes in size

- Assume that all file attributes other than data block addresses, take up negligible space in the top-level block of inode block.

- Last three entries of the top-level block of inode contain single, double, and triple indirect block addresses.

- Rest of the space in the top-level inode (between end of attributes and single-indirect block address) is used to store direct block addresses.

- Question 1: What is the largest size of a file that can be accessed through
  - direct block entries?
  - direct + single indirect block entries?
  - direct + single + double indirect block entries?
  - direct + single + double + triple indirect block entries?

- Question 2: What is the size of an inode (in bytes) for a 32GB file?