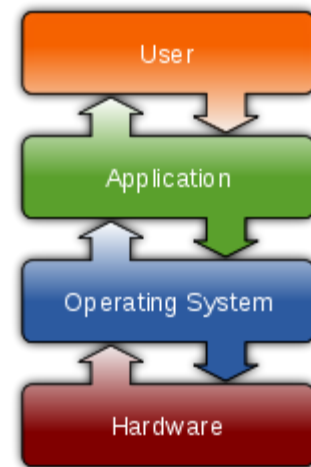# ENCE360
# Operating Systems

# File Systems

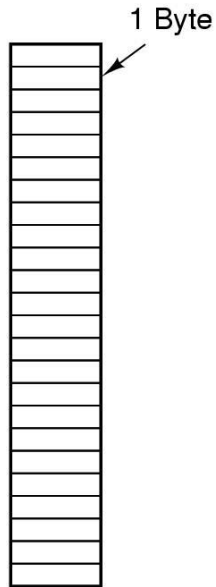- **Files**
- **Directories**
- **File system implementation**
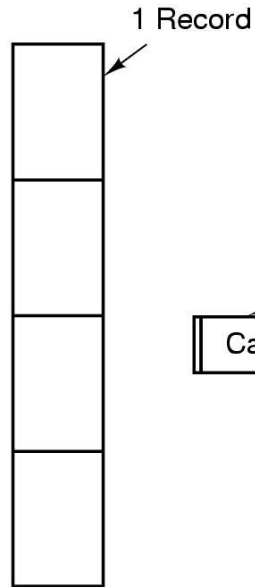- **Example file systems**

  *MOS Ch 6*

# Long-term Information Storage

1. Must store large amounts of data

2. Information stored must survive the termination of the process using it

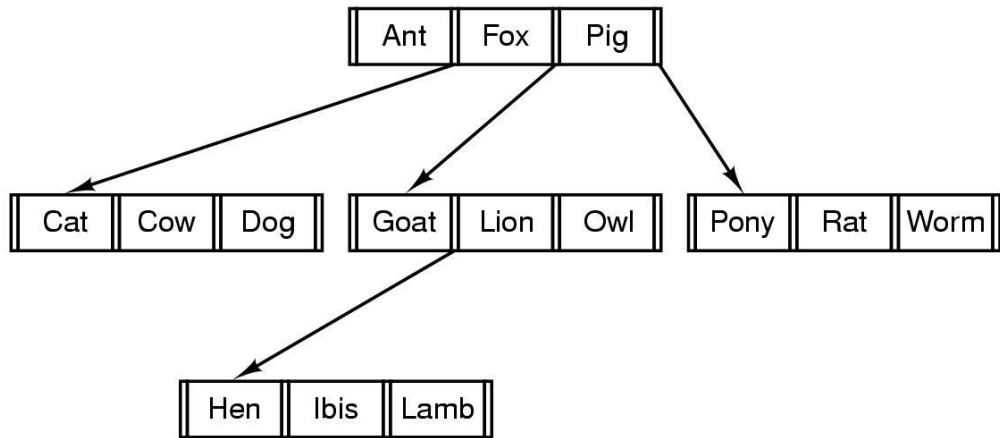3. Multiple processes must be able to access the information concurrently

# File Structure



- **Three kinds of files**
  - byte sequence
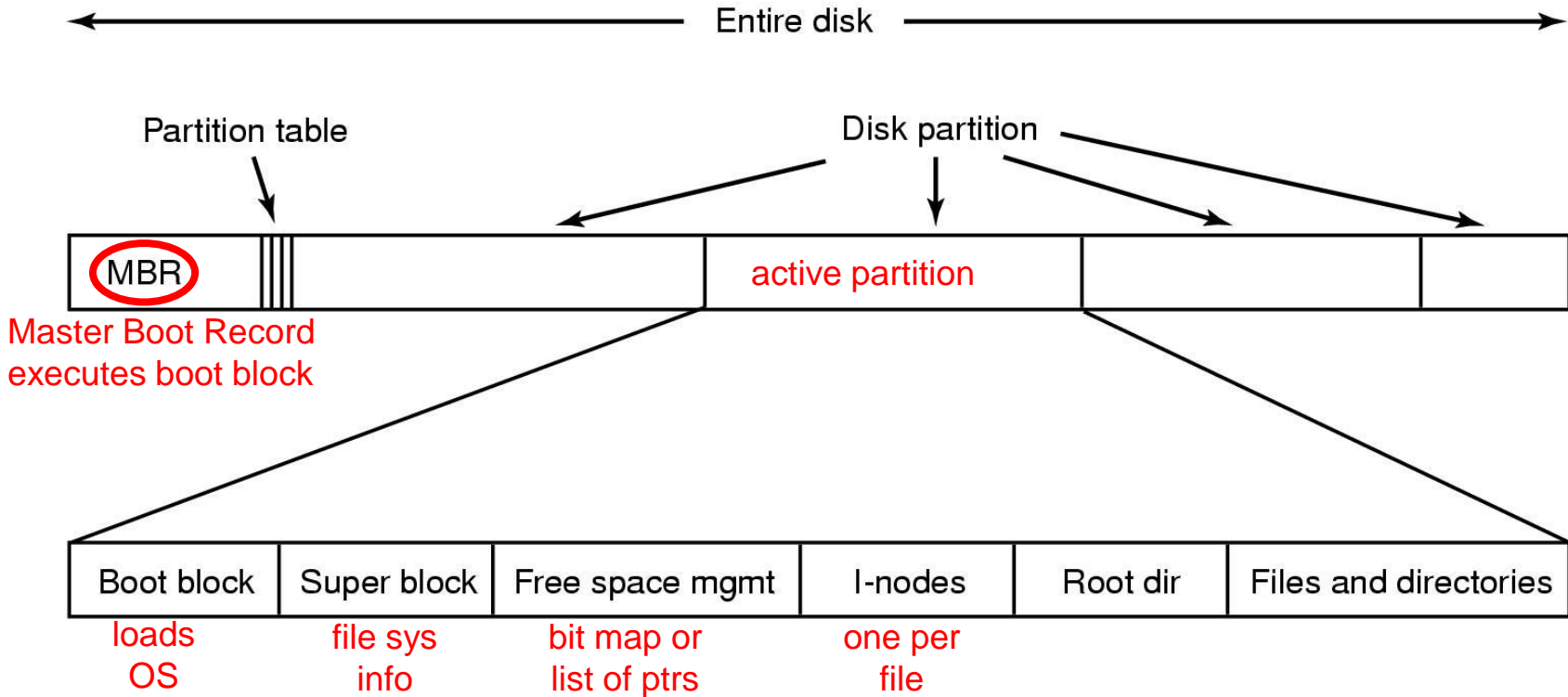  - record sequence
  - tree

3

# File Access

- Sequential access
  - read all bytes/records from the beginning
  - cannot jump around, could rewind or back up
  - convenient when medium was mag tape

- Random access
  - bytes/records read in any order
  - essential for data base systems
  - read can be …
    - move file marker (seek), then read or …
    - read and then move file marker
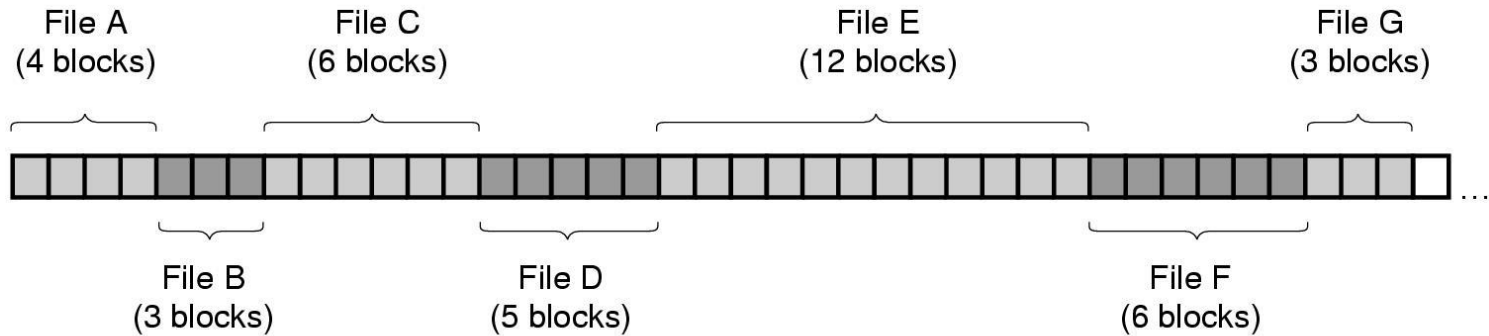
# File Operations

1. Create
2. Delete
3. Open
4. Close
5. Read
6. Write

7. Append
8. Seek
9. Get attributes
10. Set Attributes
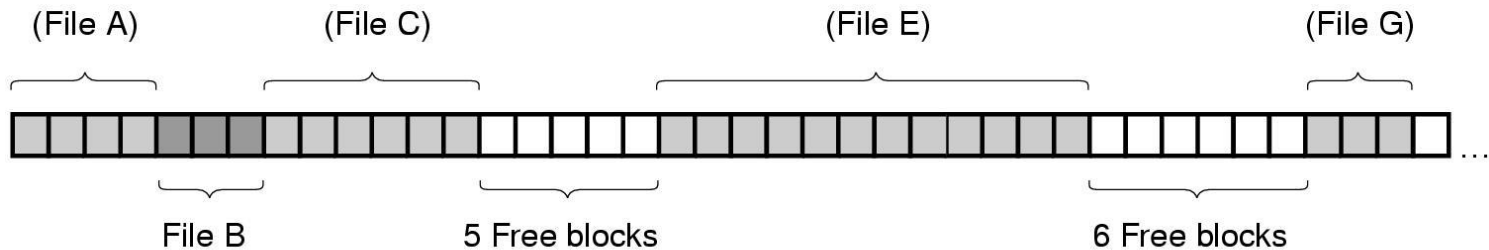11. Rename

# File System Implementation



Entire disk

Partition table          Disk partition

MBR

Master Boot Record
executes boot block

active partition

| Boot block | Super block | Free space mgmt | I-nodes | Root dir | Files and directories |
|---|---|---|---|---|---|
| loads OS | file sys info | bit map or list of ptrs | one per file | | |

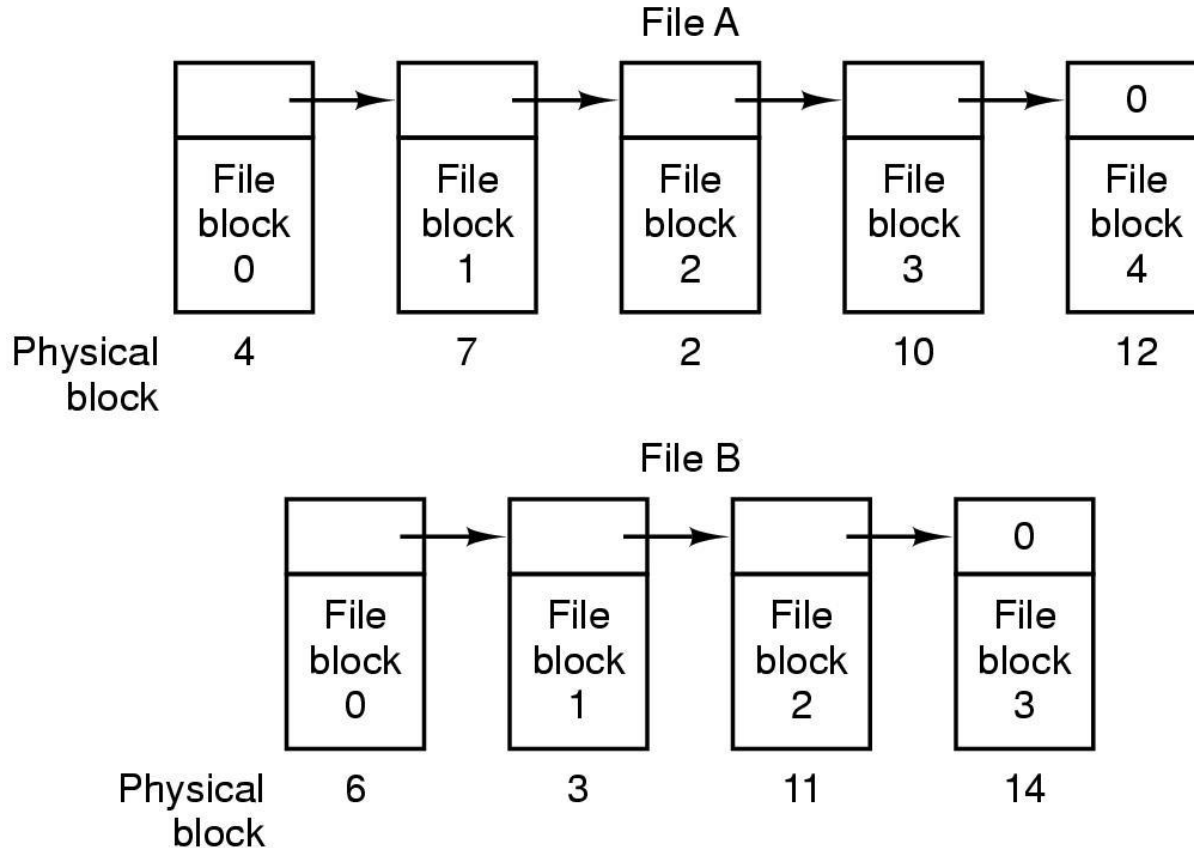## A possible file system layout

# Implementing Files (1)



(a)

(b)

(a) Contiguous allocation of disk space for 7 files
(b) Fragmened state of the disk after files *D* and *E* have
   been removed

7

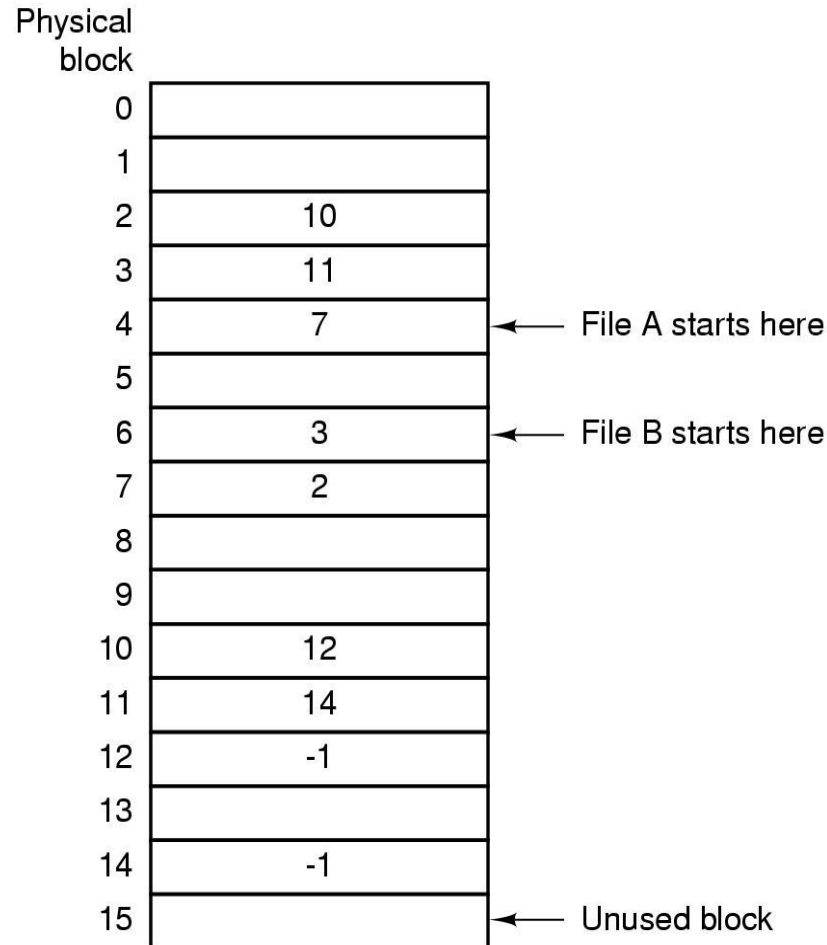# Implementing Files (2)



Storing a file as a linked list of disk blocks
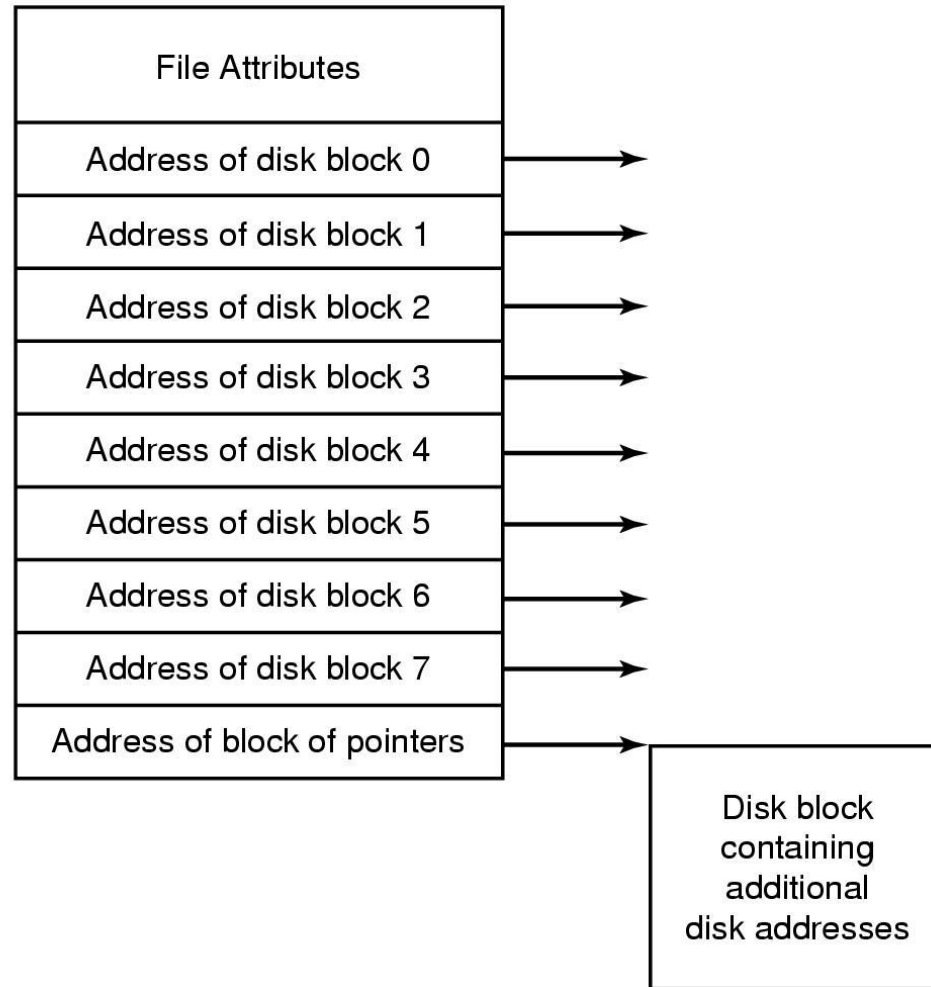But slow sequential reads and unusual block sizes

# Implementing Files (3)

Physical block

| | |
|---|---|
| 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 (FAT) in RAM
But 20GB disk with 1KB blocks with 4byte entries (pointers) = 80MB

9

# Implementing Files (4)



| File Attributes |
| --- |
| Address of disk block 0 |
| Address of disk block 1 |
| Address of disk block 2 |
| Address of disk block 3 |
| Address of disk block 4 |
| Address of disk block 5 |
| Address of disk block 6 |
| Address of disk block 7 |
| Address of block of pointers |

Disk block containing additional disk addresses

An example i-node

Only needs to be in memory when the file is open

# Implementing Directories (1)



(a)                    (b)                    Data structure
                                              containing the
                                              attributes
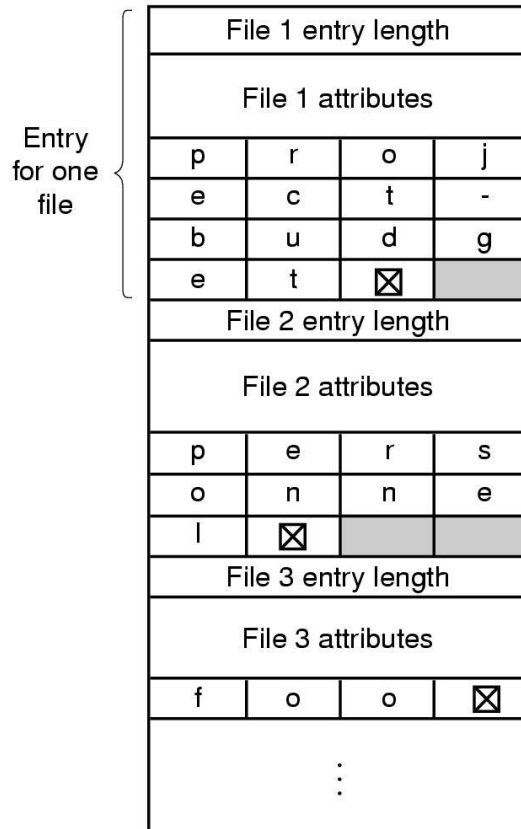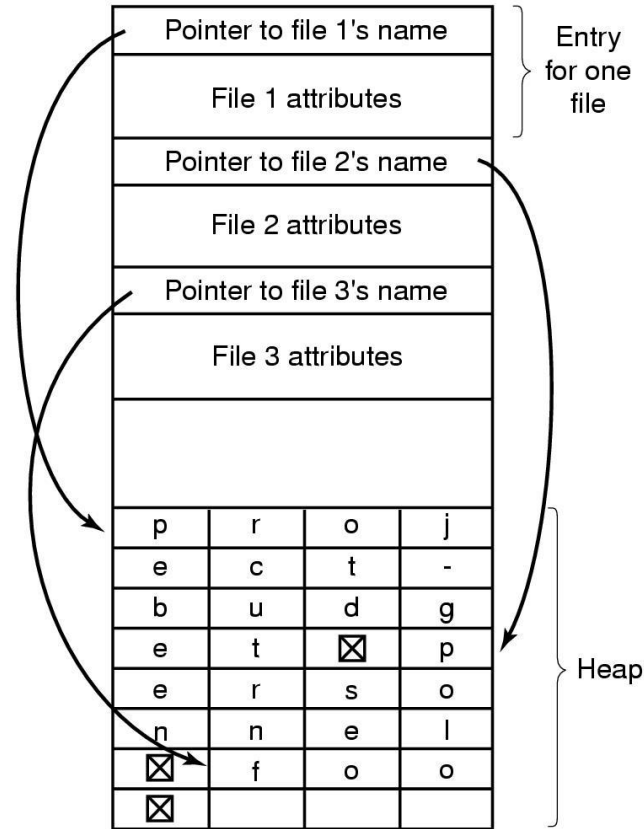
(a) A simple directory

    fixed size entries

    disk addresses and attributes in directory entry

(b) Directory in which each entry just refers to an i-node
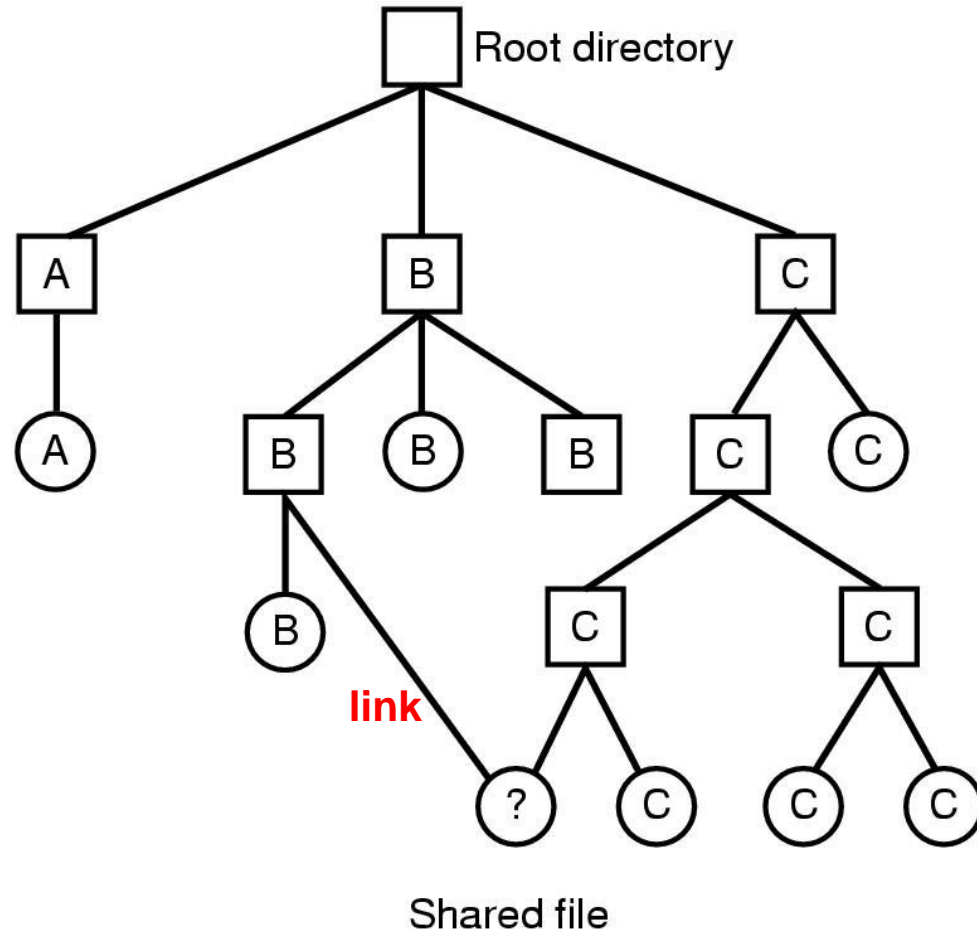
# Implementing Directories (2)

| File 1 entry length | | | |
|---|---|---|---|
| File 1 attributes | | | |
| p | r | o | j |
| e | c | t | - |
| b | u | d | g |
| e | t | ⊠ | |

| File 2 entry length | | | |
|---|---|---|---|
| File 2 attributes | | | |
| p | e | r | s |
| o | n | n | e |
| l | ⊠ | | |

| File 3 entry length | | | |
|---|---|---|---|
| File 3 attributes | | | |
| f | o | o | ⊠ |

Entry for one file

(a)

| Pointer to file 1's name |
|---|
| File 1 attributes |
| Pointer to file 2's name |
| File 2 attributes |
| Pointer to file 3's name |
| File 3 attributes |

Entry for one file

| p | r | o | j |
|---|---|---|---|
| e | c | t | - |
| b | u | d | g |
| e | t | ⊠ | p |
| e | r | s | o |
| n | n | e | l |
| ⊠ | f | o | o |
| ⊠ | | | |

Heap

(b)

- Two ways of handling long file names in directory
  - (a) In-line (becomes fragmented)   – slow linear search
  - (b) In a heap (need to manage heap) – slow linear search
- Two more ways for a faster lookup:
      (c) Hash table            (d) Cache search history

# Shared Files (1)



Shared file

File system containing a shared file
(Now a Directed Acyclic Graph (DAG) instead of a tree)

# Shared Files (2)
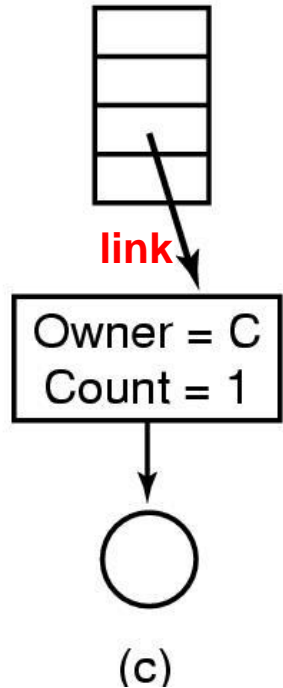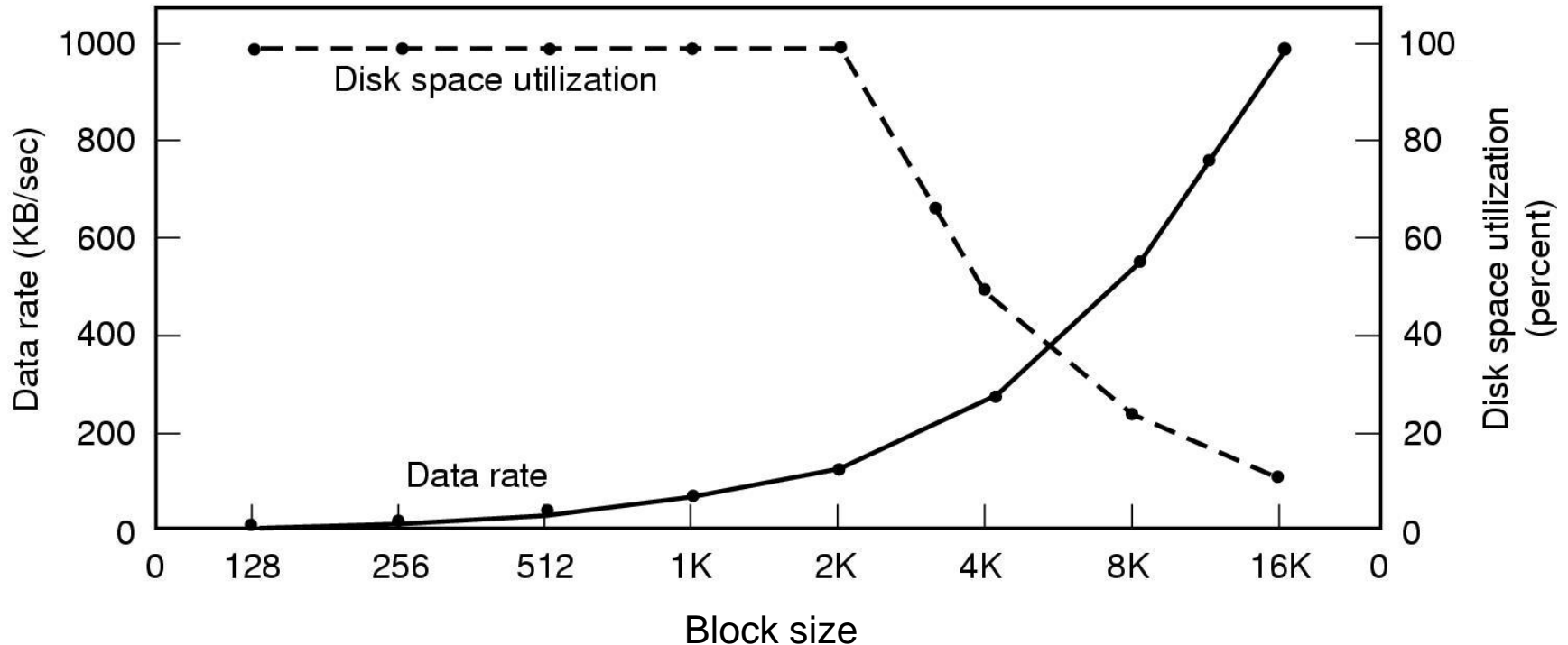


(a) Situation prior to linking
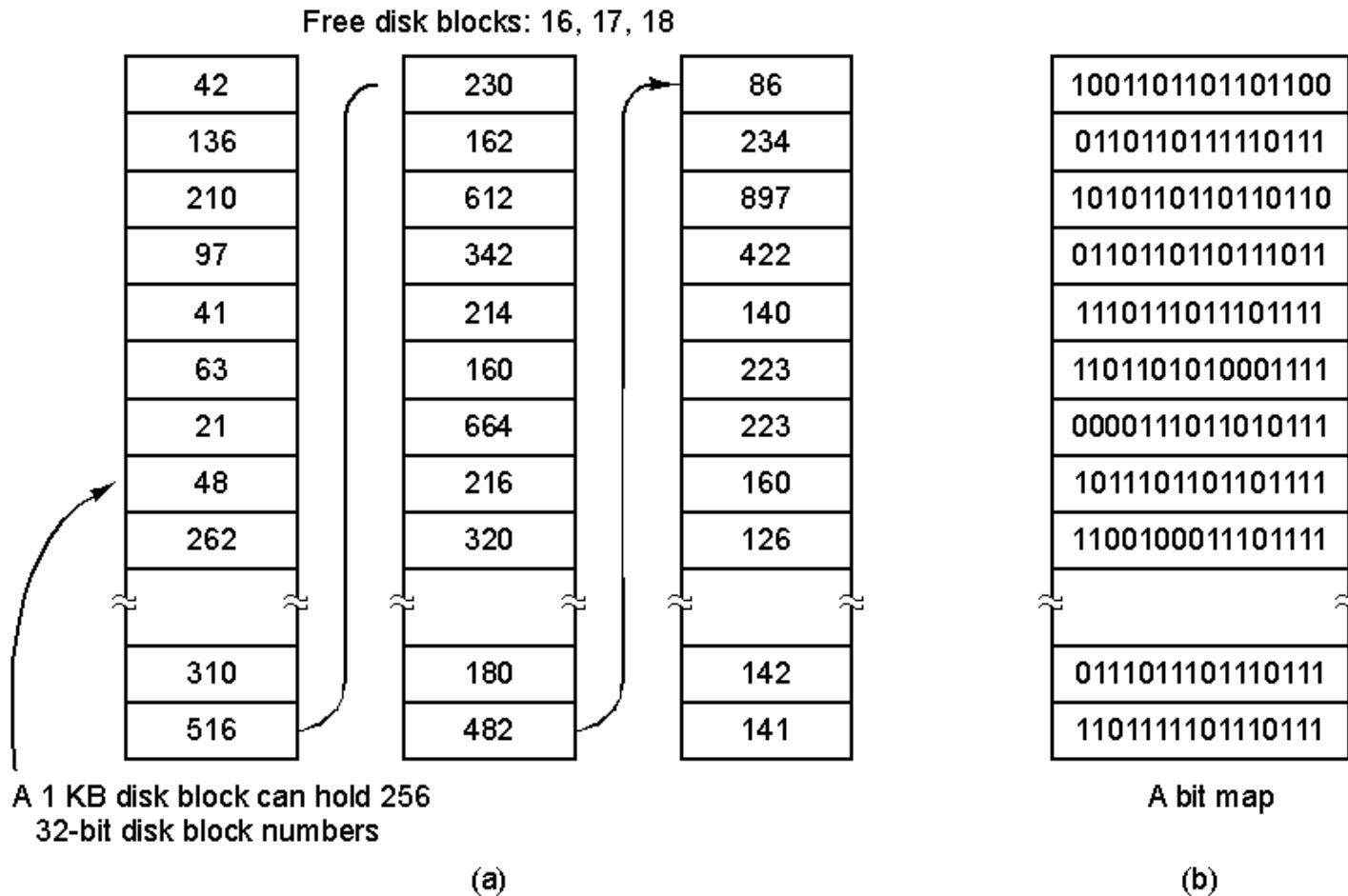
(b) After the link is created

(c) After the original owner removes the file - problem

# Disk Space Management (1)



- Dark line (left hand scale) gives data rate of a disk
- Dotted line (right hand scale) gives disk space efficiency
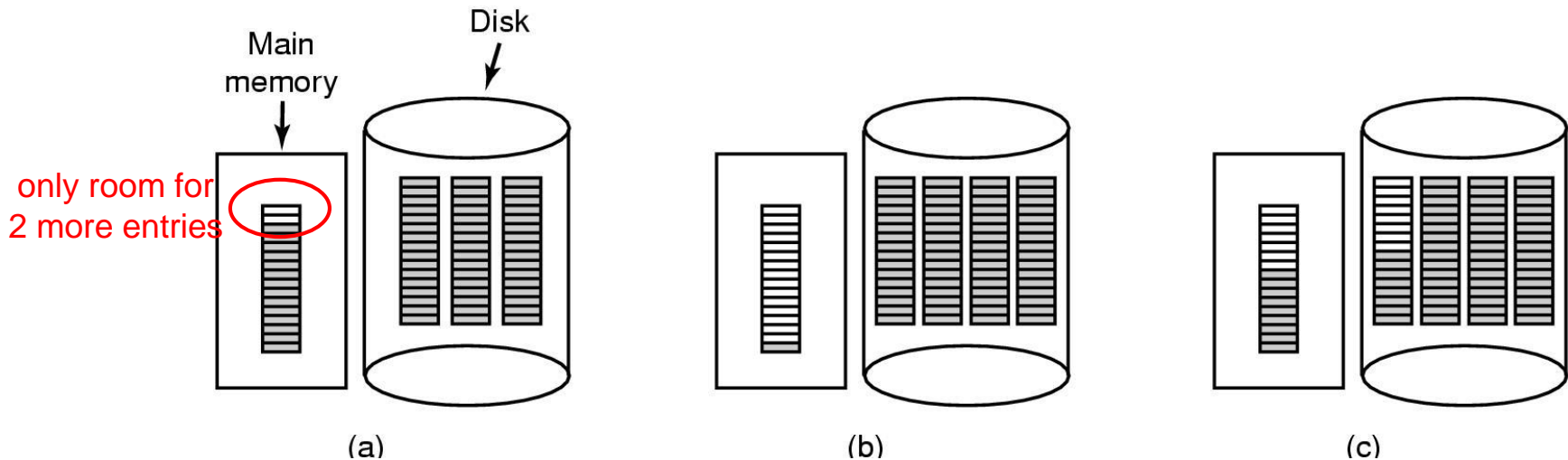- All files 2KB

# Disk Space Management (2)



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)

| 1001101101101100 |
| 0110110111110111 |
| 1010110110110110 |
| 0110110110111011 |
| 1110111011101111 |
| 1101101010001111 |
| 0000111011010111 |
| 1011101101101111 |
| 1100100011101111 |
| ~ |
| 0111011101110111 |
| 1101111101110111 |

A bit map

(b)

(a) Storing the free list on a linked list **(only better if disk nearly full)**

(b) A bit map - one bit per free block

16

# Disk Space Management (3)



only room for 2 more entries

shaded entries are pointers to free disk blocks

(a) Almost-full block of pointers to free disk blocks in RAM
- three blocks of pointers on disk

(b) Result of freeing a 3-block file

(c) Alternative strategy for handling 3 free blocks
- keep the one in memory about half full so it can handle file create and remove without disk I/O on the free list

# File System Reliability (1)



- A file system to be backed up
  - squares are directories, circles are files
  - shaded items, modified since last backup
  - each directory & file labeled by i-node number

# File System Reliability (2)

(a) | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32

mark each directory and modified file for backup

(b) | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32

unmark each directory without modified files

(c) | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32

backup each modified directory

(d) | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32

backup each modified file

Bit maps used by the logical backup algorithm

# File System Reliability (3)

Block number

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

| 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | Blocks in use |

| 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | Free blocks |

(a)

Block number

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

| 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | Blocks in use |

| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | Free blocks |

(b)

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

| 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | Blocks in use |

| 0 | 0 | 1 | 0 | 2 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | Free blocks |

(c)

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

| 1 | 1 | 0 | 1 | 0 | 2 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | Blocks in use |

| 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | Free blocks |

(d)

- File system states
  - (a) consistent
  - (b) missing block (e.g. after system crash)
  - (c) duplicate block in free list
  - (d) duplicate data block (worst case)

# File System Performance (1)



The block cache data structures
(all blocks with the same hash value are linked into a collision chain)
(cache from *cacher – to hide – French)*

# File System Performance (2)



I-nodes are located near the start of the disk

Disk is divided into cylinder groups, each with its own i-nodes

Cylinder group

(a)

(b)

- I-nodes placed at the start of the disk
- Disk divided into cylinder groups
  - each with its own blocks and i-nodes

# Log-Structured File Systems (LFS)

- Log-structured File System (LFS) Strategy structures the entire disk as a log
  - have all writes initially buffered in memory
  - periodically write these to the end of the disk log
  - when file opened, locate i-node, then find blocks

- With CPUs faster, memory larger
  - disk caches can also be larger
  - increasing number of read requests can come from cache
  - thus, most disk accesses will be writes

# Example File Systems
## CD-ROM File Systems



The ISO 9660 directory entry

# The MS-DOS File System (1)



| 8 | 3 | 1 | 10 | 2 | 2 | 2 | 4 |
|---|---|---|----|---|---|---|---|
| File name | | | | | | | Size |

Bytes

Extension  Attributes  Reserved  Time  Date  First block number

The MS-DOS directory entry

# The MS-DOS File System (2)

| Block size | FAT-12 | FAT-16 | FAT-32 |
|---|---|---|---|
| 0.5 KB | 2 MB | | |
| 1 KB | 4 MB | | |
| 2 KB | 8 MB | 128 MB | |
| 4 KB | 16 MB | 256 MB | 1 TB |
| 8 KB | | 512 MB | 2 TB |
| 16 KB | | 1024 MB | 2 TB |
| 32 KB | | 2048 MB | 2 TB |

- Maximum partition for different block sizes
- The empty boxes represent forbidden combinations

# The Windows File System (1)



The extended MOS-DOS directory entry used in Windows

# The Windows File System (2)



Bytes: 1, 10, 1, 1, 1, 12, 2, 4

| 5 characters | | 0 | | 6 characters | 0 | 2 characters |

Sequence — Attributes — Checksum

An entry for (part of) a long file name in Windows

# The Windows File System (3)

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 68 | d   o   g | A | 0 | CK | | | | | | 0 | |
| 3 | o   v   e | A | 0 | CK | t   h   e     l   a | | | | | 0 | z   y |
| 2 | w   n     f   o | A | 0 | CK | x     j   u   m   p | | | | | 0 | s |
| 1 | T   h   e     q | A | 0 | CK | u   i   c   k     b | | | | | 0 | r   o |
| T | H E Q U I ~ 1 | A | NT | S | Creation time | Last acc | Upp | | Last write | Low | Size |
| Bytes | | | | | | | | | | | |

An example of how a long name is stored in Windows

# The LINUX File System (1)



A LINUX directory entry

# The LINUX File System (2)



A LINUX i-node

# The LINUX File System (3)



| Root directory | |
|---|---|
| 1 | . |
| 1 | .. |
| 4 | bin |
| 7 | dev |
| 14 | lib |
| 9 | etc |
| 6 | usr |
| 8 | tmp |

Looking up
usr yields
i-node 6

I-node 6
is for /usr

Mode
size
times

132

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

Block 132
is /usr
directory

| 6 | • |
|---|---|
| 1 | •• |
| 19 | dick |
| 30 | erik |
| 51 | jim |
| 26 | ast |
| 45 | bal |

/usr/ast
is i-node
26

I-node 26
is for
/usr/ast

Mode
size
times

406

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

Block 406
is /usr/ast
directory

| 26 | • |
|---|---|
| 6 | •• |
| 64 | grants |
| 92 | books |
| 60 | mbox |
| 81 | minix |
| 17 | src |

/usr/ast/mbox
is i-node
60

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

# Btree – balanced tree file system

length of path from root to any leaf node should always be the same
nodes must contain a minimum number of pairs



*The leaf node's keys are ordered within the tree improving scan times, since the scan is no longer sequential. Leaf nodes are chained using pointers to each other.*

# Btree – balanced tree file system
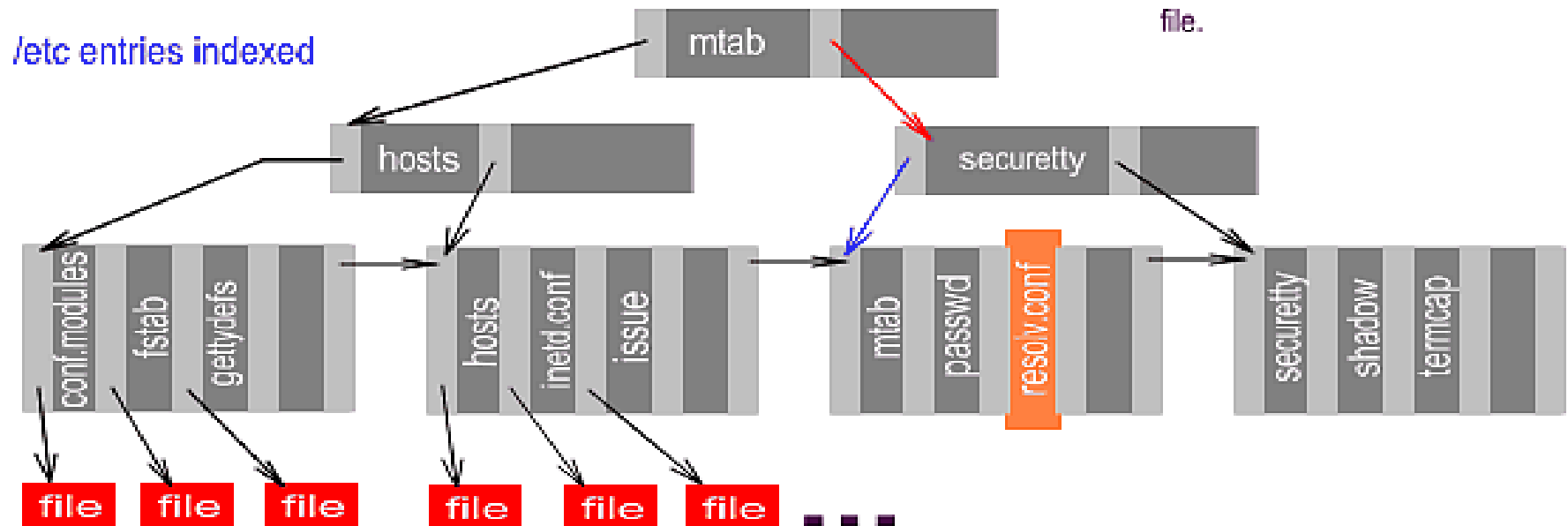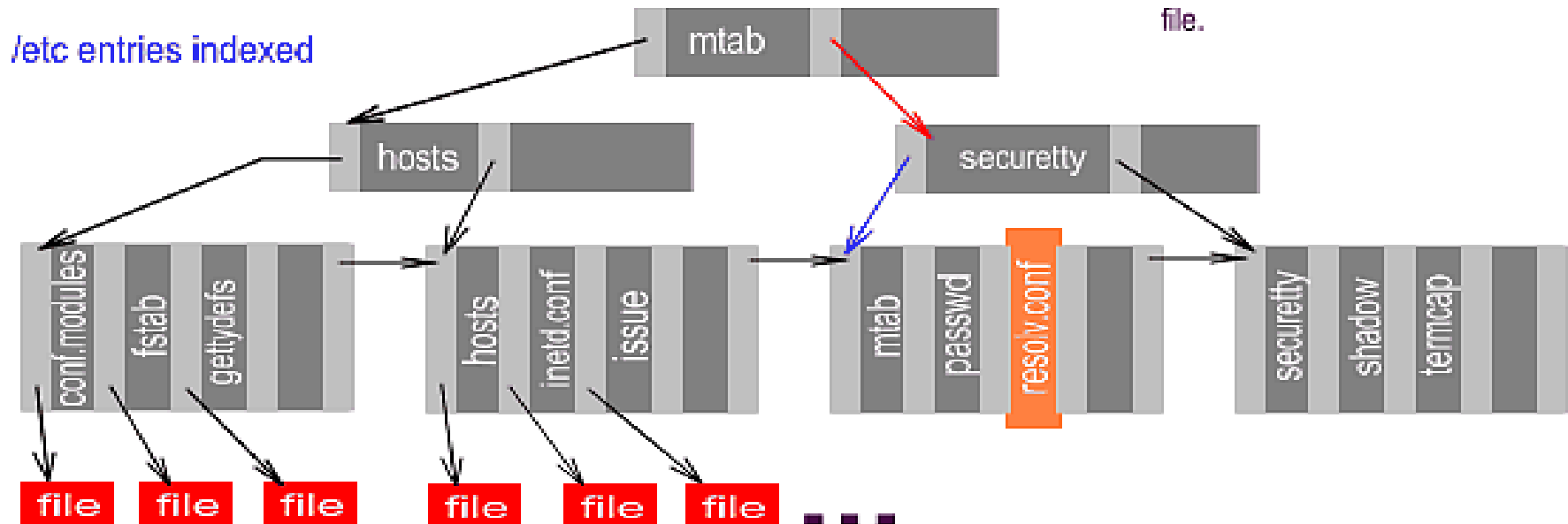
(1) to locate the file resolv.conf we begin at the tree's root. scan sequentially, and find that there is no key greater than resolv.conf, so we use the last pointer (in red)

(2) got directed to another internal node. Let's do the same. Scan through the node's keys and realise that "securetty" is a greater key than "resolv.conf". We use the accompanying pointer (in blue).

(3) we got the final leaf node. Now it's time to scan sequentially throughout the ascending ordered keys of the node. Finally, found the desired key, we should use the accompanying pointer to the "resolv.conf" named file.

/etc entries indexed

mtab

hosts

securetty

conf.modules | fstab | getdyefs

hosts | inetd.conf | issue

mtab | passwd | resolv.conf

securetty | shadow | termcap

file | file | file

file | file | file

*The keys are file names. The bottom row above the red boxes contains a key for every file in the directory: these are the leaf nodes. Above these are the internal nodes, keys that have been chosen by the system to make finding other keys faster*