# File Systems

18 January 2026
Lecture 12

Slides adapted from John Kubiatowicz (UC Berkeley)

# Concept Review

| Fragmentation | Page Table Entry | Valid/Invalid bit |
|---|---|---|
| • Internal<br>• External | | |
| Page | Page table | Multi-level page table |

# Topics for Today

- File Systems
  - Introduction to File Systems
  - Very simply file system
  - FAT
  - Inodes
  - Unix Fast File System (FFS)

# Building a File System

- File System: Layer of OS that transforms block interface of disks (or other block devices) into Files, Directories, etc.

  – File System Components
    - Disk Management: collecting disk blocks into files
    - Naming: Interface to find files by name, not by blocks
    - Protection: Layers to keep data secure
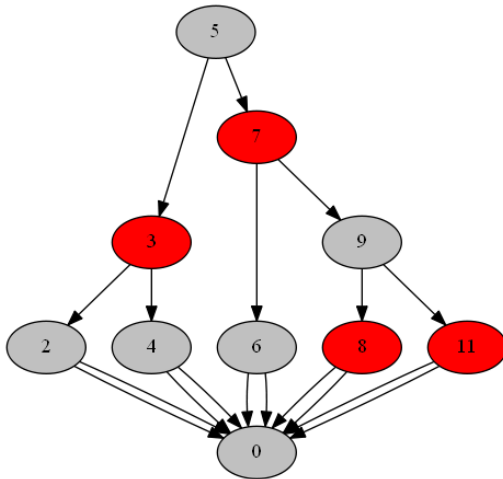    - Reliability/Durability: Keeping of files durable despite crashes, media failures, attacks, etc

# User vs. System View of a File
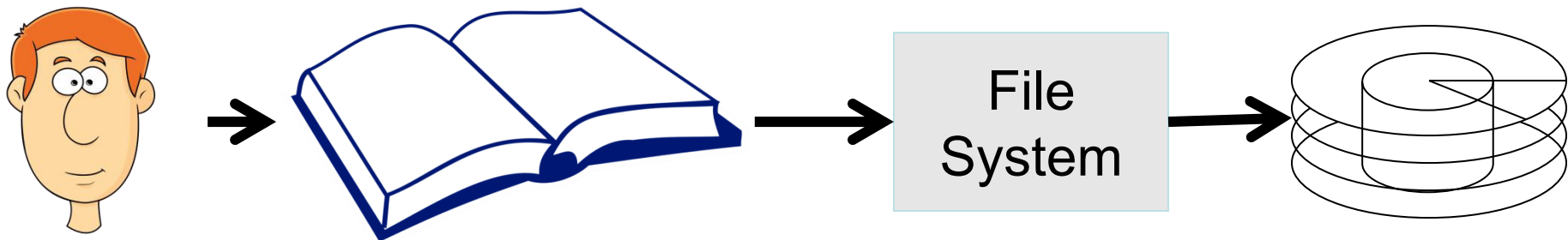
## User's view:

- Durable Data Structures



## System's View:

- System call interface:
  - Collection of Bytes (UNIX)
  - Doesn't matter to system what kind of data structures you want to store on disk!

- Inside OS:
  - Collection of blocks (a block is a logical transfer unit, while a sector is the physical transfer unit)
  - Block size $\geq$ sector size; in UNIX, block size is 4KB
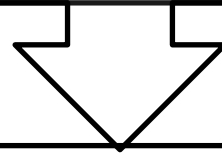
# Translating from User to System View

- What happens if user says: *Give me bytes 2-12?*
  - Fetch block corresponding to those bytes
  - Return just the correct portion of the block
- What about: *Write bytes 2-12?*
  - Fetch block, Modify portion, Write out Block
- Everything inside File System is in whole size blocks
  - For example, `getc()`, `putc()` $\Rightarrow$ buffers something like 4096 bytes, even if interface is one byte at a time
- From now on, file is a collection of blocks

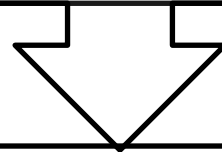# So you are going to design a file system …

What factors are critical to the design choices?

Durable data store → It's all on disk

Disk Performance!

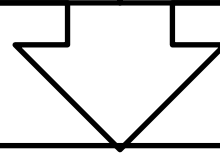Maximize sequential access, minimize seeks

Open before Read/Write

Can perform protection checks and look up where the actual file resource are, in advance

# So you are going to design a file system …
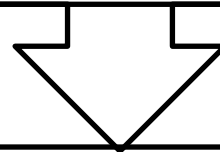
| Size is determined as files are used | |
|---|---|
| Can write (or read zeros) to expand the file | Start small and grow, need to make room |

⬇

| Organize into directories |
|---|
| What data structure (on disk) do we use for that? |

⬇

| Need to allocate and free blocks |
|---|
| Keep access efficient |

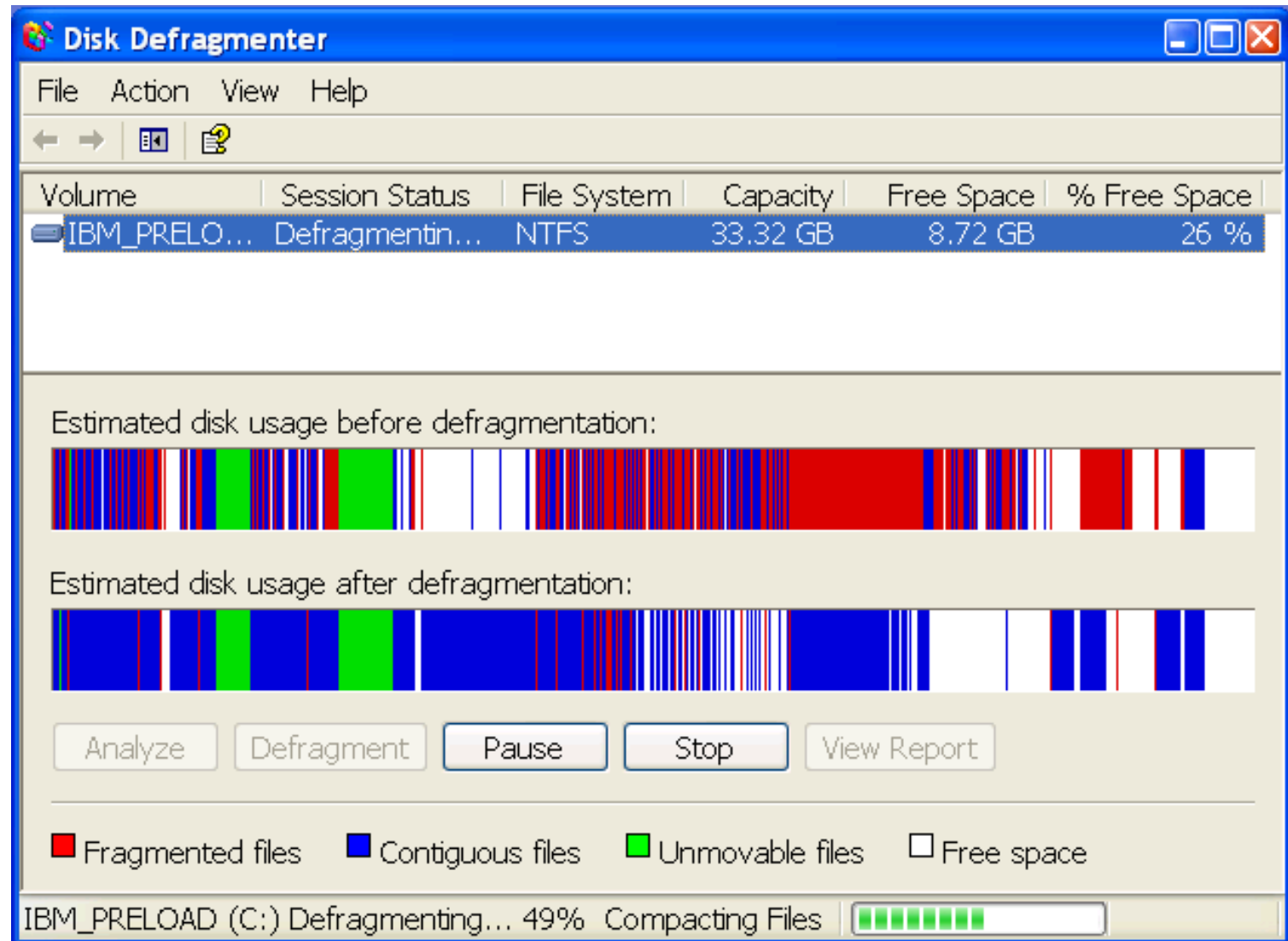# Defragmenting
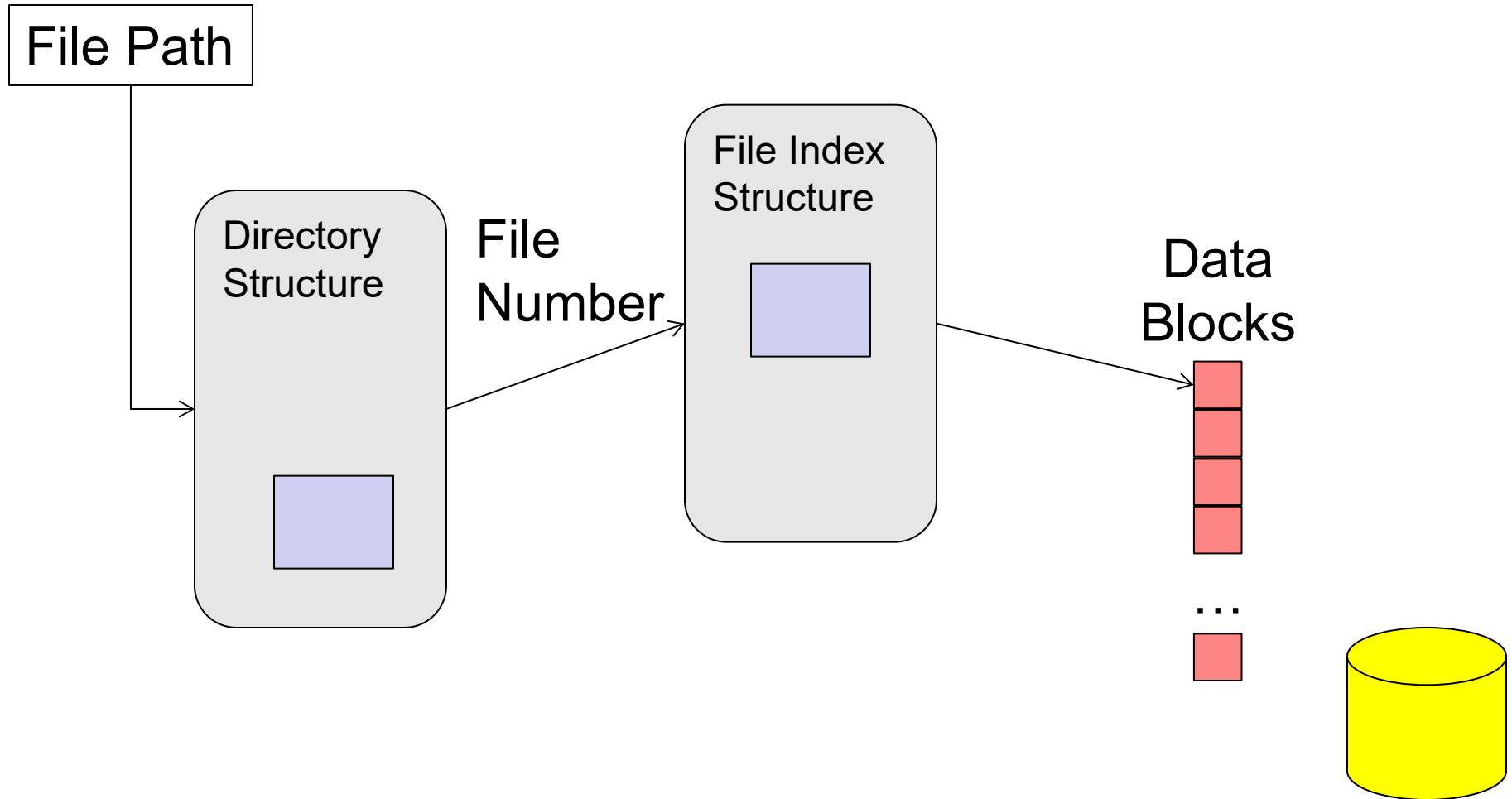


Image source: http://blog.shane-smith.com/blog:116

# Disk Management Policies

- Basic entities on a disk:
  - File: User-visible group of blocks arranged sequentially in logical space
  - Directory: user-visible index mapping names to files

- Access disk as linear array of sectors.
  Two Options:
  1. Identify sectors as vectors `[cylinder, surface, sector]`. Sort in cylinder-major order. Not used much anymore.
  2. Logical Block Addressing (LBA). Every sector has integer address from zero up to max number of sectors.

- Controller translates from address⇒physical position
  - First case: OS/BIOS must deal with bad sectors
  - Second case: Hardware shields OS from structure of disk

# Disk Management Policies

- Need way to track free disk blocks
  - Link free blocks together ⇒ too slow today
  - Use bitmap to represent free space on disk

- Need a way to structure files: File Header
  - Track which blocks belong at which offsets within the logical file structure
  - Optimize placement of files' disk blocks to match access and usage patterns

# Components of a File System

File Path

Directory Structure

File Number

File Index Structure

Data Blocks

…

# Components of a file system

- Open performs *name resolution*
  - Translates pathname into a "file number"
    - Used as an "index" to locate the blocks
  - Creates a file descriptor in PCB within kernel
  - Returns a "handle" (another `int`) to user process

- Read, Write, Seek, and Sync operate on handle
  - Mapped to **descriptor** and to **blocks**

| File Name Offset | →Directory→ | File Number Offset | →Index Structure→ | Storage Block |
|---|---|---|---|---|

# Directories

# Directory

- Basically a hierarchical structure

- Each directory entry is a collection of
  - Files
  - Directories
    - A link to another entries

- Each has a name and attributes
  - Files have data

- Links (hard links) make it a DAG, not just a tree
  - Soft links (aliases) are another name for an entry

sdcard    2013-12-10  23:51    lrwxrwxrwx    -> /mnt/sdcard

# I/O & Storage Levels

Application/Service

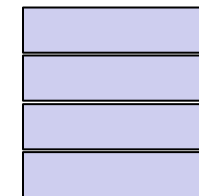| High Level I/O | Streams |
| Low Level I/O | Handles |

Registers

| Syscalls | #4 file handle |

Data Blocks

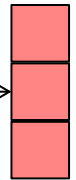| File System | Descriptors |
| I/O Driver | Commands and Data Transfers |

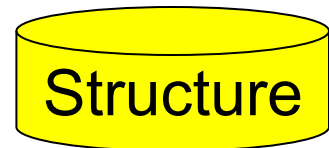Disks, Flash, Controllers, DMA
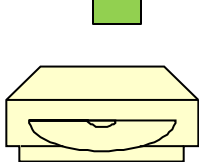
…

Directory

Structure

# File

- **Named permanent storage**

Contains

- Data
  - Blocks on disk somewhere
- Metadata (Attributes)
- Owner, size, last opened, …
- Access rights
  - R, W, X
  - Owner, Group, Other (in Unix systems)
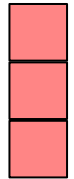  - Access control list in Windows system

Data
Blocks

File Handle
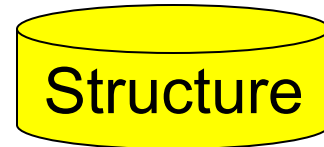_____

File Descriptor
FileObject (inode)
Position

…

Structure

# File Attributes



**info.txt Properties**

General | Security | Details | Previous Versions

info.txt

Type of file: TXT File (.txt)

Opens with: Notepad++ : a free (GNU | Change...

Location: C:\Users\Michael\Desktop

Size: 109 bytes (109 bytes)

Size on disk: 0 bytes

Created: Today, May 01, 2017, 4 hours ago

Modified: Today, May 01, 2017, 4 hours ago

Accessed: Today, May 01, 2017, 4 hours ago

Attributes: ☐ Read-only   ☐ Hidden   Advanced...

OK   Cancel   Apply

**Advanced Attributes**

Choose the settings you want for this folder.

File attributes

☑ File is ready for archiving

☑ Allow this file to have contents indexed in addition to file

Compress or Encrypt attributes

☐ Compress contents to save disk space

☐ Encrypt contents to secure data   Details

OK   Cancel

SE 317: Operating Systems

# File Attributes

SE 317: Operating Systems

# File Metadata

**alice_promo.mp4 Properties**

General | Security | **Details** | Previous Versions

| Property | Value |
|---|---|
| **Description** | |
| Title | |
| Subtitle | |
| Rating | ☆ ☆ ☆ ☆ ☆ |
| Tags | |
| Comments | |
| **Video** | |
| Length | 00:03:59 |
| Frame width | 352 |
| Frame height | 288 |
| Data rate | 245kbps |
| Total bitrate | 364kbps |
| Frame rate | 29 frames/second |
| **Audio** | |
| Bit rate | 118kbps |
| Channels | 2 (stereo) |
| Audio sample rate | 44 kHz |
| **Media** | |

Remove Properties and Personal Information

OK | Cancel

---

**alice_promo.mp4 Properties**

General | Security | **Details** | Previous Versions

| Property | Value |
|---|---|
| Mood | |
| **Part of set** | |
| Initial key | |
| Beats-per-minute | |
| Protected | No |
| **File** | |
| Name | alice_promo.mp4 |
| Item type | MP4 File |
| Folder path | C:\Users\Michael\Documents\Extern... |
| Size | 10.5 MB |
| Date created | 03-Jun-13 3:48 PM |
| Date modified | 03-Jun-13 3:56 PM |
| Attributes | A |
| Availability | Available offline |
| Offline status | |
| Shared with | Michael-Asp |
| Owner | NOV\Michael |
| Computer | NOV (this PC) |

Remove Properties and Personal Information

OK | Cancel | Apply

# In-Memory File System Structures



Open
(File Name)

Directory Structure

Directory Structure

File Control Block

User Space

Kernel Memory

Secondary Storage

- Open system call:
  - Resolves file name, finds file control block (inode)
  - Makes entries in per-process and system-wide tables
  - Returns index (called "file handle") in open-file table

# In-Memory File System Structures



- Read/write system calls:
  - Use file handle to locate inode
  - Perform appropriate reads or writes

# So Far

- File Systems
  - Introduction to File Systems
  - Very simply file system
  - FAT
  - Inodes
  - Unix Fast File System (FFS)

# A simple File System: Ingredients

## Blocks

- 4 KB
- 64 total blocks
  - Numbered 0-63

## Data space

## Metadata space

## Superblock

# Our blocks (4KB)

# Data Region



Data Region

| D D D D D D D D | D D D D D D D D | D D D D D D D D |
| 0        7 | 8      15 | 16     23 | 24     31 |

Data Region

| D D D D D D D D | D D D D D D D D | D D D D D D D D | D D D D D D D D |
| 32     39 | 40     47 | 48     55 | 56     63 |

| Blocks 8-63 will be data | 56 total | Rest will be metadata |

# Inode region – What's an inode?



Data structure holds metadata about file

- Size
- Owner
- Access rights
- Access/modify times

Typically stored in a table

- Small, say 256 bytes

# Inode region



Use 5 blocks for inodes

Each 4KB block holds 16 inodes

• 4096 / 256 = 16

Total 80 inodes

• Can't have more than 80 files or directories

# Used/Free tracking



Need to track if an inode is in used

Need to track if a data block is in use

Use a **bitmap** for each

- 0 if inode/block is free
- 1 if inode/block in use
- Example: 1011 0011

# Used/Free tracking



Inodes | Data Region

| i | d | I | I | I | I | I | DDDDDDDD | DDDDDDDD | DDDDDDDD |
0 ........ 7 8 ........ 15 16 ........ 23 24 ........ 31

Data Region

DDDDDDDD | DDDDDDDD | DDDDDDDD | DDDDDDDD
32 ........ 39 40 ........ 47 48 ........ 55 56 ........ 63

| Block 1 tracks inodes | Block 2 tracks data blocks |
|---|---|
| • Could really track 32K inodes, but let's be simple | • Could really track 32K blocks, but let's be simple |

# Superblock



Block 0 contains basic metadata

- Called Superblock

Contains info about the file system

- How many inodes
- How many blocks
- Where things are
- IDs, etc.

# Inode drill down

The Inode Table (Closeup)

|  | iblock 0 | | | | iblock 1 | | | | iblock 2 | | | | iblock 3 | | | | iblock 4 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 3 | 16 | 17 | 18 | 19 | 32 | 33 | 34 | 35 | 48 | 49 | 50 | 51 | 64 | 65 | 66 | 67 |
| | 4 | 5 | 6 | 7 | 20 | 21 | 22 | 23 | 36 | 37 | 38 | 39 | 52 | 53 | 54 | 55 | 68 | 69 | 70 | 71 |
| | 8 | 9 | 10 | 11 | 24 | 25 | 26 | 27 | 40 | 41 | 42 | 43 | 56 | 57 | 58 | 59 | 72 | 73 | 74 | 75 |
| | 12 | 13 | 14 | 15 | 28 | 29 | 30 | 31 | 44 | 45 | 46 | 47 | 60 | 61 | 62 | 63 | 76 | 77 | 78 | 79 |

Super   i-bmap   d-bmap

0KB   4KB   8KB   12KB   16KB   20KB   24KB   28KB   32KB

## Each inode has an implicit number

- i-number

## Can index the file by taking i-number X sizeof(inode)

- To read inode 32, start at byte 32 X sizeof(inode) = 8192

# What is a directory?

```
inum | reclen | strlen | name
   5      12         2      .
   2      12         3      ..
  12      12         4      foo
  13      12         4      bar
  24      36        28      foobar_is_a_pretty_longname
```

- Contains records of files and directories inside
- i-number of the file/directory
- Reclen – how many bytes in the record
- Strlen – how many bytes in the name of the file/directory
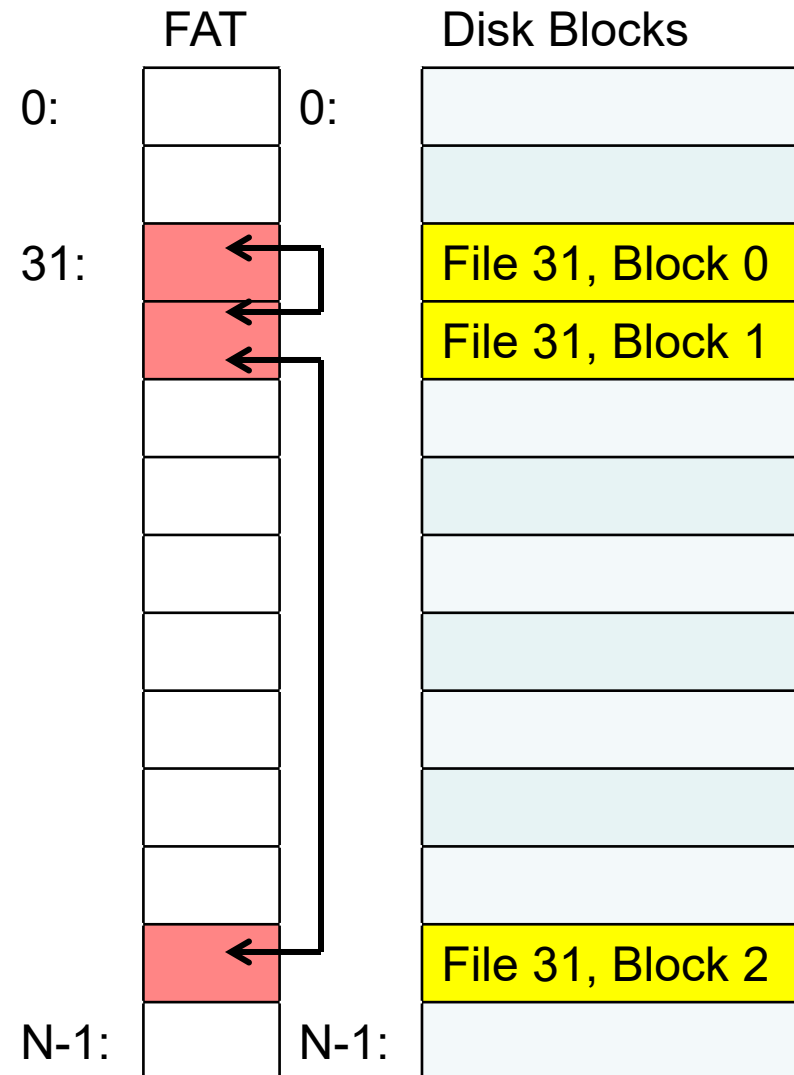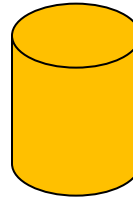- Name – the actual name ($\leq$ reclen)

# So Far

- File Systems
  - Introduction to File Systems
  - Very simply file system
  - FAT
  - Inodes
  - Unix Fast File System (FFS)

# Our first filesystem: FAT (File Allocation Table)

- **Assume we have a way to translate a path to a "file number"**
  - i.e., a directory structure
- **Disk Storage is a collection of Blocks**
  - Just hold file data
- **Ex: file_read $31, \langle 2, x \rangle$**
  - Index into FAT with file number
  - Follow linked list to block
  - Read the block from disk into mem

**FAT**

0:

31:

N-1:

**Memory**

**Disk Blocks**

0:

File 31, Block 0

File 31, Block 1

File 31, Block 2

N-1:

# FAT Properties

- File is collection of disk blocks

- FAT is linked list 1-1 with blocks

- File Number is index of root
  of block list for the file

- File offset $(o = B{:}x)$

- Follow list to get block #

- Unused blocks ⇔ FAT free list

File Number

FAT

Disk Blocks

0:

0:

31:

File 31, Block 0

File 31, Block 1

Free

File 31, Block 2

N-1:

N-1:

# Writing a File Block

Ex: `file_write(51,⟨3,y⟩)`

- Grab blocks from free list

- Linking them into file

# Create a New File

- Grow file by allocating free blocks and linking them in

- Ex: Create file, write, write

File Number

FAT

Disk Blocks

0:

0:

31:

File 31, Block 0

File 31, Block 1

Free

File 31, Block 3

File Number 2 → 63:

File 63, Block 0

File 31, Block 2

N-1:

N-1:

# Create a New File

- Grow file by allocating free blocks and linking them in

- Ex: Create file, write, write

# FAT Assessment

- Used in DOS, Windows, USB drives, …

- Where is FAT stored?
  - On disk, restore on boot, copy in memory

- What happens when you format a disk?
  - Zero the blocks, link up the FAT free-list

- Simple

FAT          Disk Blocks

0:        0:

31:       File 31, Block 0

          File 31, Block 1

          File 63, Block 1

Free

          File 31, Block 3

63:       File 63, Block 0

          File 31, Block 2

N-1:      N-1:

# FAT Assessment

- Time to find block (large files)?

- Free list usually just a bit vector (here's it's a linked list).

- Block layout for file?

- Sequential Access?

- Random Access?

- Fragmentation?

- Small files?

- Big files?

# What about the Directory?

| | File 5268830<br>"/home/mjmay" | | | | | | End of<br>File |
|---|---|---|---|---|---|---|---|
| Name | . | .. | Music | Data | | First.txt | |
| File Number | 5268830 | 88026158 | 35002320 | 85200219 | Free Space | 66212871 | Free Space |
| Next | | | | | | | |

- Essentially a file containing `<file_name:file_number>` mappings
- Free space for new entries
- In FAT: attributes kept in directory (!)
- Each directory a linked list of entries
- Where do you find root directory ( "/" )?

# Directory Structure

- How many disk accesses to resolve "`/my/book/count`"?
    - Read in file header for root (fixed spot on disk)
    - Read in first data block for `root`
        - Table of file name/index pairs.  Search linearly – ok since directories typically very small
    - Read in file header for "`my`"
    - Read in first data block for "`my`"; search for "`book`"
    - Read in file header for "`book`"
    - Read in first data block for "`book`"; search for "`count`"
    - Read in file header for "`count`"

- Current working directory: Per-address-space pointer to a directory (inode) used for resolving file names
    - Allows user to specify relative filename instead of absolute path (say CWD="`/my/book`" can resolve "`count`")

# Big FAT security holes

- FAT has no access rights
- FAT has no header in the file blocks
- Just gives an index into the FAT
  - (file number = block number)

Don't hurt me!

# So Far

- File Systems
  - Introduction to File Systems
  - Very simply file system
  - FAT
  - Inodes
  - Unix Fast File System (FFS)

# Characteristics of Files

- Most files are small

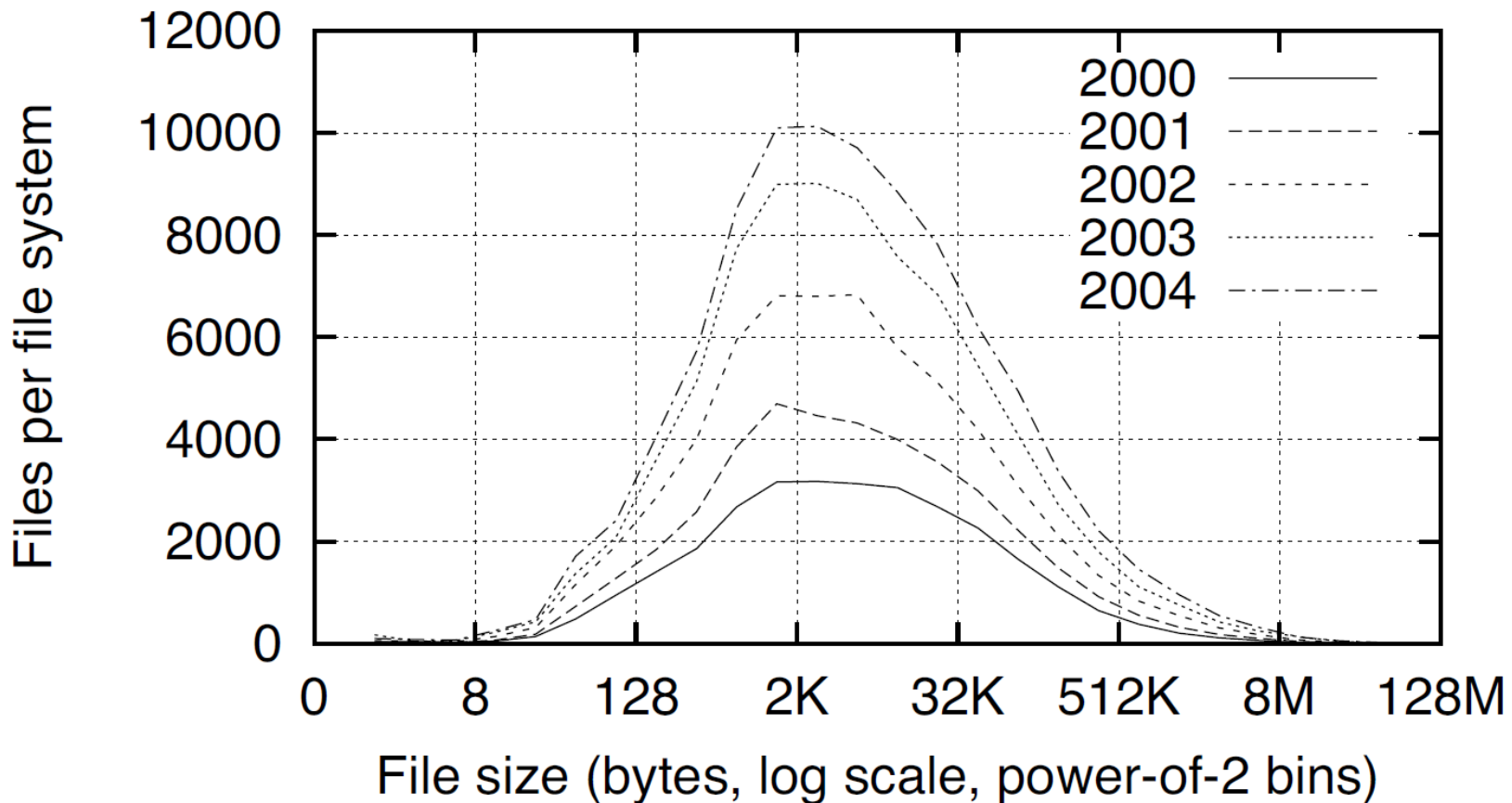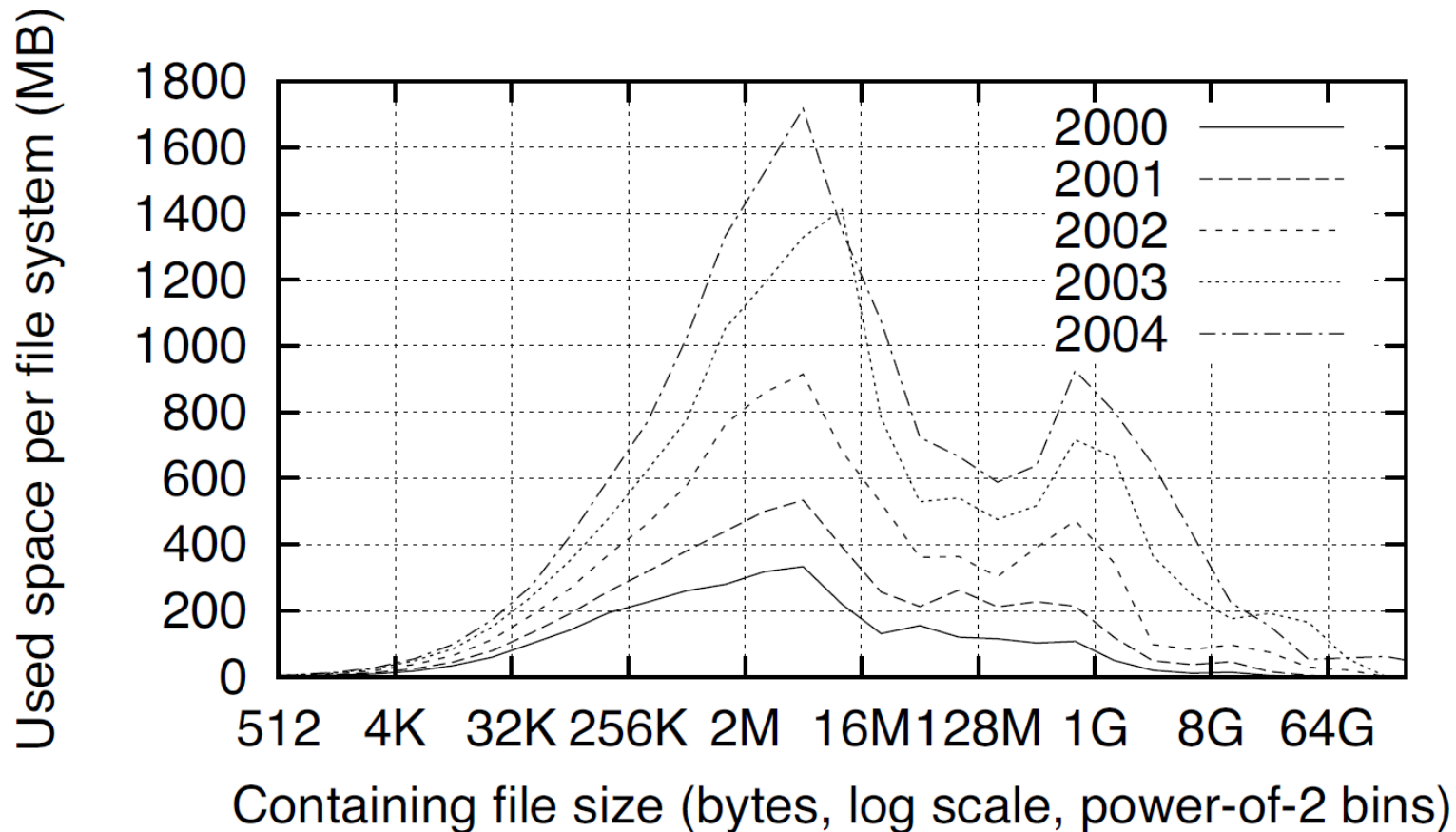- Most of the space is occupied by the rare big ones

Figure 2: Histograms of files by size

# Characteristics of Files

- Most files are small
- Most of the space is occupied by the rare big ones

Figure 4: Histograms of bytes by containing file size

18

# Meet the Inode

# Unix Fast File System

- Original inode format appeared in BSD 4.1 (1981)
  - *Berkeley Standard Distribution Unix*
  - Similar structure for Linux ext2/3
- File Number is index into inode arrays
- Multi-level index structure
  - Great for small and large files
  - Asymmetric tree with fixed sized blocks
- Metadata associated with the file
  - Rather than in the directory that points to it (FAT)
- UNIX FFS: BSD 4.2 (1983): Locality Heuristics
  - Block group placement
  - Reserve space
- Scalable directory structure

# File Attributes

Inode Array

| Inode | | Triple Indirect Blocks | Double Indirect Blocks | Indirect Blocks | Data Blocks |
|---|---|---|---|---|---|
| **Inode** | | | | | |
| File Metadata | | | | | |

User

Group

9 basic access control bits

- UGO × RWX

Setuid bit

• Execute at owner permissions, rather than user

Getgid bit

• Execute at group's permissions

12

Indirect Pointer

Double Indirect Pointer

Triple Indirect Pointer

# Data Storage

Inode Array

**Small files: 12 pointers direct to data blocks**

**Direct pointers: 4KB blocks ➔ Enough for files up to 48KB**

| Inode |
|:---:|
| File Metadata |
| 1 |
| 2 |
| 3 |
| 4 |
| 5 |
| 6 |
| 7 |
| 8 |
| 9 |
| 10 |
| 11 |
| 12 |
| Indirect Pointer |
| Double Indirect Pointer |
| Triple Indirect Pointer |

Triple Indirect Blocks    Double Indirect Blocks    Indirect Blocks    Data Blocks

Files per file system

12000
10000
8000
6000
4000
2000
0

| | 2000 |
| | 2001 |
| | 2002 |
| | 2003 |
| | 2004 |

0    8    128    2K    32K    512K    8M    128M

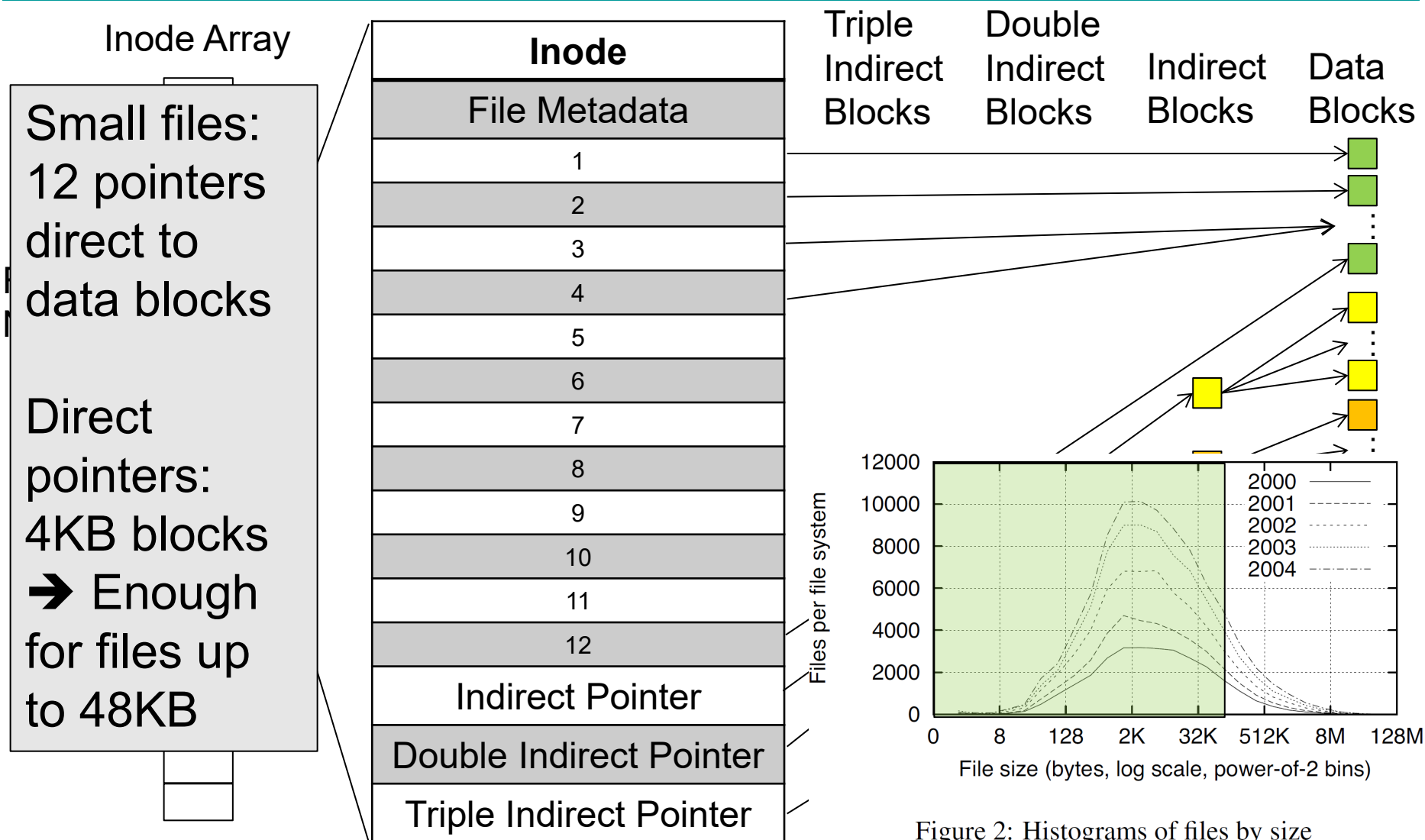File size (bytes, log scale, power-of-2 bins)

Figure 2: Histograms of files by size

# Data Storage

**Indirect pointers**
- Point to a disk block containing only pointers
- 4KB blocks
  - ➔1024 pointers
    - 4 MB at Level 2
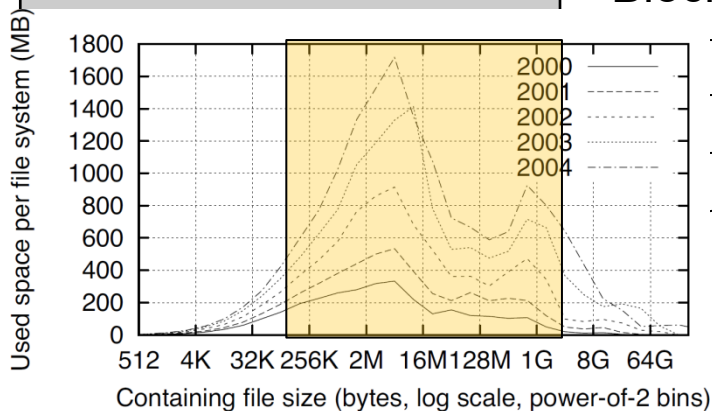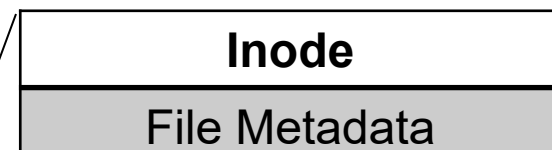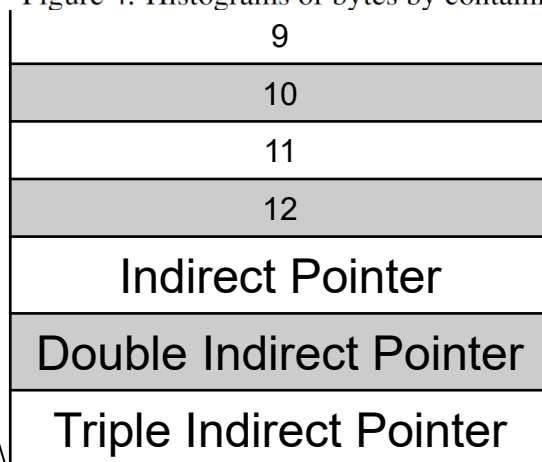    - 4GB at Level 3
    - 4TB at Level 4



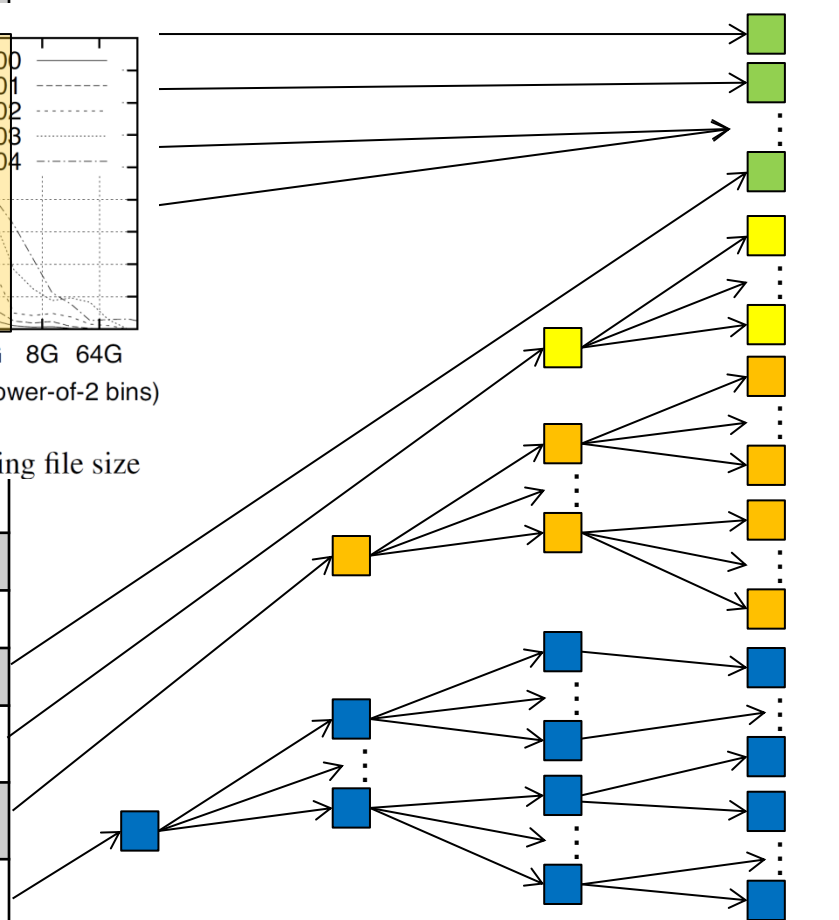Figure 4: Histograms of bytes by containing file size

| Inode |
|---|
| File Metadata |
| 9 |
| 10 |
| 11 |
| 12 |
| Indirect Pointer |
| Double Indirect Pointer |
| Triple Indirect Pointer |

Triple Indirect Blocks   Double Indirect Blocks   Indirect Blocks   Data Blocks

# So Far

- File Systems
  - Introduction to File Systems
  - Very simply file system
  - FAT
  - Inodes
  - Unix Fast File System (FFS)

# UNIX BSD 4.2 (1983)

- Same as BSD 4.1 (same file header and triply indirect blocks), except incorporated ideas from Cray DEMOS:
  - Uses bitmap allocation in place of freelist
  - Attempt to allocate files contiguously
  - 10% reserved disk space
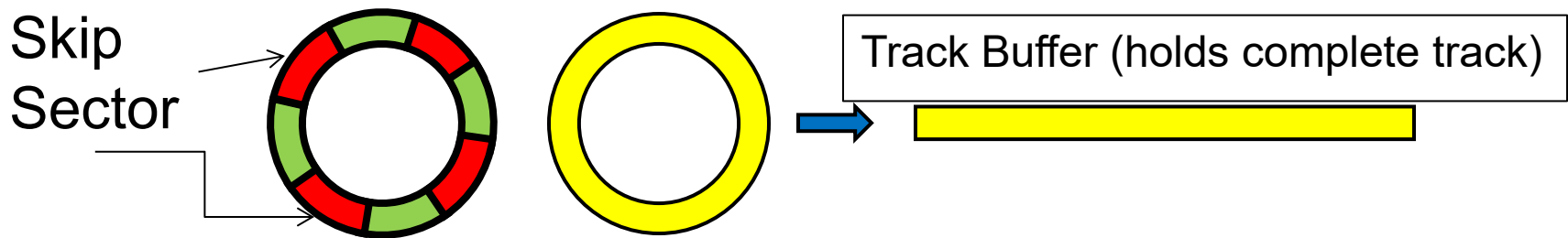  - Skip-sector positioning (soon)

# Problem 1: How big?

- When create a file, don't know how big it will become (in UNIX, most writes are by appending)
  - How much contiguous space do you allocate for a file?
  - In BSD 4.2, just find some range of free blocks
    - Put each new file at the front of different range
    - To expand a file, you first try successive blocks in bitmap, then choose new range of blocks
  - Also in BSD 4.2: store files from same directory near each other

- Fast File System (FFS)
  - Allocation and placement policies for BSD 4.2
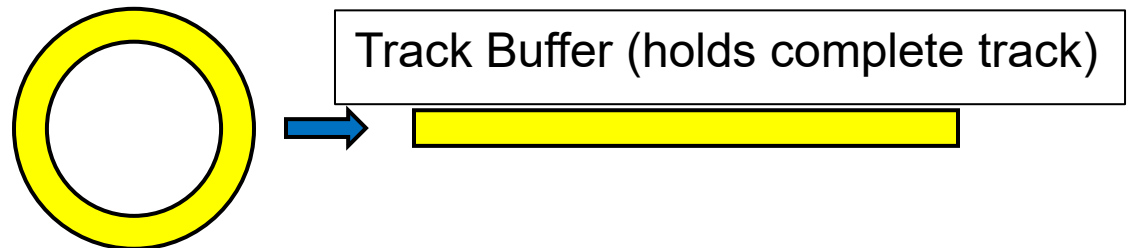
# Problem 2: Rotational Delay

- **Missing blocks due to rotational delay**
  - Issue: Read one block, do processing, and read next block.  In meantime, disk has continued turning: missed next block! Need 1 revolution/block!

Skip Sector

Track Buffer (holds complete track)

- **Solution 1**: Skip sector positioning ("interleaving")
  - Place the blocks from one file on every other block of a track: give time for processing to overlap rotation
- **Solution 2**: Read ahead: Read next block right after first, even if application hasn't asked for it yet.
  - This can be done by the OS (read ahead)  - OR -
  - By the disk itself (track buffers). Many disk controllers have internal RAM that allows them to read a complete track

# Problem 2: Rotational Delay

- Important Aside: Modern disks and controllers do many complex things "under the covers"

  - Track buffers, elevator algorithms, bad block filtering

# Conclusion

- **File Systems**
  - Introduction to File Systems
  - Very simply file system
  - FAT
  - Inodes
  - Unix Fast File System (FFS)