## CS162
## Operating Systems and Systems Programming
## Lecture 18

## File Systems

April 4th, 2016
Prof. Anthony D. Joseph
http://cs162.eecs.Berkeley.edu

---

## Quick Aside: Big Projects

- What is a big project?
  - Time/work estimation is hard
  - Programmers are eternal optimists (it will only take two days)!
    - » This is why we bug you about starting the project early
    - » Had a grad student who used to say he just needed "10 minutes" to fix something. Two hours later…
- Can a project be efficiently partitioned?
  - Partitionable task decreases in time as you add people
  - But, if you require communication:
    - » Time reaches a minimum bound
    - » With complex interactions, time increases!
  - Mythical person-month problem:
    - » You estimate how long a project will take
    - » Starts to fall behind, so you add more people
    - » Project takes even more time!

---

## Techniques for Partitioning Tasks

- Functional
  - Person A implements threads, Person B implements semaphores, Person C implements locks…
  - Problem: Lots of communication across APIs
    - » If B changes the API, A may need to make changes
    - » Story: Large airline company spent $200 million on a new scheduling and booking system. Two teams "working together." After two years, went to merge software. Failed! Interfaces had changed (documented, but no one noticed). Result: would cost another $200 million to fix.

- Task
  - Person A designs, Person B writes code, Person C tests
  - May be difficult to find right balance, but can focus on each person's strengths (Theory vs systems hacker)
  - Since Debugging is hard, Microsoft has *two* testers for *each* programmer

- Most CS162 project teams are functional, but people have had success with task-based divisions
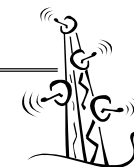
---

## Communication

- More people mean more communication
  - Changes have to be propagated to more people
  - Think about person writing code for most fundamental component of system: everyone depends on them!
  - *You should be meeting in person at least twice/week!*

- Miscommunication is common
  - "Index starts at 0? I thought you said 1!"

- Who makes decisions?
  - Individual decisions are fast but trouble
  - Group decisions take time
  - Centralized decisions require a big picture view (someone who can be the "system architect")

- Often designating someone as system architect can be a good thing
  - Better not be clueless
  - Better have good people skills
  - Better let other people do work

## Coordination

- More people ⇒ no one can make all meetings!
  - They miss decisions and associated discussion
  - Example from earlier class: one person missed meetings and did something group had rejected

- People have different work styles
  - Some people work in the morning, some at night
  - How do you decide when to meet or work together?

- What about project slippage?
  - It will happen, guaranteed!
  - Example: phase 4 of one project, everyone busy but not talking. One person way behind. No one knew until very end – too late!

- Hard to add people to existing group
  - Members have already figured out how to work together

## Recall: Device Drivers

- **Device Driver:** Device-specific code in the kernel that interacts directly with the device hardware
  - Supports a standard, internal interface
  - Same kernel I/O system can interact easily with different device drivers
  - Special device-specific configuration supported with the `ioctl()` system call

- Device Drivers typically divided into two pieces:
  - Top half: accessed in call path from system calls
    - » implements a set of standard, cross-device calls like `open()`, `close()`, `read()`, `write()`, `ioctl()`, `strategy()`
    - » This is the kernel's interface to the device driver
    - » Top half will *start* I/O to device, may put thread to sleep until finished
  - Bottom half: run as interrupt routine
    - » Gets input or transfers next block of output
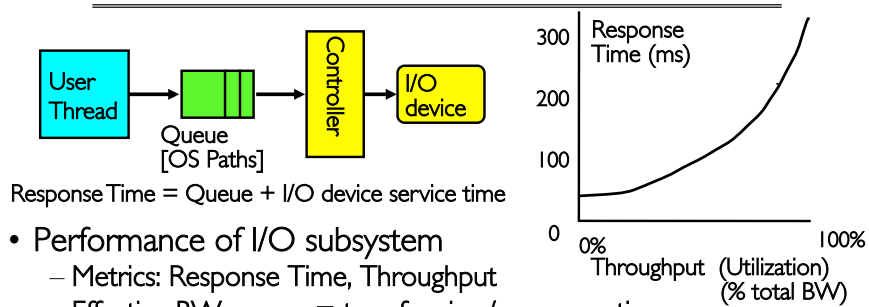    - » May wake sleeping threads if I/O now complete

## Kernel vs User-level I/O

- Both are popular and practical for different reasons:

- **Kernel-level drivers** for critical devices that must keep running (e.g., display drivers)
  - Programming is a major effort, correct operation of the rest of the kernel depends on correct driver operation

- **User-level drivers** for devices that are non-threatening (e.g., USB devices in Linux: `libusb`)
  - Provide higher-level primitives to the programmer, avoid every driver doing low-level I/O register tweaking
  - The multitude of USB devices can be supported by Less-Than-Wizard programmers
  - New drivers don't have to be compiled for each version of the OS, and loaded into the kernel

## Kernel vs User-level Programming Styles

- **Kernel-level drivers**
  - Have a much more limited set of resources available:
    - » Only a fraction of `libc` routines typically available
    - » Memory allocation (e.g., Linux `kmalloc`) much more limited in capacity and required to be physically contiguous
    - » Should avoid blocking calls
    - » Can use asynchrony with other kernel functions but tricky with user code

- **User-level drivers**
  - Similar to other application programs but:
    - » Will be called often – should do its work fast, or postpone it – or do it in the background
    - » Can use threads, blocking operations (usually much simpler) or non-blocking or asynchronous

## Recall: I/O Performance



Response Time = Queue + I/O device service time

- **Performance of I/O subsystem**
  - Metrics: Response Time, Throughput
  - Effective BW per op = transfer size / response time
    » EffectiveBW(n) = n / (S + n/B) = B / (1 + SB/n )
  - Contributing factors to latency:
    » Software paths (can be loosely modeled by a queue)
    » Hardware controller
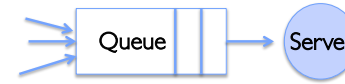    » I/O device service time
- **Queuing behavior:**
  - Can lead to big increases in latency as utilization increases
  - Solutions?

---

## Performance: Multiple Outstanding Requests



- Suppose each read takes 10 ms to service

- If a process works for 100 ms after each read, what is the utilization of the disk?
  - U = 10 ms / 110ms = 9%

- What it there are two such processes?
  - U = (10 ms + 10 ms) / 120ms ≈17%

- What if each of those processes have two such threads?

---

## Recall: How do we Hide I/O Latency?

- **Blocking Interface:** "Wait"
  - When request data (e.g., read() system call), put process to sleep until data is ready
  - When write data (e.g., write() system call), put process to sleep until device is ready for data

- **Non-blocking Interface:** "Don't Wait"
  - Returns quickly from read or write request with count of bytes successfully transferred to kernel
  - Read may return nothing, write may write nothing

- **Asynchronous Interface:** "Tell Me Later"
  - When requesting data, take pointer to user's buffer, return immediately; later kernel fills buffer and notifies user
  - When sending data, take pointer to user's buffer, return immediately; later kernel takes data and notifies user

---

## Administrivia

- HW4 – Releases on Monday 4/11 (Due 4/25)

- Project 2 code due next Friday

- Midterm II: Coming up in 2.5 weeks! (4/20)
  - 6-7:30PM (aa-eh 10 Evans, ej-oa 155 Dwinelle)
  - Covers lectures #13 to 21
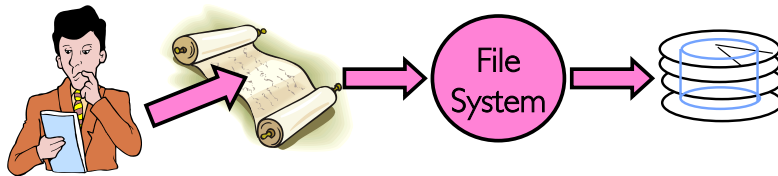  - 1 page of hand-written notes, both sides

## BREAK

## Recall: 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!
  - System's view (inside OS):
    » Collection of blocks (a block is a logical transfer unit, while a sector is the physical transfer unit)
    » Block size ≥ sector size; in UNIX, block size is 4KB

## Recall: 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()` ⇒ 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

- (Hard) Disks 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

- Size is determined as they are used !!!
  - Can write (or read zeros) to expand the file
  - Start small and grow, need to make room

- Organized into directories
  - What data structure (on disk) for that?

- Need to allocate / free blocks
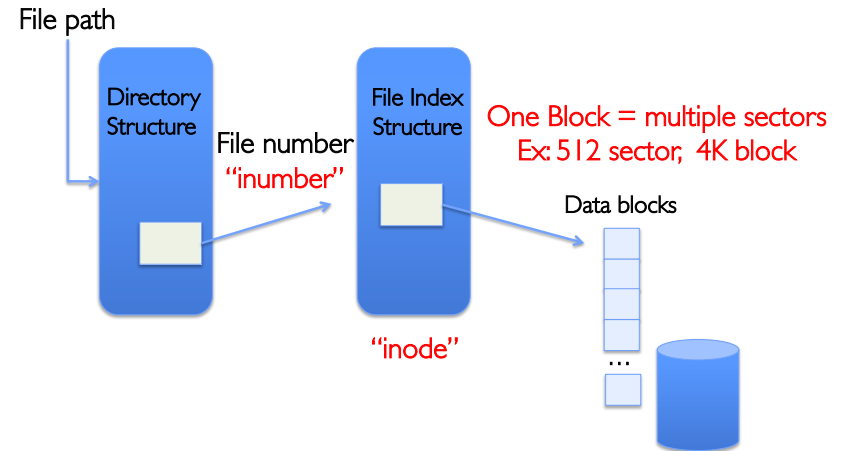  - Such that access remains efficient

## Recall: 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 (next lecture)

- Access disk as linear array of sectors. Two Options:
  - Identify sectors as vectors [cylinder, surface, sector], sort in cylinder-major order, not used anymore
  - 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

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

- Need 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 "inumber"

File Index Structure

One Block = multiple sectors
Ex: 512 sector, 4K block

Data blocks

"inode"

---

## Components of a file system

*file name* → *file number* → *Storage block*

*offset*    directory    *offset*    index structure

- 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 integer) to user process

- Read, Write, Seek, and Sync operate on handle

  - Mapped to descriptor and to blocks

---

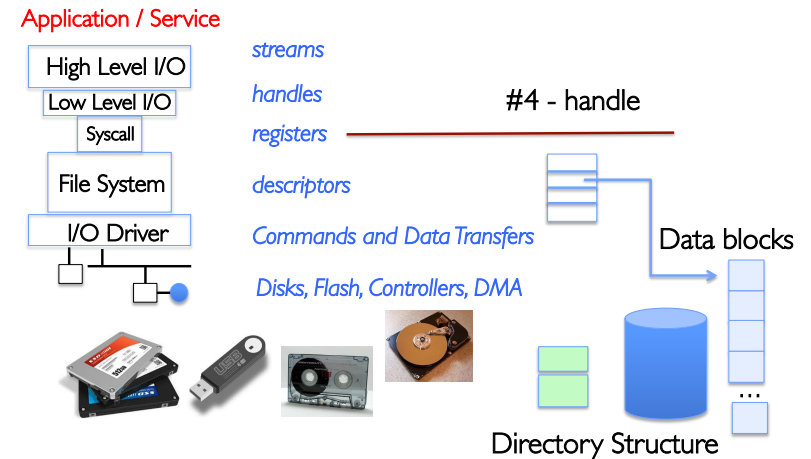## Directories

## 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
  - Softlinks (aliases) are another name for an entry

## I/O & Storage Layers

Application / Service

| High Level I/O | *streams* |
| Low Level I/O | *handles* |
| Syscall | *registers* |
| File System | *descriptors* |
| I/O Driver | *Commands and Data Transfers* |

*#4 - handle*

*Disks, Flash, Controllers, DMA*

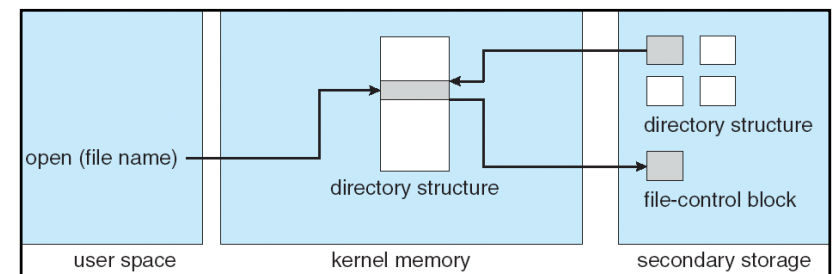Data blocks

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

## In-Memory File System Structures

open (file name)

directory structure

directory structure

file-control block

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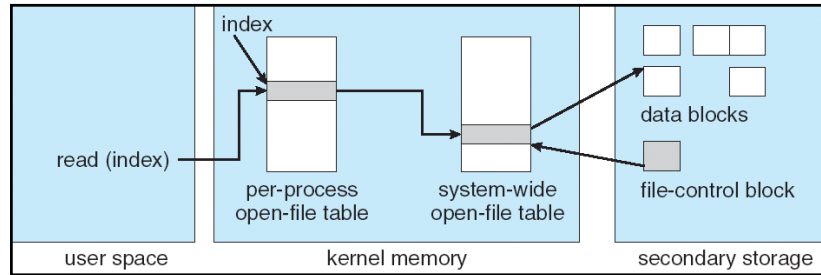
- Open system call:
  - Resolves file name, finds file control block (inode)
  - Makes entries in per-process and system-wide tables
  - Returns index (called "file handle") in open-file table

## In-Memory File System Structures



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

---

# BREAK

---

## Our first filesystem: FAT (File Allocation Table)

- The most commonly used filesystem in the world!

- Assume (for now) 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

- Example: file_read 31, < 2, x >
  - Index into FAT with file number
  - Follow linked list to block
  - Read the block from disk into memory
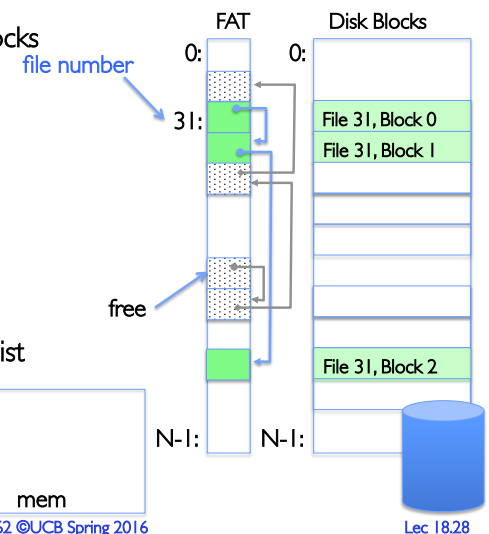
---

## 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

## FAT Properties

- File is collection of disk blocks

- FAT is linked list 1-1 with blocks

  file number

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

- File offset (o = B:x )

- Follow list to get block #

  free

- Unused blocks ⇔ FAT free list

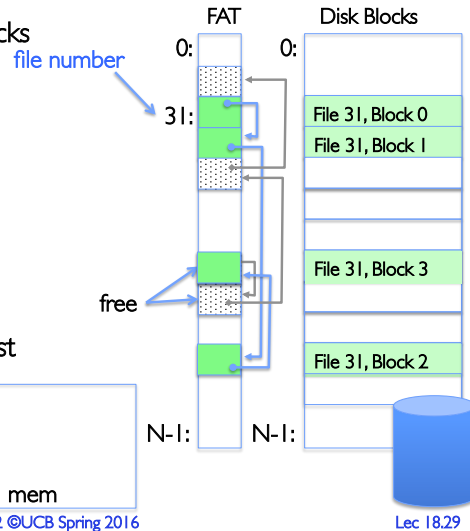- Ex: file_write(51, <3, y> )
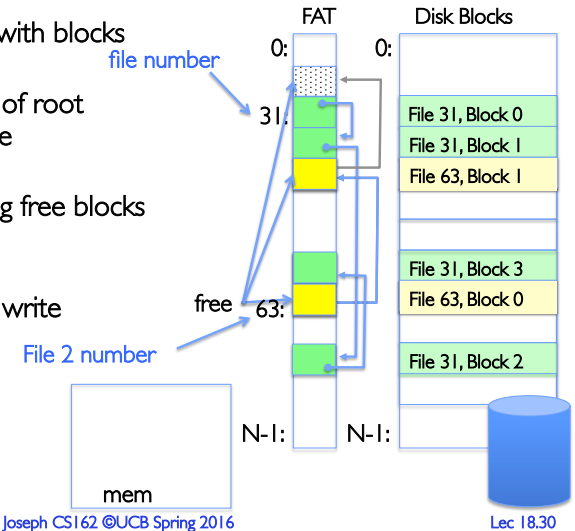  - Grab blocks from free list
  - Linking them into file

FAT    Disk Blocks

0:     0:

31:    File 31, Block 0
       File 31, Block 1

       File 31, Block 3

       File 31, Block 2

N-1:   N-1:

mem

## FAT Properties

- File is collection of disk blocks

- FAT is linked list 1-1 with blocks

  file number

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

- Grow file by allocating free blocks and linking them in
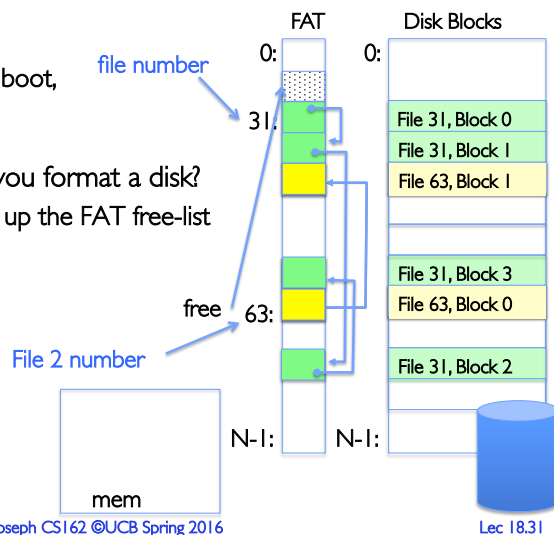
- Ex: Create file, write, write

  free   63:

  File 2 number

FAT    Disk Blocks

0:     0:

31:    File 31, Block 0
       File 31, Block 1
       File 63, Block 1

       File 31, Block 3
       File 63, Block 0

       File 31, Block 2

N-1:   N-1:

mem

## FAT Assessment

- *FAT32 (32 instead of 12 bits) used in Windows, USB drives, SD cards, …*

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

  file number

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

- Simple

  free   63:

  File 2 number

FAT    Disk Blocks

0:     0:

31:    File 31, Block 0
       File 31, Block 1
       File 63, Block 1

       File 31, Block 3
       File 63, Block 0

       File 31, Block 2

N-1:   N-1:

mem

## FAT Assessment

- Time to find block (large files) ??

- Block layout for file ???

  file number

- Sequential Access ???

- Random Access ???

- Fragmentation ???

  free   63:

- Small files ???

  File 2 number

- Big files ???

FAT    Disk Blocks

0:     0:

31:    File 31, Block 0
       File 31, Block 1
       File 63, Block 1

       File 31, Block 3
       File 63, Block 0

       File 31, Block 2

N-1:   N-1:

mem

## What about the Directory?

file 5268830
"/home/tom"

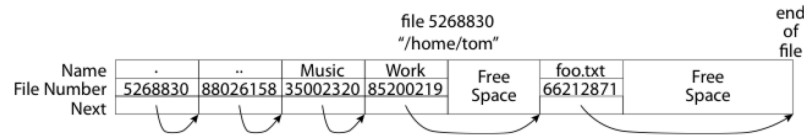| Name | . | .. | Music | Work | Free Space | foo.txt | Free Space | end of file |
| File Number | 5268830 | 88026158 | 35002320 | 85200219 | | 66212871 | | |
| Next | | | | | | | | |

- 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 (cont'd)

- 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")

---

## Many Huge 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)

---

## Characteristics of Files

- Most files are small

A Five-Year Study of File-System Metadata

NITIN AGRAWAL
University of Wisconsin, Madison
and
WILLIAM J. BOLOSKY, JOHN R. DOUCEUR, and JACOB R. LORCH
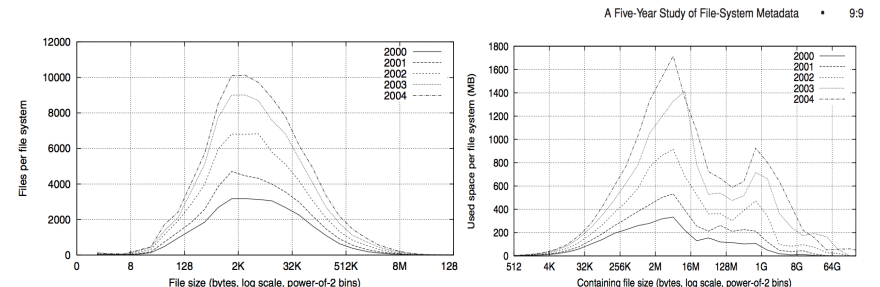Microsoft Research

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

A Five-Year Study of File-System Metadata    •    9:9



Fig. 2. Histograms of files by size.

Fig. 4. Histograms of bytes by containing file size.

## So What About a "Real" File System?

• **Meet the inode:**



Inode Array · Triple Indirect Blocks · Double Indirect Blocks · Indirect Blocks · Data Blocks · Inode · file_number · File Metadata · Direct Pointers · Indirect Pointer · Dbl. Indirect Ptr. · Tripl. Indirect Ptr.

---

## Summary

• File System:
  – Transforms blocks into Files and Directories
  – Optimize for access and usage patterns
  – Maximize sequential access, allow efficient random access

• File (and directory) defined by header, called "inode"

• File Allocation Table (FAT) Scheme
  – Linked-list approach
  – Very widely used: Cameras, USB drives, SD cards
  – Simple to implement, but poor performance

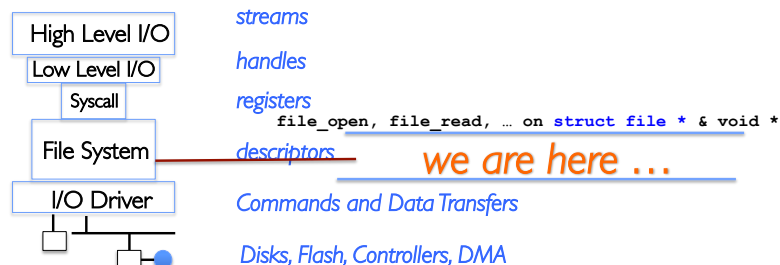---

## I/O & Storage Layers

*Operations, Entities and Interface*

Application / Service



High Level I/O — *streams*
Low Level I/O — *handles*
Syscall — *registers*
File System — *descriptors*  `file_open, file_read, … on struct file * & void *`

we are here …

I/O Driver — *Commands and Data Transfers*

*Disks, Flash, Controllers, DMA*

---

## Recall: C Low level I/O

• Operations on File Descriptors – as OS object representing the state of a file
  – User has a "handle" on the descriptor

```
#include <fcntl.h>
#include <unistd.h>
#include <sys/types.h>

int open (const char *filename, int flags [, mode_t mode])
int creat (const char *filename, mode_t mode)
int close (int filedes)
```

Bit vector of:
• Access modes (Rd, Wr, …)
• Open Flags (Create, …)
• Operating modes (Appends, …)

Bit vector of Permission Bits:
• User|Group|Other X R|W|X

http://www.gnu.org/software/libc/manual/html_node/Opening-and-Closing-Files.html

# Recall: C Low Level Operations

```
ssize_t read (int filedes, void *buffer, size_t maxsize)
 - returns bytes read, 0 => EOF, -1 => error
ssize_t write (int filedes, const void *buffer, size_t size)
 - returns bytes written

off_t lseek (int filedes, off_t offset, int whence)

int fsync (int fildes) — wait for i/o to finish
void sync (void) — wait for ALL to finish
```

- When write returns, data is on its way to disk and can be read, but it may not actually be permanent!