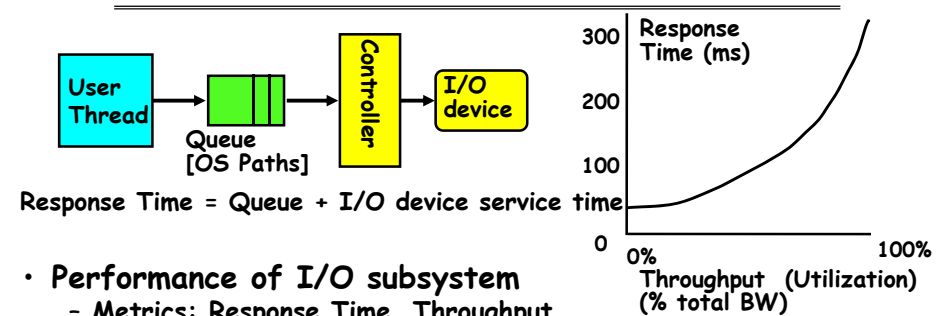


CS162 Operating Systems and Systems Programming Lecture 18

File Systems

April 6th, 2015
Prof. John Kubiatowicz
<http://cs162.eecs.Berkeley.edu>

Recall: I/O Performance



- Performance of I/O subsystem
 - Metrics: Response Time, Throughput
 - Effective BW per op = transfer size / response time
 - » $\text{EffBW}(n) = n / (S + n/B) = B / (1 + SB/n)$
 - 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 of latency as utilization increases
 - Solutions?

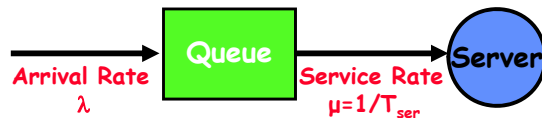
4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.2

Recall: A Little Queuing Theory: Some Results

- Assumptions:
 - System in equilibrium; No limit to the queue
 - Time between successive arrivals is random and memoryless



- Parameters that describe our system:
 - λ : mean number of arriving customers/second
 - T_{ser} : mean time to service a customer ("m1")
 - C : squared coefficient of variance = $\sigma^2/m1^2$
 - μ : service rate = $1/T_{ser}$
 - u : server utilization ($0 \leq u \leq 1$): $u = \lambda/\mu = \lambda \times T_{ser}$
- Parameters we wish to compute:
 - T_q : Time spent in queue
 - L_q : Length of queue = $\lambda \times T_q$ (by Little's law)
- Results:
 - Memoryless service distribution ($C = 1$):
 - » Called M/M/1 queue: $T_q = T_{ser} \times u/(1 - u)$
 - General service distribution (no restrictions), 1 server:
 - » Called M/G/1 queue: $T_q = T_{ser} \times \frac{1}{2}(1+C) \times u/(1 - u)$

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.3

When is the disk performance highest?

- When there are big sequential reads, or
- When there is so much work to do that they can be piggy backed (reordering queues—one moment)
- OK, to be inefficient when things are mostly idle
- Bursts are both a threat and an opportunity
- <your idea for optimization goes here>
 - Waste space for speed?
- Other techniques:
 - Reduce overhead through user level drivers
 - Reduce the impact of I/O delays by doing other useful work in the meantime

4/6/15

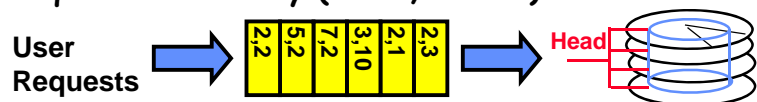
Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.4

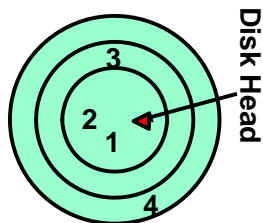
Disk Scheduling

- Disk can do only one request at a time; What order do you choose to do queued requests?

- Request denoted by (track, sector)



- Scheduling algorithms:
 - First In First Out (FIFO)
 - Shortest Seek Time First
 - SCAN
 - C-SCAN



- In our examples we ignore the sector
 - Consider only track #

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.5

FIFO: First In First Out

- Schedule requests in the order they arrive in the queue

- Example:

- Request queue:

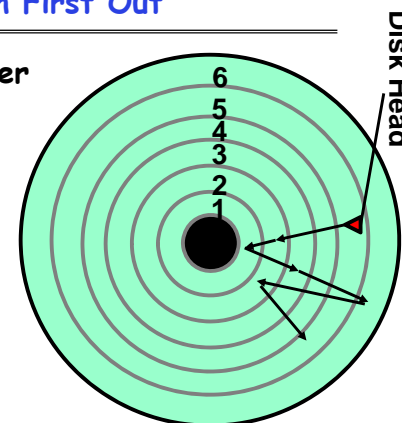
2, 1, 3, 6, 2, 5

- Scheduling order:

2, 1, 3, 6, 2, 5

- 16 tracks, 6 seeks

- Pros: Fair among requesters
- Cons: Order of arrival may be to random spots on the disk ⇒ Very long seeks



4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.6

SSTF: Shortest Seek Time First

- Pick the request that's closest to the head on the disk

- Although called SSTF, include rotational delay in calculation, as rotation can be as long as seek

- Example:

- Request queue:

2, 1, 3, 6, 2, 5

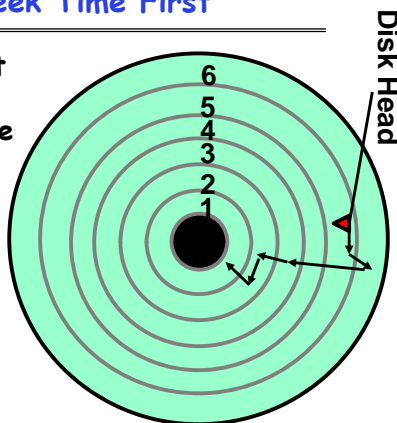
- Scheduling order:

5, 6, 3, 2, 2, 1

- 6 tracks, 4 seeks

- Pros: reduce seeks

- Cons: may lead to starvation
 - Greedy. Not optimal



4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.7

SCAN

- Implements an Elevator Algorithm: take the closest request in the direction of travel

- Example:

- Request queue:

2, 1, 3, 6, 2, 5

- Head is moving towards center

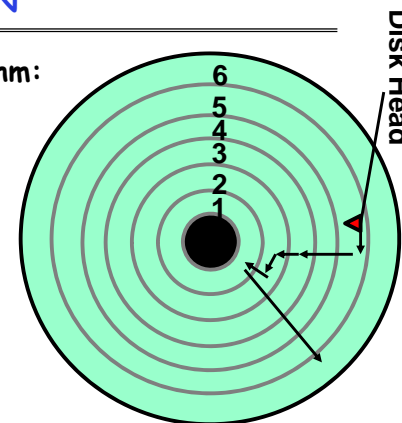
- Scheduling order:

5, 3, 2, 2, 1, 6

- 8 tracks, 4 seeks

- Pros:
 - No starvation
 - Low seek

- Cons: favors middle tracks
 - May spend time on sparse tracks while dense requests elsewhere



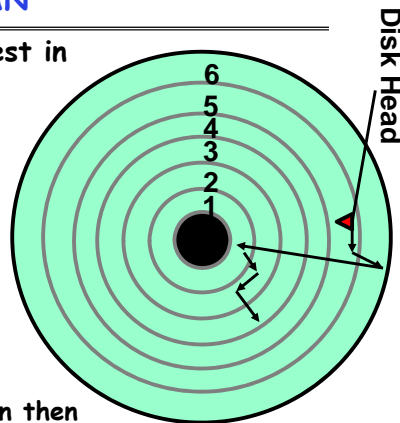
4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.8

C-SCAN

- Like SCAN but only serves request in only one direction
- Example:
 - Request queue: 2, 1, 3, 6, 2, 5
 - Head only serves request on its way from center towards edge
 - Scheduling order: 5, 6, 1, 2, 2, 3
 - 8 tracks, 5 seeks
- Pros:
 - Fairer than SCAN
 - Accumulate work in remote region then go get it
- Cons: longer seeks on the way back
- Optimization: dither to pickup nearby requests as you go



4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.9

Review: 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

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.10

Kernel vs User-level I/O

- Both are popular/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.

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.11

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.

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.12

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?
 - $U = 10 \text{ ms} / 110 \text{ ms} = 9\%$
- What if there are two such processes?
 - $U = (10 \text{ ms} + 10 \text{ ms}) / 110 \text{ ms} = 18\%$
- What if each of those processes has two such threads?

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.13

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

4/6/15

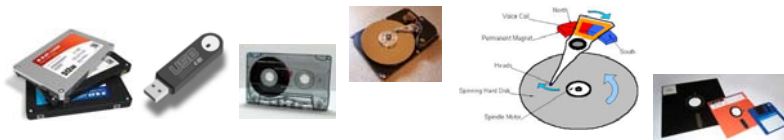
Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.14

I/O & Storage Layers

Operations, Entities and Interface

Application / Service



4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.15

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

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.16

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 filedes) - 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!

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.17

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 \geq sector size; in UNIX, block size is 4KB

4/6/15

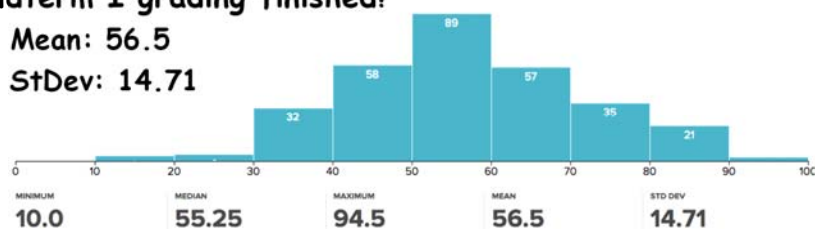
Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.18

Administrivia

- Midterm I grading finished!

- Mean: 56.5
- StDev: 14.71



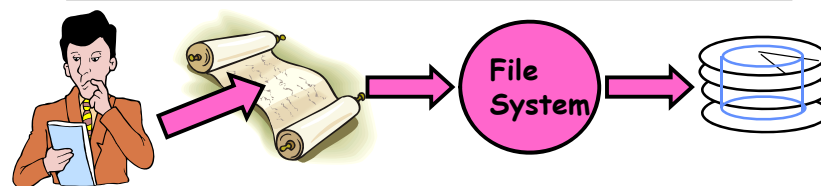
- Clearly this exam was harder than intended!
- Regrades:
 - You have until Thursday (4/9) to request a regrade
 - Be sure: If we receive a request, we may regrade whole exam (could lose points)
- Midterm II: coming up
 - April 22nd; More details upcoming

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.19

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

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.20

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
- 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

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.21

Disk Management Policies

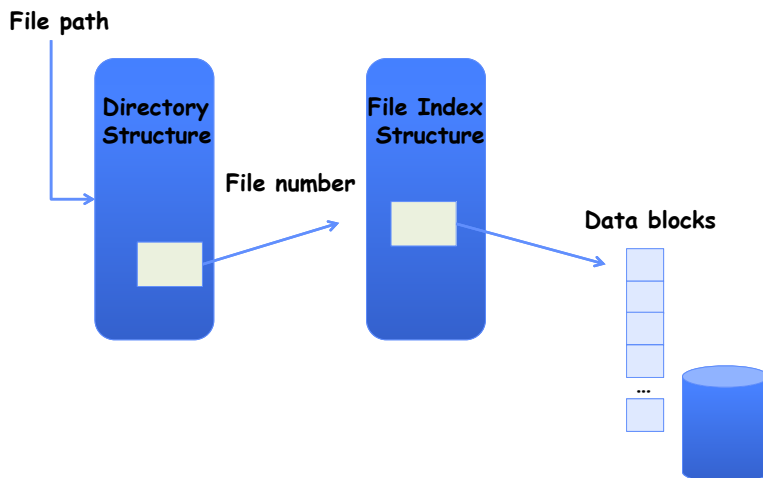
- Basic entities on a disk:
 - **File**: user-visible group of blocks arranged sequentially in logical space
 - **Directory**: user-visible index mapping names to files (next lecture)
- Access disk as linear array of sectors. Two Options:
 - Identify sectors as vectors [cylinder, surface, sector]. Sort in cylinder-major order. Not used much 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**

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.22

Components of a File System



4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.23

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

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.24

Directories

Name	Date Modified	Size	Kind
bse	Yesterday, 6:21 PM	--	Folder
Classes	Oct 13, 2014, 10:19 PM	--	Folder
AIT2008	Oct 13, 2014, 10:11 PM	--	Folder
CS-Scholars	Oct 13, 2014, 10:11 PM	--	Folder
cs61c1-f08	Oct 13, 2014, 10:17 PM	--	Folder
cs61c1-f09	Oct 13, 2014, 10:19 PM	--	Folder
cs162	Today, 8:36 AM	--	Folder
AndersonDahlin	Oct 13, 2014, 10:11 PM	--	Folder
fa14	Today, 8:36 AM	--	Folder
162prereqcheckSept8.xlsx	Sep 10, 2014, 3:20 PM	36 KB	Micros...kbook
coursecomparison.xlsx	Aug 6, 2014, 7:50 AM	31 KB	Micros...kbook
CS 162 apps.xlsx	Jun 29, 2014, 6:35 AM	53 KB	Micros...kbook
cs162git	Sep 23, 2014, 11:33 AM	--	Folder
devel	Oct 15, 2014, 11:40 AM	--	Folder
exams	Oct 13, 2014, 10:12 PM	--	Folder
gitprojects	Oct 8, 2014, 4:52 PM	--	Folder
group0	Today, 8:35 AM	--	Folder
pintos	Today, 8:35 AM	--	Folder
src	Today, 8:35 AM	--	Folder
gradesheet.xls	Sep 19, 2014, 4:48 PM	68 KB	Micros...kbook
CSI Section Coverage.xlsx	Aug 22, 2014, 1:29 PM	11 KB	Micros...kbook
Lectures	Today, 8:22 AM	--	Folder
pintos-notes.txt	Sep 14, 2014, 2:10 PM	1 KB	Plain Text
pintos.pdf	Jul 21, 2014, 10:17 AM	549 KB	PDF Document
roster-9-13.xls	Sep 13, 2014, 5:12 PM	83 KB	Micros...kbook
roster-9-19.xls	Sep 19, 2014, 4:39 PM	84 KB	Micros...kbook
staff.xlsx	Aug 6, 2014, 7:14 AM	34 KB	Micros...kbook
student	Oct 13, 2014, 10:12 PM	--	Folder
studentsExcelFile-10-20	Yesterday, 9:53 AM	84 KB	Micros...kbook
syllabus-fa14.xlsx	Sep 12, 2014, 10:00 AM	38 KB	Micros...kbook
tmp	Oct 13, 2014, 10:12 PM	--	Folder
pintos	Aug 8, 2014, 6:06 AM	--	Folder
sp14	May 14, 2014, 9:02 PM	--	Folder
cs194	Oct 13, 2014, 10:16 PM	--	Folder
cs262b	Aug 7, 2013, 7:55 AM	--	Folder

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.25

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

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.26

I/O & Storage Layers



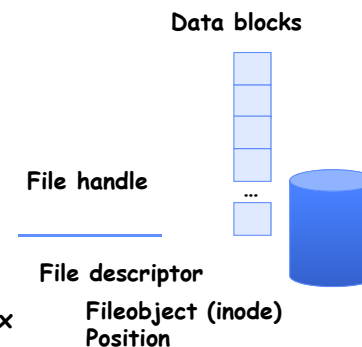
4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.27

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



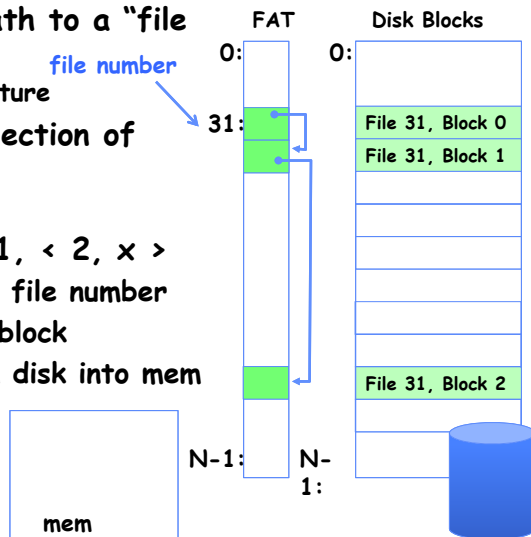
4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.28

FAT (File Allocation Table)

- 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 mem



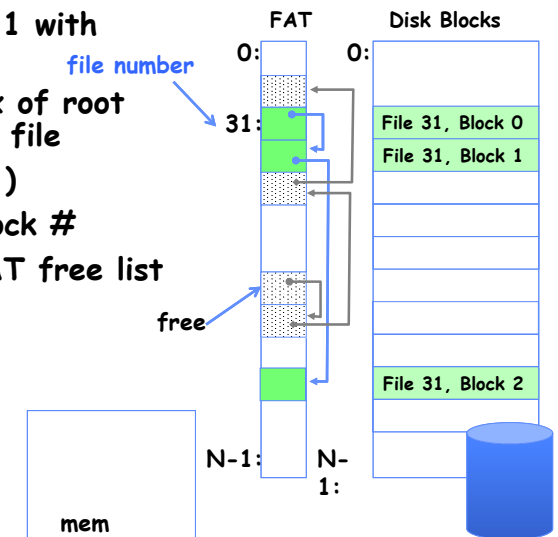
4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.29

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 \leftrightarrow FAT free list



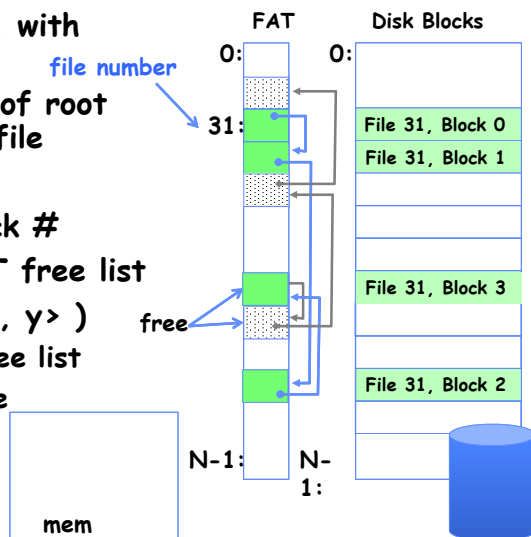
4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.30

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



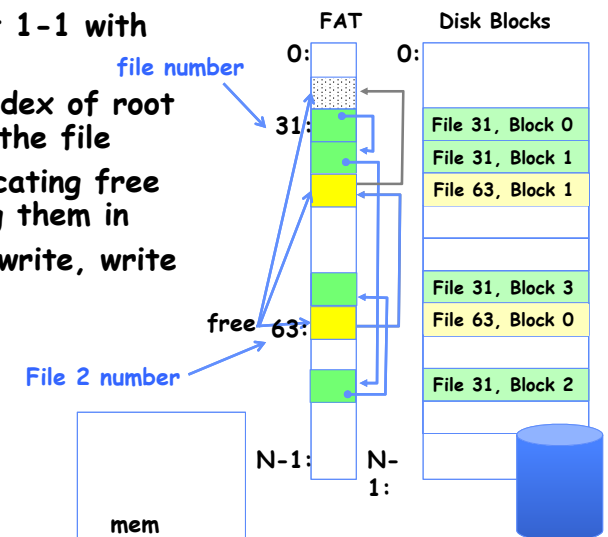
4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.31

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



4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.32

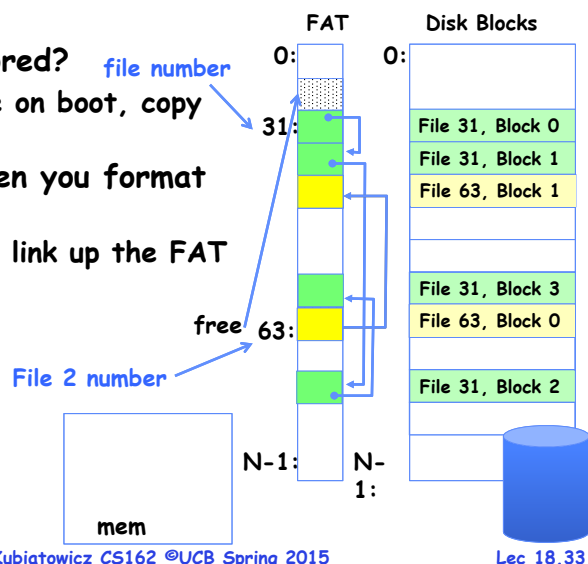
FAT Assessment

- Used in DOS, Windows, thumb 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



4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.33

FAT Assessment

- Time to find block (large files) ??

- Free list usually just a bit vector

- Next fit algorithm

- Block layout for file ???

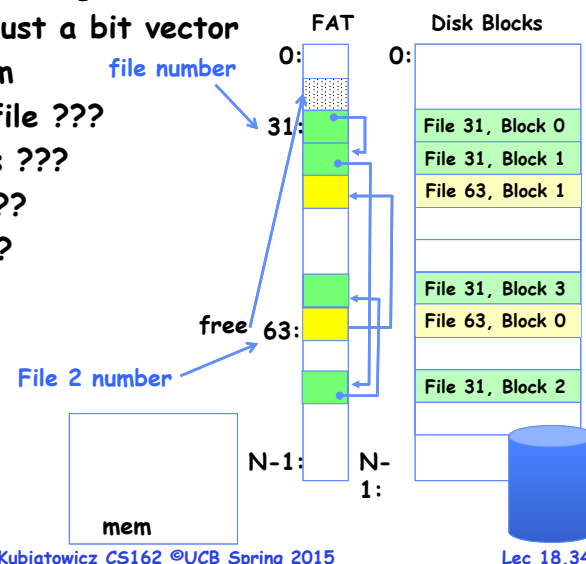
- Sequential Access ???

- Random Access ???

- Fragmentation ???

- Small files ???

- Big files ???

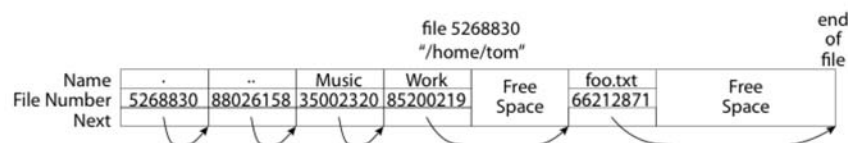


4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.34

What about the Directory?



- 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 ("/")?

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.35

Directory Structure (Con't)

- 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")

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.36

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)

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.37

Characteristics of Files

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

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

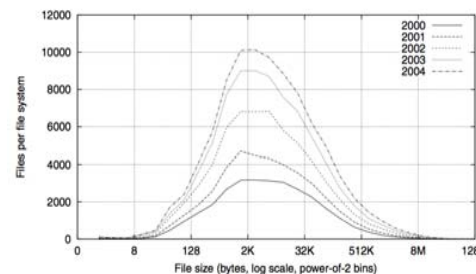


Fig. 2. Histograms of files by size.

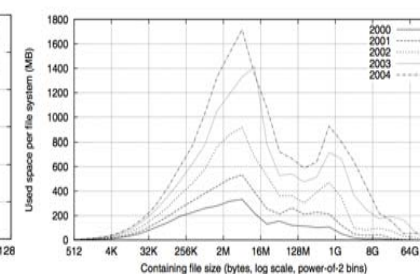


Fig. 4. Histograms of bytes by containing file size.

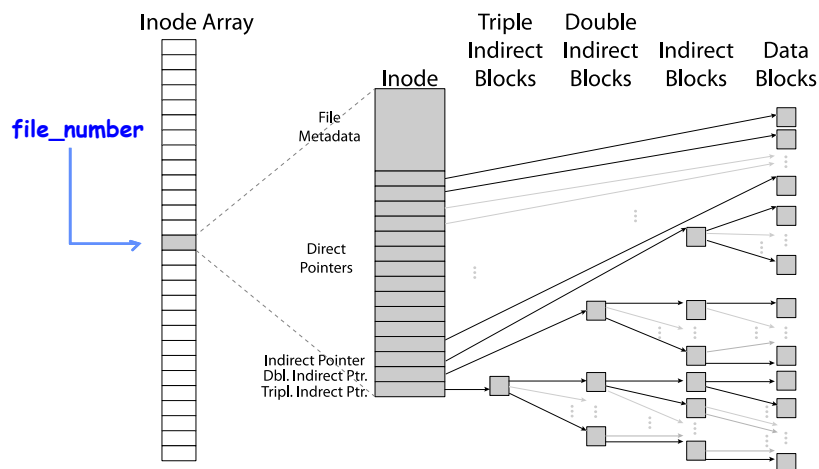
4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.38

So what about a "real" file system

- Meet the inode



4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.39

Unix Fast File System (Optimization on Unix Filesystem)

- Original inode format appeared in BSD 4.1
 - Berkeley Standard Distribution Unix
 - Part of your heritage!
- File Number is index into inode arrays
- Multi-level index structure
 - Great for little to large files
 - Asymmetric tree with fixed sized blocks
- Metadata associated with the file
 - Rather than in the directory that points to it
- UNIX FFS: BSD 4.2: Locality Heuristics
 - Block group placement
 - Reserve space
- Scalable directory structure

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.40

An "almost real" file system

- Pintos: src/filesys/file.c, inode.c

```

/* An open file. */
struct file
{
    struct inode *inode; /* File's inode. */
    off_t pos; /* Current position. */
    bool deny_write; /* Has file_deny_write() been called? */
};

/* In-memory inode. */
struct inode
{
    struct list_elem elem; /* Element in inode list. */
    block_sector_t sector; /* Sector number of disk location. */
    int open_cnt; /* Number of openers. */
    bool removed; /* True if deleted, false otherwise. */
    int deny_write_cnt; /* 0: writes ok, >0: deny writes. */
    struct inode_disk data; /* Inode content. */
};

/* On-disk inode.
   Must be exactly BLOCK_SECTOR_SIZE bytes long. */
struct inode_disk
{
    block_sector_t start; /* First data sector. */
    off_t length; /* File size in bytes. */
    unsigned magic; /* Magic number. */
    uint32_t unused[125]; /* Not used. */
};
    
```

Diagram showing the mapping of file metadata to disk blocks:

- File Array** (vertical stack of boxes) points to **File Metadata** (box).
- File Metadata** points to **Inode** (box).
- Inode** points to **Direct Blocks** (vertical stack of boxes).
- Inode** points to **Triple Indirect Blocks** (vertical stack of boxes).
- Inode** points to **Double Indirect Blocks** (vertical stack of boxes).
- Inode** points to **Single Indirect Blocks** (vertical stack of boxes).
- Inode** points to **Data Blocks** (vertical stack of boxes).

4/6/15

FFS: File Attributes

- Inode metadata

Inode Array (vertical stack of boxes) points to **File Metadata** (box).

File Metadata points to **Inode** (box).

Inode points to **Direct Blocks** (vertical stack of boxes).

Inode points to **Triple Indirect Blocks** (vertical stack of boxes).

Inode points to **Double Indirect Blocks** (vertical stack of boxes).

Inode points to **Single Indirect Blocks** (vertical stack of boxes).

Inode points to **Data Blocks** (vertical stack of boxes).

User Group (box) contains:

- 9 basic access control bits
 - UGO x RWX
- Setuid bit
 - execute at owner permissions
 - rather than user
- Getgid bit
 - execute at group's permissions

4/6/15

Kubiatowicz CS162 @UCB Spring 2015

Lec 18.42

FFS: Data Storage

- Small files: 12 pointers direct to data blocks

Direct pointers (box): 4kB blocks ⇒ sufficient For files up to 48KB

Inode (box) points to **File Metadata** (box).

File Metadata points to **Inode** (box).

Inode points to **Direct Pointers** (vertical stack of boxes).

Inode points to **Triple Indirect Blocks** (vertical stack of boxes).

Inode points to **Double Indirect Blocks** (vertical stack of boxes).

Inode points to **Single Indirect Blocks** (vertical stack of boxes).

Inode points to **Data Blocks** (vertical stack of boxes).

Indirect Pointer (box) contains: Dbl. Indirect Ptr., Tripl. Indirect Ptr.

Fig. 2. Histograms of files by size. (Graph showing Files per file system vs File size (bytes, log scale, power-of-2 bins) for years 2000-2004. The x-axis ranges from 0 to 128M, and the y-axis ranges from 0 to 12000. The distribution shows a peak around 32K-64K bytes.

4/6/15

Kubiatowicz CS162 @UCB Spring 2015

Fig. 2. Histograms of files by size.

FFS: Data Storage

- Large files: 1,2,3 level indirect pointers

Indirect pointers (box):

- point to a disk block containing only pointers
- 4 kB blocks => 1024 ptrs
- => 4 MB @ level 2
- => 4 GB @ level 3
- => 4 TB @ level 4

Inode (box) points to **File Metadata** (box).

File Metadata points to **Inode** (box).

Inode points to **Direct Pointers** (vertical stack of boxes).

Inode points to **Triple Indirect Blocks** (vertical stack of boxes).

Inode points to **Double Indirect Blocks** (vertical stack of boxes).

Inode points to **Single Indirect Blocks** (vertical stack of boxes).

Inode points to **Data Blocks** (vertical stack of boxes).

Pointers (box) contains: Direct Pointers, Indirect Ptr., Triple Indirect Ptr.

Fig. 4. Histograms of bytes by containing file size. (Graph showing Used space per file system (MB) vs Containing file size (bytes, log scale, power-of-2 bins) for years 2000-2004. The x-axis ranges from 0 to 64G, and the y-axis ranges from 0 to 1800. The distribution shows a peak around 4K-8K bytes.

4/6/15

Kubiatowicz CS162 @UCB Spring 2015

Lec 18.44

Freespace Management

- Bit vector with a bit per storage block
- Stored at a fixed location within the file system

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.45

Where are inodes stored?

- In early UNIX and DOS/Windows' FAT file system, headers stored in special array in outermost cylinders
 - Header not stored anywhere near the data blocks. To read a small file, seek to get header, seek back to data.
 - Fixed size, set when disk is formatted. At formatting time, a fixed number of inodes were created (They were each given a unique number, called an "inumber")

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.46

Where are inodes stored?

- Later versions of UNIX moved the header information to be closer to the data blocks
 - Often, inode for file stored in same "cylinder group" as parent directory of the file (makes an ls of that directory run fast).
 - Pros:
 - » UNIX BSD 4.2 puts a portion of the file header array on each of many cylinders. For small directories, can fit all data, file headers, etc. in same cylinder ⇒ no seeks!
 - » File headers much smaller than whole block (a few hundred bytes), so multiple headers fetched from disk at same time
 - » Reliability: whatever happens to the disk, you can find many of the files (even if directories disconnected)
 - Part of the Fast File System (FFS)
 - » General optimization to avoid seeks

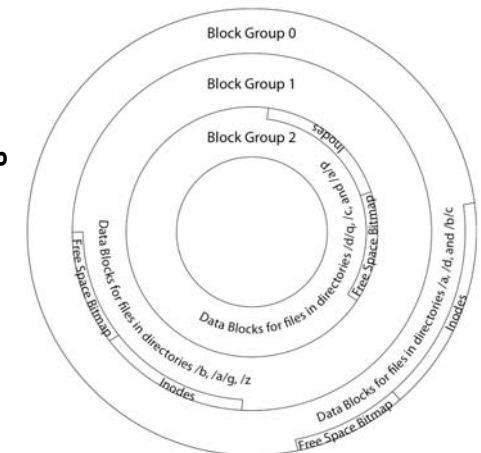
4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.47

Locality: Block Groups

- File system volume is divided into a set of block groups
 - Close set of tracks
- File data blocks, metadata, and free space are interleaved within block group
 - Avoid huge seeks between user data and system structure
- Put directory and its files in common block group
- First-Free allocation of new file block
 - Few little holes at start, big sequential runs at end of group
 - Avoids fragmentation
 - Sequential layout for big
- Reserve space in the BG

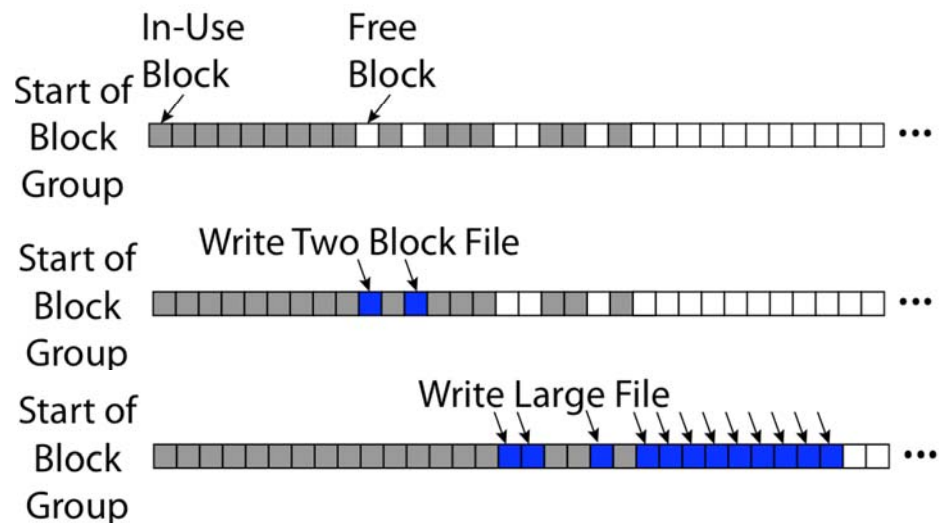


4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.48

FFS First Fit Block Allocation



4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.49

FFS

- **Pros**
 - Efficient storage for both small and large files
 - Locality for both small and large files
 - Locality for metadata and data
- **Cons**
 - Inefficient for tiny files (a 1 byte file requires both an inode and a data block)
 - Inefficient encoding when file is mostly contiguous on disk (no equivalent to superpages)
 - Need to reserve 10-20% of free space to prevent fragmentation

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.50

File System 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"
- **Multilevel Indexed Scheme**
 - Inode contains file info, direct pointers to blocks,
 - indirect blocks, doubly indirect, etc..
- **4.2 BSD Multilevel index files**
 - Inode contains pointers to actual blocks, indirect blocks, double indirect blocks, etc.
 - Optimizations for sequential access: start new files in open ranges of free blocks, rotational Optimization
- **Naming: act of translating from user-visible names to actual system resources**
 - Directories used for naming for local file systems

4/6/15

Kubiatowicz CS162 ©UCB Spring 2015

Lec 18.51