

CS162 Operating Systems and Systems Programming Lecture 18

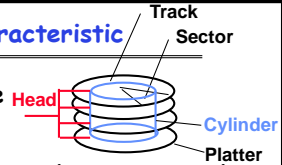
File Systems, Naming, and Directories

April 7, 2008

Prof. Anthony D. Joseph

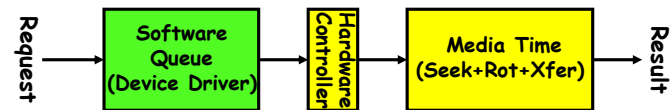
<http://inst.eecs.berkeley.edu/~cs162>

Review: Magnetic Disk Characteristic



- Cylinder: all the tracks under the head at a given point on all surface
- Read/write data is a three-stage process:
 - Seek time: position the head/arm over the proper track (into proper cylinder)
 - Rotational latency: wait for the desired sector to rotate under the read/write head
 - Transfer time: transfer a block of bits (sector) under the read-write head

• **Disk Latency = Queueing Time + Controller time + Seek Time + Rotation Time + Xfer Time**



- **Highest Bandwidth:**
 - transfer large group of blocks sequentially from one track

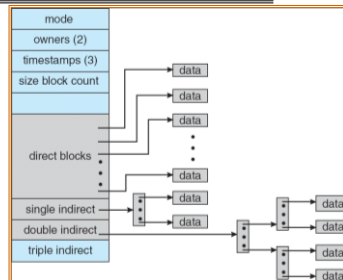
4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.2

Review: Multilevel Indexed Files (UNIX 4.1)

- **Multilevel Indexed Files:** Like multilevel address translation (from UNIX 4.1 BSD)
 - Key idea: efficient for small files, but still allow big files



- File hdr contains 13 pointers
 - Fixed size table, pointers not all equivalent
 - This header is called an "inode" in UNIX
- File Header format:
 - First 10 pointers are to data blocks
 - Ptr 11 points to "indirect block" containing 256 block ptrs
 - Pointer 12 points to "doubly indirect block" containing 256 indirect block ptrs for total of 64K blocks
 - Pointer 13 points to a triply indirect block (16M blocks)

4/7/08

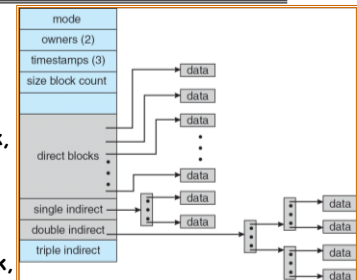
Joseph CS162 ©UCB Spring 2008

Lec 18.3

Review: Example of Multilevel Indexed Files

- **Sample file in multilevel indexed format:**

- How many accesses for block #23? (assume file header accessed on open)
 - » Two: One for indirect block, one for data
- How about block #5?
 - » One: One for data
- Block #340?
 - » Three: double indirect block, indirect block, and data



- UNIX 4.1 Pros and cons
 - Pros: Simple (more or less)
 - Files can easily expand (up to a point)
 - Small files particularly cheap and easy
 - Cons: Lots of seeks
 - Very large files must read many indirect blocks (four I/Os per block!)

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.4

Review: Multilevel Indexed Files (UNIX 4.1)

- Basic technique places an upper limit on file size that is approximately 16Gbytes
 - Designers thought this was bigger than anything anyone would need. Much bigger than a disk at the time...
 - Fallacy: today, EOS producing 2TB of data per day
- Pointers get filled in dynamically: need to allocate indirect block only when file grows > 10 blocks
 - On small files, no indirection needed

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.5

Goals for Today

- File Systems
 - Structure, Naming, Directories
- Caching in File Systems

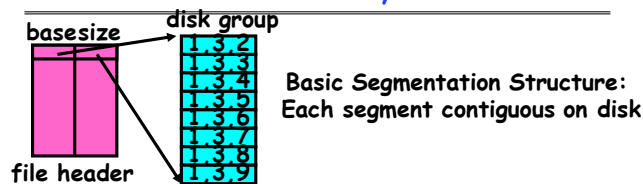
Note: Some slides and/or pictures in the following are adapted from slides ©2005 Silberschatz, Galvin, and Gagne. Many slides generated from my lecture notes by Kubiawicz.

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.6

File Allocation for Cray-1 DEMOS



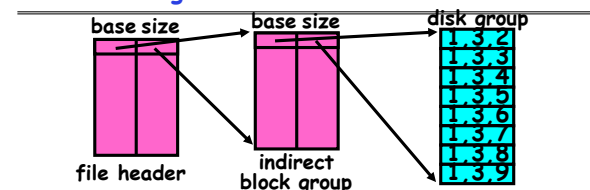
- DEMOS: File system structure similar to segmentation
 - Idea: reduce disk seeks by
 - » using contiguous allocation in normal case
 - » but allow flexibility to have non-contiguous allocation
 - Cray-1 had 12ns cycle time, so CPU:disk speed ratio about the same as today (a few million instructions per seek)
- Header: table of base & size (10 "block group" pointers)
 - Each block chunk is a contiguous group of disk blocks
 - Sequential reads within a block chunk can proceed at high speed - similar to continuous allocation
- How do you find an available block group?
 - Use freelist bitmap to find block of 0's.

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.7

Large File Version of DEMOS



- What if need much bigger files?
 - If need more than 10 groups, set flag in header: BIGFILE
 - » Each table entry now points to an indirect block group
 - Suppose 1000 blocks in a block group ⇒ 80GB max file
 - » Assuming 8KB blocks, 8byte entries ⇒
(10 ptrs × 1024 groups / ptr × 1000 blocks / group) × 8K = 80GB
- Discussion of DEMOS scheme
 - Pros: Fast sequential access, Free areas merge simply, Easy to find free block groups (when disk not full)
 - Cons: Disk full ⇒ No long runs of blocks (fragmentation), so high overhead allocation/access
 - Full disk ⇒ worst of 4.1BSD (lots of seeks) with worst of continuous allocation (lots of recompaction needed)

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.8

How to keep DEMOS performing well?

- In many systems, disks are always full
 - CS department growth: 300 GB to 1TB in a year
 - » That's 2GB/day! (Now at 65+50 TB!)
 - How to fix? Announce that disk space is getting low, so please delete files?
 - » Don't really work: people try to store their data faster
 - Sidebar: Perhaps we are getting out of this mode with new disks... However, let's assume disks full for now
- Solution:
 - Don't let disks get completely full: reserve portion
 - » Free count = # blocks free in bitmap
 - » Scheme: Don't allocate data if count < reserve
 - How much reserve do you need?
 - » In practice, 10% seems like enough
 - Tradeoff: pay for more disk, get contiguous allocation
 - » Since seeks so expensive for performance, this is a very good tradeoff

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.9

Administrivia

- Plan Ahead: this month will be difficult!!
 - Project deadlines every week
- Project #3 design doc due today at 11:59pm
- Midterm #2 is next Wednesday (April 16th)
 - 6-7:30pm in 10 Evans
 - All material from projects 1-3, lectures #9 (2/25) to #19 (4/9)
 - » OS History, Services, and Structure; CPU Scheduling; Kernel and Address Spaces; Address Translation, Caching and TLBs; Demand Paging; I/O Systems; Filesystems, Disk Management, Naming, and Directories; Distributed Systems
 - Email cs162@cory with conflicts
- Projects have a grading standard

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.10

UNIX BSD 4.2

- Same as BSD 4.1 (same file header and triply indirect blocks), except incorporated ideas from DEMOS:
 - Uses bitmap allocation in place of freelist
 - Attempt to allocate files contiguously
 - 10% reserved disk space
 - Skip-sector positioning (mentioned next slide)
- Problem: 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 Demos, power of 2 growth: once it grows past 1MB, allocate 2MB, etc
 - 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

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.11

Attack of the Rotational Delay

- Problem 2: 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!



- Solution1: Skip sector positioning ("interleaving")
 - » Place the blocks from one file on every other block of a track: give time for processing to overlap rotation
- Solution2: Read ahead: read next block right after first, even if application hasn't asked for it yet.
 - » This can be done either by OS (read ahead)
 - » By disk itself (track buffers). Many disk controllers have internal RAM that allows them to read a complete track
- Important Aside: Modern disks+controllers do many complex things "under the covers"
 - Track buffers, elevator algorithms, bad block filtering

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.12

How do we actually access files?

- All information about a file contained in its file header
 - UNIX calls this an "inode"
 - » Inodes are global resources identified by index ("inumber")
 - Once you load the header structure, all the other blocks of the file are locatable
- Question: how does the user ask for a particular file?
 - One option: user specifies an inode by a number (index).
 - » Imagine: `open("14553344")`
 - Better option: specify by textual name
 - » Have to map name→inumber
 - Another option: Icon
 - » This is how Apple made its money. Graphical user interfaces. Point to a file and click.
- **Naming:** The process by which a system translates from user-visible names to system resources
 - In the case of files, need to translate from strings (textual names) or icons to innumbers/inodes
 - For global file systems, data may be spread over globe→need to translate from strings or icons to some combination of physical server location and inumber

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.13

Directories

- **Directory:** a relation used for naming
 - Just a table of (file name, inumber) pairs
- How are directories constructed?
 - Directories often stored in files
 - » Reuse of existing mechanism
 - » Directory named by inode/inumber like other files
 - Needs to be quickly searchable
 - » Options: Simple list or Hashtable
 - » Can be cached into memory in easier form to search
- How are directories modified?
 - Originally, direct read/write of special file
 - System calls for manipulation: `mkdir`, `rmdir`
 - Ties to file creation/destruction
 - » On creating a file by name, new inode grabbed and associated with new file in particular directory

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.14

Directory Organization

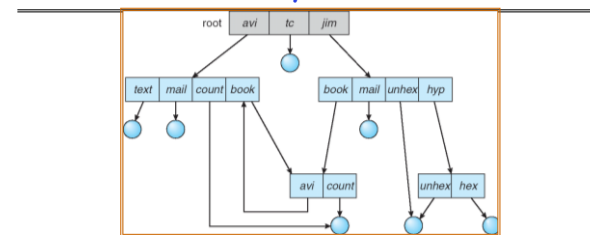
- Directories organized into a hierarchical structure
 - Seems standard, but in early 70's it wasn't
 - Permits much easier organization of data structures
- Entries in directory can be either files or directories
- Files named by ordered set (e.g., `/programs/p/list`)

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.15

Directory Structure



- Not really a hierarchy!
 - Many systems allow directory structure to be organized as an acyclic graph or even a (potentially) cyclic graph
 - Hard Links: different names for the same file
 - » Multiple directory entries point at the same file
 - Soft Links: "shortcut" pointers to other files
 - » Implemented by storing the logical name of actual file
- **Name Resolution:** The process of converting a logical name into a physical resource (like a file)
 - Traverse succession of directories until reach target file
 - Global file system: May be spread across the network

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.16

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")

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.17

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, see 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")

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.18

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 cylinder. 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

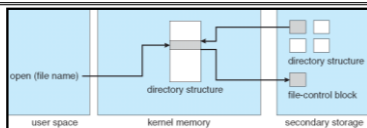
4/7/08

Joseph CS162 ©UCB Spring 2008

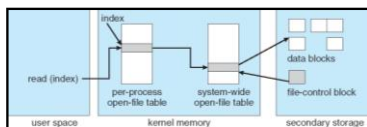
Lec 18.19

BREAK

In-Memory File System Structures



- **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



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

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.21

File System Caching

- **Key Idea:** Exploit locality by caching data in memory
 - Name translations: Mapping from paths→inodes
 - Disk blocks: Mapping from block address→disk content
- **Buffer Cache:** Memory used to cache kernel resources, including disk blocks and name translations
 - Can contain "dirty" blocks (blocks yet on disk)
- **Replacement policy?** LRU
 - Can afford overhead of timestamps for each disk block
 - Advantages:
 - » Works very well for name translation
 - » Works well in general as long as memory is big enough to accommodate a host's working set of files.
 - Disadvantages:
 - » Fails when some application scans through file system, thereby flushing the cache with data used only once
 - » Example: `find . -exec grep foo {} \;`
- **Other Replacement Policies?**
 - Some systems allow applications to request other policies
 - Example, 'Use Once':
 - » File system can discard blocks as soon as they are used

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.22

File System Caching (con't)

- **Cache Size:** How much memory should the OS allocate to the buffer cache vs virtual memory?
 - Too much memory to the file system cache ⇒ won't be able to run many applications at once
 - Too little memory to file system cache ⇒ many applications may run slowly (disk caching not effective)
 - Solution: adjust boundary dynamically so that the disk access rates for paging and file access are balanced
- **Read Ahead Prefetching:** fetch sequential blocks early
 - Key Idea: exploit fact that most common file access is sequential by prefetching subsequent disk blocks ahead of current read request (if they are not already in memory)
 - Elevator algorithm can efficiently interleave groups of prefetches from concurrent applications
 - How much to prefetch?
 - » Too many imposes delays on requests by other applications
 - » Too few causes many seeks (and rotational delays) among concurrent file requests

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.23

File System Caching (con't)

- **Delayed Writes:** Writes to files not immediately sent out to disk
 - Instead, `write()` copies data from user space buffer to kernel buffer (in cache)
 - » Enabled by presence of buffer cache: can leave written file blocks in cache for a while
 - » If some other application tries to read data before written to disk, file system will read from cache
 - Flushed to disk periodically (e.g. in UNIX, every 30 sec)
 - Advantages:
 - » Disk scheduler can efficiently order lots of requests
 - » Disk allocation algorithm can be run with correct size value for a file
 - » Some files need never get written to disk! (e.g. temporary scratch files written /tmp often don't exist for 30 sec)
 - Disadvantages:
 - » What if system crashes before file has been written out?
 - » Worse yet, what if system crashes before a directory file has been written out? (lose pointer to inode!)

4/7/08

Joseph CS162 ©UCB Spring 2008

Lec 18.24

Summary

- **Cray DEMOS: optimization for sequential access**
 - Inode holds set of disk ranges, similar to segmentation
- **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
- **Buffer cache used to increase performance**
 - Read Ahead Prefetching and Delayed Writes