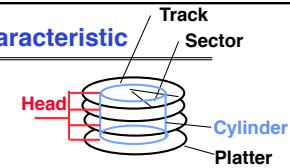


CS162 Operating Systems and Systems Programming Lecture 13

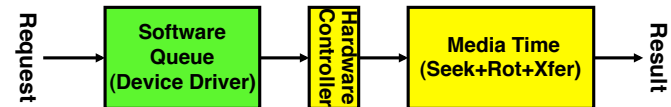
File Systems, Naming, Directories, and Caching

March 5, 2012
Anthony D. Joseph and Ion Stoica
<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 = Queuing Time + Controller time + Seek Time + Rotation Time + Xfer Time**



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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.2

Goals for Today

- Finish SSD discussion
- Important System Properties
- File Systems
 - Structure, Naming, Directories, Caching

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.

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.3

Review: Solid State Disks (SSDs)

- NAND Flash
 - Sector addressable, *but stores 4-64 “sectors” per memory page*
 - No moving parts (no rotate/seek motors)
 - Very low power and lightweight
- Reading data similar to memory read (25µs)
 - No seek or rotational latency
 - Transfer time: transfer a block of bits (sector)
 - » Limited by controller and disk interface (SATA: 300-600MB/s)
 - **Disk Latency = Queuing Time + Controller time + Xfer Time**
 - **Highest Bandwidth:** Sequential OR Random reads

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.4

Review: SSD Architecture – Writes

- Writing data is complex! (~200µs – 1.7ms)
 - No seek or rotational latency, Xfer time: transfer a sector
- But, can only write empty pages (erase takes ~1.5ms!)
 - Controller maintains pool of empty pages by coalescing used sectors (read, erase, write), also reserve some % of capacity
- Typical steady state behavior when SSD is almost full
 - One erase every 64 or 128 writes (depending on page size)
- Write and erase cycles require “high” voltage
 - Damages memory cells, limits SSD lifespan
 - Controller uses ECC, performs wear leveling
 - OS may provide TRIM information about “deleted” sectors
- Result is very workload dependent performance
 - **Disk Latency = Queuing Time + Controller time (Find Free Block) + Xfer Time**

Rule of thumb: writes 10x more expensive than reads, and erases 10x more expensive than writes

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 @UCB Spring 2012

13.7

Drive Health: SMART

- Self-Monitoring, Analysis and Reporting Technology
 - Drive reports on its own health



3/5/2012

Anthony D. Joseph and Ion Stoica CS162 @UCB Spring 2012

13.6

Storage Performance & Price

	Bandwidth (sequential R/W)	Cost/GB	Size
HDD	50-100 MB/s	\$0.05-0.1/GB	2-4 TB
SSD ¹	200-500 MB/s (SATA) 6 GB/s (PCI)	\$1.5-5/GB	200GB-1TB
DRAM	10-16 GB/s	\$5-10/GB	64GB-256GB

¹<http://www.fastestssd.com/featured/ssd-rankings-the-fastest-solid-state-drives/>

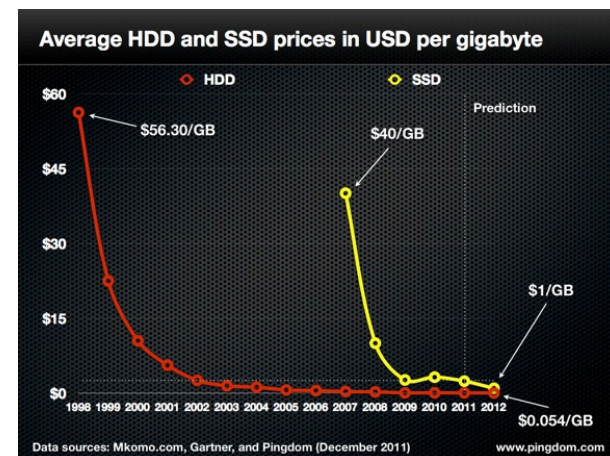
BW: SSD up to x10 than HDD, DRAM > x10 than SSD
Price: HDD x30 less than SSD, SSD x4 less than DRAM

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 @UCB Spring 2012

13.7

Is 2012 the Tipping Point for SSDs?



3/5/2012

Anthony D. Joseph and Ion Stoica CS162 @UCB Spring 2012

13.8

SSD Summary

- Pros (vs. magnetic disk drives):
 - Low latency, high throughput (eliminate seek/rotational delay)
 - No moving parts:
 - » Very light weight, low power, silent, very shock insensitive
 - Read at memory speeds (limited by controller and I/O bus)
- Cons
 - Small storage (0.1-0.5x disk), very expensive (30x disk)
 - » Hybrid alternative: combine small SSD with large HDD
 - Asymmetric block write performance: read pg/erase/write pg
 - » Controller GC algorithms have major effect on performance
 - » Sequential write performance may be worse than HDD
 - Limited drive lifetime (NOR is higher, more expensive)
 - » 50-100K writes/page for SLC, 1-10K writes/page for MLC

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.9

Important “ilities”

- **Availability**: the probability that the system can accept and process requests
 - Often measured in “nines” of probability. So, a 99.9% probability is considered “3-nines of availability”
 - Key idea here is independence of failures
- **Durability**: the ability of a system to recover data despite faults
 - This idea is fault tolerance applied to data
- **Durability doesn't imply Availability**
 - Information on pyramids was very durable, but could not be accessed until discovery of Rosetta Stone
- **Reliability**: the ability of a system or component to perform its required functions under stated conditions for a specified period of time (IEEE definition)
 - Usually stronger than simply availability: means that the system is not only “up”, but also working correctly
 - Includes availability, security, fault tolerance/durability
 - Must make sure data survives system crashes, disk crashes, other problems

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.10

Building a File System

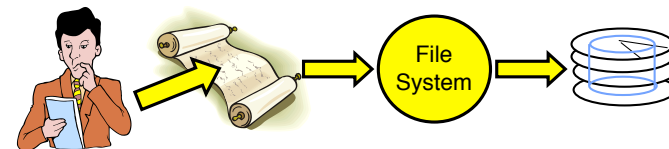
- **File System**: Layer of OS that transforms block interface of disks (or other block devices) into Files, Directories, etc.
- File System Components
 - Disk Management: collecting disk blocks into files
 - Naming: Interface to find files by name, not by blocks
 - Protection: Layers to keep data secure
 - 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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.11

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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.12

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.
 - **Logical Block Addressing (LBA)**: Every sector has integer address from zero up to max number of sectors.
 - Controller translates from address \Rightarrow 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 \Rightarrow 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**

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.13

Designing the File System: Access Patterns

- How do users access files?
 - Need to know type of access patterns user is likely to throw at system
- Sequential Access: bytes read in order (“give me the next X bytes, then give me next, etc.”)
 - Almost all file access are of this flavor
- Random Access: read/write element out of middle of array (“give me bytes i–j”)
 - Less frequent, but still important. For example, virtual memory backing file: page of memory stored in file
 - Want this to be fast – don't want to have to read all bytes to get to the middle of the file
- Content-based Access: (“find me 100 bytes starting with JOSEPH”)
 - Example: employee records – once you find the bytes, increase my salary by a factor of 2
 - Many systems don't provide this; instead, build DBs on top of disk access to index content (requires efficient random access)
 - Example: Mac OSX Spotlight search (do we need directories?)

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.14

Designing the File System: Usage Patterns

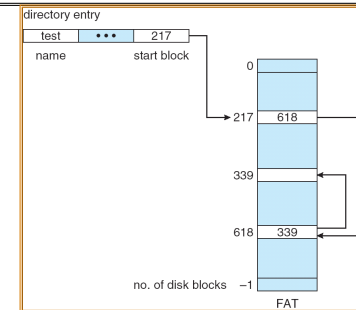
- Most files are small (for example, .login, .c, .java files)
 - A few files are big – executables, .jar, core files, etc.; the .jar is as big as all of your .class files combined
 - However, most files are small – .class's, .o's, .c's, etc.
- Large files use up most of the disk space and bandwidth to/ from disk
 - May seem contradictory, but a few enormous files are equivalent to an immense # of small files
- Although we will use these observations, beware usage patterns:
 - Good idea to look at usage patterns: beat competitors by optimizing for frequent patterns
 - Except: changes in performance or cost can alter usage patterns. Maybe UNIX has lots of small files because big files are really inefficient?
- **File System Goals:**
 - Maximize sequential performance
 - Easy random access to file
 - Easy management of file (growth, truncation, etc)

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.15

Linked Allocation: File-Allocation Table (FAT)



- MSDOS links pages together to create a file
 - Links not in pages, but in the File Allocation Table (FAT)
 - » FAT contains an entry for each block on the disk
 - » FAT Entries corresponding to blocks of file linked together
 - Access properties:
 - » Sequential access expensive unless FAT cached in memory
 - » Random access expensive always, but *really* expensive if FAT not cached in memory

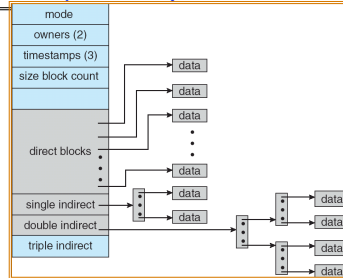
3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.16

Multilevel Indexed Files (UNIX 4.1)

- Multilevel Indexed Files: (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)

3/5/2012 Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012 13.17

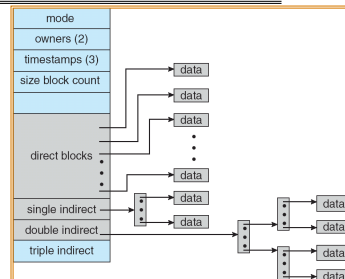
Multilevel Indexed Files (UNIX 4.1): Discussion

- 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

3/5/2012 Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012 13.18

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/O's per block!)

3/5/2012 Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012 13.19

Administrivia

- Midterm Wednesday 3/7 at 5-6:30PM in **10 Evans**
 - Closed-book, 1 double-sided page of handwritten notes
 - Covers lectures/readings #1-12 (Wed 3/1) and project one
- Midterm review session **today** 7-9PM in 141 McCone

3/5/2012 Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012 13.20

5min Break

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 @UCB Spring 2012

13.21

UNIX BSD 4.2

- Same as BSD 4.1 (same file header and triply indirect blocks), except incorporated ideas from Cray-1 DEMOS:
 - Uses bitmap allocation in place of freelist
 - Attempt to allocate files contiguously
 - 10% reserved disk space (mentioned next slide)
 - Skip-sector positioning (mentioned in two slides)
- 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 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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 @UCB Spring 2012

13.22

How to Deal with Full Disks?

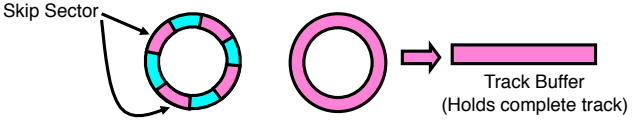
- In many systems, disks are always full
 - EECS department growth: 300 GB to 1TB in a year
 - » That's 2GB/day! (Now at 65+50 TB!)
 - How to fix? Announce that disk space is getting low, so please delete files?
 - » Don't really work: people try to store their data faster
 - Sidebar: Perhaps we are getting out of this mode with new disks... However, let's assume disks are 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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 @UCB Spring 2012

13.23

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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 @UCB Spring 2012

13.24

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 inumbers/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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.25

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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.26

Directory Organization

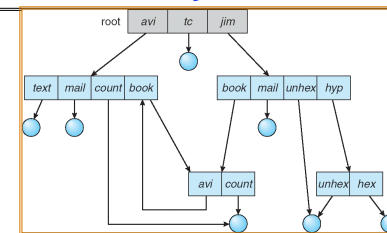
- Directories organized into a hierarchical structure
 - Seems standard, but in early 70's it wasn't
 - Permits much easier organization of data structures
- Entries in directory can be either files or directories
- Files named by ordered set (e.g., `/programs/p/list`)

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.27

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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.28

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"`)

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.29

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")

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.30

Where are inodes stored?

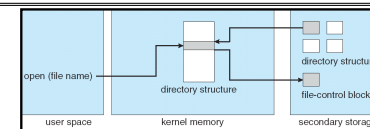
- Later versions of UNIX moved the header information to be closer to the data blocks
 - Often, inode for file stored in same "cylinder group" as parent directory of the file (makes an `ls` of that directory run fast).
 - Pros:
 - » UNIX BSD 4.2 puts a portion of the file header array on each cylinder. 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

3/5/2012

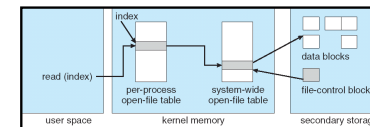
Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.31

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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.32

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? Least Recently Used (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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.33

File System Caching (cont’d)

- 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

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.34

File System Caching (cont’d)

- **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!)

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.35

How to make file system durable?

- Disk blocks contain Reed-Solomon error correcting codes (ECC) to deal with small defects in disk drive
 - Can allow recovery of data from small media defects
- Make sure writes survive in short term
 - Either abandon delayed writes or
 - use special, battery-backed RAM (called non-volatile RAM or **NVRAM**) for dirty blocks in buffer cache.
- Make sure that data survives in long term
 - Need to replicate! More than one copy of data!
 - Important element: **independence of failure**
 - » Could put copies on one disk, but if disk head fails...
 - » Could put copies on different disks, but if server fails...
 - » Could put copies on different servers, but if building is struck by lightning....
 - » Could put copies on servers in different continents...
- **RAID:** Redundant Arrays of Inexpensive Disks
 - Data stored on multiple disks (redundancy)
 - Either in software or hardware
 - » In hardware case, done by disk controller; file system may not even know that there is more than one disk in use

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.36

Log Structured and Journalled File Systems

- Better reliability through use of log
 - All changes are treated as *transactions*
 - A transaction is *committed* once it is written to the log
 - » Data forced to disk for reliability
 - » Process can be accelerated with NVRAM
 - Although File system may not be updated immediately, data preserved in the log
- Difference between “Log Structured” and “Journalled”
 - In a Log Structured file system, data stays in log form
 - In a Journalled file system, Log used for recovery
- For Journalled system:
 - Log used to asynchronously update filesystem
 - » Log entries removed after used
 - After crash:
 - » Remaining transactions in the log performed (“Redo”)
 - » Modifications done in way that can survive crashes
- Examples of Journalled File Systems:
 - Ext3 (Linux), XFS (Unix), HDFS (Mac), NTFS (Windows), etc.

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.37

Summary (1/2)

- Important system properties
 - Availability: how often is the resource available?
 - Durability: how well is data preserved against faults?
 - Reliability: how often is resource performing correctly?
- 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” with index called “inumber”
- Multilevel Indexed Scheme
 - Inode contains file info, direct pointers to blocks,
 - indirect blocks, doubly indirect, etc..

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.38

Summary (2/2)

- 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 file system performance
 - Read Ahead Prefetching and Delayed Writes

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.39

How to organize files on disk

- Goals:
 - Maximize sequential performance
 - Easy random access to file
 - Easy management of file (growth, truncation, etc)
- First Technique: Continuous Allocation
 - Use continuous range of blocks in logical block space
 - » Analogous to base+bounds in virtual memory
 - » User says in advance how big file will be (disadvantage)
 - Search bit-map for space using best fit/first fit
 - » What if not enough contiguous space for new file?
 - File Header Contains:
 - » First block/LBA in file
 - » File size (# of blocks)
 - Pros: Fast Sequential Access, Easy Random access
 - Cons: External Fragmentation/Hard to grow files
 - » Free holes get smaller and smaller
 - » Could compact space, but that would be *really* expensive
- Continuous Allocation used by IBM 360
 - Result of allocation and management cost: People would create a big file, put their file at the start

3/5/2012

Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012

13.40

File Allocation for Cray-1 DEMOS

Basic Segmentation Structure:
Each segment contiguous on disk

- 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.

3/5/2012 Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012 13.41

Large File Version of DEMOS

indirect block group

Basic Segmentation Structure:
Each segment contiguous on disk

- 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 \Rightarrow 80GB max file
 - » Assuming 8KB blocks, 8byte entries \Rightarrow
 $(10 \text{ ptrs} \times 1024 \text{ groups/ptr} \times 1000 \text{ blocks/group}) \times 8K = 80GB$
- 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 \Rightarrow No long runs of blocks (fragmentation), so high overhead allocation/access
 - Full disk \Rightarrow worst of 4.1BSD (lots of seeks) with worst of continuous allocation (lots of recompactation needed)

3/5/2012 Anthony D. Joseph and Ion Stoica CS162 ©UCB Spring 2012 13.42