

CS194-24
Advanced Operating Systems
Structures and Implementation
Lecture 13

File Systems (Con't)
RAID/Journaling/VFS

March 9th, 2014
Prof. John Kubiatowicz
<http://inst.eecs.berkeley.edu/~cs194-24>

Goals for Today

- File Systems (Continued)
- Durability
 - RAID
 - Log-structured File System, Journaling
- VFS

Interactive is important!

Ask Questions!

Note: Some slides and/or pictures in the following are adapted from slides ©2013

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.2

Recall: Implementing LRU

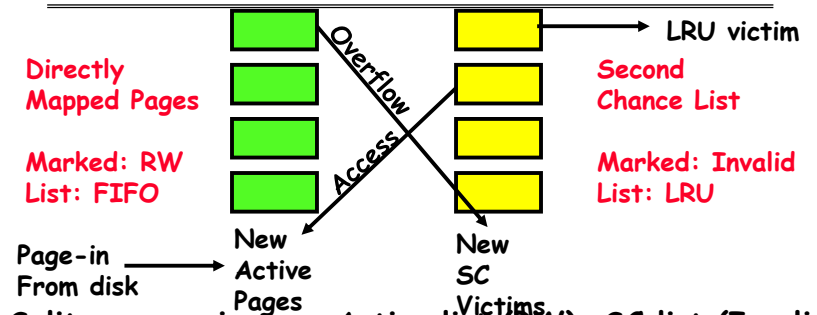
- Perfect:
 - Timestamp page on each reference
 - Keep list of pages ordered by time of reference
 - Too expensive to implement in reality for many reasons
- **Clock Algorithm:** Arrange physical pages in circle with single clock hand
 - Approximate LRU (approx to approx to MIN)
 - Replace **an** old page, not **the oldest** page
- Details:
 - Hardware "use" bit per physical page:
 - » Hardware sets use bit on each reference
 - » If use bit isn't set, means not referenced in a long time
 - » Nachos hardware sets use bit in the TLB; you have to copy this back to page table when TLB entry gets replaced
 - On page fault:
 - » Advance clock hand (not real time)
 - » Check use bit: 1→used recently; clear and leave alone
0→selected candidate for replacement
 - Will always find a page or loop forever?
 - » Even if all use bits set, will eventually loop around⇒FIFO

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.3

Recall: Second-Chance List Algorithm (VAX/VMS)



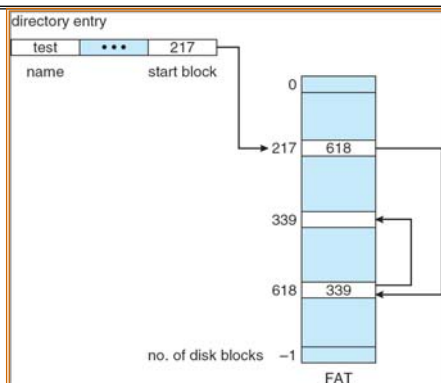
- Split memory in two: Active list (RW), SC list (Invalid)
- Access pages in Active list at full speed
- Otherwise, Page Fault
 - Always move overflow page from end of Active list to front of Second-chance list (SC) and mark invalid
 - Desired Page On SC List: move to front of Active list, mark RW
 - Not on SC list: page in to front of Active list, mark RW; page out LRU victim at end of SC list

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.4

Recall: 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

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.5

Multilevel Indexed Files (UNIX BSD 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 header 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
 - » Block 11 points to "indirect block" containing 256 blocks
 - » Block 12 points to "doubly indirect block" containing 256 indirect blocks for total of 64K blocks
 - » Block 13 points to a triply indirect block (16M blocks)
- **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

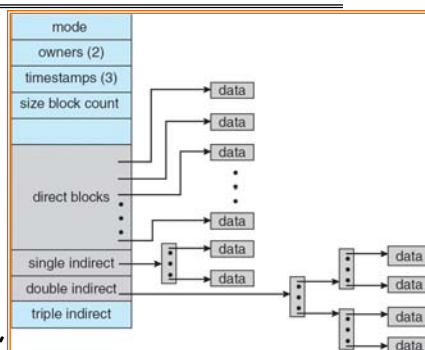
3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.6

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 block (four I/Os per block!)

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.7

Administrivia

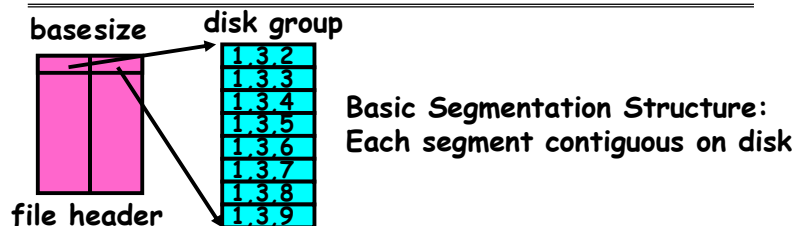
- **Midterm I: Wednesday (3/12)!**
 - Intention is a 1.5 hour exam over 3 hours
 - No class on day of exam!
- **Midterm Timing:**
 - 7:00-10:00PM in 306 Soda Hall (Here!)
- **Topics: everything up to Monday**
 - OS Structure, BDD, Process support, Synchronization, Memory Management, File systems
 - Labs/Papers
- **Notes: 1 page both sides, hand written**

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.8

File Allocation for Cray-1 DEMOS



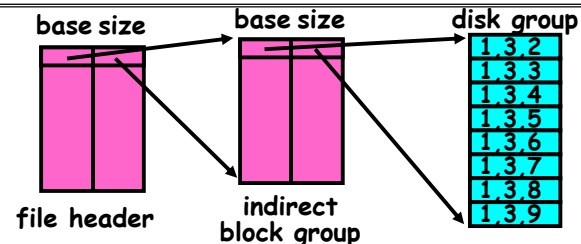
- DEMOS: File system structure similar to segmentation
 - Idea: reduce disk seeks by
 - » using contiguous allocation in normal case
 - » but allow flexibility to have non-contiguous allocation
 - Cray-1 had 12ns cycle time, so CPU:disk speed ratio about the same as today (a few million instructions 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.

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.9

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 \Rightarrow 80GB max file
 - » Assuming 8KB blocks, 8byte entries \Rightarrow (10 ptrs \times 1024 groups/ptr \times 1000 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 recompaction needed)

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.10

How to keep DEMOS performing well?

- In many systems, disks are always full
 - 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

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.11

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)
- Fast File System (FFS)
 - Allocation and placement policies for BSD 4.2
- 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: store files from same directory near each other
- Problem: Block size (512) too small
 - Increase block size to 4096 or 8192
 - Bitmask to allow files to consume partial fragment

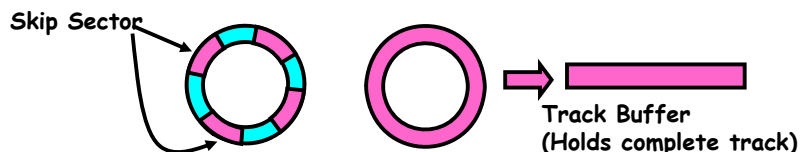
3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.12

Attack of the Rotational Delay

- **Problem: 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**

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.13

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

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.14

Directories

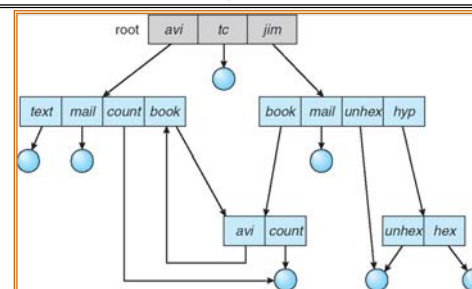
- **Directory: a relation used for naming**
 - Just a table of (file name, inumber) pairs
- **How are directories constructed?**
 - Directories often stored in files
 - » Reuse of existing mechanism
 - » Directory named by inode/inumber like other files
 - Needs to be quickly searchable
 - » Options: 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
- **Directories organized into a hierarchical structure**
 - Entries in directory can be either files or directories
 - Files named by "path" through directory structure: e.g. `/usr/homes/george/data.txt`

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.15

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

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.16

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"`)

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.17

Where are inodes stored?

- In early UNIX and DOS/Windows' FAT file system, headers stored in special array in outermost cylinders
 - Header not stored 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")

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.18

Where are inodes stored?

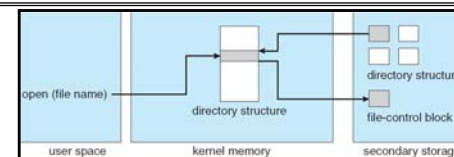
- Later versions of UNIX moved the header information to be closer to the data blocks
 - Often, inode for file stored in same "cylinder group" as parent directory of the file (makes an `ls` of that directory run fast).
 - Pros:
 - » UNIX BSD 4.2 puts a portion of the file header array on each cylinder.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

3/9/14

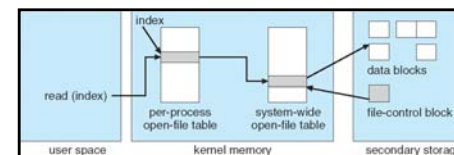
Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.19

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

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.20

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

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.21

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

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.22

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!)

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.23

The Linux Page Cache

- **Goal:** Minimize disk I/O by storing data in memory
 - Caches *any* page-based objects, including files and memory mappings
- **Cache consists of many `address_space` objects**
 - Really misnamed should be a "page_cache_entity" or "physical_pages_of_a_file"
 - Can hold one or more pages from a file or swapper
 - Often associated with an inode, but doesn't have to be
- **Searching an `address_space` object for given page is done by `offset`**
 - Contains a set of `address_space` operations
 - » For reading/writing pages from disk
 - An `address_space` contains a radix tree of all pages
- **Flusher Threads**
 - Start pushing dirty blocks back to disk when free memory shrinks below a specified hreshold

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.24

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
 - Doesn't necessarily 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

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.25

How to make file system durable?

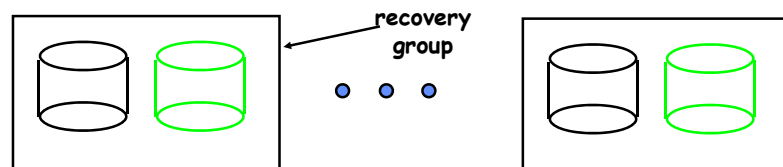
- Disk blocks contain Reed-Solomon error correcting codes (ECC) to deal with small defects in disk drive
 - Can allow recovery of data from small media defects
- Make sure writes survive in short term
 - Either abandon delayed writes or
 - use special, battery-backed RAM (called non-volatile RAM or **NVRAM**) 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

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.26

Redundant Arrays of Disks RAID 1: Disk Mirroring/Shadowing



- Each disk is fully duplicated onto its "shadow"
Very high availability can be achieved
- Bandwidth sacrifice on write:
Logical write = two physical writes
- Reads may be optimized
- Most expensive solution: 100% capacity overhead

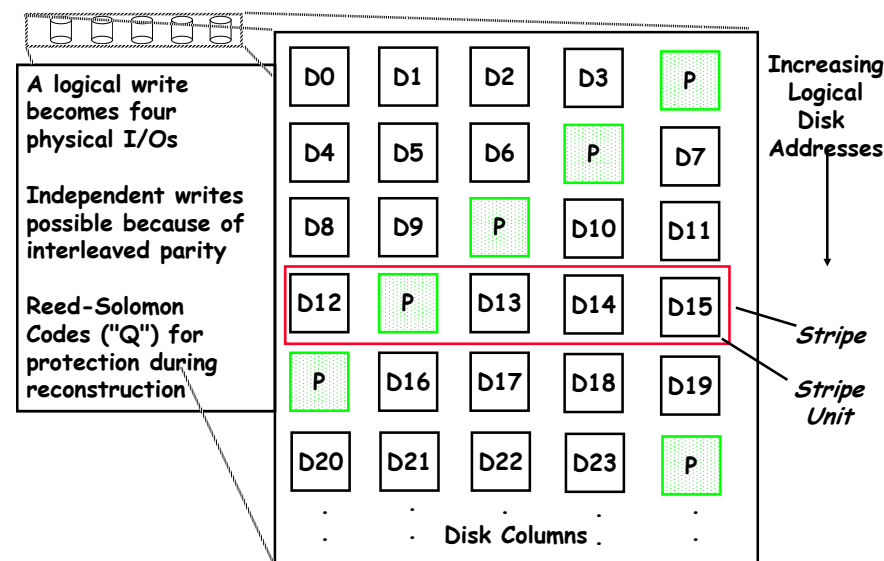
Targeted for high I/O rate, high availability environments

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.27

Redundant Arrays of Disks RAID 5+: High I/O Rate Parity



3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.28

Log Structured and Journalled File Systems

- Better (write) performance through use of log
 - Optimized for writes to disk (log is contiguous)
 - Assume that reads handled through page cache
- 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 filesystem, data stays in log form
 - In a Journalled filesystem, 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), etc.

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.29

Log-Structured File System - Motivation

- Radically different file system design
- Technology motivations:
 - CPUs outpacing disks: I/O becoming more-and-more of a bottleneck
 - Large RAM: file caches work well, making most disk traffic writes
- Problems with (then) current file systems:
 - Lots of little writes
 - Synchronous: wait for disk in too many places - makes it hard to win much from RAIDs, too little concurrency
 - 5 seeks to create a new file: (rough order)
 1. file i-node (create)
 2. file data
 3. directory entry
 4. file i-node (finalize)
 5. directory i-node (modification time)

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.30

LFS Basic Idea

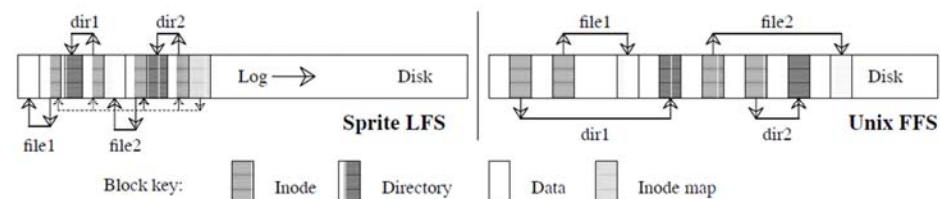
- Log all data and metadata with efficient, large, sequential writes
- Treat the log as the truth, but keep an index on its contents
- Rely on a large memory to provide fast access through caching
- Data layout on disk has "temporal locality" (good for writing), rather than "logical locality" (good for reading)
 - Why is this a better? Because caching helps reads but not writes!
- Two potential problems:
 - Log retrieval on cache misses
 - Wrap-around: what happens when end of disk is reached?
 - » No longer any big, empty runs available
 - » How to prevent fragmentation?

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.31

Comparison of LFS and FFS



- Comparison of LFS (Sprite) with FFS (Unix)
 - Creation of two single-block files named dir1/file1 and dir2/file2
 - Each writes new blocks and inodes for file 1 and file 2
 - Each writes new data blocks and inodes for directories
- FFS Traffic:
 - Ten non-sequential writes for new information
 - Inodes each written twice to ease recovery from crashes
- LFS Traffic:
 - Single large write
- For both when reading back: same number of disk accesses

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.32

LFS Log Retrieval

- Keep same basic file structure as UNIX (inode, indirect blocks, data)
- Retrieval is just a question of finding a file's inode
- UNIX inodes kept in one or a few big arrays, LFS inodes must float to avoid update-in-place
- Solution: an *inode map* that tells where each inode is (Also keeps other stuff: version number, last access time, free/allocated)
- inode map gets written to log like everything else
- Map of inode map gets written in special checkpoint location on disk; used in crash recovery

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.33

LFS Disk Wrap-Around

- Compact live info to open up large runs of free space
 - Problem: long-lived information gets copied over-and-over
- Thread log through free spaces
 - Problem: disk fragments, causing I/O to become inefficient again
- Solution: *segmented log*
 - Divide disk into large, fixed-size segments
 - Do compaction within a segment; thread between segments
 - When writing, use only clean segments (i.e. no live data)
 - Occasionally clean segments: read in several, write out live data in compacted form, leaving some fragments free
 - Try to collect long-lived info into segments that never need to be cleaned
 - Note there is not free list or bit map (as in FFS), only a list of clean segments

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.34

LFS Segment Cleaning

- Which segments to clean?
 - Keep estimate of free space in each segment to help find segments with lowest utilization
 - Always start by looking for segment with utilization=0, since those are trivial to clean...
 - If utilization of segments being cleaned is U :
 - » write cost =
(total bytes read & written)/(new data written) = $2/(1-U)$
(unless U is 0)
 - » write cost increases as U increases: $U = .9 \Rightarrow \text{cost} = 20!$
 - » Need a cost of less than 4 to 10; $\Rightarrow U$ of less than .75 to .45
- How to clean a segment?
 - Segment summary block contains map of the segment
 - Must list every i-node and file block
 - For file blocks you need {i-number, block #}

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.35

Analysis and Evolution of Journaling File Systems

- Write-ahead logging: commit data by writing it to log, synchronously and sequentially
- Unlike LFS, then later moved data to its normal (FFS-like) location - this write is called *checkpointing* and like segment cleaning, it makes room in the (circular) journal
- Better for random writes, slightly worse for big sequential writes
- All reads go to the fixed location blocks, not the journal, which is only read for crash recovery and checkpointing
- Much better than FFS (fsck) for crash recovery (covered below) because it is much faster
- Ext3/ReiserFS/Ext4 filesystems are the main ones in Linux

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.36

Three modes for a JFS

- **Writeback mode:**
 - Journal only metadata
 - Write back data and metadata independently
 - Metadata may thus have dangling references after a crash (if metadata written before the data with a crash in between)
- **Ordered mode:**
 - Journal only metadata, but always write data blocks before their referring metadata is journaled
 - This mode generally makes the most sense and is used by Windows NTFS and IBM's JFS
- **Data journaling mode:**
 - Write both data and metadata to the journal
 - Huge increase in journal traffic; plus have to write most blocks twice, once to the journal and once for checkpointing

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.37

What about remote file systems?

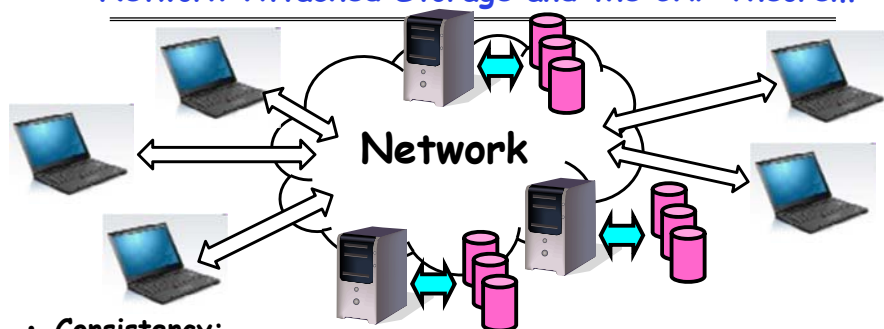
- **Remote File System:**
 - Storage available over a network
 - Local storage is only a cache on permanent storage
- **Advantages?**
 - Someone else worries about keeping data safe
 - Data Accessible from multiple physical locations
- **Disadvantages?**
 - Performance - may take one or more network roundtrips to fetch data
 - Privacy: your data is available over the network, others can possibly see your data
 - Integrity: without sufficient protections, others can overwrite/delete your data

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.38

Network-Attached Storage and the CAP Theorem



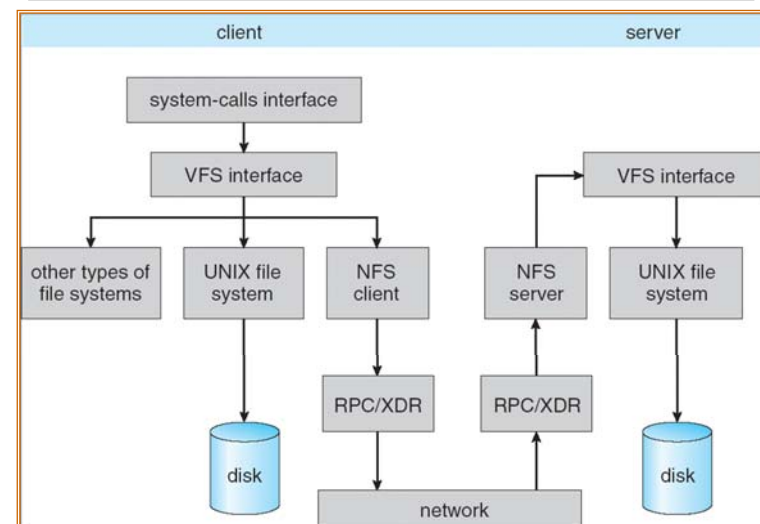
- **Consistency:**
 - Changes appear to everyone in the same serial order
- **Availability:**
 - Can get a result at any time
- **Partition-Tolerance**
 - System continues to work even when network becomes partitioned
- **Consistency, Availability, Partition-Tolerance (CAP) Theorem:**
Cannot have all three at same time
 - Otherwise known as "Brewer's Theorem"

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.39

Schematic View of NFS Architecture



3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.40

Virtual Filesystem Switch (VFS)



```

inf = open("/floppy/TEST", O_RDONLY, 0);
outf = open("/tmp/test",
            O_WRONLY|O_CREAT|O_TRUNC, 0600);
do {
    i = read(inf, buf, 4096);
    write(outf, buf, i);
} while (i);
close(outf);
close(inf);
    
```

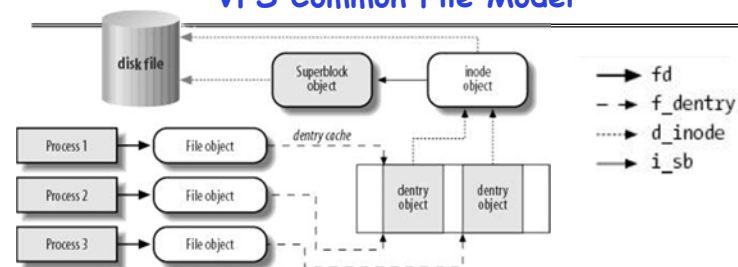
- **VFS: Virtual abstraction similar to local file system**
 - Provides virtual superblocks, inodes, files, etc
 - Compatible with a variety of local and remote file systems
 - » provides object-oriented way of implementing file systems
- **VFS allows the same system call interface (the API) to be used for different types of file systems**
 - The API is to the VFS interface, rather than any specific type of file system

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.41

VFS Common File Model



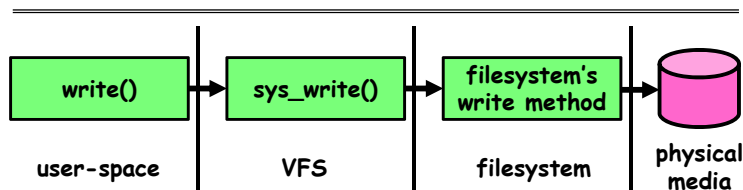
- **Four primary object types for VFS:**
 - superblock object: represents a specific mounted filesystem
 - inode object: represents a specific file
 - dentry object: represents a directory entry
 - file object: represents open file associated with process
- **There is no specific directory object (VFS treats directories as files)**
- **May need to fit the model by faking it**
 - Example: make it look like directories are files
 - Example: make it look like have inodes, superblocks, etc.

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.42

Linux VFS



- **An operations object is contained within each primary object type to set operations of specific filesystems**
 - "super_operations": methods that kernel can invoke on a specific filesystem, i.e. write_inode() and sync_fs().
 - "inode_operations": methods that kernel can invoke on a specific file, such as create() and link()
 - "dentry_operations": methods that kernel can invoke on a specific directory entry, such as d_compare() or d_delete()
 - "file_operations": methods that process can invoke on an open file, such as read() and write()
- **There are a lot of operations**
 - You need to read Bovet Chapter 12 and Love Chapter 13

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.43

Summary (1/2)

- **Multilevel Indexed Scheme**
 - Inode contains file info, direct pointers to blocks, indirect blocks, doubly indirect, etc..
- **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
- **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?

3/9/14

Kubiatowicz CS194-24 ©UCB Fall 2014

Lec 13.44

Summary (2/2)

- **Log Structured File system (LFS)**
 - The Log *is* the file system
 - All updates written sequentially in the log
 - Inode map tracks where inodes lie in the log
- **Journaling File System (JFS, Ext3, ...)**
 - Use of log to help durability
 - Primary storage in read-optimized format
- **Distributed File System:**
 - Transparent access to files stored on a remote disk
 - Caching for performance
- **VFS: Virtual File System layer**
 - Provides mechanism which gives same system call interface for different types of file systems
 - More Next Time!