

**CSE 451: Operating Systems**  
**Spring 2009**

**Lecture 14**  
**FFS and LFS**

**Steve Gribble**

# File System Implementations

- We've looked at disks and file systems generically
  - now it's time to bridge the gap by talking about specific file system implementations
- We'll focus on two:
  - BSD Unix FFS
    - what's at the heart of most UNIX file systems
  - LFS
    - a research file system originally from Berkeley

# BSD UNIX FFS

- FFS = “Fast File System”
  - original (i.e. 1970’s) file system was very simple and straightforwardly implemented
    - but had very poor disk bandwidth utilization
    - why? far too many disk seeks on average
- BSD UNIX folks did a redesign in the mid ’80’s
  - FFS: improved disk utilization, decreased response time
  - McKusick, Joy, Fabry, and Leffler
  - basic idea is FFS is aware of disk structure
    - i.e., place related things on nearby cylinders to reduce seeks

# File System “data structure” reminder

- How does the FS use the disk to store files?
- FS defines a block size (e.g., 4KB)
  - disk space allocated in granularity of blocks
- A “Master Block” defines the location of root directory
  - always at a well-known location
  - usually replicated for reliability
- A “free map” lists which blocks are free vs. allocated
  - usually a bitmap, one bit per block on the disk
  - also stored on disk, and cached in memory for performance
- Remaining disk blocks are used to store files/dirs
  - how this is done is the essence of FFS

# Inodes and Path Search

- Unix Inodes are NOT directories
  - they describe where on disk the blocks for a file are placed
    - directories are just files, so each directory also has an inode that describes where the blocks for the directory is placed
- Directory entries map file names to inodes
  - to open “/one”, use master block to find inode for “/” on disk
    - open “/”, look for entry for “one”
    - this gives the disk block number for inode of “one”
  - read the inode for “one” into memory
    - this inode says where the first data block is on disk
    - read that data block into memory to access the data in the file

# Data and Inode placement

- Original (non-FFS) unix FS had two major problems:
  - 1. data blocks are allocated randomly in aging file systems
    - blocks for the same file allocated sequentially when FS is new
    - as FS “ages” and fills, need to allocate blocks freed up when other files are deleted
      - problem: deleted files are essentially randomly placed
      - so, blocks for new files become scattered across the disk!
  - 2. inodes are allocated far from blocks
    - all inodes at beginning of disk, far from data
    - traversing file name paths, manipulating files, directories requires going back and forth from inodes to data blocks
  - BOTH of these generate many long seeks!

# Cylinder groups

- FFS addressed these problems using notion of a cylinder group
  - disk partitioned into groups of cylinders
  - data blocks from a file all placed in same cylinder group
  - files in same directory placed in same cylinder group
  - inode for file in same cylinder group as file's data
- Introduces a free space requirement
  - to be able to allocate according to cylinder group, the disk must have free space scattered across all cylinders
  - in FFS, 10% of the disk is reserved just for this purpose!
    - good insight: keep disk partially free at all times!
    - this is why it may be possible for df to report >100%

# File Buffer Cache

- Exploit locality by caching file blocks in memory
  - the cache is system wide, shared by all processes
  - even a small (4MB) cache can be very effective (why?)
  - many FS's "read-ahead" into buffer cache
- What about writes?
  - some apps assume data is on disk after write
    - either "write-through" the buffer cache
    - or "write-behind"
      - maintain queue of uncommitted blocks, periodically flush. Unreliable!
      - NVRAM: write into battery-backed RAM. Expensive!
      - LFS: we'll talk about this soon!
- Buffer cache issues:
  - competes with VM for physical frames
    - integrated VM/buffer cache?
  - need replacement algorithms here
    - LRU usually



# Other FFS innovations

- Small blocks (1KB) caused two problems:
  - low bandwidth utilization
  - small max file size (function of block size)
  - FFS fixes by using a larger block (4KB)
    - allows for large files (1MB only uses 2 level indirect)
    - but, introduces internal fragmentation
      - there are many small files (i.e., <4KB)
    - fix: introduce “fragments”
      - 1KB pieces of a block. More bookkeeping, but more flexible.
- Old FS was unaware of disk parameters
  - FFS: parameterize FS according to disk and CPU characteristics
    - e.g.: account for CPU interrupt and processing time to layout sequential blocks

# Log-Structured File System (LFS)

- LFS was designed in response to two trends in workload and disk technology:
  - 1. Disk bandwidth scaling significantly (40% a year)
    - but, latency is not
  - 2. Large main memories in machines
    - therefore, large buffer caches
      - absorb large fraction of read requests in caches
    - can use for writes as well
      - coalesce small writes into large writes
- LFS takes advantage of both to increase FS performance
  - Rosenblum and Ousterhout (Berkeley, '91)
    - note: Rosenblum went on to become Stanford prof, and to co-found VMware, inc!

# FFS problems that LFS solves

- FFS: placement improved, but can still have many small seeks
  - possibly related files are physically separated
  - inodes separated from files (small seeks)
  - directory entries separate from inodes
- FFS: metadata required synchronous writes
  - with small files, most writes are to metadata
  - synchronous writes are very slow!

# LFS: The Basic Idea

- Treat the entire disk as a single log for appending
  - collect writes in the disk buffer cache, and write out the entire collection of writes in one large request
    - leverages disk bandwidth with large sequential write
    - no seeks at all! (assuming head at end of log)
  - all info written to disk is appended to log
    - data blocks, attributes, inodes, directories, .etc.
- Sounds simple!
  - but it's really complicated under the covers

# LFS Challenges

- There are two main challenges with LFS:
  - 1. locating data written in the log
    - FFS places files in a well-known location, LFS writes data “at the end of the log”
  - 2. managing free space on the disk
    - disk is finite, and therefore log must be finite
    - cannot always append to log!
      - need to recover deleted blocks in old part of log
      - need to fill holes created by recovered blocks

# LFS: locating data

- FFS uses inodes to locate data blocks
  - inodes preallocated in each cylinder group
  - directories contain locations of inodes
- LFS appends inodes to end of log, just like data
  - makes them hard to find
- Solution:
  - use another level of indirection: **inode maps**
  - inode maps map inode #s to inode location
  - so how do you find the inode map?
    - after all, changes to it must be appended to the log
    - location of inode map blocks are kept in a **checkpoint region**
    - checkpoint region has a fixed location
  - cache inode maps in memory for performance

# LFS: free space management

- LFS: append-only quickly eats up all disk space
  - need to recover deleted blocks
- Solution:
  - fragment log into segments
  - thread segments on disk
    - segments can be anywhere
  - reclaim space by cleaning segments
    - read segment
    - copy live data to end of log
    - now have free segment you can reuse!
  - cleaning is a big problem
    - costly overhead, when do you do it?
      - “idleness is not sloth”

# Other ideas...take 551 for more!

- **File system consistency issues**
  - any given file operation involves touching data blocks as well as metadata (inodes, directory entries, free lists/bitmaps)
  - in what order should you touch these blocks? (what if a crash happens partway through?)
- **Journaling file systems**
  - halfway in-between FFS and LFS: write operations to a log region (the journal) before updating the file system itself
  - after a crash, replay recent operations from the log
    - simpler to maintain consistency! and, no more fsck!
- **Other kinds of file systems**
  - indexed file systems (e.g., Spotlight on Mac)
  - file systems for large-scale clusters (e.g., GFS, bigtable)
  - SSD-optimized file systems