

CS 202: Advanced Operating Systems

Log-Structured File System

Log-Structured/Journaling File System



- 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)
 - 6. (not to mention bitmap updates)

LFS Basic Idea



- > Log all data and metadata with efficient, large, sequential writes
 - > Do not update blocks in place just write new versions in the log
- > Treat the log as the truth, but keep an index on its contents
- > Not necessarily good for reads, but trends help
 - > 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!

Basic idea





- We buffer all updates, and write them together in one big sequential write
 - Good for the disk
 - Example above, writes to two different files were written together (along with the new version of i-node) in one write
 - > How much should we buffer?
 - > What happens if too much? If too little?
- > But how do we find a file??
 - > All problems in CS solved with another level of indirection ©

Devil is in the details



- > Two potential problems:
 - Log retrieval on cache misses how do we find the data?
 - Wrap-around: what happens when end of disk is reached?
 - > No longer any big, empty runs available
 - How to prevent fragmentation?



LFS vs. UFS



i-node map





- > A map keeping track of the location of i-nodes
- Anytime an i-node is written to disk, the imap is updated
 - > But is that any better? In a second
- Most of the time the imap is in memory, so access is fast
- > Updated imap is saved as part of the log!
 - but how do we find <u>it!</u>

Final piece to the solution





- Checkpoint region is written to point to the location of the imap
 - Also serves as an indicator of a stable point in the file system for crash recovery
- > So, to read a file from LFS:
 - > Read the CR, use it to read and cache the imap
 - > After that, it is identical to FFS
 - > Are reads fast?

What about directories?





- > When a file is updated, its inode changes (new copy)
 - > We need to update the directory inode (also creating a copy)
 - We need to update its parent directory
- > Ugh....what to do?
 - Inode map helps with that too just keep track of inode number and resolve it through inode map

LFS Disk Wrap-Around/Garbage collection





- > 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

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 => cost = 20!
 - Need a cost of less than 4 to 10; => 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 #}
 - Through i-map you check if this block is still being used for the (inumber, block #)

Evaluation Results



Figure 8 — Small-file performance under Sprite LFS and SunOS.

Figure (a) measures a benchmark that created 10000 one-kilobyte files, then read them back in the same order as created, then deleted them. Speed is measured by the number of files per second for each operation on the two file systems. The logging approach in Sprite LFS provides an order-of-magnitude speedup for creation and deletion. Figure (b) estimates the performance of each system for creating files on faster computers with the same disk. In SunOS the disk was 85% saturated in (a), so faster processors will not improve performance much. In Sprite LFS the disk was only 17% saturated in (a) while the CPU was 100% utilized; as a consequence I/O performance will scale with CPU speed.



Figure 9 — Large-file performance under Sprite LFS and SunOS.

The figure shows the speed of a benchmark that creates a 100-Mbyte file with sequential writes, then reads the file back sequentially, then writes 100 Mbytes randomly to the existing file, then reads 100 Mbytes randomly from the file, and finally reads the file sequentially again. The bandwidth of each of the five phases is shown separately. Sprite LFS has a higher write bandwidth and the same read bandwidth as SunOS with the exception of sequential reading of a file that was written randomly.

Is this a good paper?



- > What were the authors' goals?
- > What about the evaluation/metrics?
- Did they convince you that this was a good system/approach?
- Does the system/approach meet the "Test of Time" challenge?
- > How would you review this paper today?