

Lecture Notes, 9/2/98

UNIX Fast File System Measurements of a Distributed File System

I. A Fast File System for UNIX

Original UNIX file system was simple and elegant, but slow.

Could only achieve about 20 KB/sec/arm; ~2% of 1982 disk bandwidth

Problems:

- o blocks too small
- o consecutive blocks of files not close together
- o i-nodes far from data
- o i-nodes of directory not close together

Aspects of new file system:

- o 4096 or 8192 byte block size (why not larger?)
- o large blocks and small fragments
- o disk divided into cylinder groups
- o each contains superblock, i-nodes, bitmap of free blocks, usage summary info
- o keeps i-node near file, i-nodes of a directory together
- o cylinder groups ~ 16 cylinders, or 7.5 MB
- o cylinder headers spread around so not all on one platter

Two techniques for locality:

- o don't let disk fill up in any one area
- o paradox: to achieve locality, must spread unrelated things far apart
- o note: new file system got 175KB/sec because free list contained sequential blocks (it did generate locality), but an old system has randomly ordered blocks and only got 30 KB/sec

Specific application of these techniques:

- o goal: keep directory within a cylinder group, spread out different directories
- o goal: allocate runs of blocks within a cylinder group, every once in a while switch to a

Lecture 2

new cylinder group (jump at 1MB).

- o layout policy: global and local
- o global policy allocates files & directories to cylinder groups. Picks “optimal” next block for block allocation.
- o local allocation routines handle specific block requests. Select from a sequence of alternative if need to.

Results:

- o 20-40% of disk bandwidth for large reads/writes.
- o 10-20x original UNIX speeds.
- o Size: 3800 lines of code vs. 2700 in old system.
- o 10% of total disk space unusable (except at 50% perf. price)

Could have done more; later versions do.

Enhancements made to system interface: (really a second mini-paper)

- o long file names (14 -> 255)
- o advisory file locks
- o symbolic links (contrast to hard links)
- o atomic rename capability
- o disk quotas

3 key features of paper:

- o parameterize FS implementation for the hardware it's running on.
- o measurement-driven design decisions
- o locality “wins”

A major flaws:

- o measurements derived from a single installation.
- o ignored technology trends

A lesson for the future: don't ignore underlying hardware characteristics.

Contrasting research approaches: improve what you've got vs. design something new.

Lecture 2

II. Measurements of a DFS

History:

- o two cs262 projects yield original BSD study in 1985
- o follow-on cs262 project in 1990,
- o Baker *et al.* work.

Goals: analysis paper

- o Understand how access patterns have changed (it at all) since 1985.
- o Understand behavior of distributed caches in network file systems.

BSD comparison results: *Which of these do you believe will hold in the future?*

- o Throughput 20x higher; process migration can add another 6x
Why has throughput gone up 20x?
 - bigger binaries
 - bigger projects; more files to link, grep, etc.
 - studied user community does “big file” things.
- o Why did they get unexpected process migration results?
 - really just remote execution
 - only compute-intensive things get migrated
- o Accesses still mostly sequential
- o Files short (80% < 10KB), but significant numbers of large files (50% of all bytes are transferred to files > 500 KB).
- o File lifetimes: 30-50% of bytes die within 30 sec.
- o Note Hartmann letter correcting paper results.

Caching results:

- o Caches not as effective as predicted in BSD study (40% read misses with 5-10 MB caches): why?
 - large files tend to flush caches
 - per-client caches less effective
- o Paging is a significant portion of server traffic (35%)
- o Client caches are big enough:
 - Block lifetimes in client caches ~ 1 hr or more
 - All write-backs for integrity, not for block replacement (70% 30-sec. delay, 16% fsync, 12% recall)
- o Write sharing doesn't happen often, but happens often enough to affect many users.

Lecture 2

- o Simplest consistency mechanism seems good enough.

Good figure technique: individual data points explained as examples.

An interesting assertion: local disks considered “harmful”! Why? Network is faster than disk.

Some key features of the paper:

- o Validated (most of) the original BSD study results for 1990 timeframe.
- o Data caching works well, but not as well as originally predicted.
- o Consistency for write sharing must be supported, but a simple implementation is good enough.

A flaw:

- o How general is the measured user community and the workload it generates?

A lesson: Understanding the workload of your system is important.