PERSISTENCE: LOG-STRUCTURED FILESYSTEM

Shivaram Venkataraman CS 537, Spring 2020

ADMINISTRIVIA

Discussion today: Prep for final exam

AGENDA / LEARNING OUTCOMES

How to design a filesystem that performs better for small writes?

What are some similarities or differences with FFS?

Very Simple FS -> FFS Log Structured FS

RECAP







HOW CAN FILE SYSTEM FIX INCONSISTENCIES?

Solution #1:

FSCK = file system checker

Strategy:

After crash, scan whole disk for contradictions and "fix" if needed Keep file system off-line until FSCK completes

For example, how to tell if data bitmap block is consistent?

Read every valid inode+indirect block

If pointer to data block, the corresponding bit should be 1; else bit is 0









Write 8, 9, 10, 11, 12, 13 Barrier Write 2,4,5,6

Txn Start 4 blocks TxEnd In-place updates

LOG STRUCTURED FILE SYSTEM (LFS)

LFS PERFORMANCE GOAL

litmap

write operation

seeles

Motivation:

- Growing gap between sequential and random I/O performance
- RAID-5 especially bad with small random writes

Idea: use disk purely sequentially

Design for writes to use disk sequentially – how?

WHERE DO INODES GO?



LFS STRATEGY

File system buffers writes in main memory until "enough" data

- How much is enough?
- Enough to get good sequential bandwidth from disk (MB)

Write buffered data sequentially to new **segment** on disk

Never overwrite old info: old copies left behind



BUFFERED WRITES



AO A1 A2 A3 <u>Inode[j]</u> A5 Inode[k] 4 blocks to file 1 1 block to file 2

WHAT ELSE IS DIFFERENT FROM FFS?

What data structures has LFS removed?

allocation structs: data + inode bitmaps

How to do reads?

Inodes are no longer at fixed offset

Use (imap) structure to map: inode number => inode location on disk image boots 28 7ead $4/ \rightarrow 1$ 15 23 $6/ \rightarrow 1$ 2 23 7ead $4/ \rightarrow 1$ 2 23 7ead $4/ \rightarrow 1$ 2 23 7ead $4/ \rightarrow 1$ 7ead $4/ \rightarrow 1$ 7ead $4/ \rightarrow 1$ 7ead 7ea 7ead 7ead 7ead 7ead 7ead 7ead7ea









WHAT TO DO WITH OLD DATA? old inde

١

T. MAG

new ina

Old versions of files \rightarrow garbage

Approach I: garbage is a feature!

- Keep old versions in case user wants to revert files later
- Versioning file systems
- Example: Dropbox

Approach 2: garbage collection

GARBAGE COLLECTION

garfale

pdit

rode

fore

inad

ma

0

Valio

Need to reclaim space:

I.When no more references (any file system)

2. After newer copy is created (COW file system)

LFS reclaims segments (not individual inodes and data blocks)

Valia

- Want future overwites to be to sequential areas
- Tricky, since segments are usually partly valid



GARBAGE COLLECTION



When moving data blocks, copy new inode to point to it When move inode, update imap to point to it

GARBAGE COLLECTION

General operation:

Pick M segments, compact into N (where N < M).

Mechanism:

How does LFS know whether data in segments is valid?



Policy: Which/segments to compact?

GARBAGE COLLECTION MECHANISM

Is an inode the latest version?

- Check imap to see if this inode is pointed to
- Fast!
- Is a data block the latest version?
 - Scan ALL inodes to see if any point to this data
 - Very slow!

How to track information more efficiently?

Segment summary lists inode and data offset corresponding to each data block in segment (reverse pointers)





60% garlege 40% parkey GARBAGE COLLECTION

2011

Pick most empty,

hot segment or cold segment

General operation:

251. garbage Pick M segments, compact into N (where M < M). n604. gainy Mechanism:

Use segment summary, imap to determine liveness / Cold

Policy:

Which segments to compact?

- clean most empty first
- clean coldest (ones undergoing least change)

60% garbage

more complex heuristics...

CRASH RECOVERY

What data needs to be recovered after a crash?

Need imap (lost in volatile memory)

Better approach?

- Occasionally save to checkpoint region the pointers to imap pieces

How often to checkpoint?

- Checkpoint often: random I/O
- Checkpoint rarely: lose more data, recovery takes longer
- Example: checkpoint every 30 secs

CRASH RECOVERY



CHECKPOINT SUMMARY

Checkpoint occasionally (e.g., every 30s)

Upon recovery:

- read checkpoint to find most imap pointers and segment tail
- find rest of imap pointers by reading past tail

What if crash <u>during</u> checkpoint?

CHECKPOINT STRATEGY

Have two checkpoint regions

Only overwrite one checkpoint at a time

Use checksum/timestamps to identify newest checkpoint



LFS SUMMARY

Journaling: / FF-S

Put final location of data wherever file system chooses (usually in a place optimized for future reads)

LFS:

Puts data where it's fastest to write, assume future reads cached in memory

Other COW file systems: WAFL, ZFS, btrfs

OUIZ 29 https://tinyurl.com/cs537-sp20-quiz29

```
block 100: [("." 0), (".." 0), ("foo" 1)] // a data block
block 101: [size=1,ptr=100,type=d] // an inode
block 102: [size=0,ptr=-,type=r] // an inode
block 103: [imap: 0->101,1->102] // a piece of the imap
```



// a data block block 104: [SOME DATA] block 105: [SOME DATA] // a data block block 106: [size=2,ptr=104,ptr=105,type=r] // an inode block 107: [imap: 0->101,1->106] // a piece of the imap If we fill out 100 segments in a newly created LFS, how long does it take to complete a write?

If we read this file (reads do not hit cache), how long does it take to read the entire file?

If we now read this file backwards, one segment at a time (and reads do not hit in cache), how long does this backwards read take?

NEXT STEPS

Project 5 is one week away!

Discussion: Final practice quiz

