PERSISTENCE: FAST FILE SYSTEM

Shivaram Venkataraman CS 537, Spring 2020

ADMINISTRIVIA

2 weeks from now

P5 release tonight, due April 23rd, 10pm

Attend the discussion section for more details!

Optional project after PS?

AGENDA / LEARNING OUTCOMES

How does file system represent files, directories?

What steps must reads/writes take?

Very Sinple file System

How does FFS improve performance?

RECAP

FILE API WITH FILE DESCRIPTORS

int fd = open(char *path, int flag, mode t mode) read(int fd, void *buf, size t nbyte) write(int fd, void *buf, size t nbyte) close(int fd)

advantages:

- string names

/root / test. txt

- hierarchical
- traverse once
- offsets precisely defined

Permittence La perices Hard dists RAID File API to access data on disks

FILE API SUMMARY

Using multiple types of name provides convenience and efficiency

flushing) atomically moving data to a file from dist one path to another

Hard and soft link features provide flexibility.

Special calls (fsync, rename) let developers communicate requirements to file system



Data blocks

D	D	D	D	D	D		D
8							15
D	D	D	D	D	D	D	D
24							31
D	D	D	D	D	D	D	D
40							47
D	D	D	D	D	D	D	D
56							63

INODE

type (file or dir?)
uid (owner)
rwx (permissions)
size (in bytes)
Blocks
time (access)
ctime (create)
links_count (# paths)
addrs[N] (N data blocks)





How to handle even larger files?

SIMPLE DIRECTORY LIST EXAMPLE



FS OPERATIONS

FILE API

- open
- read
- close
- create file
- write

open /foo/bar

TIME

(already exists)



read /foo/bar – assume opened

TIME



close /foo/bar





write to /foo/bar (assume file exists and has been opened) TIME



EFFICIENCY

How can we avoid this excessive I/O for basic ops?

Cache for:

- reads
- write buffering

Periodi cally or m

Sync



WRITE BUFFERING

Overwrites, deletes, scheduling

Shared structs (e.g., bitmaps+dirs) often overwritten.

Tradeoffs: how much to buffer, how long to buffer

Configurable leg users





QUIZ 26 https://tinyurl.com/cs537-sp20-quiz26
inode bitmap 11100000 /
$$(d_{a:1})$$
 r:2] [d a:2 r:2] [] [] [] [] []
data bitmap 11100000
data [(.,0) (..,0) [d,1]) (w,2)] [????] [(.,2) (..,0)] [] [] [] [] []
mk dir (/d) mk dir (/w)
[(.,1) (..,0)] = /d
inode bitmap 11000000
inodes [d a:0 r:2] [f a:1 r:2] [] [] [] [] [] [] [] dir
data bitmap ????????
data [(.,0) (..,0) (c,1) (m,1)] [foofoofoo] [] [] [] [] [] []
() C Am refer to same inothe refer do same inothe refer dir foofoofoo]
bitmap foofoofoo



FILE LAYOUT IMPORTANCE



Layout is not disk-aware!

DISK-AWARE FILE SYSTEM

How to make the disk use more efficient?

Where to place meta-data and data on disk?



Use groups across disks;

Strategy: allocate inodes and data blocks in same group.

PLACEMENT TECHNIQUE: GROUPS

In FFS, groups were ranges of cylinders called cylinder group

In ext2, ext3, ext4 groups are ranges of blocks called block group $\leq SSTP$





REPLICATED SUPER BLOCKS





top platter damage?

solution: for each group, store super-block at different offset

SMART POLICY



Where should new inodes and data blocks go?

Policy

PLACEMENT STRATEGY

Put related pieces of data near each other.

Rules:

I. Put directory entries near directory inodes.

2. Put inodes near directory entries.

3. Put data blocks near inodes. for files

Problem: File system is one big tree

All directories and files have a common root.

All data in same FS is related in some way

Trying to put everything near everything else doesn't make any choices!

REVISED STRATEGY Compilie code Src/ac 1 a.h 1 b.c

Put more-related pieces of data near each other group inodes Group data Put less-related pieces of data far micdir (1) micdir (2) ls-R/ acde---- accddee /a/b € bf---bff 1 Create (/a/e) /a/c ર /a/d create (lalc) /b/f 5 6 mkdir ((16)

POLICY SUMMARY

File inodes: allocate in <u>same</u> group with dir

Dir inodes: allocate in <u>new</u> group with fewer used inodes than average group | load belowing

First data block: allocate near inode

Other data blocks: allocate near previous block



one seek for each of many small files



12 direct 4 KB

Define "large" as requiring an indirect block

Starting at indirect (e.g., after 48 KB) put blocks in a new block group. 4KB block oddr 4 bytes

Each chunk corresponds to one indirect block Block size 4KB, 4 byte per address => 1024 address per_indirect 1024*4KB = 4MB contiguous "chunk"

POLICY SUMMARY

File inodes: allocate in same group with dir

Dir inodes: allocate in new group with fewer used inodes than average group

First data block: allocate near inode Other data blocks: allocate near previous block

Large file data blocks: after 48KB, go to new group. Move to another group (w/ fewer than avg blocks) every subsequent MB.





OTHER FFS FEATURES

FFS also introduced several new features:

- large blocks (with libc buffering / fragments)
- long file names
- atomic rename
- symbolic links

usability

virtual file system VFS FFS Impl

FFS SUMMARY

First disk-aware file system

- Bitmaps
- Locality groups
- Rotated superblocks
- Smart allocation policy

Inspired modern files systems, including ext2 and ext3

NEXT STEPS

P5 will be released later today Details in the discussion section

Next class: Filesystem consistency

inde addr 32 33 34 77 42

internal fragmentation

