

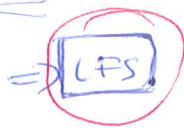
LECTURE

LFS

Perf: FFS (main idea?)
block groups, etc.

Crash consistency: Journaling
main idea?

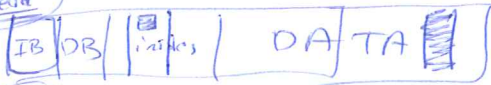
Today: Perf + Crash Consistency



Problem: FFS perf good, not great
e.g. create a 7-block file
(what structures are read, written?)

read dir inode
read dir data (see if filename is unique)

① all the blocks you must read



② all the blocks you must write

block group
read inode bitmap (find free)
write inode bitmap (mark allocated)
write dir data (name → new inode num)
write dir inode (update length)
read inode block (mark allocated)
write inode block (now initialized)

read data bitmap
write data too

what happens w/ cache?

(as it gets larger?)

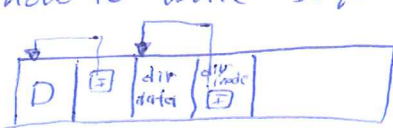
cache is central shaper of traffic
reads: avoided writes: delayed + scheduled

conclude: Perf of FFS under lots of file creates: localized, but still lots of (seeks/rotates)

real goal: use disk (purely sequentially)

hard for reads → why? if file x, y are not near each other...
easy for writes → why? can always write to unallocated space

let's pretend: all reads are in cache
how to write sequentially?



(in memory)

segment (say 1MB)

Protocol: attempt #1:

- ① put data, inodes into in-memory segment
- ② given full, write to disk
disk: just a bunch of empty segments

Problem: open ("/foo", O_RDONLY)
typical fs:

read root inode
read root data (find foo)
read foo inode

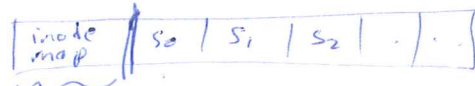
given inode #, how to find? (easy: calculate)
inode

LFS: how to find inode?

- 1) scan disk
- 2) indirection: inode map

array: indexed by inode #
content: addr of latest copy of inode

Protocol #2: inode map @ fixed location

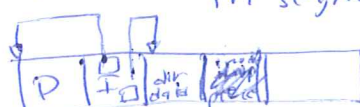


millions of entries (4B each)
① fill segment (in memory)
update inode map (in memory)

- ② write segment
write dirty parts of inode map

but: problem? writes not that sequential anymore

Protocol #3: inode map pieces in segments!



- ① fill ② write ⇒ purely sequential
- problem? can't find inode map

Protocol #4: checkpoint region



includes pointers to pieces of inode map

- ① fill blocks, update inode map (in memory)
- ② write segment (seq.)
- ③ occasionally (every 30secs) checkpoint by updating CR w/ inode map

what if crash?

- problem between CR update
- ② during CR update

Review: how to open/read a file?

assume: all on disk

CR \rightarrow inode map \rightarrow inode \rightarrow {file...}

once here, save

really: on boot, read in inode map (circle it)
 \Rightarrow no slower

problem

influence:

NetApp

(wall)

ZFS

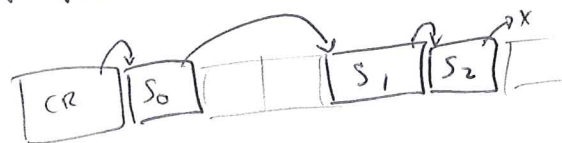
(cow)

new problem: disk gets full



what to do? is disk really full?

what about crash recovery?

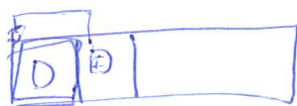


cleaner (garbage collector)

\rightarrow goal: find dead blocks, make available for reuse

\Rightarrow roll forward

mechanism: how to determine if dead?
cleaner



① read in segment

need:

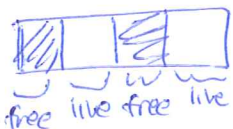
{ inode # of data block
offset in file }

② how to determine if inode is live? (easy: map)

how to determine if data is live? (easy: check live inode's map)

then: Check if live inode \Rightarrow this data block
pts to

result



\Rightarrow problem: can't write sequentially

solution:

read in N segments ~~blocks~~, write M like ~~blocks~~ segments
($M < N$)

free N old segments

problem: which blocks to clean?
random: }
hot/cold: } policy

problem: how big to make segments

