

filesystems 4 / distribution (start)

last time

(double/triple-)indirect blocks

block groups

typical file sizes

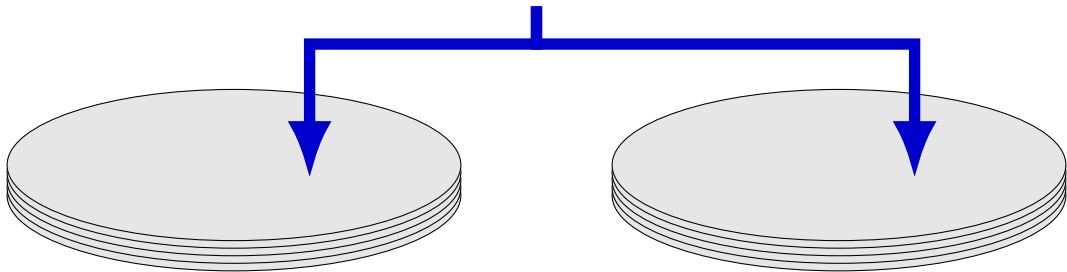
extents

trees on disk

mirroring whole disks

alternate strategy: write everything to **two disks**

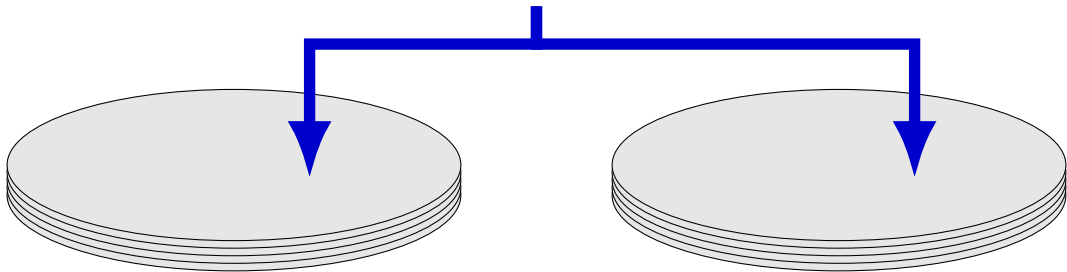
always write to both



mirroring whole disks

alternate strategy: write everything to **two disks**

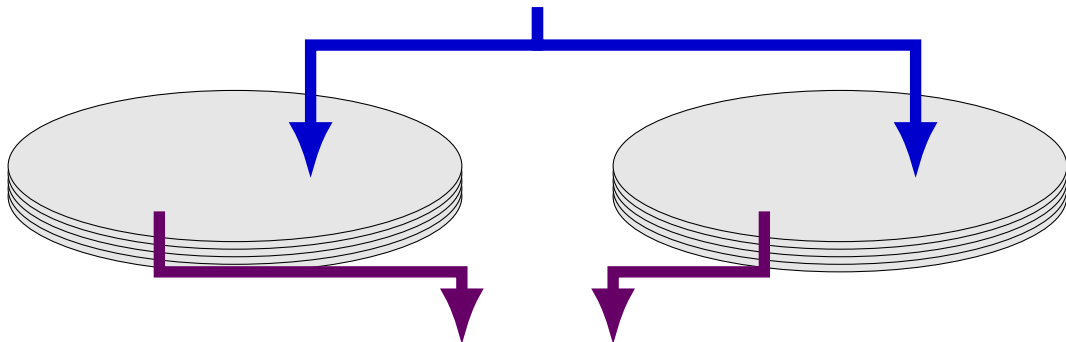
always **write to both**



mirroring whole disks

alternate strategy: write everything to **two disks**

always write to both



read from either
(or different parts of both – **faster!**)

RAID 4 parity

\oplus — bitwise xor

disk 1

disk 2

disk 3

A_1 : sector 0	A_2 : sector 1	A_p : $A_1 \oplus A_2$
B_1 : sector 2	B_2 : sector 3	B_p : $B_1 \oplus B_2$

...

...

...

RAID 4 parity

\oplus — bitwise xor

disk 1

disk 2

disk 3

A_1 : sector 0	A_2 : sector 1	A_p : $A_1 \oplus A_2$
B_1 : sector 2	B_2 : sector 3	B_p : $B_1 \oplus B_2$

...

...

...

$$A_p = A_1 \oplus A_2$$

$$A_1 = A_p \oplus A_2$$

$$A_2 = A_1 \oplus A_p$$

can compute contents of any disk!

RAID 4 parity

\oplus — bitwise xor

disk 1	disk 2	disk 3
A_1 : sector 0	A_2 : sector 1	A_p : $A_1 \oplus A_2$
B_1 : sector 2	B_2 : sector 3	B_p : $B_1 \oplus B_2$
...

exercise: how to replace sector 3 (B_2) with new value?
how many writes? how many reads?

RAID 4 parity (more disks)

disk 1	disk 2	disk 3	disk 4
A_1 : sector 0	A_2 : sector 1	A_3 sector 2	A_p : $A_1 \oplus A_2 \oplus A_3$
B_1 : sector 3	B_2 : sector 4	B_3 : sector 5	B_p : $B_1 \oplus B_2 \oplus B_3$
...	

RAID 4 parity (more disks)

disk 1	disk 2	disk 3	disk 4
A_1 : sector 0	A_2 : sector 1	A_3 : sector 2	A_p : $A_1 \oplus A_2 \oplus A_3$
B_1 : sector 3	B_2 : sector 4	B_3 : sector 5	B_p : $B_1 \oplus B_2 \oplus B_3$
...	

$$A_p = A_1 \oplus A_2 \oplus A_3$$

$$A_1 = A_p \oplus A_2 \oplus A_3$$

$$A_2 = A_1 \oplus A_p \oplus A_3$$

$$A_3 = A_1 \oplus A_2 \oplus A_p$$

can still compute contents of any disk!

RAID 4 parity (more disks)

disk 1	disk 2	disk 3	disk 4
A_1 : sector 0	A_2 : sector 1	A_3 sector 2	A_p : $A_1 \oplus A_2 \oplus A_3$
B_1 : sector 3	B_2 : sector 4	B_3 : sector 5	B_p : $B_1 \oplus B_2 \oplus B_3$
...	

exercise: how to replace sector 3 (B_1) with new value now?
how many writes? how many reads?

RAID 5 parity

disk 1	disk 2	disk 3	disk 4
A_1 : sector 0	A_2 : sector 1	A_3 : sector 2	A_p : $A_1 \oplus A_2 \oplus A_3$
B_1 : sector 3	B_2 : sector 4	B_p : $B_1 \oplus B_2 \oplus B_3$	B_3 : sector 5
C_1 : sector 6	C_p : $C_1 \oplus C_2 \oplus C_3$	C_2 : sector 7	C_3 : sector 8
...	

RAID 5 parity

disk 1	disk 2	disk 3	disk 4
A_1 : sector 0	A_2 : sector 1	A_3 : sector 2	A_p : $A_1 \oplus A_2 \oplus A_3$
B_1 : sector 3	B_2 : sector 4	B_p : $B_1 \oplus B_2 \oplus B_3$	B_3 : sector 5
C_1 : sector 6	C_p : $C_1 \oplus C_2 \oplus C_3$	C_2 : sector 7	C_3 : sector 8
...	

spread out parity updates across disks
so each disk has about same amount of work

more general schemes

RAID 6: tolerate loss of any two disks

can generalize to 3 or more failures

justification: takes days/weeks to replace data on missing disk
...giving time for more disks to fail

probably more in CS 4434?

but none of this addresses consistency

RAID-like redundancy

usually appears to filesystem as 'more reliable disk'

hardware or software layers to implement extra copies/parity

some filesystems (e.g. ZFS) implement this themselves

more flexibility — e.g. change redundancy file-by-file

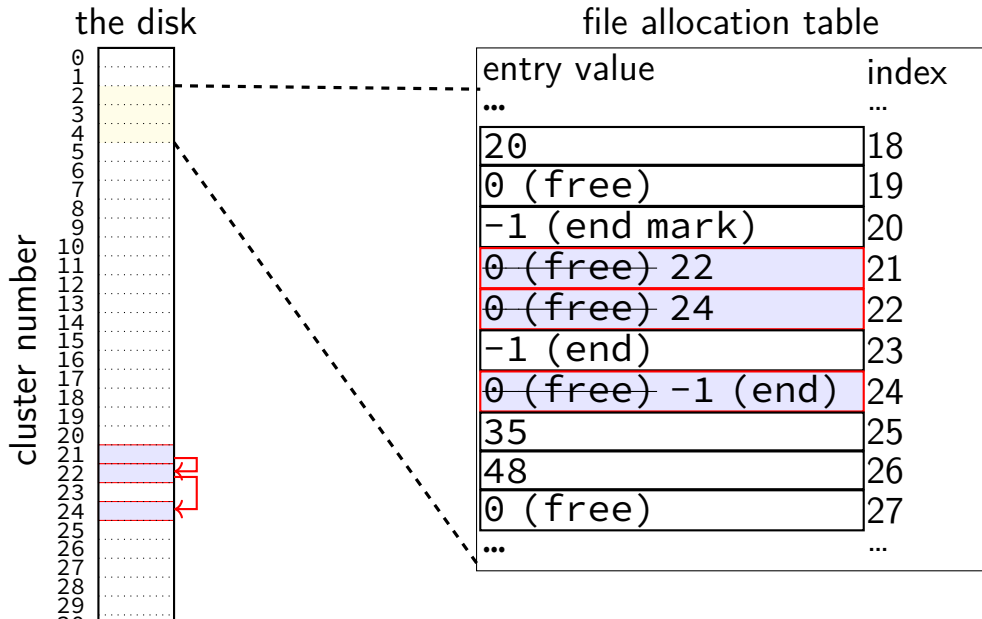
ZFS combines with its own checksums — don't trust disks!

RAID: missing piece

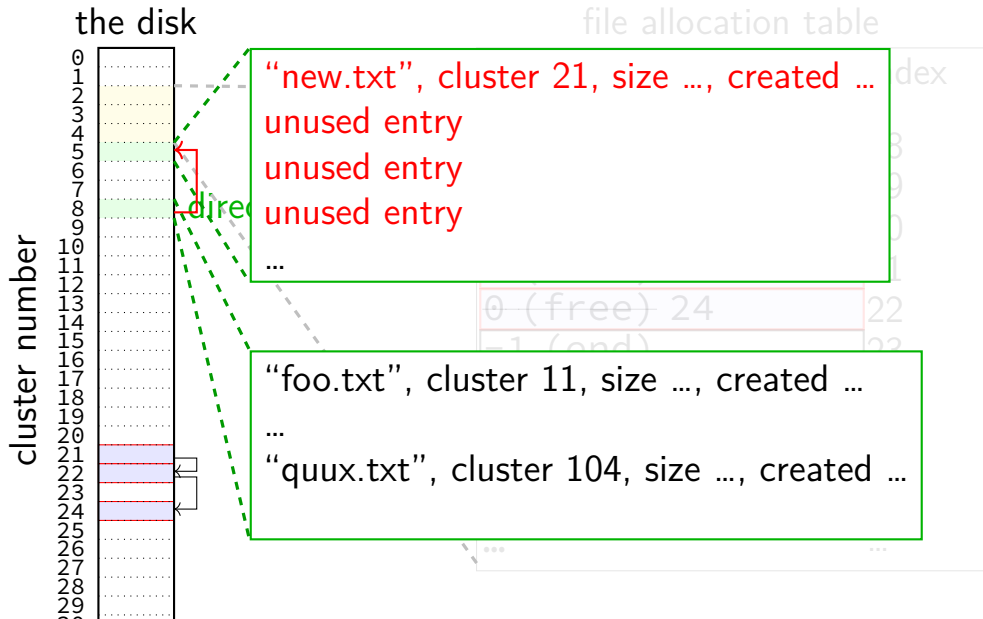
what about losing data while blocks being updated

very tricky/failure-prone part of RAID implementations

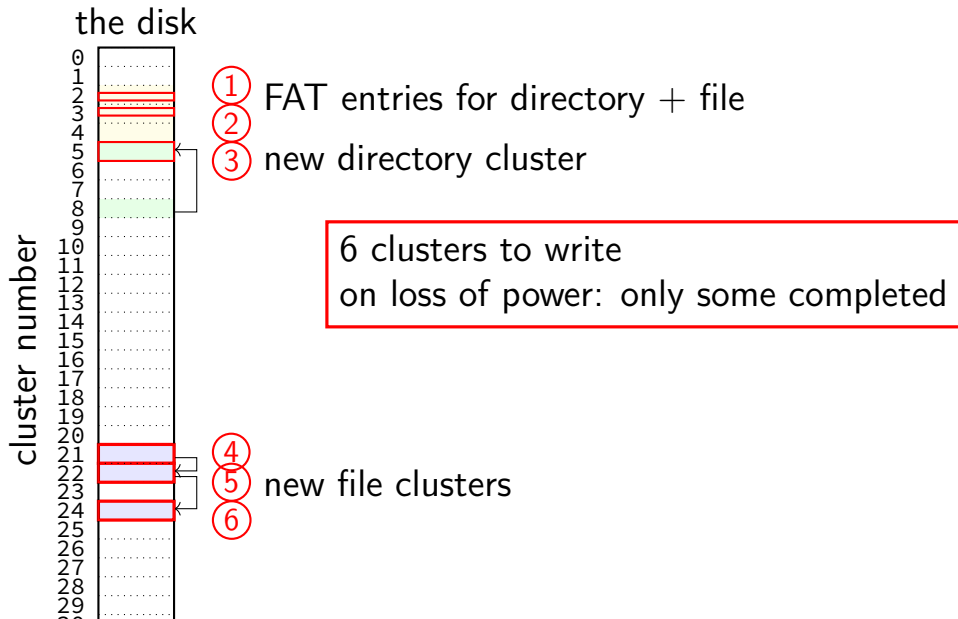
recall: FAT: file creation (1)



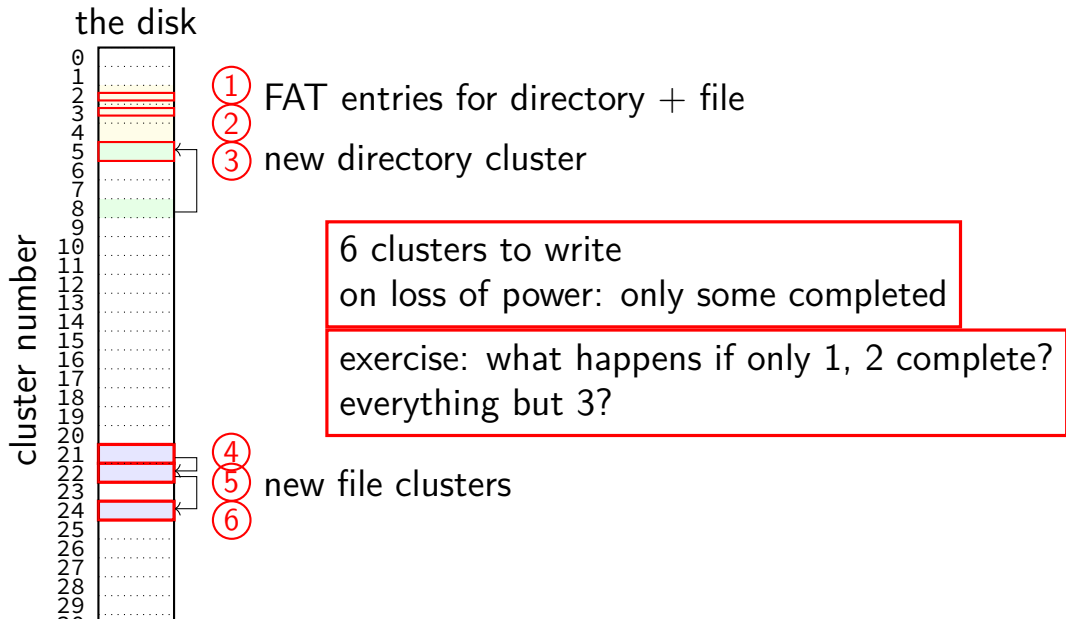
recall: FAT: file creation (2)



exercise: FAT file creation



exercise: FAT file creation



exercise: FAT ordering

(creating a file that needs new cluster of direntries)

1. FAT entry for extra directory cluster
2. FAT entry for new file clusters
3. file clusters
4. file's directory entry (in new directory cluster)

what ordering is best if a crash happens in the middle?

- A. 1, 2, 3, 4
- B. 4, 3, 1, 2
- C. 1, 3, 4, 2
- D. 3, 4, 2, 1
- E. 3, 1, 4, 2

exercise: xv6 FS ordering

(creating a file that needs new block of direntries)

1. free block map for new directory block
2. free block map for new file block
3. directory inode
4. new file inode
5. new directory entry for file (in new directory block)
6. file data blocks

what ordering is best if a crash happens in the middle?

- A. 1, 2, 3, 4, 5, 6
- B. 6, 5, 4, 3, 2, 1
- C. 1, 2, 6, 5, 4, 3
- D. 2, 6, 4, 1, 5, 3
- E. 3, 4, 1, 2, 5, 6

inode-based FS: careful ordering

mark blocks as allocated before referring to them from directories

write data blocks before writing pointers to them from inodes

write inodes before directory entries pointing to it

remove inode from directory before marking inode as free
or decreasing link count, if there's another hard link

idea: better to waste space than point to bad data

inode-based FS: creating a file

normal operation

allocate data block

write data block

update free block map

update file inode

update directory entry

filename+inode number

update directory inode

modification time

inode-based FS: creating a file

normal operation

allocate data block
write data block
update free block map
update file inode
update directory entry
 filename+inode number
update directory inode
 modification time

general rule:

better to waste space
than point to bad data

mark blocks/inodes used before writing

inode-based FS: creating a file

normal operation

- allocate data block
- write data block
- update free block map
- update file inode
- update directory entry
filename+inode number
- update directory inode
modification time

recovery (fsck)

- read all directory entries
- scan all inodes
 - free unused inodes
unused = not in directory
- free unused data blocks
 - unused = not in inode lists
- scan directories for missing
update/access times

inode-based FS: exercise: unlink

what order to remove a hard link (= directory entry) for file?

1. overwrite directory entry for file
2. decrement link count in inode (but link count still > 1 so don't remove)

assume not the last hard link

inode-based FS: exercise: unlink

what order to remove a hard link (= directory entry) for file?

1. overwrite directory entry for file
2. decrement link count in inode (but link count still > 1 so don't remove)

assume not the last hard link

what does recovery operation do?

inode-based FS: exercise: unlink last

what order to remove a hard link (= directory entry) for file?

1. overwrite last directory entry for file
2. mark inode as free (link count = 0 now)
3. mark inode's data blocks as free

assume **is the last hard link**

inode-based FS: exercise: unlink last

what order to remove a hard link (= directory entry) for file?

1. overwrite last directory entry for file
2. mark inode as free (link count = 0 now)
3. mark inode's data blocks as free

assume **is the last hard link**

what does recovery operation do?

fsck

Unix typically has an `fsck` utility

Windows equivalent: `chkdsk`

checks for *filesystem consistency*

is a data block marked as used that no inodes uses?

is a data block referred to by two different inodes?

is an inode marked as used that no directory references?

is the link count for each inode = number of directories referencing it?

...

assuming careful ordering, can fix errors after a crash without loss

maybe can fix other errors, too

fsck costs

my desktop's filesystem:

2.4M used inodes; 379.9M of 472.4M used blocks

recall: check for data block marked as used that no inode uses:

- read blocks containing all of the 2.4M used inodes

- add each block pointer to a list of used blocks

- if they have indirect block pointers, read those blocks, too

- get list of all used blocks (via direct or indirect pointers)

- compare list of used blocks to actual free block bitmap

pretty expensive and slow

running fsck automatically

common to have “clean” bit in superblock

last thing written (to set) on shutdown

first thing written (to clear) on startup

on boot: if clean bit clear, run fsck first

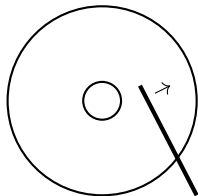
ordering and disk performance

recall: seek times

would like to **order writes based on locations on disk**

write many things in one pass of disk head

write many things in cylinder in one rotation

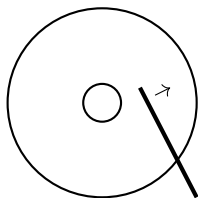


ordering and disk performance

recall: seek times

would like to **order writes based on locations on disk**

write many things in one pass of disk head
write many things in cylinder in one rotation



ordering constraints make this hard:

free block map for file (start), then file blocks (middle), then...

file inode (start), then directory (middle), ...

beyond ordering

recall: updating a sector is atomic
happens entirely or doesn't

can we make filesystem updates work this way?

beyond ordering

recall: updating a sector is atomic
happens entirely or doesn't

can we make filesystem updates work this way?

yes — 'just' make updating one sector do the update

concept: transaction

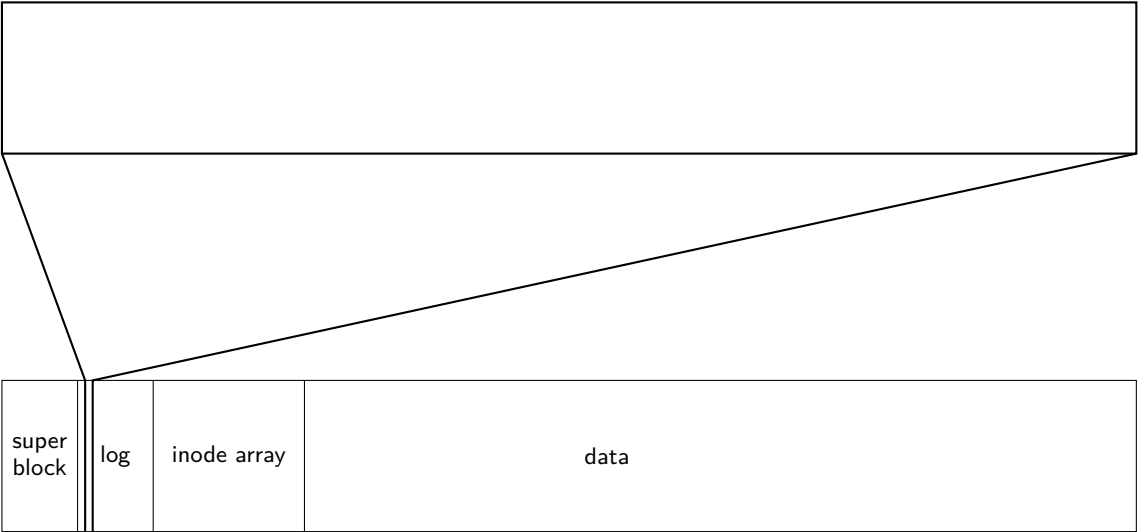
transaction: bunch of updates that happen all at once

implementation trick: one update means transaction “commits”

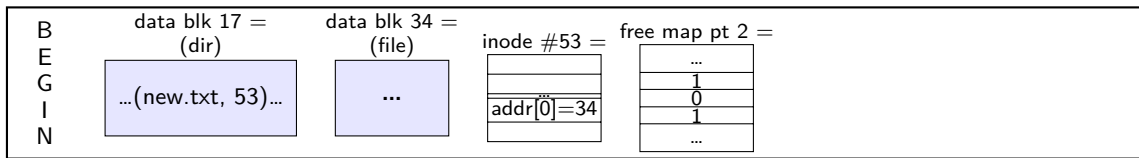
- update done — whole transaction happened

- update not done — whole transaction did not happen

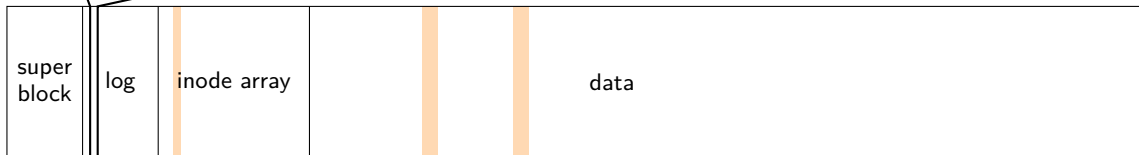
redo logging: file creation



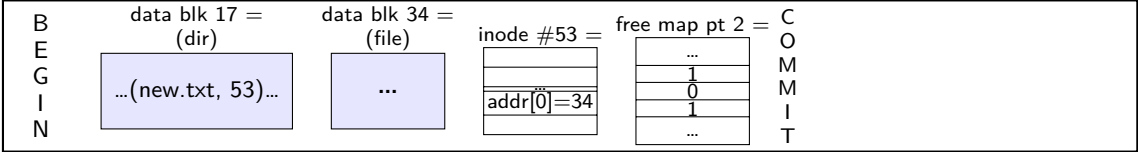
redo logging: file creation



write log entries with **intended operations**



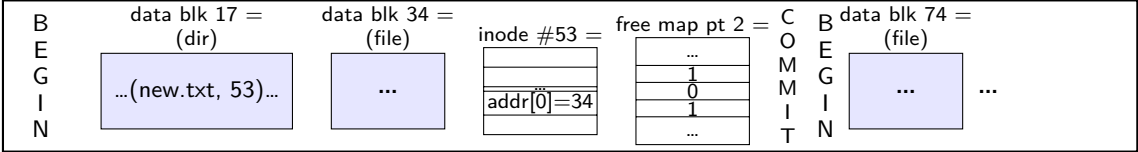
redo logging: file creation



write commit message to log



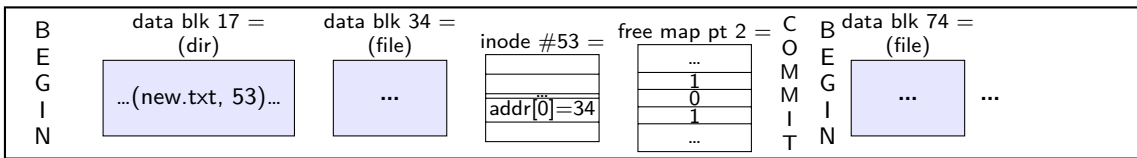
redo logging: file creation



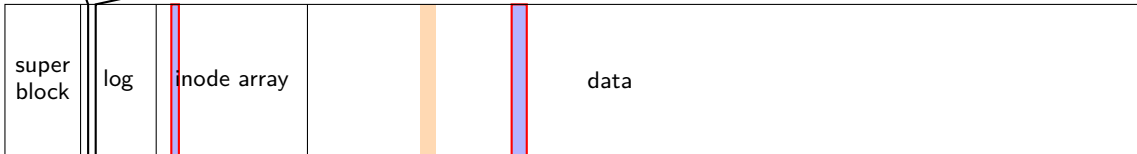
...and start more transactions



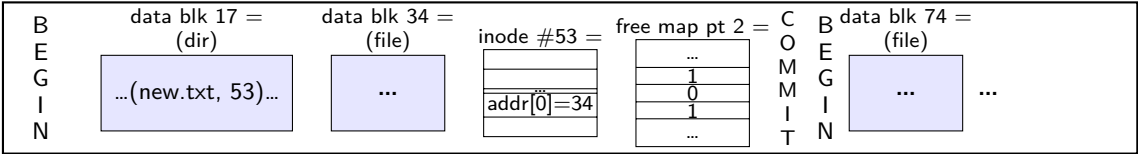
redo logging: file creation



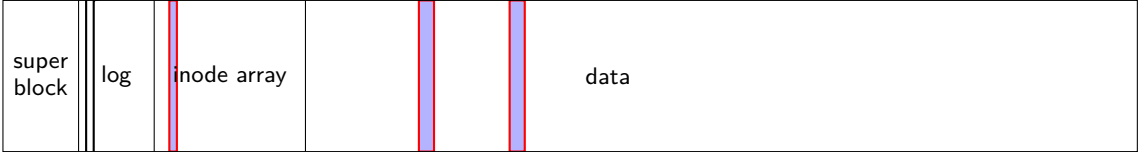
later, start applying results to actual disk



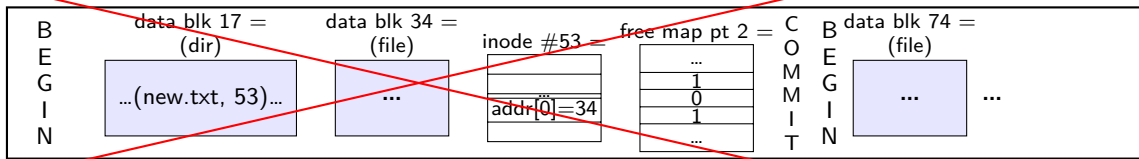
redo logging: file creation



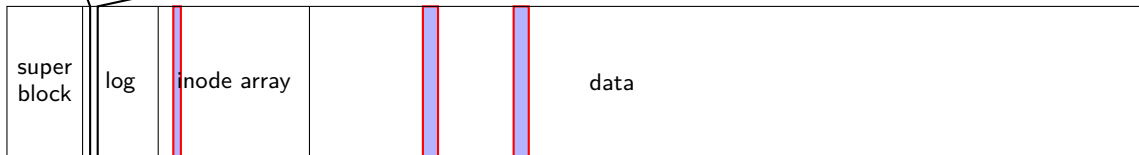
when everything is written, can overwrite log



redo logging: file creation



when everything is written, can overwrite log



redo logging: file creation

normal operation

write to log transaction steps:

- data blocks to create
- directory entry, inode to write
- directory inode (size, time)
- update

write to log “commit transaction”

in any order:

- update file data blocks
- update directory entry
- update file inode
- update directory inode

reclaim space in log

redo logging: file creation

normal operation

write to log transaction steps:

- data blocks to create
- directory entry, inode to write
- directory inode (size, time)
- update

write to log “commit transaction”

in any order:

- update file data blocks
- update directory entry
- update file inode
- update directory inode

reclaim space in log

crash before *commit*?

file not created

no partial operation to real data

redo logging: file creation

normal operation

write to log transaction steps:

- data blocks to create
- directory entry, inode to write
- directory inode (size, time)
- update

write to log “commit transaction”

in any order:

- update file data blocks
- update directory entry
- update file inode
- update directory inode

reclaim space in log

crash after *commit*?

file created

promise: **will perform logged updates**
(after system reboots/recovers)

redo logging: file creation

normal operation

write to log transaction steps:

- data blocks to create
- directory entry, inode to write
- directory inode (size, time)
- update

write to log “commit transaction”

in any order:

- update file data blocks
- update directory entry
- update file inode
- update directory inode

reclaim space in log

redo logging: file creation

normal operation

write to log transaction steps:

- data blocks to create
- directory entry, inode to write
- directory inode (size, time)
- update

write to log “commit transaction”

in any order:

- update file data blocks
- update directory entry
- update file inode
- update directory inode

reclaim space in log

recovery

read log and...

ignore any operation with no
“commit”

redo any operation with
“commit”

- already done? — okay, setting
inode twice

reclaim space in log

idempotency

logged operations should be *okay to do twice* = *idempotent*

good example: set inode link count to 4

bad example: increment inode link count

good example: overwrite inode number X with new value
as long as last committed inode value in log is right...

bad example: allocate new inode with particular contents

good example: overwrite data block with new value

bad example: append data to last used block of file

redo logging summary

write intended operation to the log

before ever touching 'real' data
in format that's safe to do twice

write marker to commit to the log

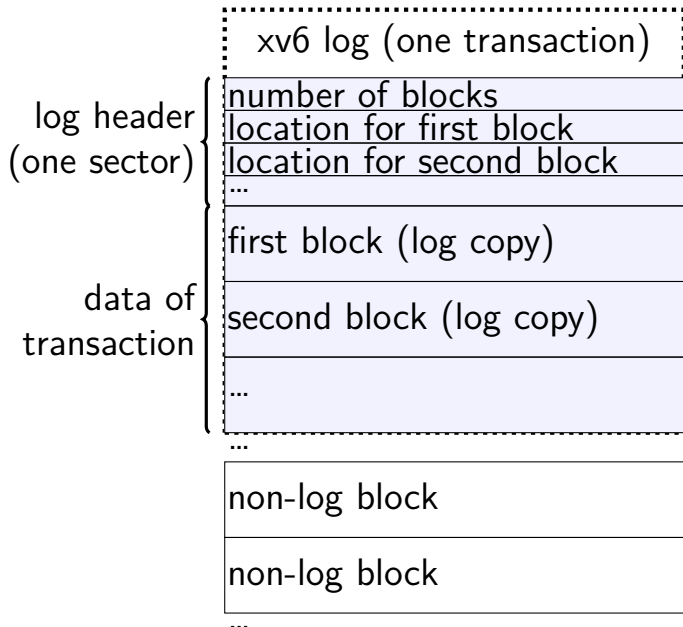
if exists, the operation *will be done eventually*

actually update the real data

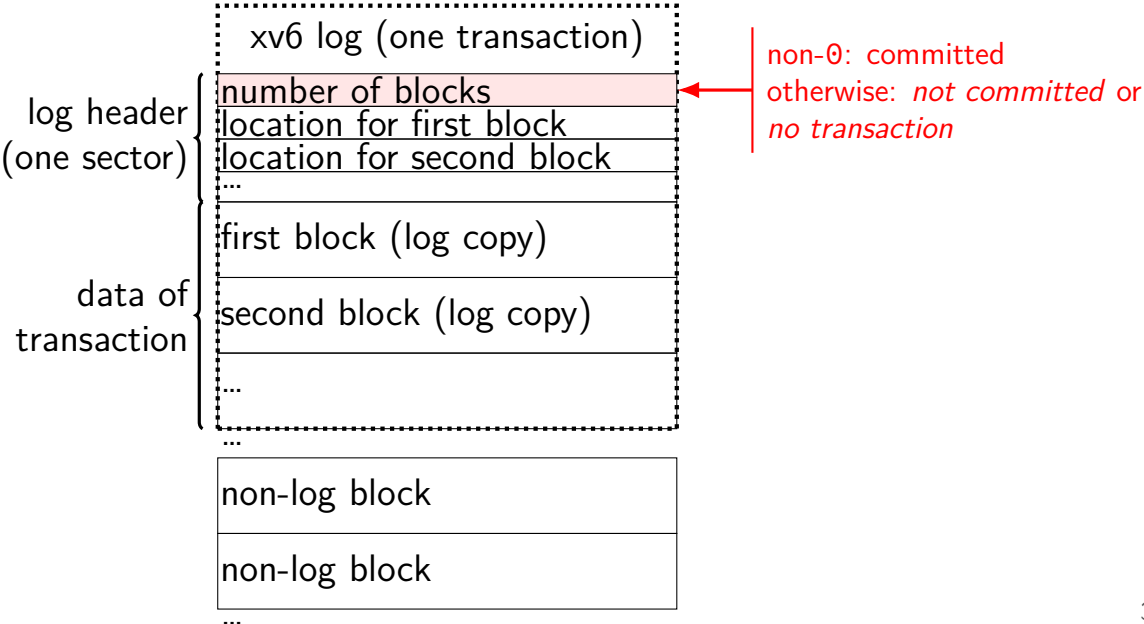
redo logging and filesystems

filesystems that do redo logging are called *journaling filesystems*

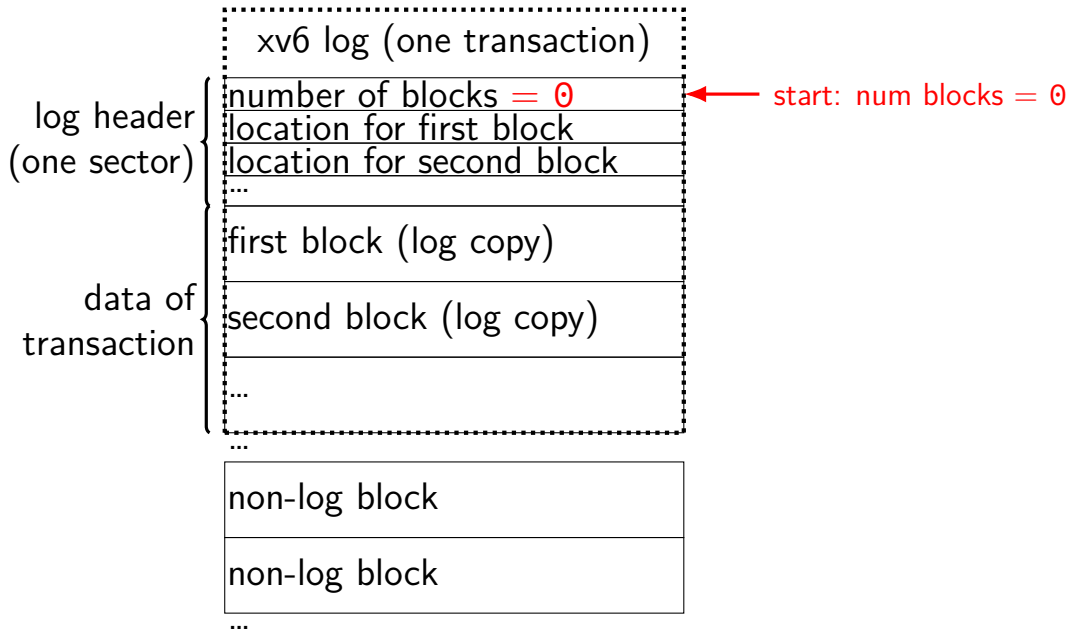
the xv6 journal



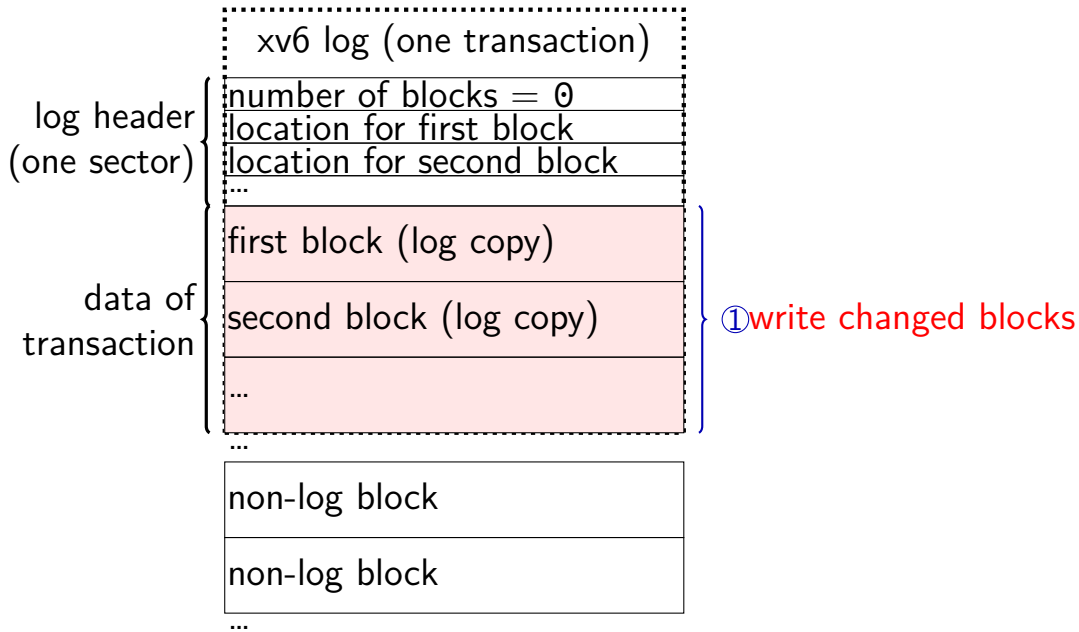
the xv6 journal



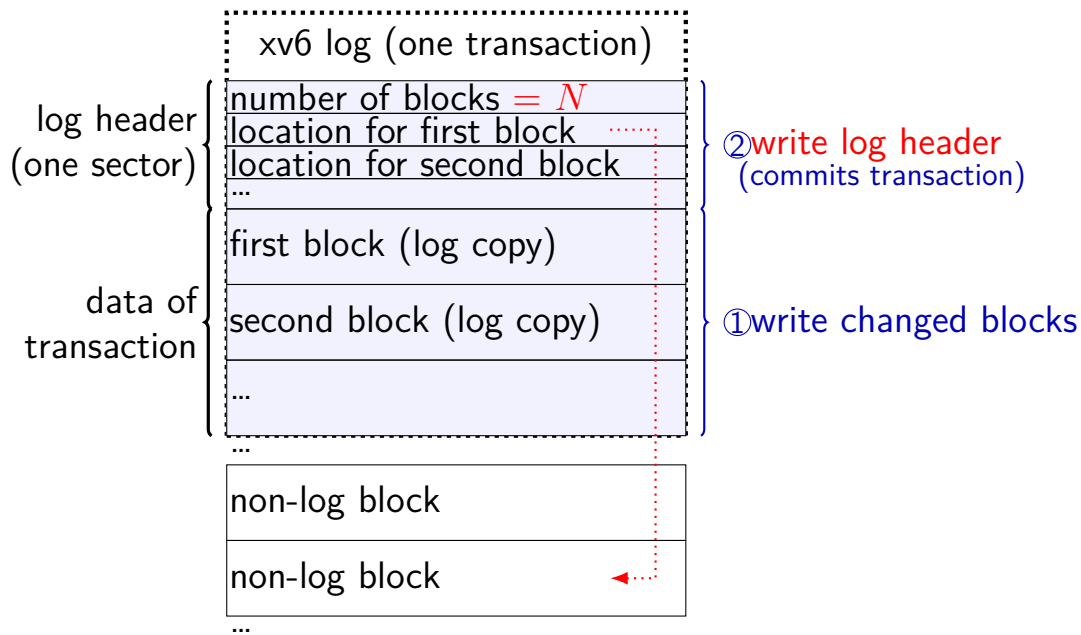
the xv6 journal



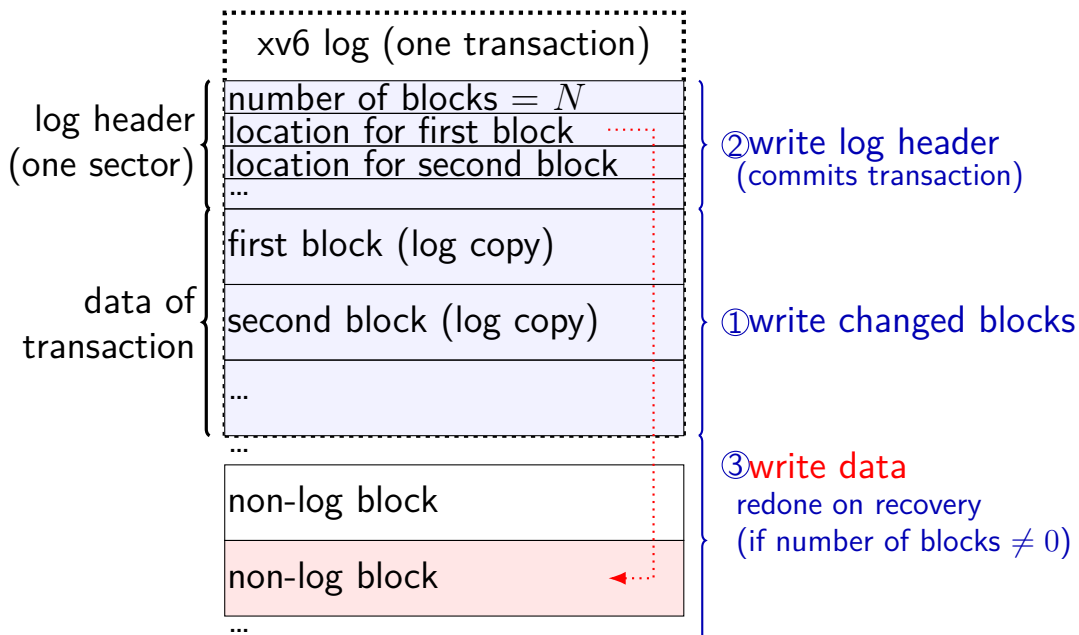
the xv6 journal



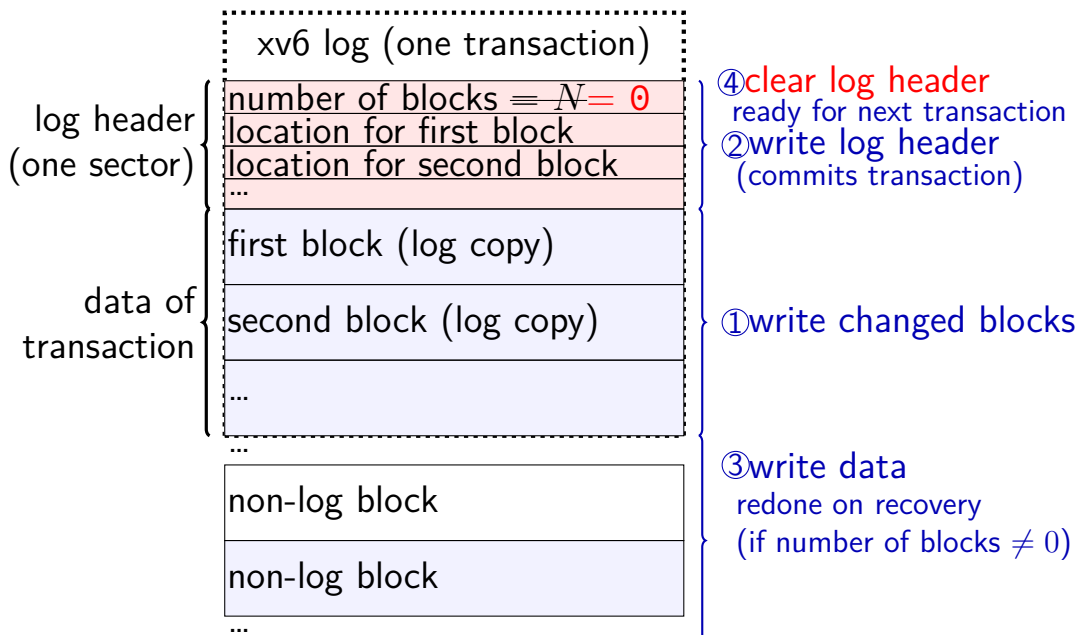
the xv6 journal



the xv6 journal



the xv6 journal



what is a transaction?

so far: each file update?

faster to do batch of updates together

- one log write finishes lots of things
- don't wait to write

xv6 solution: combine lots of updates into one transaction

only commit when...

- no active file operation, *or*
- not enough room left in log for more operations

what is a transaction?

so far: each file update?

faster to do **batch of updates together**

- one log write finishes lots of things
- don't wait to write

xv6 solution: combine lots of updates into one transaction

only commit when...

- no active file operation, *or*
- not enough room left in log for more operations

redo logging problems

doesn't the log get infinitely big?

writing everything twice?

redo logging problems

doesn't the log get infinitely big?

writing everything twice?

limiting log size

once transaction is written to real data, can discard

sometimes called “garbage collecting” the log

may sometimes need to block to free up log space

perform logged updates before adding more to log

hope: usually log cleanup happens “in the background”

redo logging problems

doesn't the log get infinitely big?

writing everything twice?

lots of writing? (1)

entire log can be **written sequentially**

- ideal for hard disk performance
- also pretty good for SSDs

multiple updates can be done **in any order**

- can reorder to minimize seek time/rotational latency/etc.
- can interleave updates that make up multiple transactions

no waiting for 'real' updates

- application can proceed while updates are happening
- files will be updated even if system crashes

often better for performance!

lots of writing? (2)

updating 1000 files?

with redo logging — 2 big seeks

- write all updates to log in order

- write all updates to file/inode/directory data in order

lots of writing? (2)

updating 1000 files?

with redo logging — 2 big seeks

- write all updates to log in order

- write all updates to file/inode/directory data in order

careful ordering — lots of seeks?

- write to free block map

- seek + write to inode

- seek + write to directory entry

- repeat 1000x

maybe could also combine file updates with careful ordering??

- but sure starts to get complicated to track order requirements

- redo logging is probably simpler?

degrees of consistency

not all journalling filesystem use redo logging for everything

some use it *only for metadata operations*

some use it *for both metadata and user data*

only metadata: avoids lots of duplicate writing

metadata+user data: integrity of user data guaranteed

snapshots

filesystem snapshots

idea: filesystem keeps old versions of files around

accidental deletion? old version still there

eventually discard some old versions

can access *snapshot* of files at prior time

snapshots

filesystem snapshots

idea: filesystem keeps old versions of files around
accidental deletion? old version still there
eventually discard some old versions

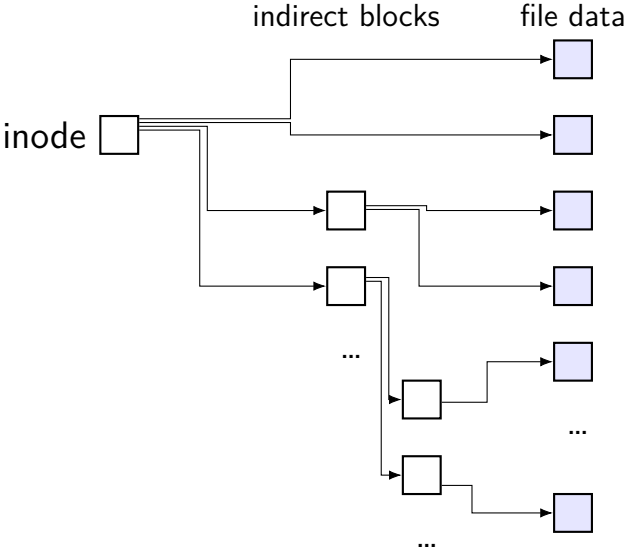
can access *snapshot* of files at prior time

mechanism: **copy-on-write**

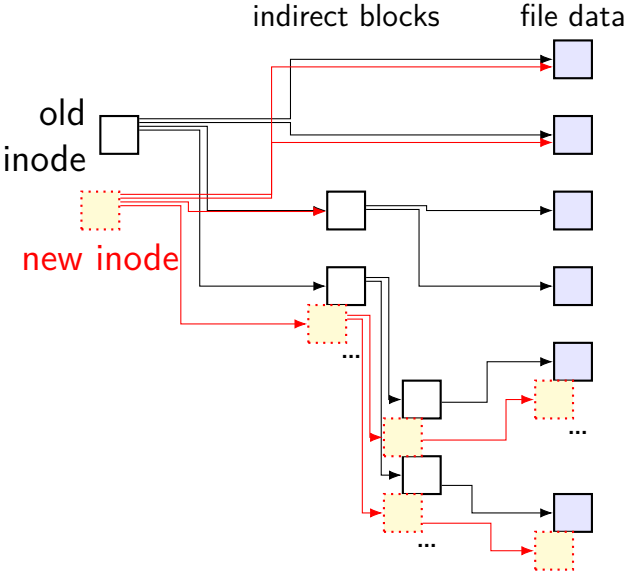
changing file makes **new copy** of filesystem

common parts shared between versions

inode and copy-on-write



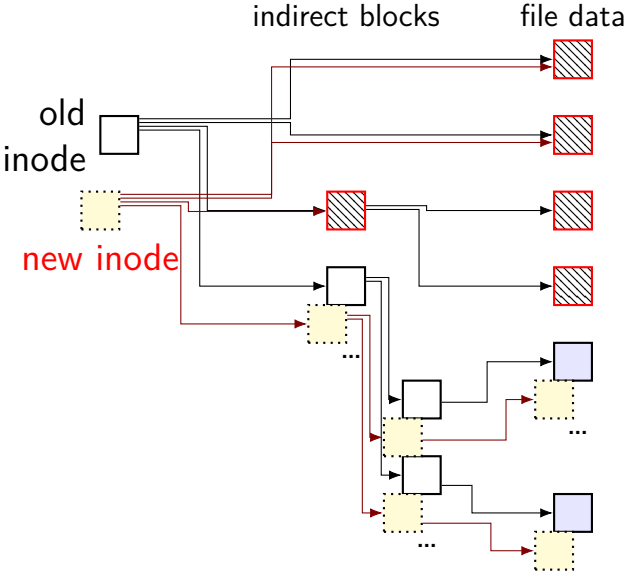
inode and copy-on-write



update: new data blocks
+ new indirect blocks
+ new inode

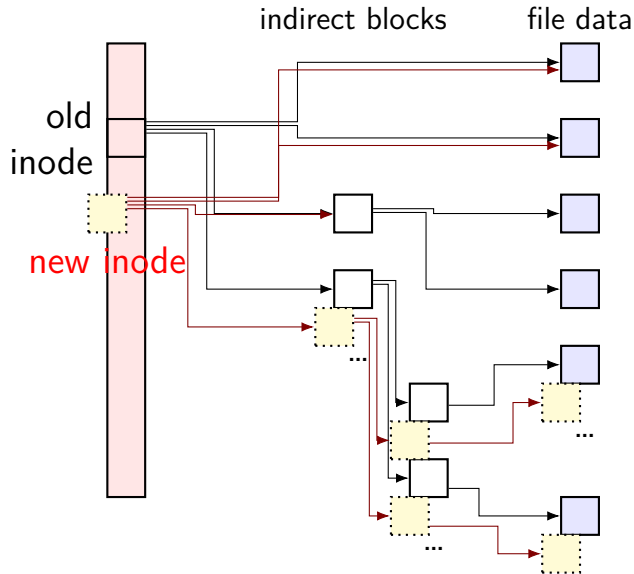
both old+new inode valid

inode and copy-on-write



unchanged parts of file shared

inode and copy-on-write

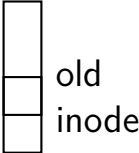


challenge: FFS/xv6/ext2 design
has big array of inodes

don't want to write new copy
of *entire inode array*

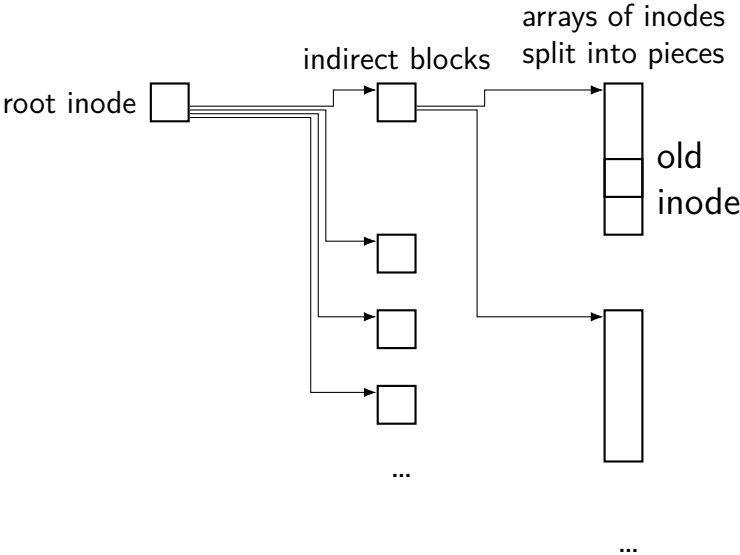
extra indirection for inode array

arrays of inodes
split into pieces

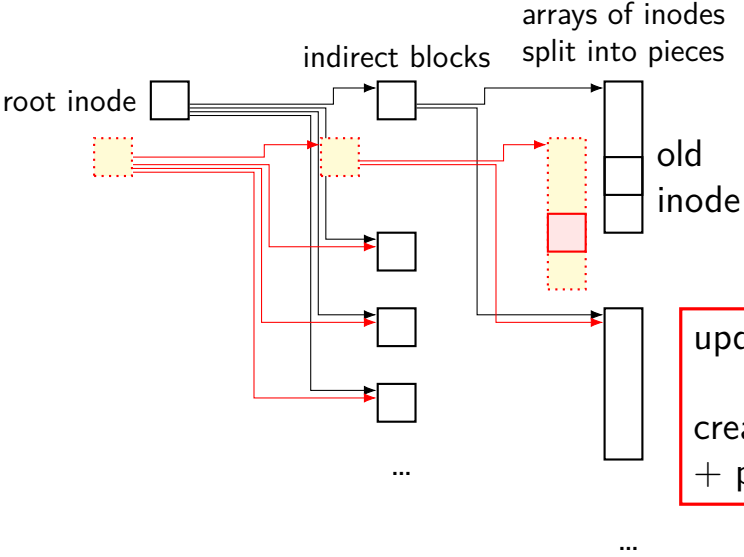


...

extra indirection for inode array

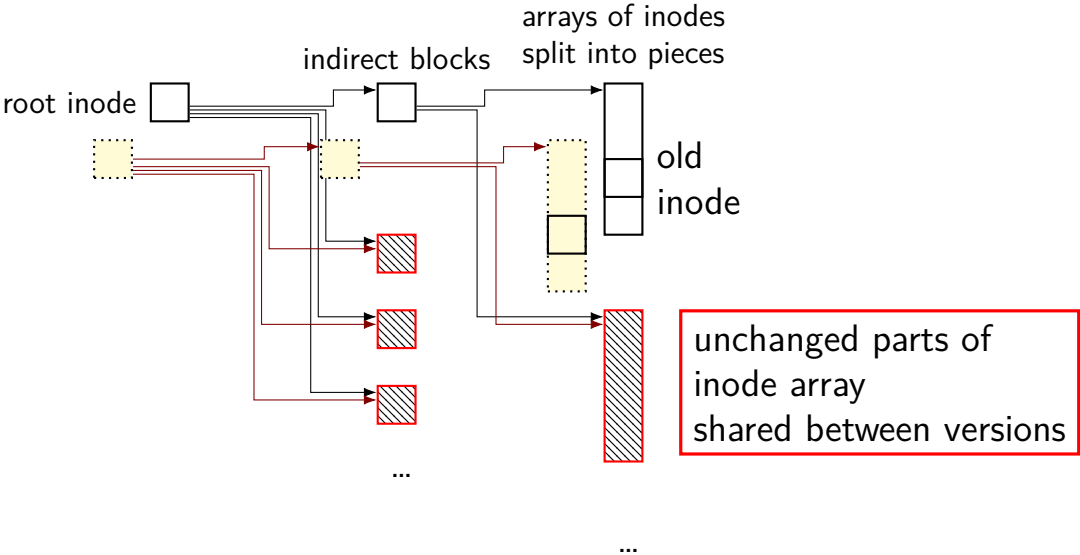


extra indirection for inode array

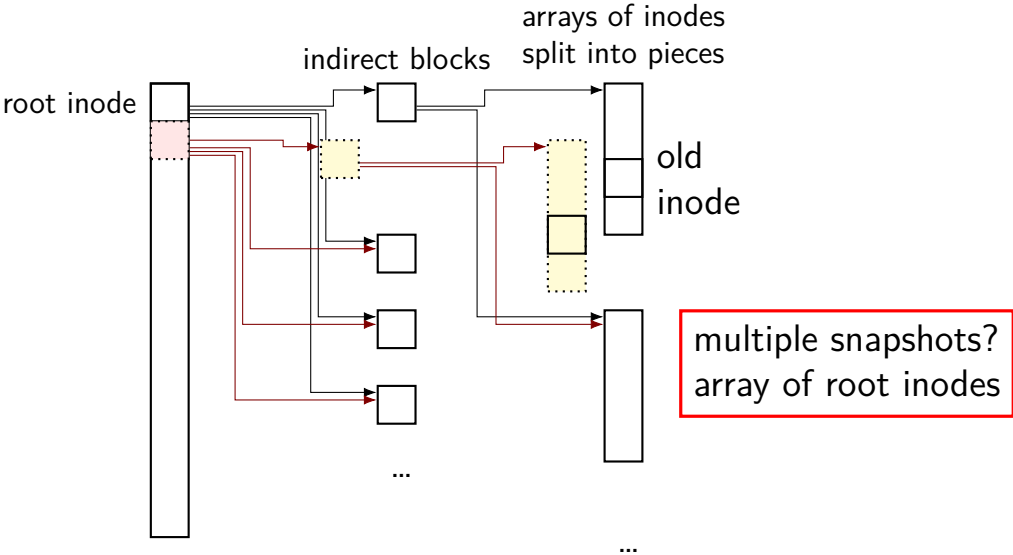


update one inode?
create new root inode
+ pointers

extra indirection for inode array



extra indirection for inode array



copy-on-write indirection

file update = replace with new version

array of **versions of entire filesystem**

only copy modified parts

keep reference counts, like for paging assignment

lots of pointers — only change pointers where modifications happen

snapshots in practice

ZFS (used on department machines) implements this

example: `.zfs/snapshots/11.11.18-06` pseudo-directory

contains contents of files at 11 November 2018 6AM

copy-on-write and logging

copy-on-write is a nice solution to duplicate writes

before (data journalling)

- write new data to journal

- copy new data to real location

after (copy-on-write)

- write new data to new location

- update pointer to point to new location

useful even without snapshots

- but maybe not keeping file data in best place?

aside: fsync

filesystem can order things carefully

filesystem can make sure data on disk before proceeding

what if I, non-OS programmer want to do that?

POSIX mechanism: `fsync`

“please actually write this file to disk now — I’ll wait”

some stories of broken implementations of `fsync`

nasty problem — how do you test it???

some varying interpretations

some only send to disk, but *don't wait for disk to finish writing*

does not guarantee updating file's directory entry

mounting filesystems

Unix-like system

root filesystem appears as /

other filesystems *appear as directory*

e.g. lab machines: my home dir is in filesystem at /net/zf15

directories that are filesystems look like normal directories

/net/zf15/.. is /net (even though in different filesystems)

mounts on a dept. machine

```
/dev/sda1 on / type ext4 (rw,errors=remount-ro)
proc on /proc type proc (rw,noexec,nosuid,nodev)
...
udev on /dev type devtmpfs (rw,mode=0755)
devpts on /dev/pts type devpts (rw,noexec,nosuid,gid=5,mode=0620)
tmpfs on /run type tmpfs (rw,noexec,nosuid,size=10%,mode=0755)
...
/dev/sda3 on /localtmp type ext4 (rw)
...
zfs1:/zf2 on /net/zf2 type nfs (rw,hard,intr,proto=udp,nfsvers=3,
                               noacl,sloppy,addr=128.143.136.9)
zfs3:/zf19 on /net/zf19 type nfs (rw,hard,intr,proto=udp,nfsvers=3,
                                   noacl,sloppy,addr=128.143.67.236)
zfs4:/sw on /net/sw type nfs (rw,hard,intr,proto=udp,nfsvers=3,
                               noacl,sloppy,addr=128.143.136.9)
zfs3:/zf14 on /net/zf14 type nfs (rw,hard,intr,proto=udp,nfsvers=3,
                                   noacl,sloppy,addr=128.143.67.236)
...
```

kernel FS abstractions

Linux: *virtual file system* API

object-oriented, based on FFS-style filesystem

to implement a filesystem, create object types for:

- superblock (represents “header”)

- inode (represents file)

- dentry (represents cached directory entry)

- file (represents *open file*)

common code handles directory traversal
and caches directory traversals

common code handles file descriptors, etc.

linux VFS operations

superblock: write_inodez, sync_fs, ...

inode: create, link, unlink, mkdir, open ...
most just for inodes which are directories

dentry: compare, delete ...
more commonly argument to inode operation
can be created for non-yet-existing files

file: read, write, ...

linux VFS operations example

```
struct inode_operations {
    struct dentry * (*lookup) (struct inode *,struct dentry *, unsigned
    ...
    int (*create) (struct inode *,struct dentry *, umode_t, bool);
    int (*link) (struct dentry *,struct inode *,struct dentry *);
    int (*unlink) (struct inode *,struct dentry *);
    int (*symlink) (struct inode *,struct dentry *,const char *);
    int (*mkdir) (struct inode *,struct dentry *,umode_t);
    int (*rmdir) (struct inode *,struct dentry *);
    int (*mknod) (struct inode *,struct dentry *,umode_t,dev_t);
    int (*rename) (struct inode *, struct dentry *,
                    struct inode *, struct dentry *, unsigned int);
    ...
    int (*update_time)(struct inode *, struct timespec64 *, int);
    int (*atomic_open)(struct inode *, struct dentry *,
                       struct file *, unsigned open_flag,
                       umode_t create_mode);
    ..
}
```

FS abstractions and awkward FSes

example: inode object for FAT?

fake it: point to directory entry?

backup/if time slides

log-structured filesystems

logging is a great access pattern for hard drives and SSDs

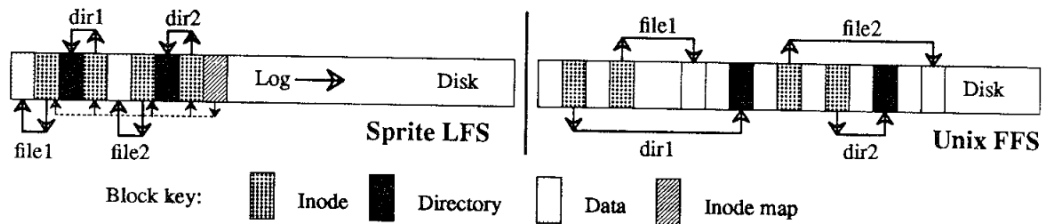
sequential

right for SSDs — write everything once before writing again

how about designing a filesystem around it!

idea: log-structured filesystems

log-structured filesystem



log-structured filesystem ideas

write inodes + data + free map + etc. to log instead of disk

problem: scanning log to find latest version of inode?

periodically write *inode maps* to log
 computed latest location of inodes

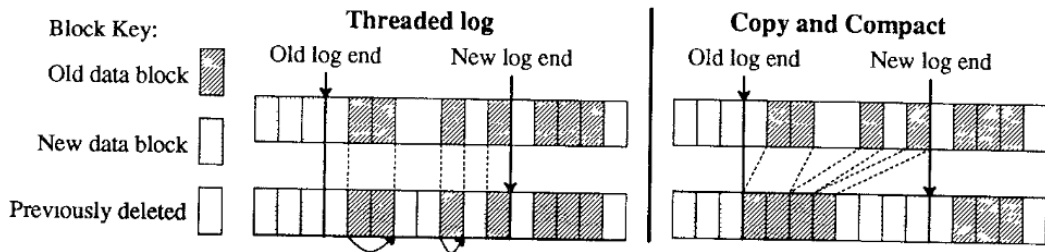
searching limited to last inode map

log-structured FS garbage collection

challenge: what happens when log gets to the end of the disk?

want to start from beginning of disk again...

either: copy data to free space or 'thread' log around used space:



log-structured filesystems in practice

the kind of ideas you'd use to implement an SSD

used for some filesystems that work directly with Flash chips

changing file atomically?

often applications want to update a file all at once

changing file atomically?

often applications want to update a file all at once

on Unix, one way to do this:

create a new file with a hard-to-guess name in the same directory

rename the new file to replace the old file

overwrites that directory entry

no one will ever read partially written file