

Changelog

12 April 2022: change xv6 file sizes slide to correct block pointer size

last time

FAT filesystem

header — global/general information

divide disk into clusters (possibly larger than sector)

files+directory 'data' stored in one+ whole clusters

linked list identifies clusters making up file/directory

file allocation table — one number per (potential data) cluster

- next pointers for linked list

- indicate which clusters are free

directory entries

- 'data' for directories

- starting location, name, etc. about file/dir in directory

on debugging issues

having better tests for paging/protection helped less than I hoped

too many students not getting to the point of looking at tests

tests not good enough at diagnosing certain types of memory corruption

- e.g. freeing non-heap pages incorrectly

probably should have discouraged students from modifying `kfree()`

- more elegant/less code, but harder to debug than other options

on office hour queues

too much OH time per student given number waiting

we can't spend 15 uninterrupted minutes/student with one TA + 15 students waiting

TAs/I sometimes have trouble switching away from students (e.g. while they're gathering debugging info, when they get to a new problem)
adjusting the queue ordering can't really fix this issue

we can be more useful with better problems

we're less useful when student has done less investigation of what ran before crash/etc.

need faster switching between students

ideally can help other students while waiting for student to add debugging/etc.

in practice: especially on Discord, doesn't happen as much as I hoped

xv6 filesystem

xv6's filesystem similar to modern Unix filesystems

better at doing contiguous reads than FAT

better at handling crashes

supports *hard links*

divides disk into *blocks* instead of clusters

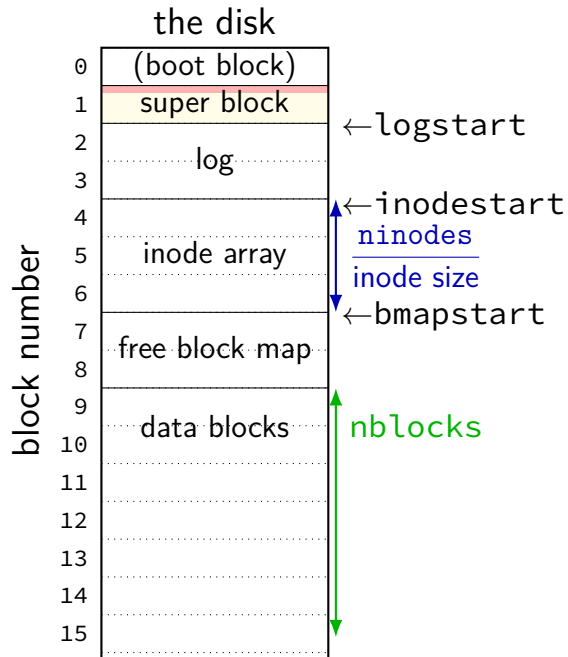
file block numbers, free blocks, etc. in different tables

xv6 disk layout

the disk

0	(boot block)
1	super block
2	log
3	
4	inode array
5	
6	
7	free block map
8	
9	data blocks
10	
11	
12	
13	
14	
15	

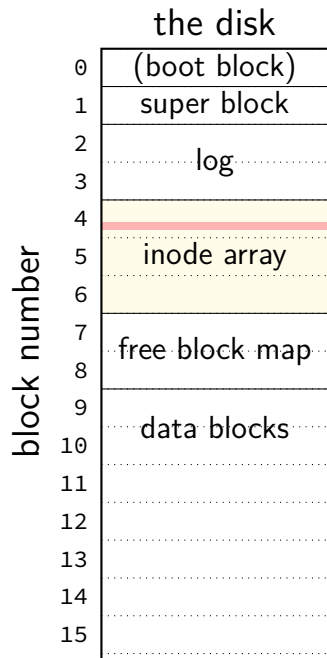
xv6 disk layout



superblock — “header”

```
struct superblock {  
    uint size;  
        // Size of file system image (b  
    uint nblocks;  
        // # of data blocks  
    uint ninodes;  
        // # of inodes  
    uint nlog;  
        // # of log blocks  
    uint logstart;  
        // block # of first log block  
    uint inodestart;  
        // block # of first inode block  
    uint bmapstart;  
        // block # of first free map bl  
};
```

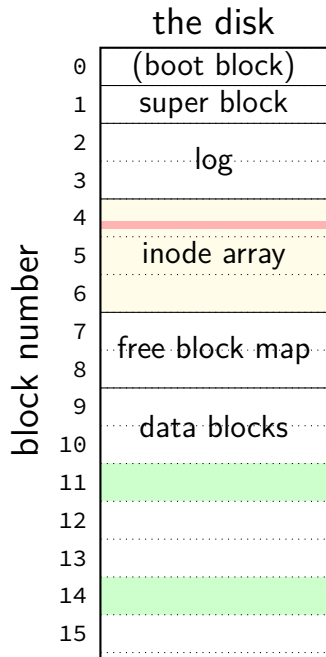

xv6 disk layout



inode — file information

```
struct dinode {  
    short type; // File type  
              // T_DIR, T_FILE, T_DEV  
  
    short major; short minor; // T_DEV only  
  
    short nlink;  
    // Number of links to inode in file syst  
    uint size; // Size of file (bytes)  
    uint addrs[NDIRECT+1];  
    // Data block addresses  
};
```

xv6 disk layout



inode — file information

```
struct dinode {  
    short type; // File type  
                // T_DIR, T_FILE, T_DEV  
  
    short major; short minor; // T_DEV only  
  
    short nlink;  
    // Number of links to inode in file syst  
    uint size; // Size of file (bytes)  
    uint addrs[NDIRECT+1];  
    // Data block addresses  
};
```

location of data as block numbers:
e.g. `addrs[0] = 11; addrs[1] = 14;`
special case for larger files

xv6 disk layout

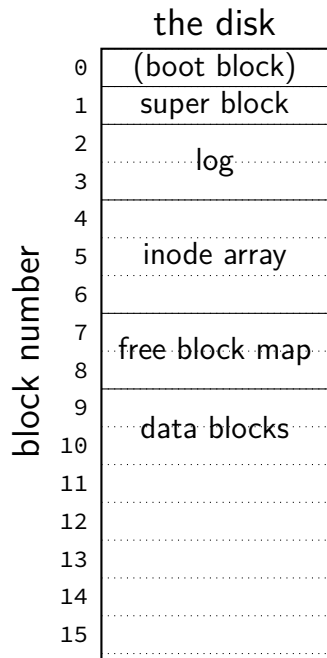
the disk

0	(boot block)
1	super block
2	log
3	
4	inode array
5	
6	free block map
7	
8	data blocks
9	
10	
11	
12	
13	
14	
15	

free block map — 1 bit per data block
1 if available, 0 if used

allocating blocks: scan for 1 bits
contiguous 1s — contiguous blocks

xv6 disk layout



what about finding free inodes

xv6 solution: scan for type = 0

typical Unix solution: separate free inode map

xv6 directory entries

```
struct dirent {  
    ushort inum;  
    char name[DIRSIZ];  
};
```

inum — index into inode array on disk

name — name of file or directory

each directory reference to inode called a *hard link*
multiple hard links to file allowed!

xv6 allocating inodes/blocks

need new inode or data block: linear search

simplest solution: xv6 always takes the first one that's free

xv6 FS pros versus FAT

- support for reliability — log
 - more on this later

- possibly easier to scan for free blocks
 - more compact free block map

- easier to find location of k th block of file
 - element of addrs array

- file type/size information held with block locations
 - inode number = everything about open file
 - easier to read/modify file info all at once?

missing pieces

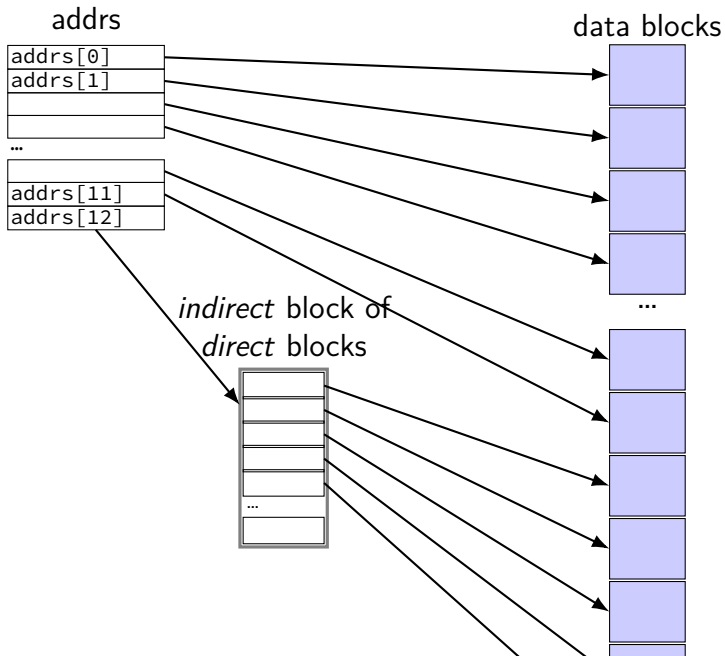
what's the log? (more on that later)

other file metadata?

creation times, etc. — xv6 doesn't have it

not good at taking advantage of HDD architecture

xv6 inode: direct and indirect blocks



xv6 file sizes

512 byte blocks

2-byte block pointers: 256 block pointers in the indirect block

256 blocks = 131072 bytes of data referenced

12 direct blocks @ 512 bytes each = 6144 bytes

1 indirect block @ 131072 bytes each = 131072 bytes

maximum file size = 6144 + 131072 bytes

Linux ext2 inode

```
struct ext2_inode {  
    __le16 i_mode;           /* File mode */  
    __le16 i_uid;            /* Low 16 bits of Owner Uid */  
    __le32 i_size;           /* Size in bytes */  
    __le32 i_atime;          /* Access time */  
    __le32 i_ctime;          /* Creation time */  
    __le32 i_mtime;          /* Modification time */  
    __le32 i_dtime;          /* Deletion Time */  
    __le16 i_gid;            /* Low 16 bits of Group Id */  
    __le16 i_links_count;     /* Links count */  
    __le32 i_blocks;         /* Blocks count */  
    __le32 i_flags;          /* File flags */  
    ...  
    __le32 i_block[EXT2_N_BLOCKS]; /* Pointers to blocks */  
    ...  
};
```

Linux ext2 inode

```
struct ext2_inode {
    __le16 i_mode;           /* File mode */
    __le16 i_uid;            /* Low 16 bits of Owner Uid */
    __le32 i_size;           /* Size in bytes */
    __le32 i_atime;          /* Access time */
    __le32 i_ctime;          /* Creation time */
    -- type (regular, directory, device)
    -- and permissions (read/write/execute for owner/group/others)
    __le16 i_links_count;     /* Links count */
    __le32 i_blocks;         /* Blocks count */
    __le32 i_flags;          /* File flags */
    ...
    __le32 i_block[EXT2_N_BLOCKS]; /* Pointers to blocks */
    ...
};
```

Linux ext2 inode

```
struct ext2_inode {  
    __le16 i_mode;           /* File mode */  
    __le16 i_uid;            /* Low 16 bits owner and group */  
    __le32 i_size;           /* Size in bytes */  
    __le32 i_atime;          /* Access time */  
    __le32 i_ctime;          /* Creation time */  
    __le32 i_mtime;          /* Modification time */  
    __le32 i_dtime;          /* Deletion Time */  
    __le16 i_gid;            /* Low 16 bits of Group Id */  
    __le16 i_links_count;     /* Links count */  
    __le32 i_blocks;         /* Blocks count */  
    __le32 i_flags;          /* File flags */  
    ...  
    __le32 i_block[EXT2_N_BLOCKS]; /* Pointers to blocks */  
    ...  
};
```

Linux ext2 inode

```
struct ext2_inode {  
    __le16 i_mode;           /* File mode */  
    __le16 i_uid;            /* Low 16 bits of user id of file  
                               whole bunch of times  
    __le32 i_size;           /* Size in bytes */  
    __le32 i_atime;          /* Access time */  
    __le32 i_ctime;          /* Creation time */  
    __le32 i_mtime;          /* Modification time */  
    __le32 i_dtime;          /* Deletion Time */  
    __le16 i_gid;            /* Low 16 bits of Group Id */  
    __le16 i_links_count;    /* Links count */  
    __le32 i_blocks;         /* Blocks count */  
    __le32 i_flags;          /* File flags */  
    ...  
    __le32 i_block[EXT2_N_BLOCKS]; /* Pointers to blocks */  
    ...  
};
```

Linux ext2 inode

```
struct ext2_inode {  
    __le16 i_mod;   
    __le16 i_uid;   
    __le32 i_size;   
    __le32 i_atime;   
    __le32 i_ctime;   
    __le32 i_mtime;   
    __le32 i_dtime;   
    __le16 i_gid;   
    __le16 i_links_count;   
    __le32 i_blocks;   
    __le32 i_flags;   
    ...  
    __le32 i_block[EXT2_N_BLOCKS];   
    ...  
};
```

similar pointers like xv6 FS — but more indirection

/ Size in bytes */*

/ Access time */*

/ Creation time */*

/ Modification time */*

/ Deletion Time */*

/ Low 16 bits of Group Id */*

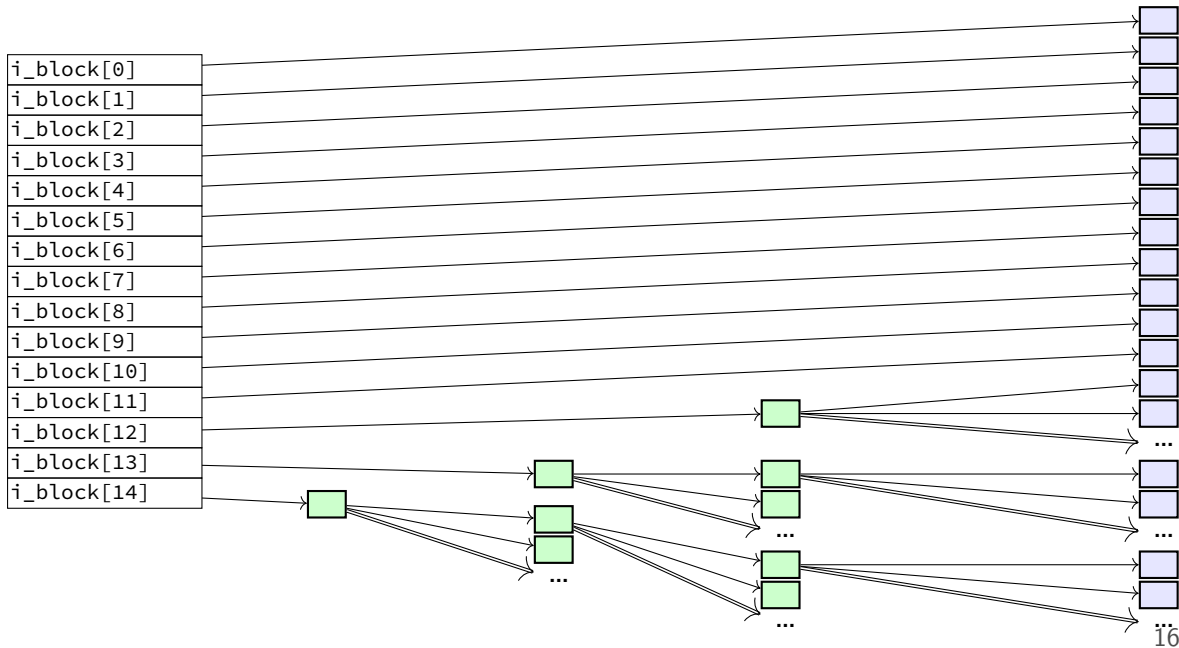
/ Links count */*

/ Blocks count */*

/ File flags */*

/ Pointers to blocks */*

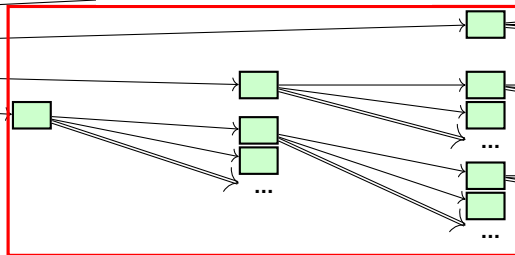
double/triple indirect



double/triple indirect

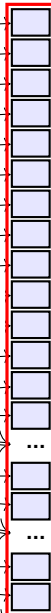
block pointers

i_block[0]
i_block[1]
i_block[2]
i_block[3]
i_block[4]
i_block[5]
i_block[6]
i_block[7]
i_block[8]
i_block[9]
i_block[10]
i_block[11]
i_block[12]
i_block[13]
i_block[14]



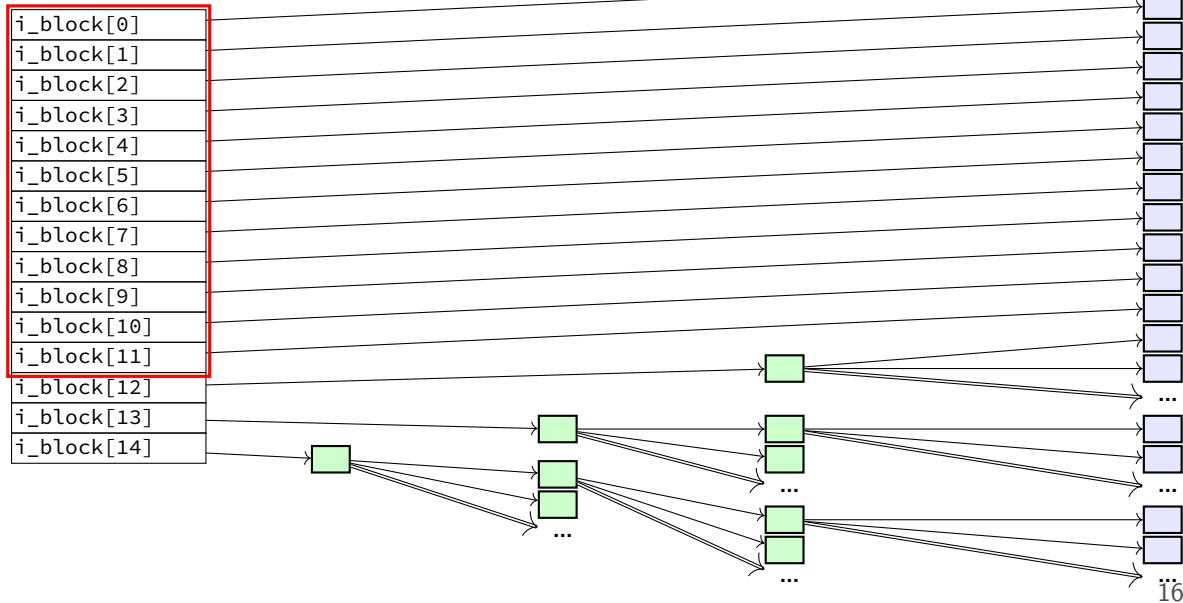
blocks of block pointers

data blocks

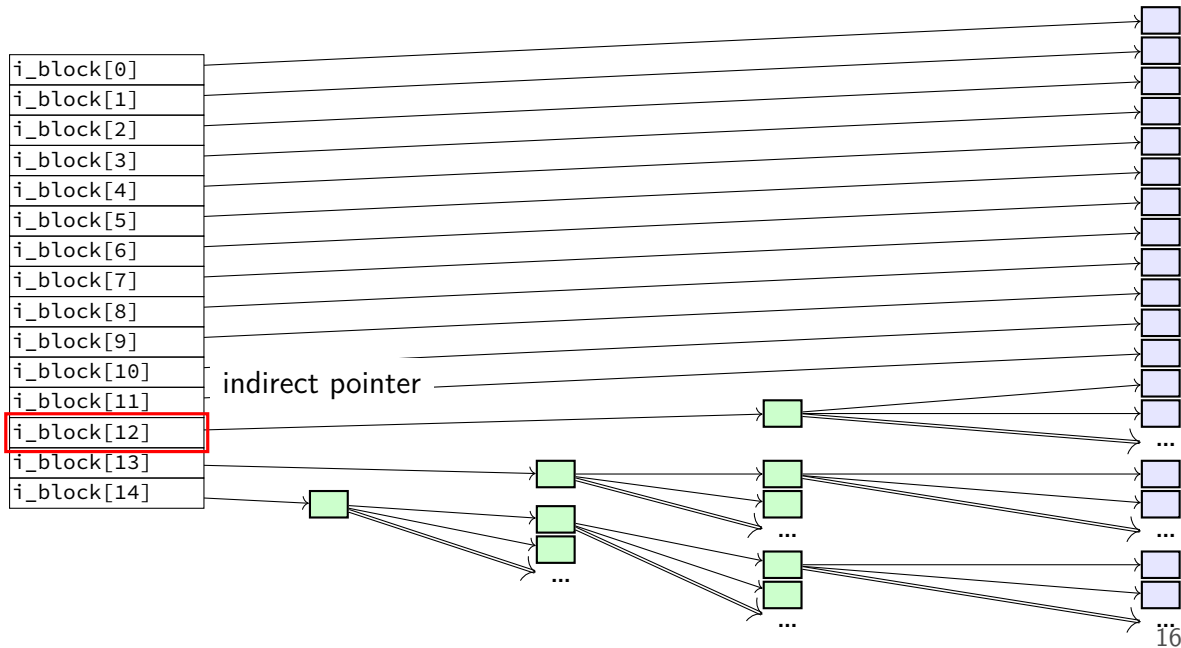


double/triple indirect

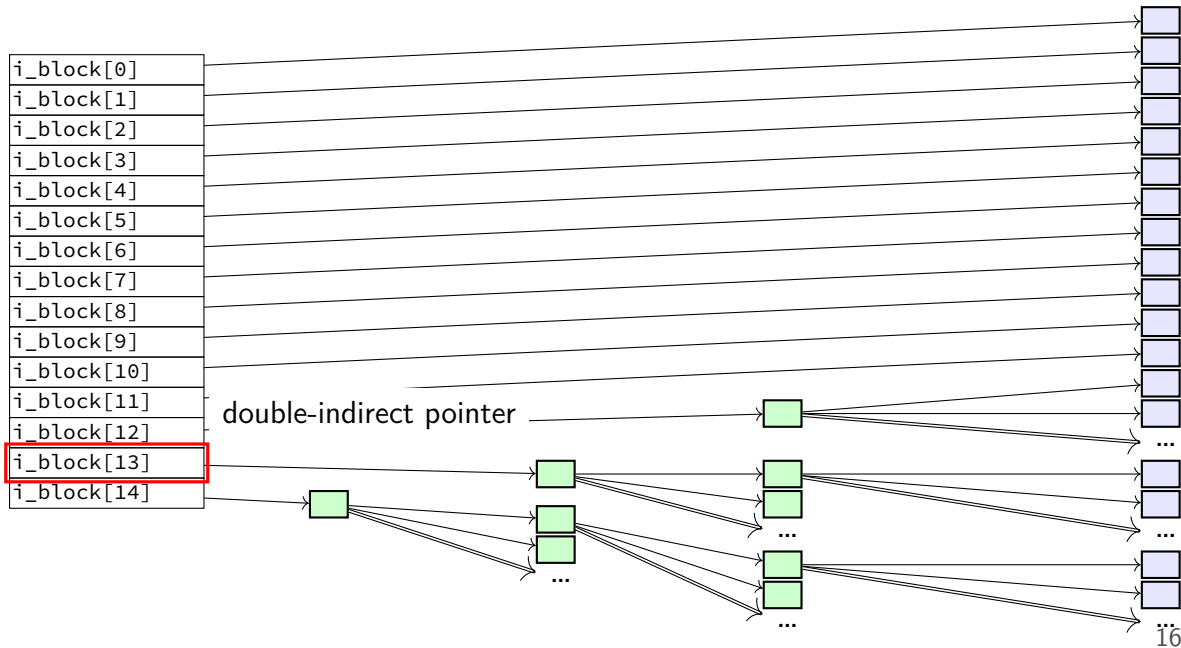
12 direct pointers



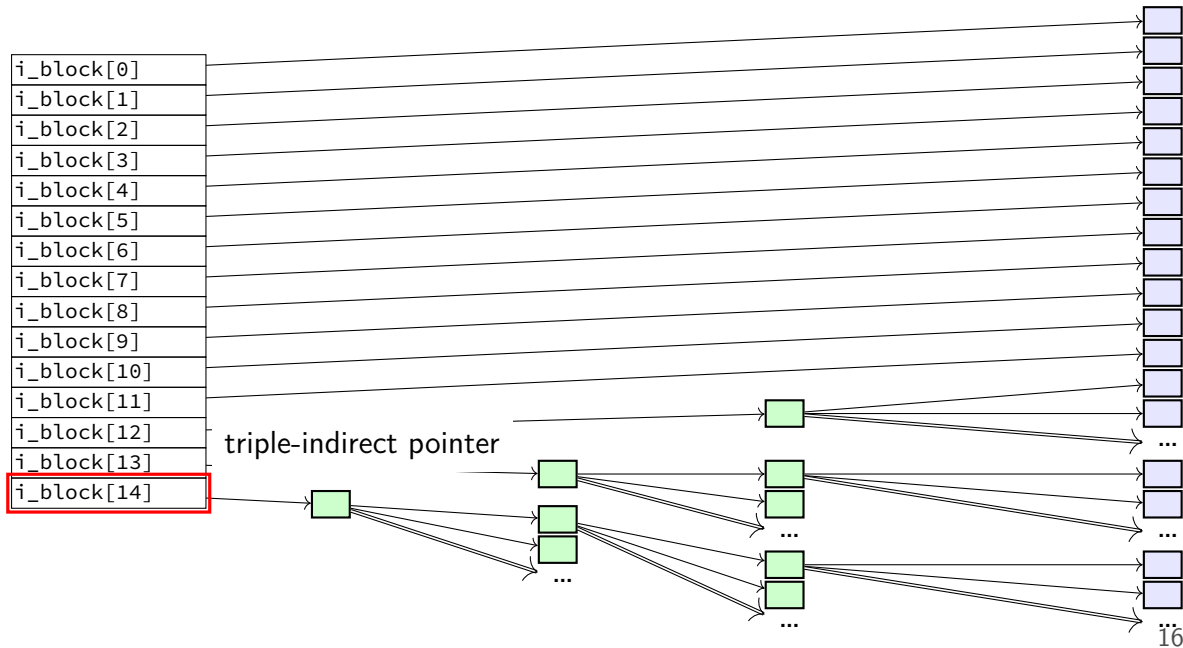
double/triple indirect



double/triple indirect



double/triple indirect



ext2 indirect blocks (1)

12 direct block pointers

1 indirect block pointer

pointer to block containing more direct block pointers

1 double indirect block pointer

pointer to block containing more indirect block pointers

1 triple indirect block pointer

pointer to block containing more double indirect block pointers

ext2 indirect blocks (1)

12 direct block pointers

1 indirect block pointer

pointer to block containing more direct block pointers

1 double indirect block pointer

pointer to block containing more indirect block pointers

1 triple indirect block pointer

pointer to block containing more double indirect block pointers

exercise: if 1K blocks, 4 byte block pointers, how big can a file be?

ext2 indirect blocks (solution)

12 direct pointers: first 1K (block size) \times 12 bytes of data

1 indirect pointer:

points to block with $1\text{K (block size)} / 4 \text{ byte (pointer size)} = 256$ pointers

256 pointers point to 1K blocks

next 256KB of data

1 double indirect pointer

points to block with $1\text{K (block size)} / 4 \text{ byte (pointer size)} = 256$ pointers

256 pointers point to pointers that each are like an indirect pointer

256KB per indirect pointer \rightarrow next $256 \cdot 256$ KB of data

1 triple indirect

next $256 \cdot 256 \cdot 256$ KB of data

total size: $12 + 256 + 256^2 + 256^3$ KB $= 16843020$ KB $\approx 16\text{GB}$

ext2 indirect blocks (2)

12 direct block pointers

1 indirect block pointer

1 double indirect block pointer

1 triple indirect block pointer

exercise: if 1K (2^{10} byte) blocks, 4 byte block pointers,
how does OS find byte 2^{15} of the file?

- (1) using indirect pointer or double-indirect pointer in inode?
- (2) what index of block pointer array pointed to by pointer in inode?

ext2 indirect blocks (2) (solution)

byte $2^{15} = 32\text{KB}$ into file

12 direct pointers: first 1K (block size) \times 12 bytes of data

1 indirect pointer:

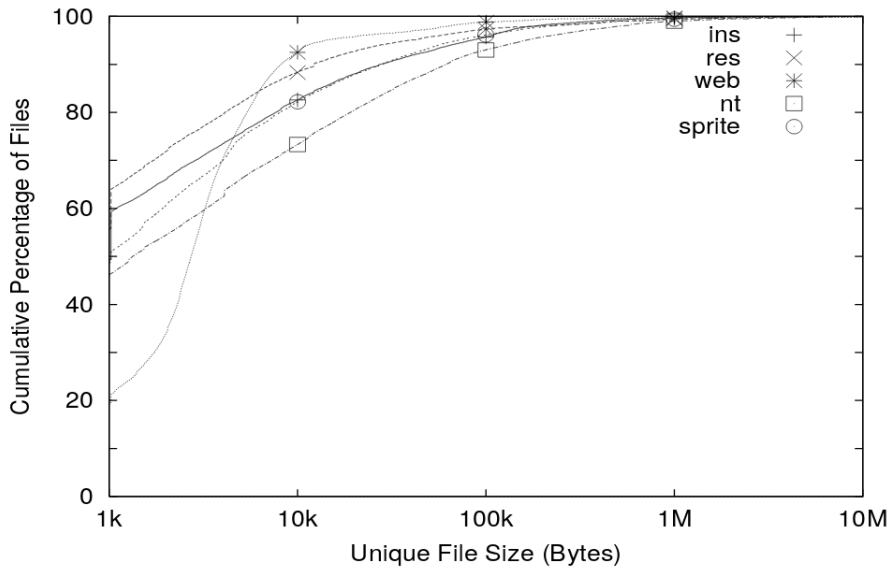
points to block with $1\text{K (block size)} / 4 \text{ byte (pointer size)} = 256$ pointers

256 pointers point to 1K blocks

next 256KB of data

going to be (32 - 12)th element

empirical file sizes



typical file sizes

most files are small

- sometimes 50+% less than 1kbyte

- often 80-95% less than 10kbyte

doesn't mean large files are unimportant

- still take up most of the space

- biggest performance problems

extents

large file? lists of many thousands of blocks is awkward
...and requires multiple reads from disk to get

solution: store **extents**: (start disk block, size)
replaces or supplements block list

Linux's ext4 and Windows's NTFS both use this

allocating extents

challenge: finding contiguous sets of free blocks

NTFS: scan block map for “best fit”

- look for big enough chunk of free blocks

- choose smallest among all the candidates

don't find any? okay: use more than one extent

seeking with extents

challenge: finding byte X of the file

with block pointers: can compute index

with extents: need to scan list?

filesystem reliability

a crash happens — what's the state of my filesystem?

hard disk atomicity

interrupt a hard drive write?

write whole disk sector or corrupt it

hard drive/SSD stores checksum for each sector

write interrupted? — checksum mismatch

hard drive/SSD returns read error

reliability issues

is the filesystem in a consistent state?

- do we know what blocks are free?

- do we know what files exist?

- is the data for files actually what was written?

also important topics, but won't spend much time on these:

what data will I lose if storage fails?

- mirroring, erasure coding (e.g. RAID) — using multiple storage devices

- idea: if one storage device fails, other(s) still have data

what data will I lose if I make a mistake?

- filesystem can store *multiple versions*

- “snapshots” of what was previously there

several bad options (1)

suppose we're moving a file from one directory to another on xv6
steps:

A: write new directory entry

B: overwrite (remove) old directory entry

several bad options (1)

suppose we're moving a file from one directory to another on xv6
steps:

A: write new directory entry

B: overwrite (remove) old directory entry

if we do A before B and crash happens after A:

- can have extra pointer of file

- problem: if old directory entry removed later, will get confused and free the file!

several bad options (1)

suppose we're moving a file from one directory to another on xv6
steps:

A: write new directory entry

B: overwrite (remove) old directory entry

if we do A before B and crash happens after A:

- can have extra pointer of file

- problem: if old directory entry removed later, will get confused and free the file!

if we do B before A and crash happens after B:

- the file disappeared entirely!

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

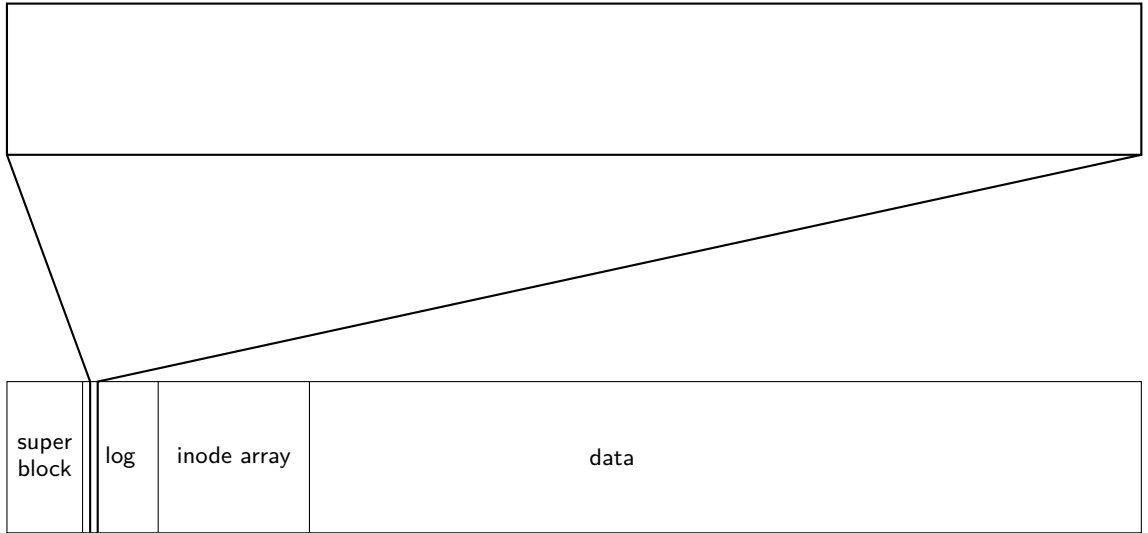
transaction: bunch of updates that happen all at once

implementation trick: one update means transaction “commits”

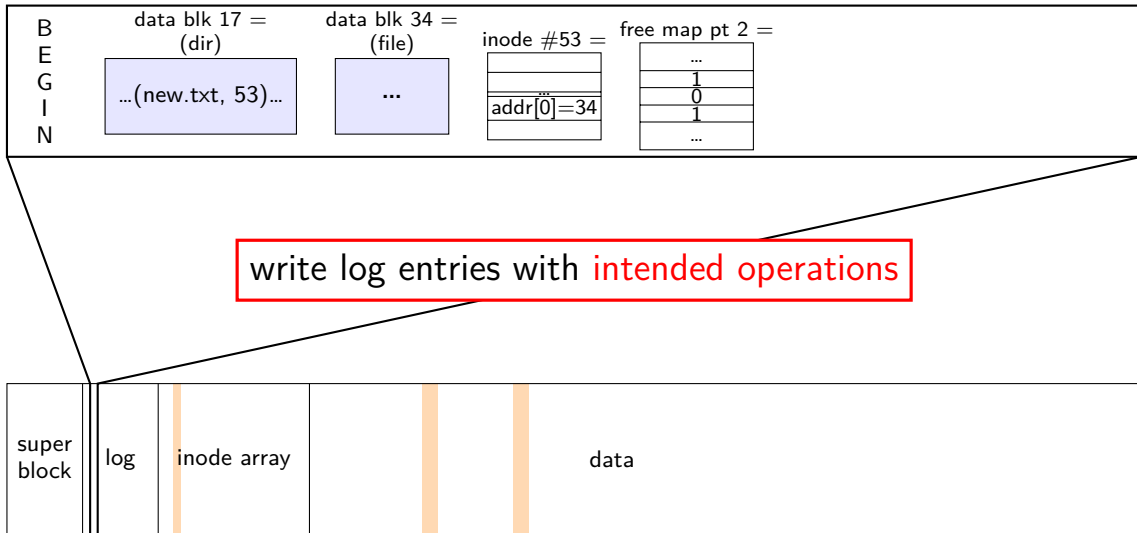
- update done — whole transaction happened

- update not done — whole transaction did not happen

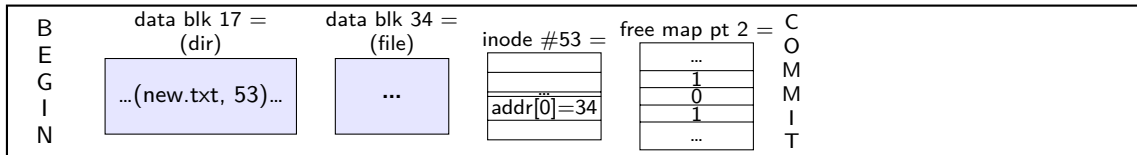
redo logging: file creation



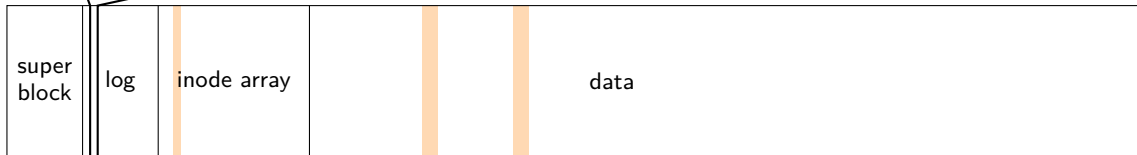
redo logging: file creation



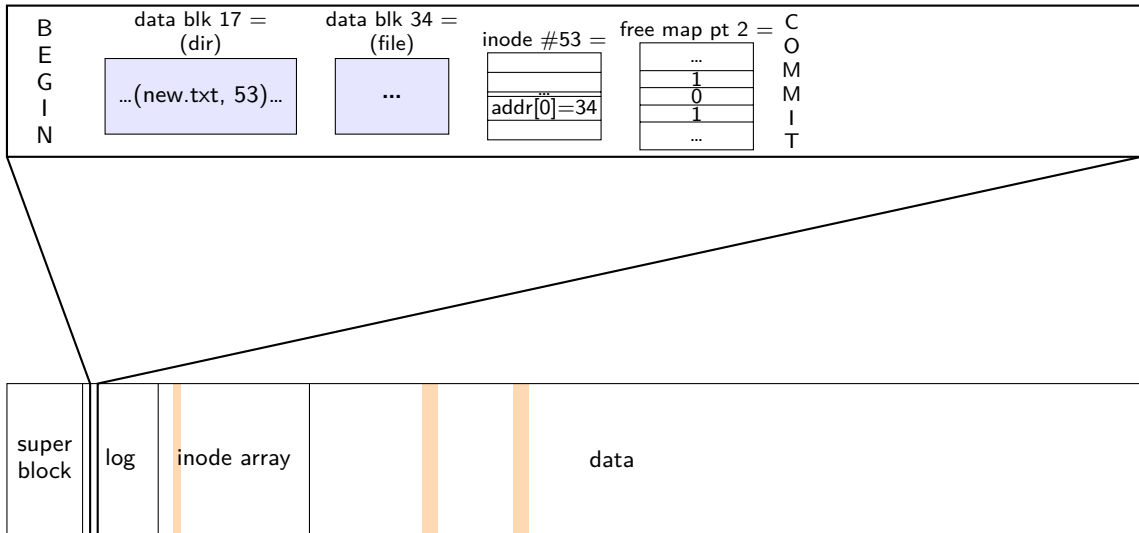
redo logging: file creation



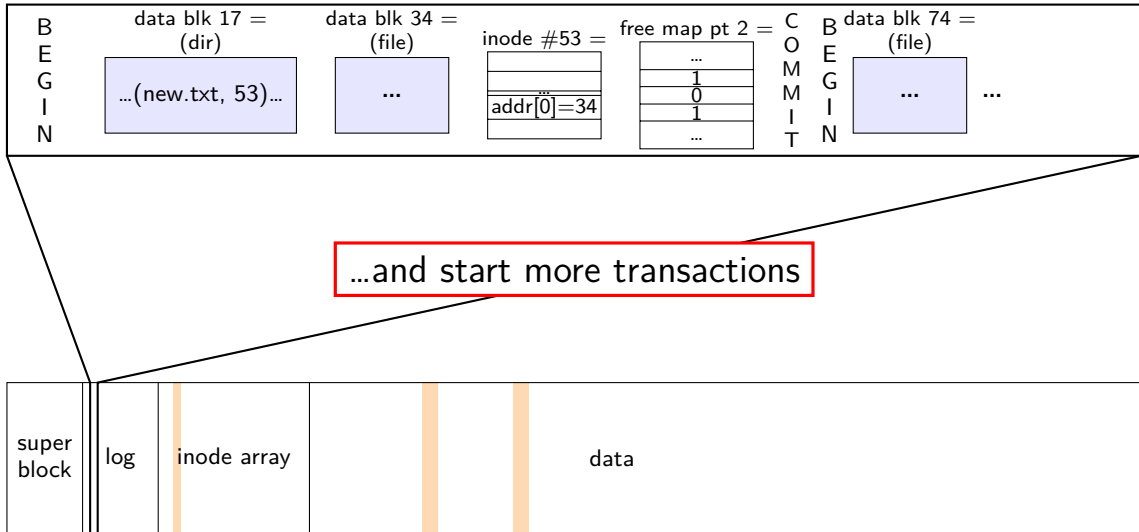
filesystem needs to ensure that committed updates **will definitely happen!**
mechanism: check this log for commit messages later, and **redo them** (just in case)



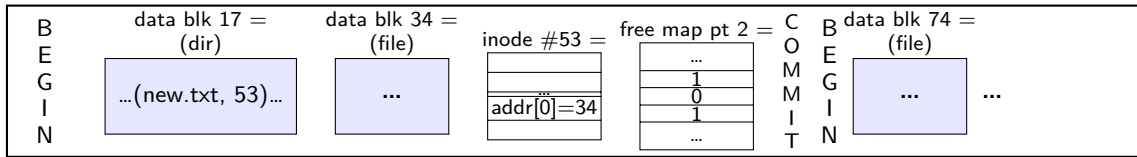
redo logging: file creation



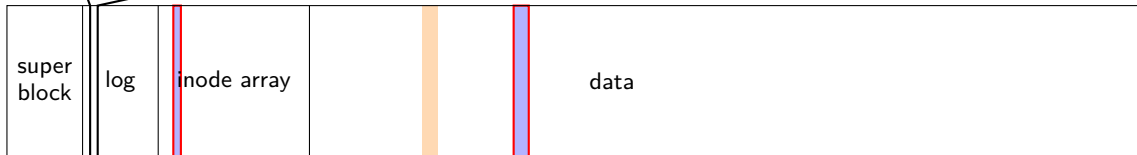
redo logging: file creation



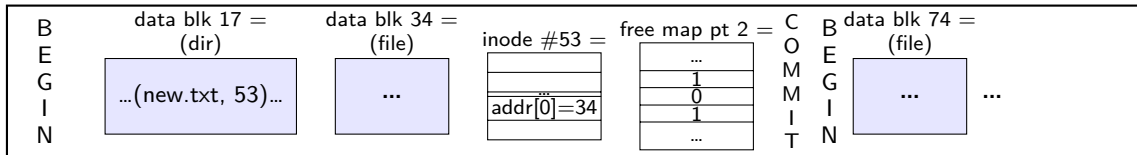
redo logging: file creation



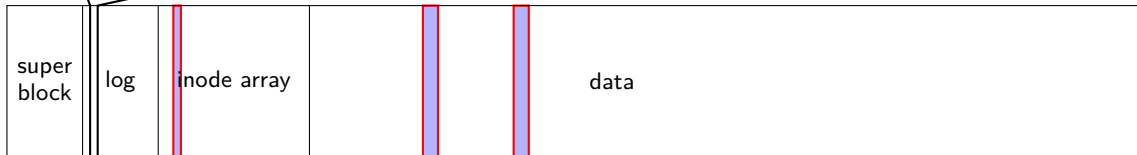
later, start applying results to actual disk



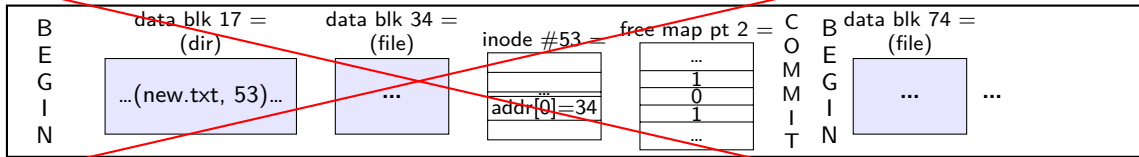
redo logging: file creation



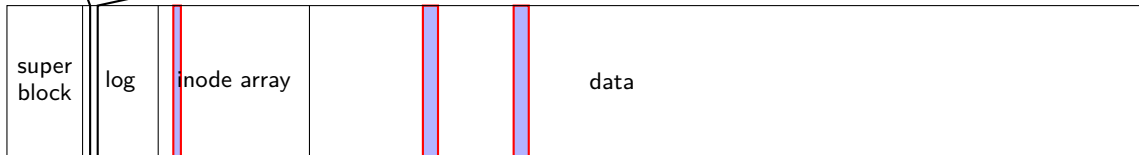
when everything is written, can overwrite log



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

“garbage collection”

redo logging: file creation

normal operation

write to log transaction steps:

- data blocks to create
- direcotry 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

“garbage collection”

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
- direcotry 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

“garbage collection”

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

“garbage collection”

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

“garbage collection”

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*

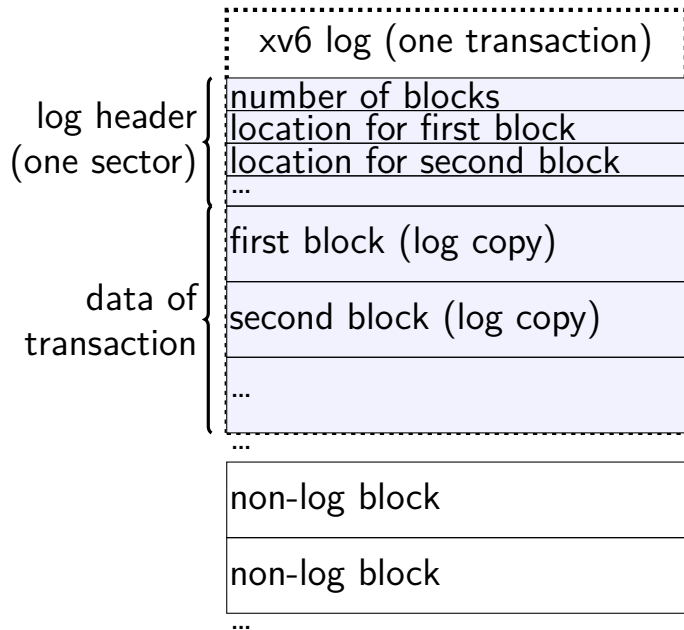
actually update the real data

redo logging and filesystems

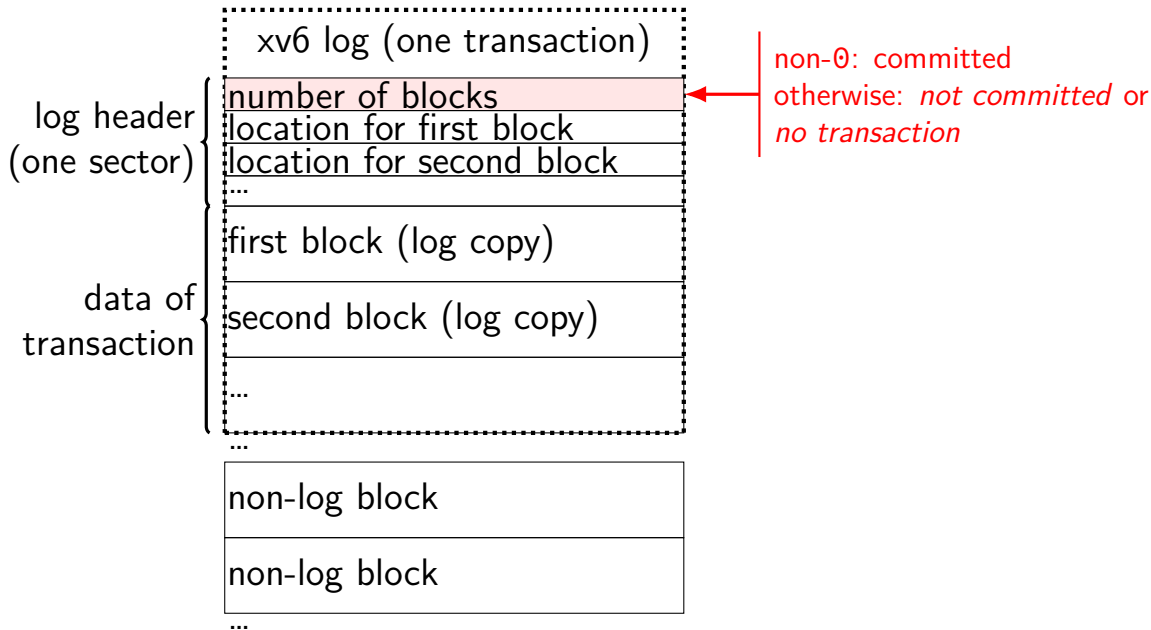
filesystems that do redo logging are called *journalling filesystems*

backup slides

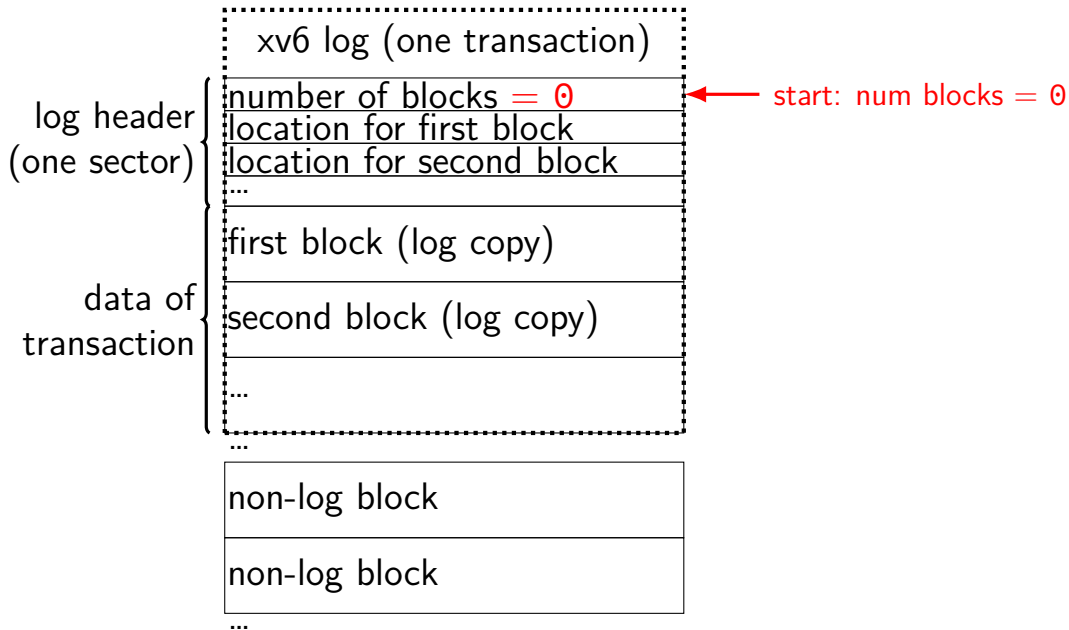
the xv6 journal



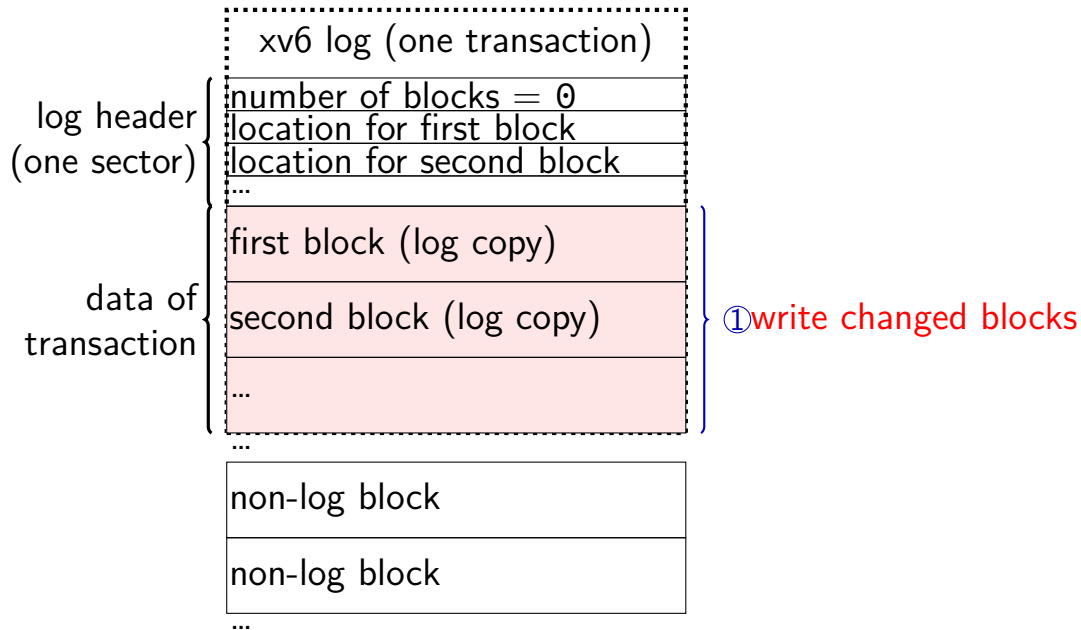
the xv6 journal



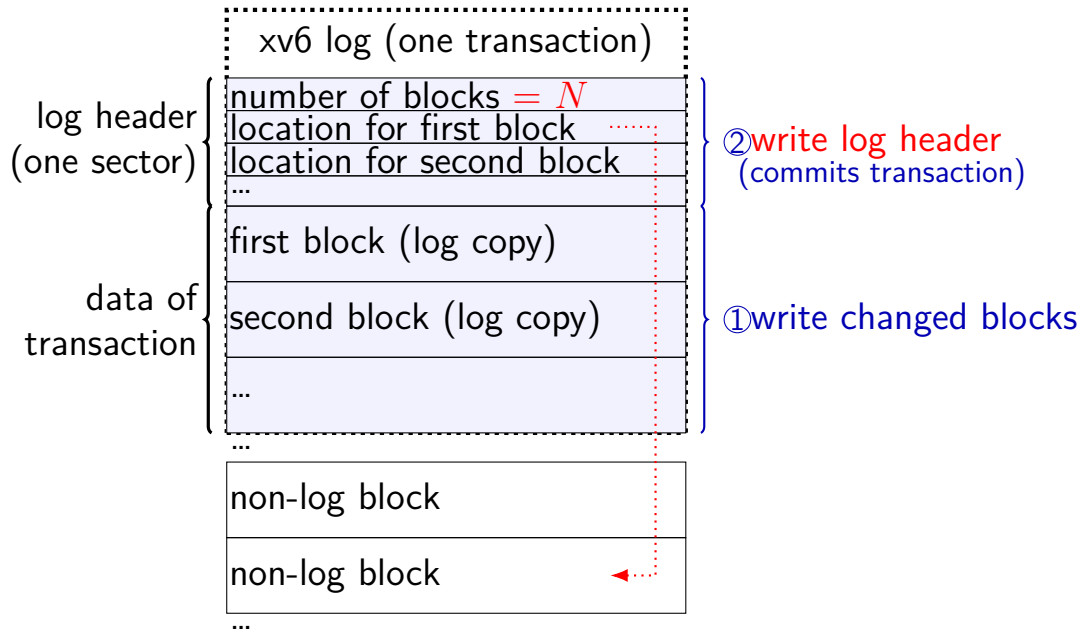
the xv6 journal



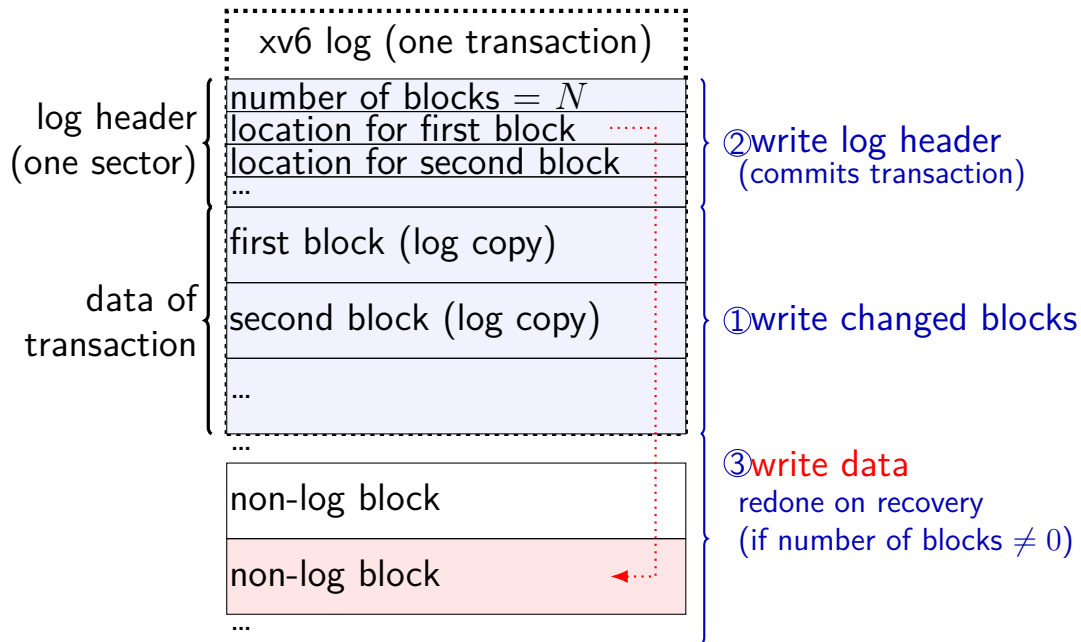
the xv6 journal



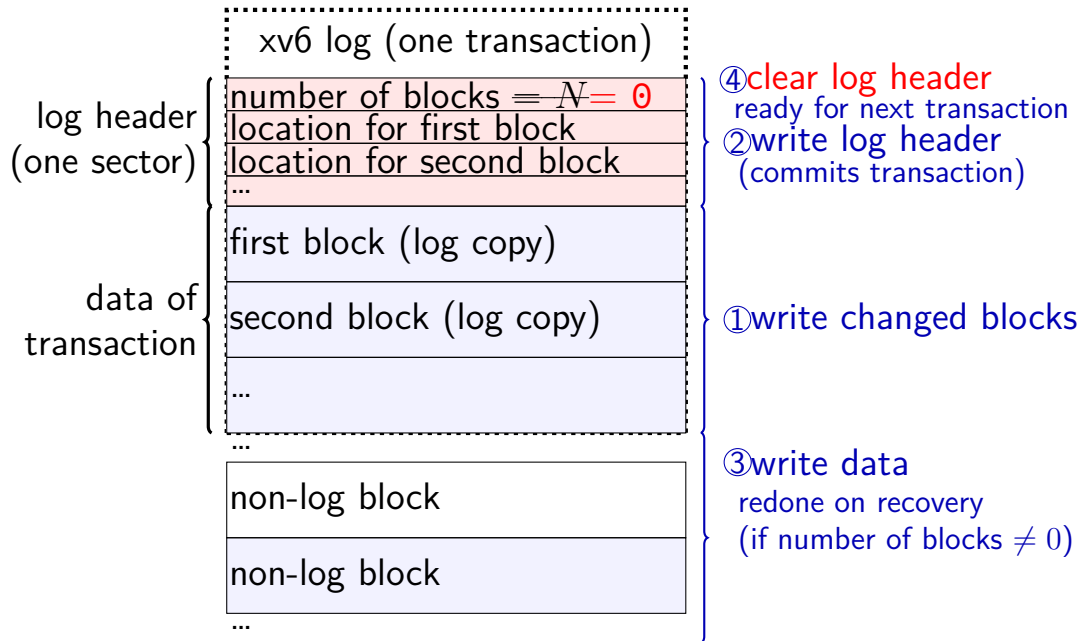
the xv6 journal



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

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.

exercise

say xv6 filesystem with:

- 64-byte inodes (12 direct + 1 indirect pointer)

- 16-byte directory entries

- 512 byte blocks

- 2-byte block pointers

how many blocks (not storing inodes) is used to store a directory of 200 30464B ($29 \cdot 1024 + 256$ byte) files?

remember: blocks could include blocks storing data or block pointers or directory enties

how many blocks is used to store a directory of 2000 3KB files?

fragments

Linux FS: a file's last block can be a *fragment* — only part of a block

each block split into approx. 4 fragments

each fragment has its own index

extra field in inode indicates that last block is fragment

allows one block to store data for several small files

beyond mirroring

mirroring seems to waste a lot of space

10 disks of data? mirroring \rightarrow 20 disks

10 disks of data? how good can we do with 15 disks?

best possible: lose 5 disks, still okay

can't do better or it wasn't really 10 disks of data

schemes that do this based on *erasure codes*

erasure code: encode data in way that handles parts missing (being erased)

erasure code example

store 2 disks of data on 3 disks

recompute original 2 disks of data from any 2 of the 3 disks

extra disk of data: some formula based on the original disks

common choice: bitwise XOR

common set of schemes like this: RAID

Redundant Array of Independent Disks

snapshots

filesystem snapshots

idea: filesystem keeps old versions of files around

- accidental deletion? old version stil 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

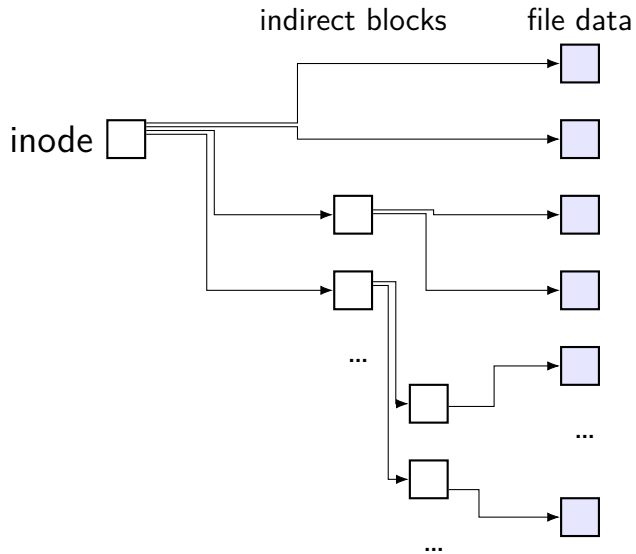
can access *snapshot* of files at prior time

mechanism: **copy-on-write**

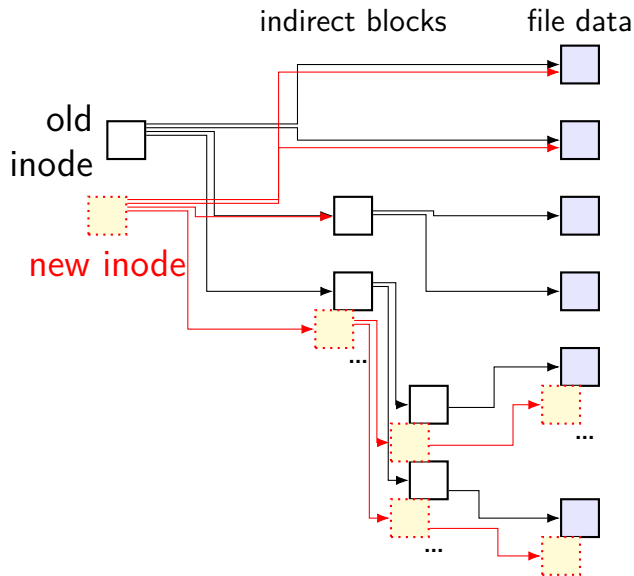
changing file makes **new copy** of filesystem

common parts shared between versions

inode and copy-on-write



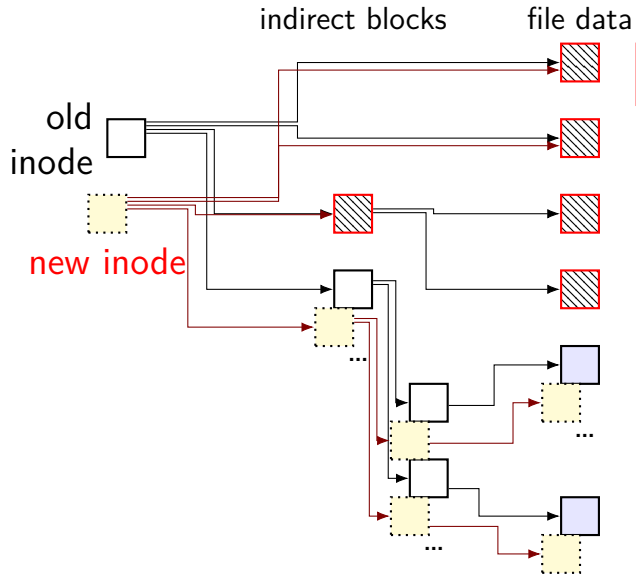
inode and copy-on-write



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

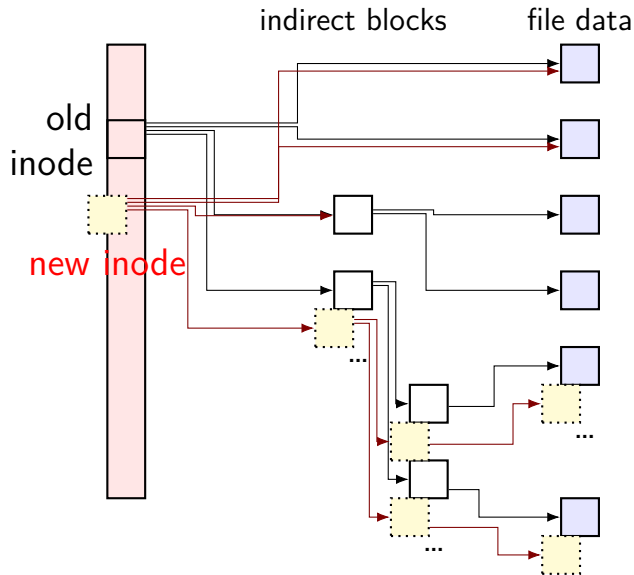
both old + new inode valid

inode and copy-on-write



unchanged parts of file shared

inode and copy-on-write

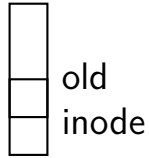


challenge: FFS/xv6/ext2 design
has big array of inodes

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

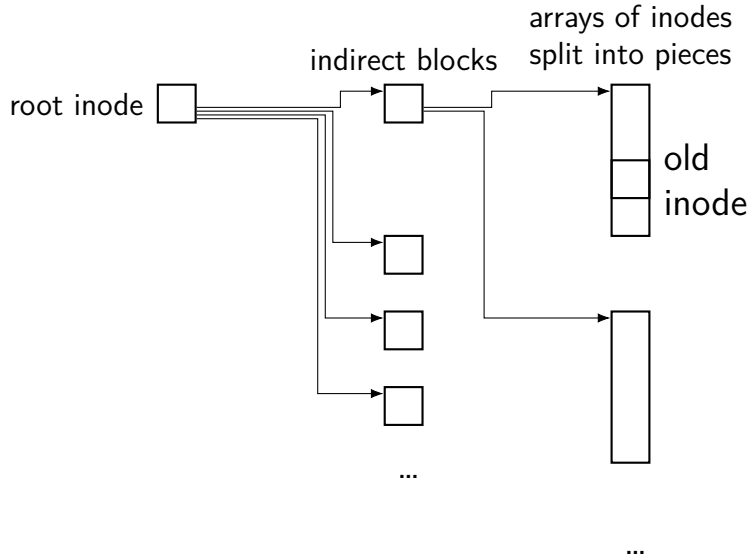
extra indirection for inode array

arrays of inodes
split into pieces

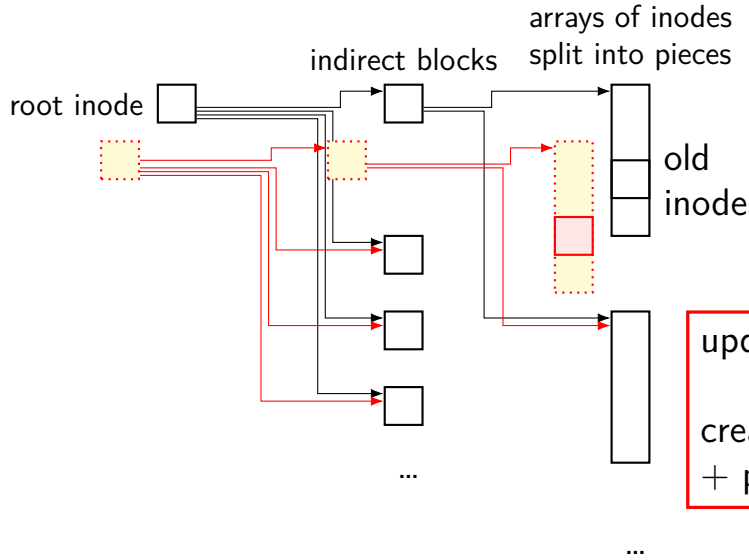


...

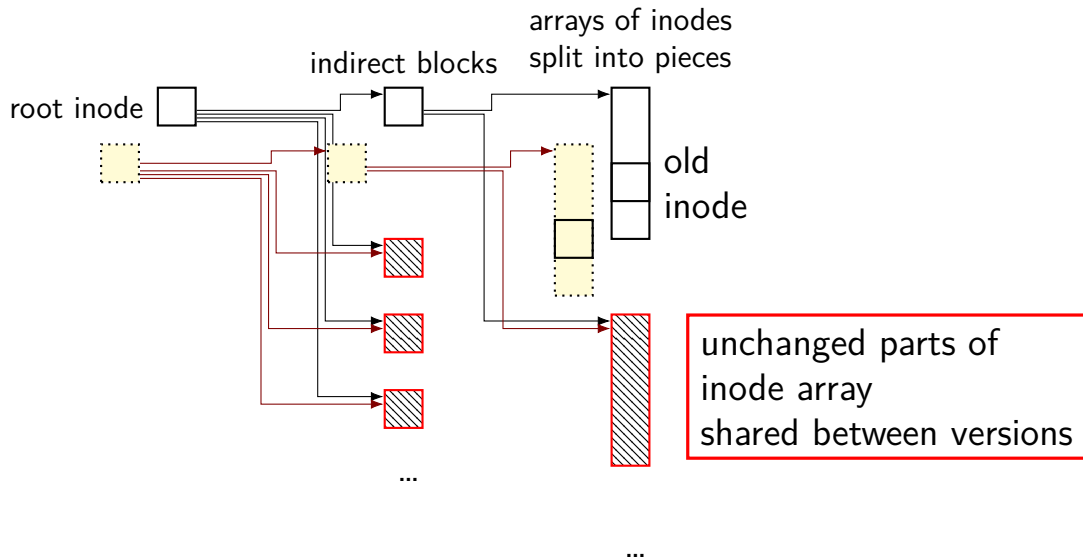
extra indirection for inode array



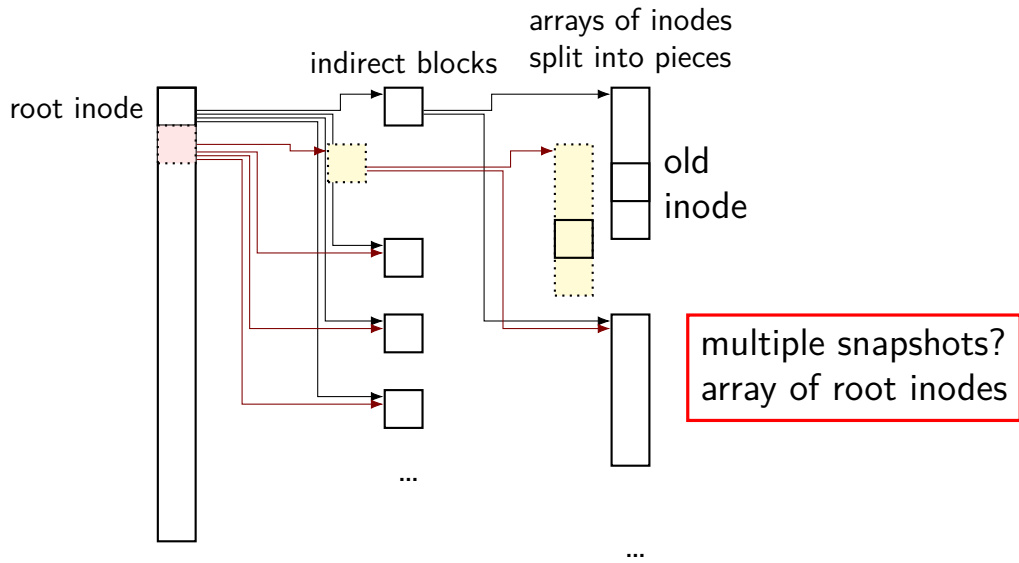
extra indirection for inode array



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 supports this (if turned on)

example: `.zfs/snapshots/11.11.18-06` pseudo-directory
contains contents of files at 11 November 2018 6AM

multiple copies

FAT: multiple copies of file allocation table and header

in inode-based filesystems: often multiple copies of superblocks

if part of disk's data is lost, have an extra copy

- always update both copies

- hope: disk failure to small group of sectors

hope: enough to recover most files on disk failure

- extra copy of metadata that is important for all files

- but won't recover specific files/directories whose data was lost

aside: FAT date encoding

seperate date and time fields (16 bits, little-endian integers)

bits 0-4: seconds (divided by 2), 5-10: minute, 11-15: hour

bits 0-4: day, 5-8: month, 9-15: year (minus 1980)

sometimes extra field for 100s(?) of a second

Fast File System

the Berkeley Fast File System (FFS) 'solved' some of these problems

McKusick et al, "A Fast File System for UNIX" <https://people.eecs.berkeley.edu/~brewer/cs262/FFS.pdf>
avoids long seek times, wasting space for tiny files

Linux's ext2 filesystem based on FFS

some other notable newer solutions (beyond what FFS/ext2 do)

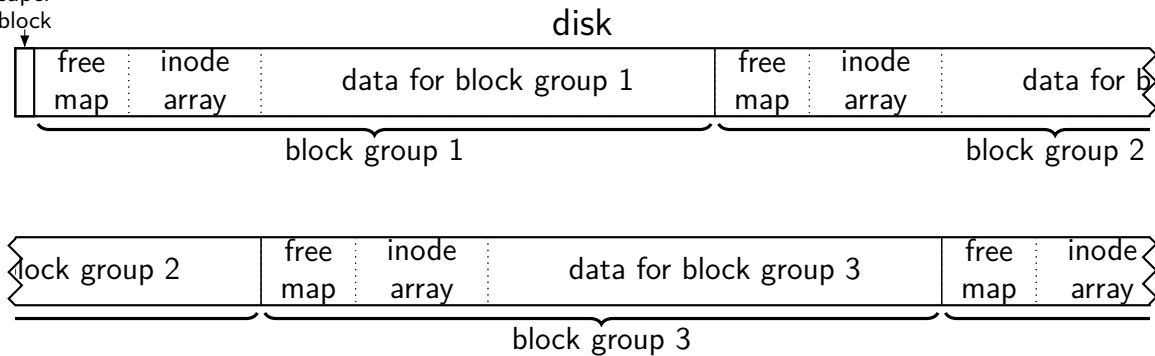
- better handling of very large files
- avoiding linear directory searches

block groups

(AKA cluster groups)

super

block



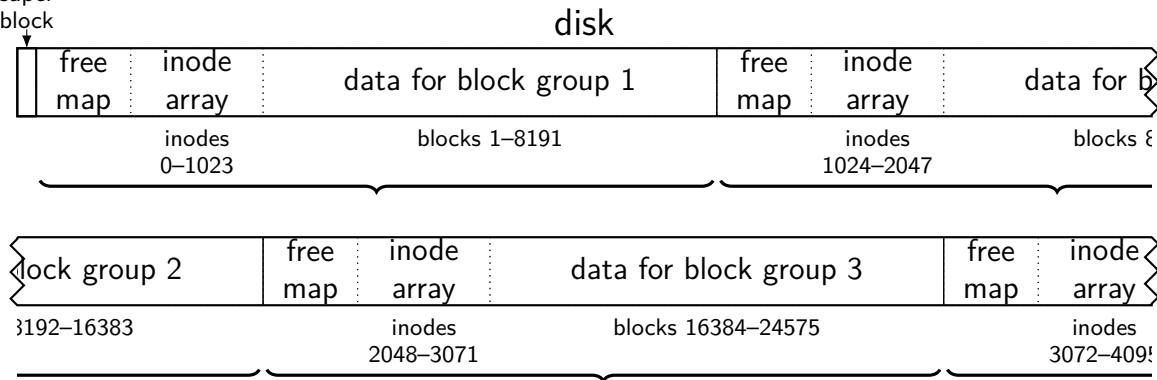
split disk into block groups
each block group like a mini-filesystem

block groups

(AKA cluster groups)

super

block



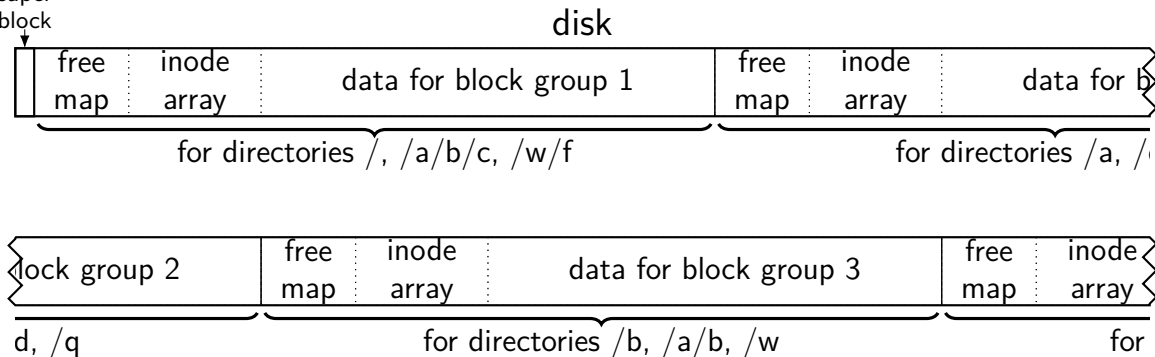
split block + inode numbers across the groups
inode in one block group can reference blocks in another
(but would rather not)

block groups

(AKA cluster groups)

super

block



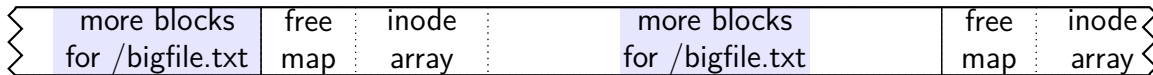
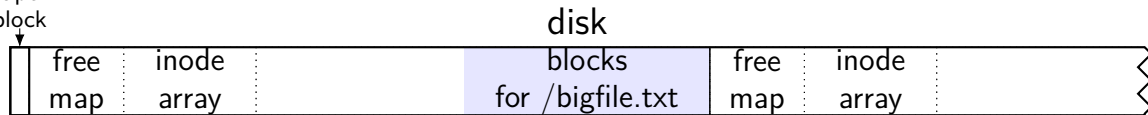
goal: *most data* for each directory within a block group
directory entries + inodes + file data close on disk
lower seek times!

block groups

(AKA cluster groups)

super

block

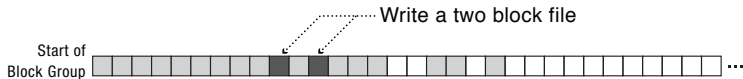


large files might need to be split across block groups

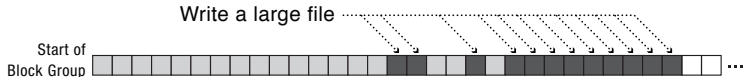
allocation within block groups



Expected typical arrangement.



Small files fill holes near start of block group.



Large files fill holes near start of block group and then write most data to sequential range blocks.

FFS block groups

making a subdirectory: new block group
for inode + data (entries) in different

writing a file: same block group as directory, first free block
intuition: non-small files get contiguous groups at end of block
FFS keeps disk deliberately underutilized (e.g. 10% free) to ensure this
can wait until dirty file data flushed from cache to allocate blocks
makes it easier to allocate contiguous ranges of blocks

several bad options (2)

suppose we're creating a new file

A: mark blocks as used in free block map

B: write inode for file

C: write directory entry for file

several bad options (2)

suppose we're creating a new file

A: mark blocks as used in free block map

B: write inode for file

C: write directory entry for file

if we do A before B+C and crash happens after A:

have blocks we can't use (not free), but which are unused

several bad options (2)

suppose we're creating a new file

A: mark blocks as used in free block map

B: write inode for file

C: write directory entry for file

if we do A before B+C and crash happens after A:

have blocks we can't use (not free), but which are unused

if we do B before A+C and crash happens after B:

have inode we can't use (not free), but which is not really used

several bad options (2)

suppose we're creating a new file

A: mark blocks as used in free block map

B: write inode for file

C: write directory entry for file

if we do A before B+C and crash happens after A:

have blocks we can't use (not free), but which are unused

if we do B before A+C and crash happens after B:

have inode we can't use (not free), but which is not really used

if we do C before A+B and crash happens after C:

have directory entry that points to junk — will behave weirdly

xv6 filesystem performance issues

inode, block map stored far away from file data

long seek times for reading files

unintelligent choice of file/directory data blocks

xv6 finds *first free block/inode*

result: files/directory entries scattered about

blocks are pretty small — needs lots of space for metadata

could change size? but waste space for small files

large files have giant lists of blocks

linear searches of directory entries to resolve paths

xv6 filesystem performance issues

inode, block map stored far away from file data

long seek times for reading files

unintelligent choice of file/directory data blocks

xv6 finds *first free block/inode*

result: files/directory entries scattered about

blocks are pretty small — needs lots of space for metadata

could change size? but waste space for small files

large files have giant lists of blocks

linear searches of directory entries to resolve paths

xv6 filesystem performance issues

inode, block map stored far away from file data

long seek times for reading files

unintelligent choice of file/directory data blocks

xv6 finds *first free block/inode*

result: files/directory entries scattered about

blocks are pretty small — needs lots of space for metadata

could change size? but waste space for small files

large files have giant lists of blocks

linear searches of directory entries to resolve paths

xv6 filesystem performance issues

inode, block map stored far away from file data

long seek times for reading files

unintelligent choice of file/directory data blocks

xv6 finds *first free block/inode*

result: files/directory entries scattered about

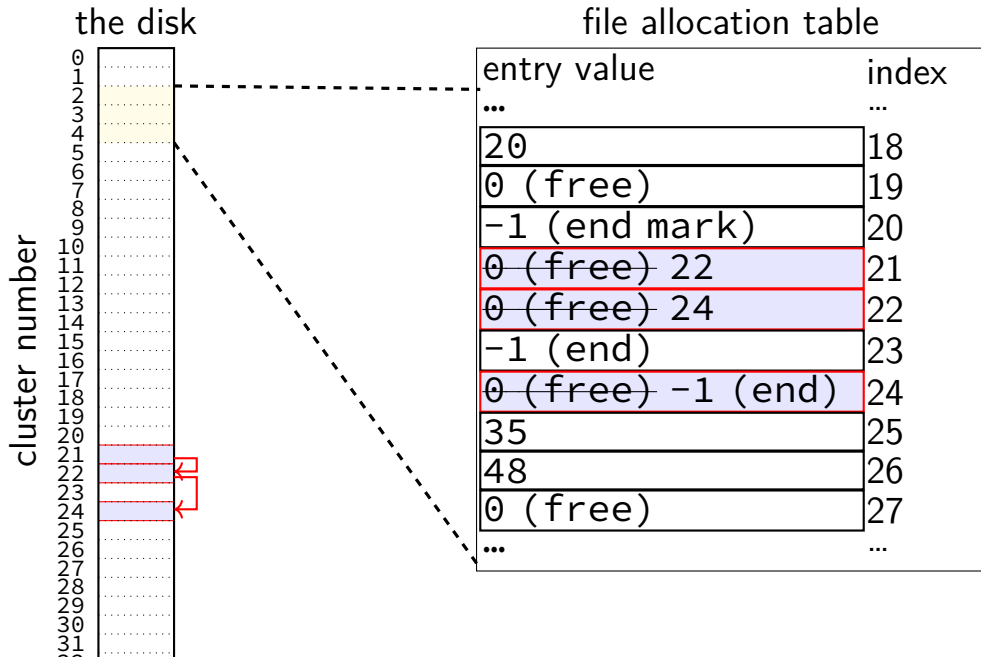
blocks are pretty small — needs lots of space for metadata

could change size? but waste space for small files

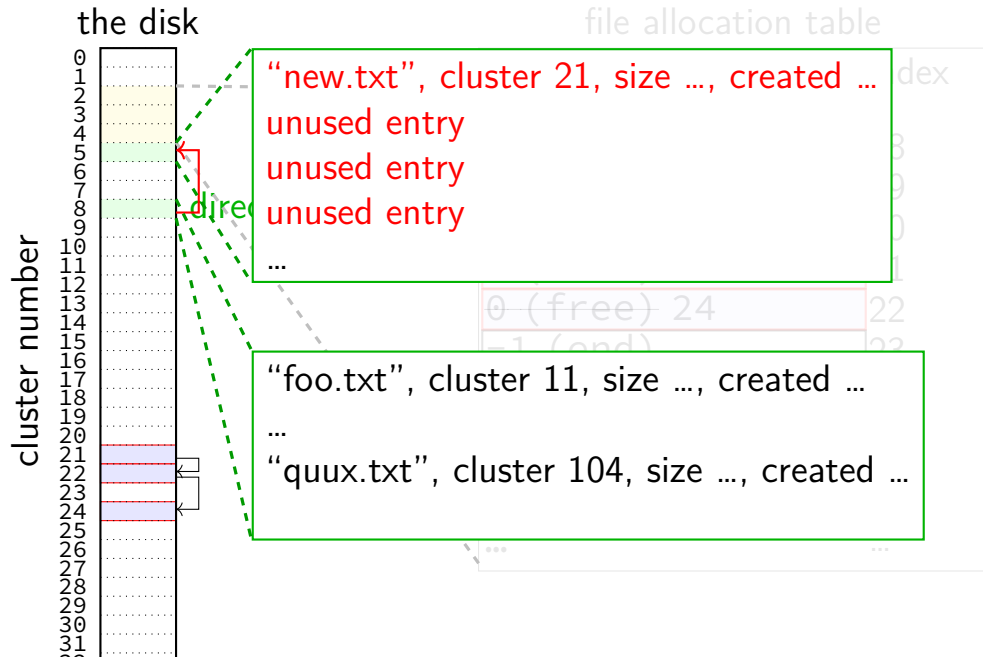
large files have giant lists of blocks

linear searches of directory entries to resolve paths

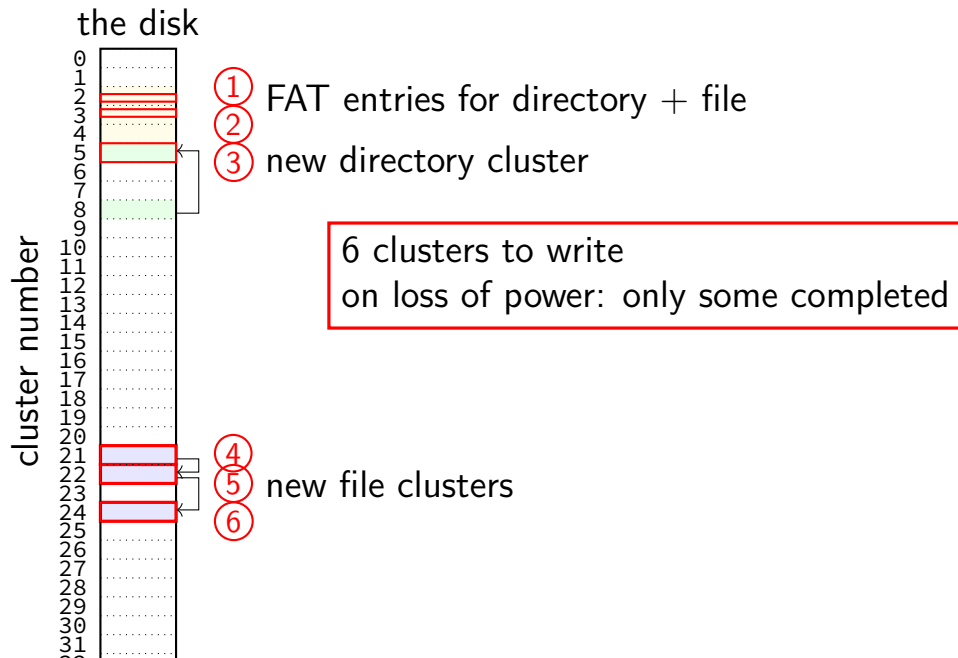
recall: FAT: file creation (1)



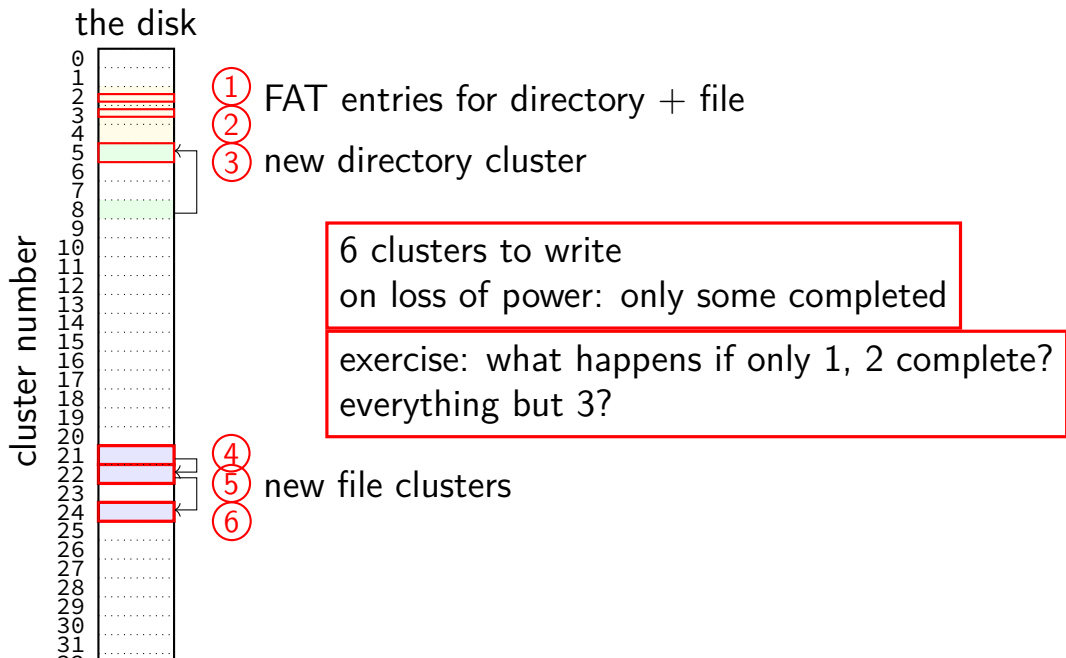
recall: FAT: file creation (2)



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

recovery with careful ordering

avoiding data loss → can 'fix' inconsistencies

programs like `fsck` (filesystem check), `chkdsk` (check disk)
run manually or periodically or after abnormal shutdown

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 direcotry 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 direcotry 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 a 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

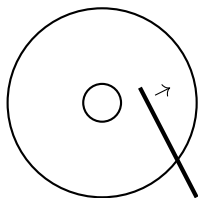
ordering and disk performance

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



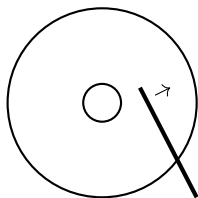
ordering and disk performance

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:

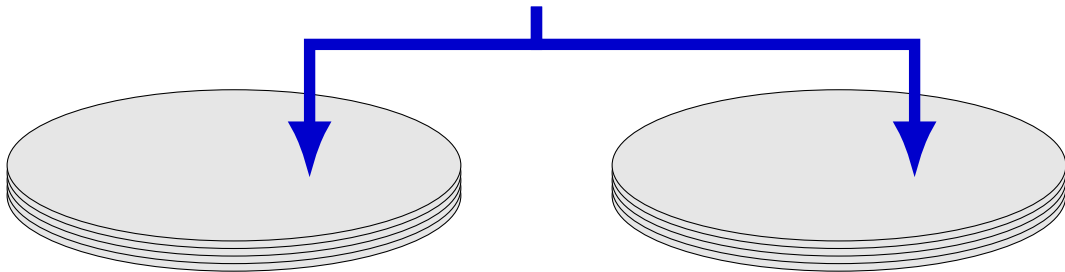
free block map for file (start), then file blocks (middle), then...

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

mirroring whole disks

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

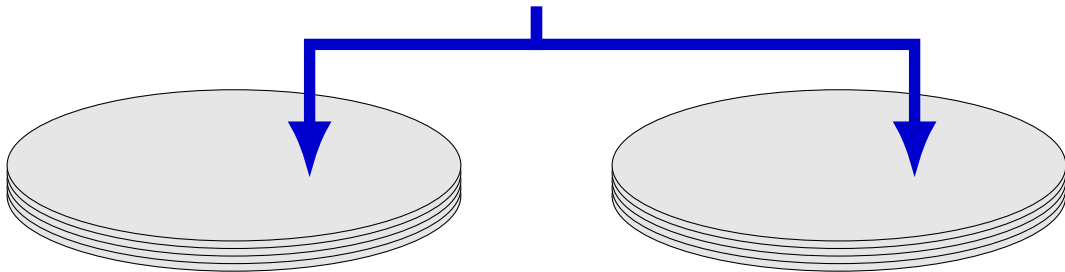
always write to both



mirroring whole disks

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

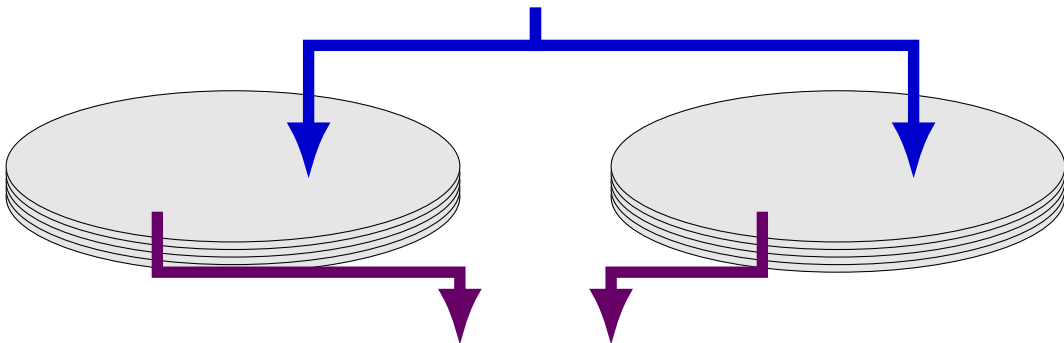
always write to both



mirroring whole disks

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

always write to both



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

beyond mirroring

mirroring seems to waste a lot of space

10 disks of data? mirroring \rightarrow 20 disks

10 disks of data? how good can we do with 15 disks?

best possible: lose 5 disks, still okay

can't do better or it wasn't really 10 disks of data

schemes that do this based on *erasure codes*

erasure code: encode data in way that handles parts missing (being erased)

erasure code example

store 2 disks of data on 3 disks

recompute original 2 disks of data from any 2 of the 3 disks

extra disk of data: some formula based on the original disks

common choice: bitwise XOR

common set of schemes like this: RAID

Redundant Array of Independent Disks

exercice

filesystem has:

- root directory with 2 subdirectories

- each subdirectory contains 3 512B files, 2 4MB files

- (1MB = 1024KB; 1KB = 1024B)

- 32B directory entries

- 4B block pointers

- 4KB blocks

- inode: 12 direct pointers, 1 indirect pointer, 1 double-indirect, 1 triple-indirect

(a) how many inodes used?

(b) how many blocks (outside of inodes) with 1KB fragments?
[minimum w/partial blocks]

(c) how many blocks (outside of inodes) with block pointers
replaced by 8B extents (no fragments)? [compute minimum]

inodes used

per each of 2 subdirectories: 5 files + 1 inode for subdirectory = 6

plus 1 for root directory itself

$$= 12 + 1 = 13$$

blocks with fragments

each of 6 512B files uses a single 1KB fragment
wastes 512Bs of it

each of 2 subdirectory needs $32B \cdot 5 \ll 1KB$ (1 fragment)
(5 directory entries; probably also additional entries for ..)

root directory needs $32B \cdot 2 \ll 1KB$ (1 fragment)

9 1KB fragments \rightarrow minimum 3 (4KB) blocks

each of 4 4MB file uses 1024 data blocks

1 indirect block for blocks 13-(1024+13) [last 12 pointers unused]

= 4096 blocks (4MB files data) + 4 (4MB file indirects) + 3 (for fragments)

= 4103 blocks

blocks with extents

each of 6 512B files uses a single 4KB block
extent specifying block

each of 2 subdirectory needs $32B \cdot 5 \ll 4KB$ (1 block)

root directory needs $32B \cdot 2 \ll 4KB$ (1 block)

each of 2 4MB file uses 2048 data blocks

no indirect blocks assuming 2048 data blocks are contiguous (one extent in inode)

$= 4096 \text{ blocks (4MB files data)} + 6 \text{ (small files)} + 3 \text{ (directory entries)} = 4105 \text{ blocks}$

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?

entire log can be **written sequentially**

- ideal for hard disk performance

- also pretty good for SSDs

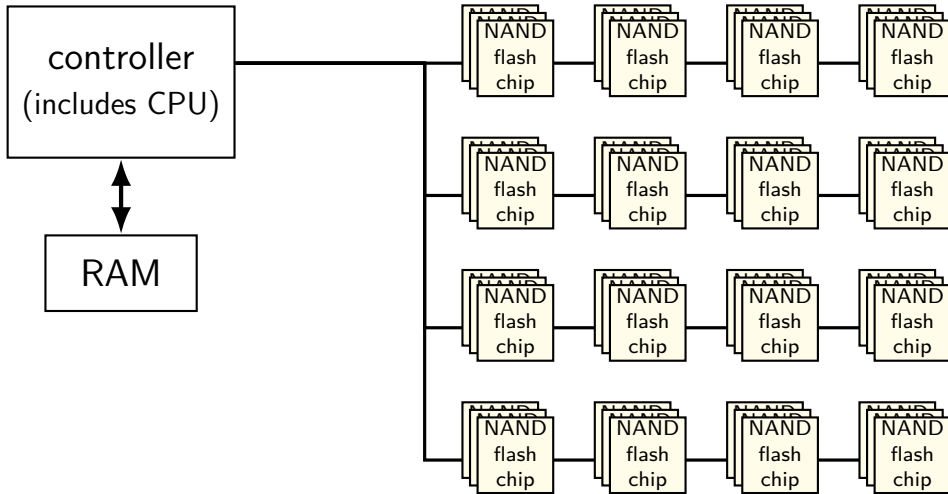
no waiting for 'real' updates

- application can proceed while updates are happening

- files will be updated even if system crashes

often better for performance!

solid state disk architecture



flash

- no moving parts

 - no seek time, rotational latency

- can read in sector-like sizes (“pages”) (e.g. 4KB or 16KB)

- write once between erasures

- erasure only in large *erasure blocks* (often 256KB to megabytes!)

- can only rewrite blocks order tens of thousands of times

 - after that, flash starts failing

SSDs: flash as disk

SSDs: implement hard disk interface for NAND flash

- read/**write** sectors at a time

- sectors much smaller than erasure blocks

- sectors sometimes smaller than flash 'pages'

- read/write with use sector numbers, not addresses

- queue of read/writes

need to hide **erasure blocks**

- trick: block remapping — move where sectors are in flash

need to hide limit on number of erases

- trick: wear leveling — spread writes out

block remapping

Flash
Translation
Layer
remapping table

logical	physical
0	93
1	260
...	...
31	74
32	75
...	...

OS sector numbers

flash locations

block remapping

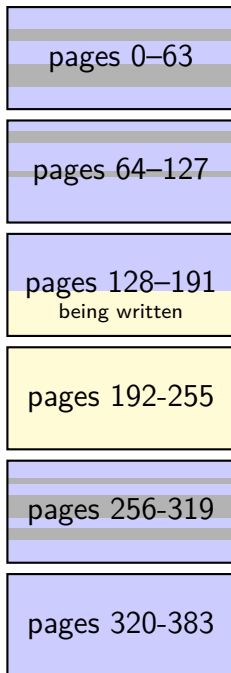
Flash
Translation
Layer
remapping table

logical	physical
0	93
1	260
...	...
31	74
32	75
...	...

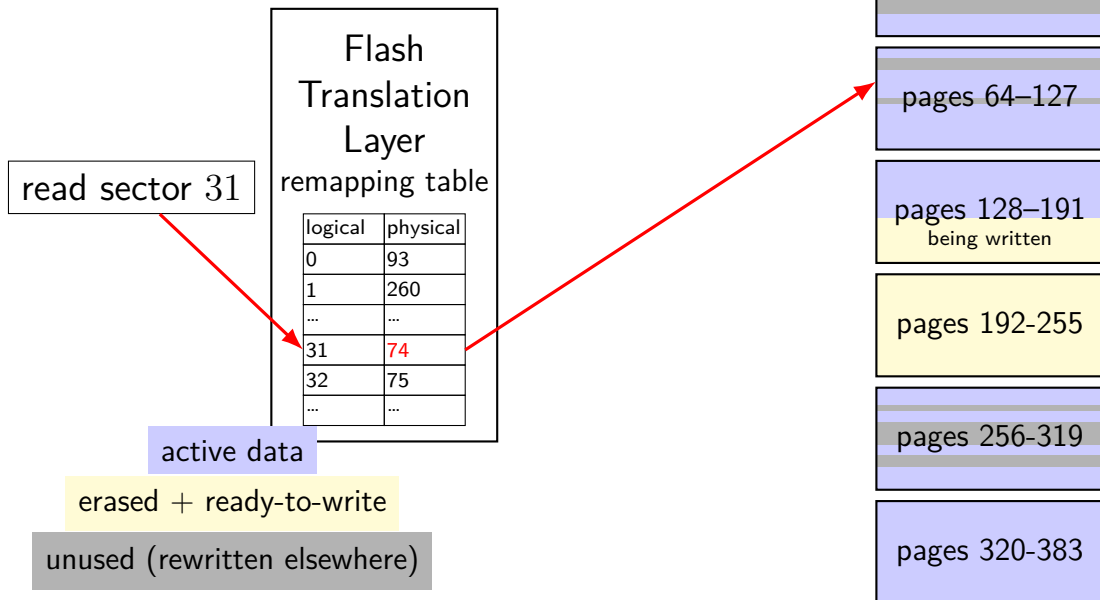
active data

erased + ready-to-write

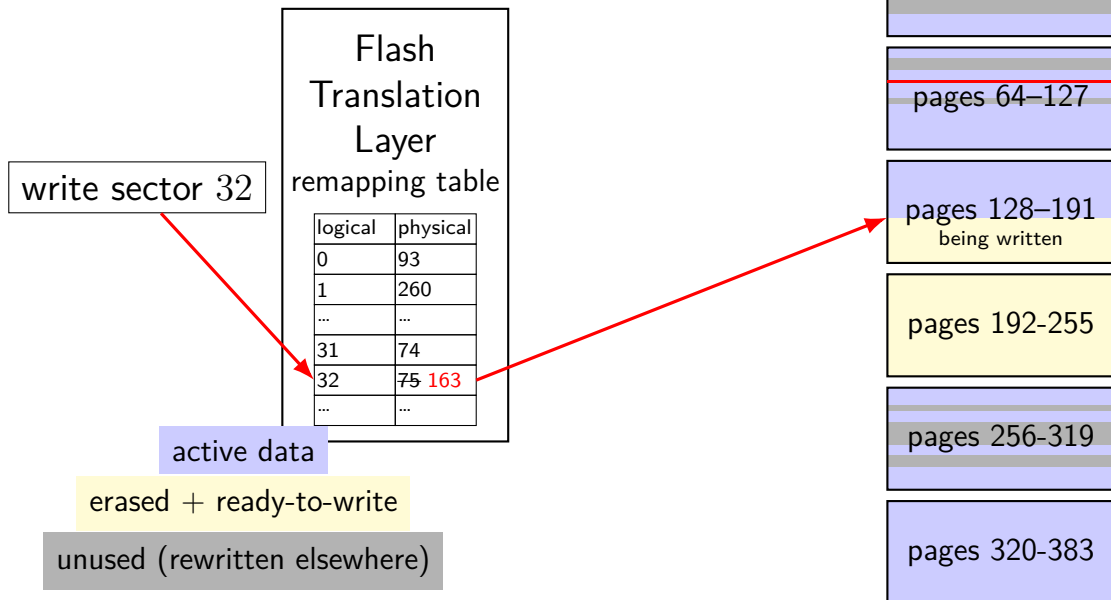
unused (rewritten elsewhere)



block remapping



block remapping



block remapping

Flash
Translation
Layer
remapping table

logical	physical
0	93
1	260 187
...	...
31	74
32	75 163
...	...

active data

erased + ready-to-write

unused (rewritten elsewhere)

“garbage collection”
(free up new space)

pages 128–191

copied from erased

pages 192–255

pages 256–319
erased block

can only erase
whole “erasure block”

pages 0–63

pages 64–127

pages 128–191
being written

pages 192–255

pages 256–319

pages 320–383

block remapping

controller contains mapping: sector \rightarrow location in flash

on write: write sector to *new location*

eventually do *garbage collection* of sectors

- if erasure block contains some replaced sectors and some current sectors...
copy current blocks to new location to reclaim space from replaced sectors

doing this efficiently is very complicated

SSDs sometimes have a 'real' processor for this purpose

exercise

Assuming a FAT-like filesystem on an SSD, which of the following are likely to be stored in the same (or very small number of) erasure block?

- [a] the clusters of a set of log file all in one directory written continuously over months by a server and assigned a contiguous range of cluster numbers
- [b] the data clusters of a set of images, copied all at once from a camera and assigned a variety of cluster numbers
- [c] all the entries of the FAT (assume the OS only rewrites a sector of the FAT if it is changed)

SSD performance

reads/writes: sub-millisecond

contiguous blocks don't really matter

can depend a lot on the controller

- faster/slower ways to handle block remapping

writing can be slower, especially when almost full

- controller may need to move data around to free up erasure blocks

- erasing an erasure block is pretty slow (milliseconds?)

extra SSD operations

SSDs sometimes implement non-HDD operations

on operation: TRIM

way for OS to mark sectors as unused/erase them

SSD can remove sectors from block map

- more efficient than zeroing blocks

- freed up more space for writing new blocks

aside: future storage

emerging non-volatile memories...

slower than DRAM (“normal memory”)

faster than SSDs

read/write interface like DRAM but persistent

capacities similar to/larger than DRAM

backup slides