

- ☑ 1. Last time
 - ☑ 2. Finish directories
 - ☑ 3. Performance
 - ☐ 4. Crash recovery
 - intro
 - ad hoc
 - copy on write
 - journaling
-

2. Finish directories

Example: /a/foo.c
/b/c/essay.txt inode #7 (by assumption)

We can give a file another name; for example
/a/hello.txt

This is because the file itself does not know its own name. Looking only at the file, the "name" is just a i-node number.

Human-usable names, such as /a/foo.c
for directory entries, for example:

Come from direct...

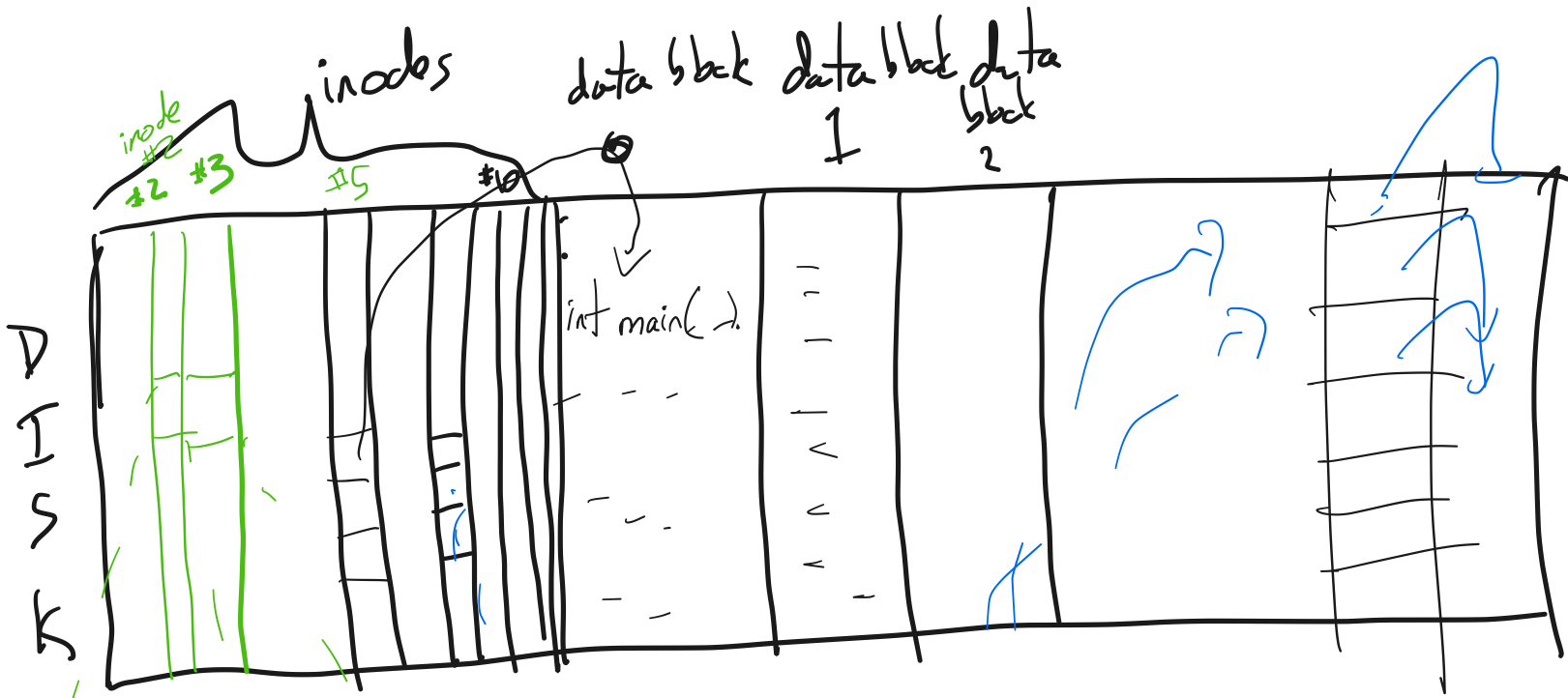
dirent

← readdir()

lab 2 d-name
d-no

< a, 3 >
< b, 4 >

which means "let us all call inode 3 by the segment "/a", and similarly for inode 4 and "/b".



< a, 3 >

< b, 4 >

< foo.c, 5 >

< mylink, 10 >

< hello.txt, 7 >

\$ ln /b/c/essay.txt /a/hello.txt

"link"

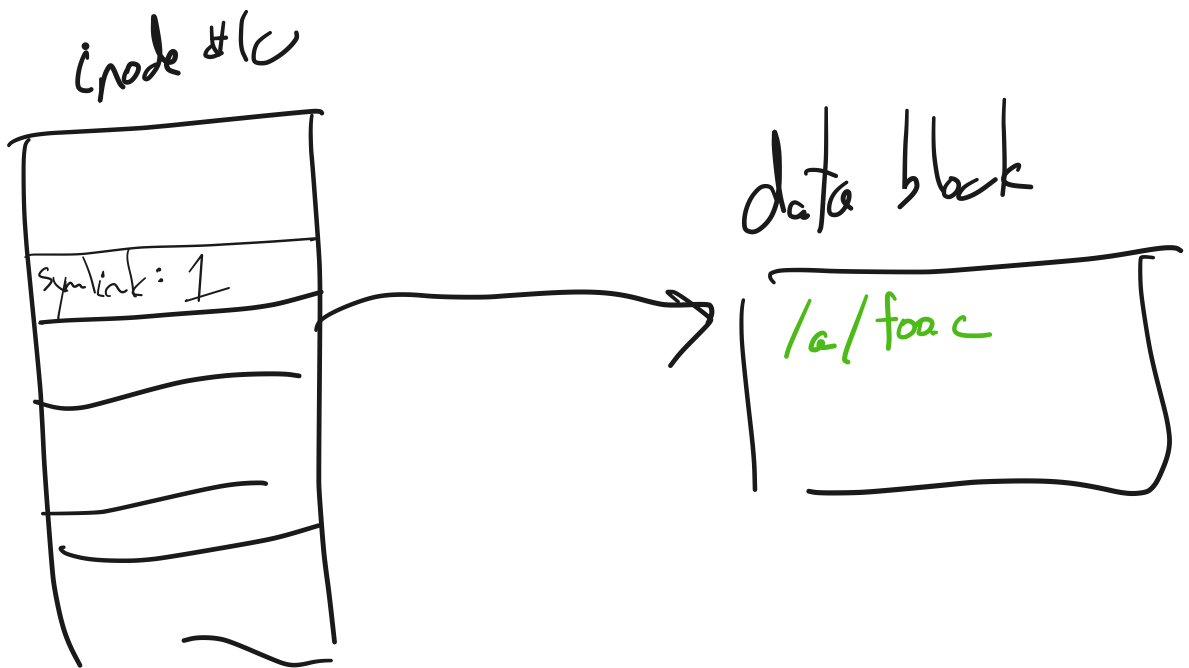
syscall: "unlink();"

soft links vs. hard links

10

\$ ln -s /a/foo.c /mylink

creates a new inode ($\neq 10$)



symlink bit set 1.



/a/b/c'

\$ ln ^{dir name}

mount points

\$ cat /a/foo.c



3. Performance

case study: FFS (1984)

problems w/ the original:

- blocks too small (512)
- inode array was at beginning of disk
- free blocks were stored in a linked list on the disk

- poor clustering of related objects

consecutive file blocks

inodes relative to disk blocks

directories

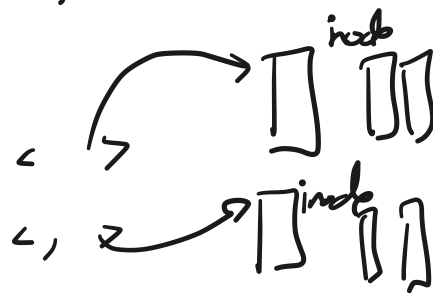
inodes for a given directory

result:

ls
grep <path> *.c } ⇒ slow

Improvements?

- Make data blocks / inodes close to each other



- Cluster files in the same directory
- Make data blocks bigger (4KB, 8KB, 16KB)
- Free blocks: store separately (bitmap)

110
1

$$\frac{1 \text{ block}}{1 \text{ bit}} \times \frac{500 \text{ GB disk}}{4 \text{ KB blocks}} = 12$$

$$2^{39} / 2^{12} = 2^{27} \approx \boxed{\text{disk: } 128,000,000 \text{ bits}}$$

12 2 } 120, 1

≈ 16 MB ← store in memory

$$2^{27} \text{ bits} \times \frac{1 \text{ byte}}{2^3 \text{ bits}} = 2^{24} \text{ bytes} = 16 \text{ MB}$$

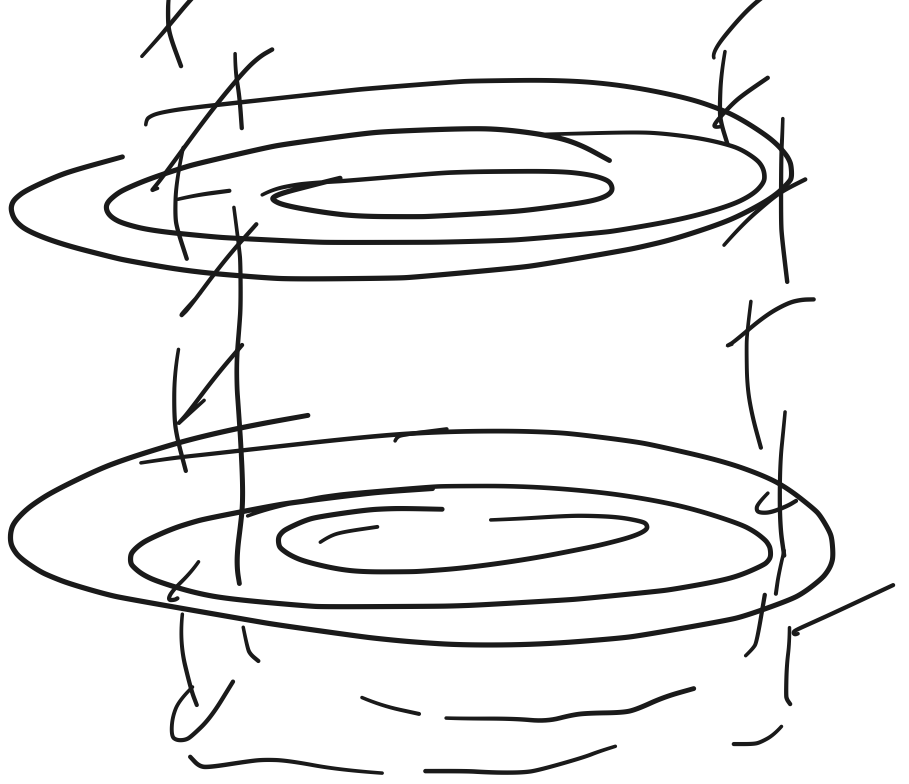
⇒ reserve space (lie to the user abt free space)

- symbolic links

- atomic "rename" = (\$ mv abc.txt def.txt)

- cylinder groups (for clustering)





[superblock | bookkeeping | inodes | bitmap | data blocks
 (512 bytes each)]

attempt: put inodes + their data blocks in the same cyl. grp.

attempt: put inodes of ~~the~~ files in the same dir. in the same cyl. group.

new directory: choose a cyl. group w/ higher than avg. # of free inodes.

as a file grows, after it crosses 48KB, spill to next cyl. group, and do likewise for every 1MB thereafter.



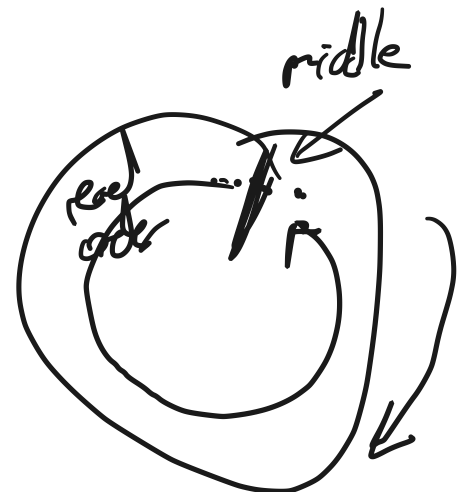
total perf:

20-40% of disk B/w for large files

10-20x improvement on the predecessor.

Other things they did:

- buffer cache
- read entire track



4. Crash recovery

- intro
- ad-hoc
- copy-on-write
- journaling

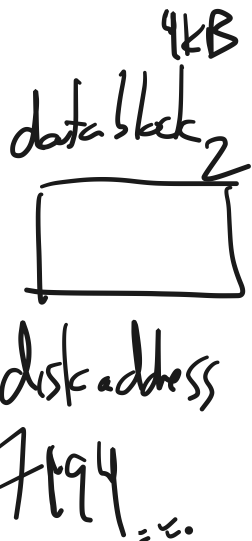
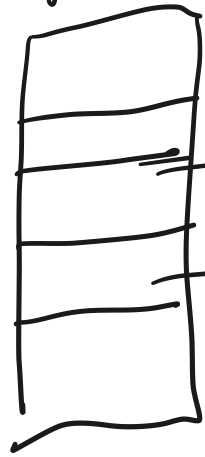


Ex.

binary or
C prog.

```
fd = open();  
seek(fd);  
write(fd, buf, 256);
```

inode



- add to inode
- write data to the block

• update the bitmap

What happens if there is a crash anywhere in this list of operations?

CS202 Handout 13

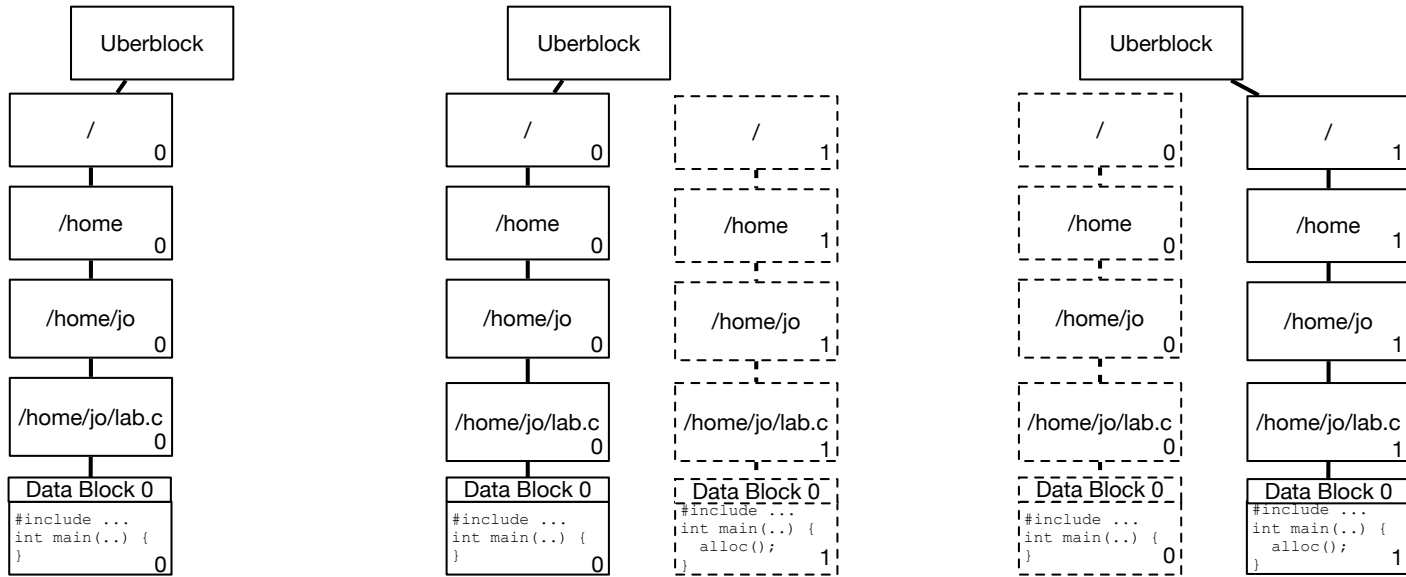
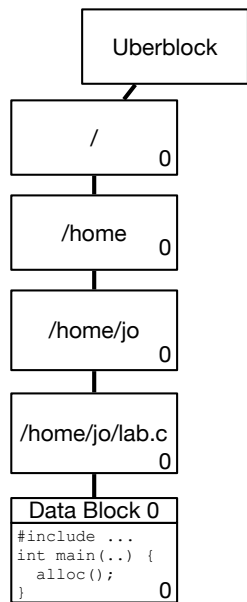
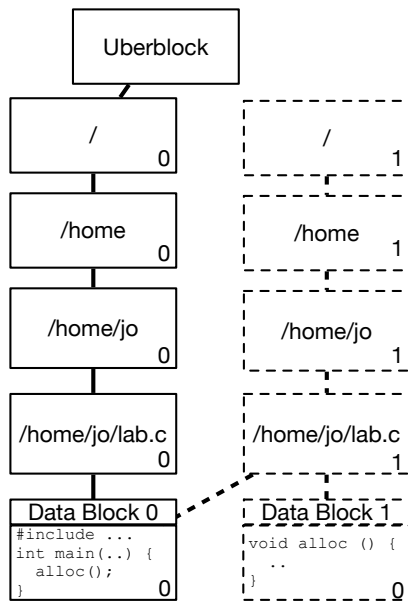


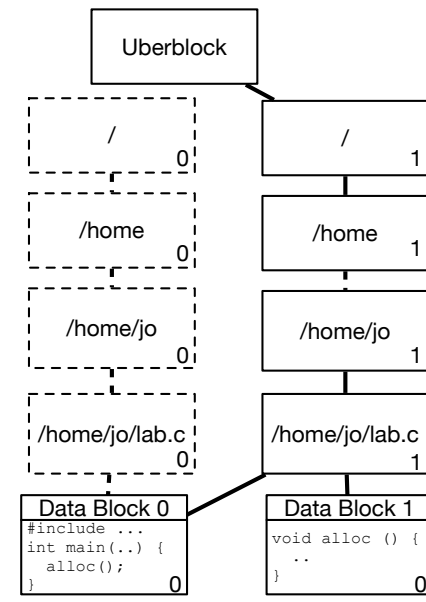
Figure 1: Copy-on-write filesystem: modifying a data block



(a) Initial State

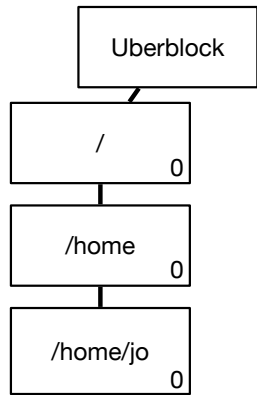


(b) System allocates and creates new versions of all modified blocks.

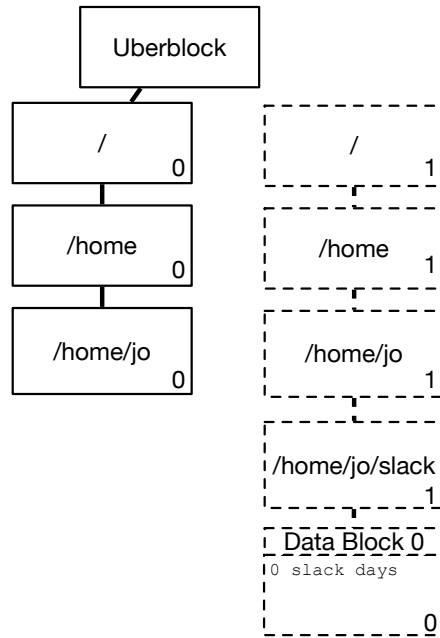


(c) System updates Uberblock to point to new version of blocks.

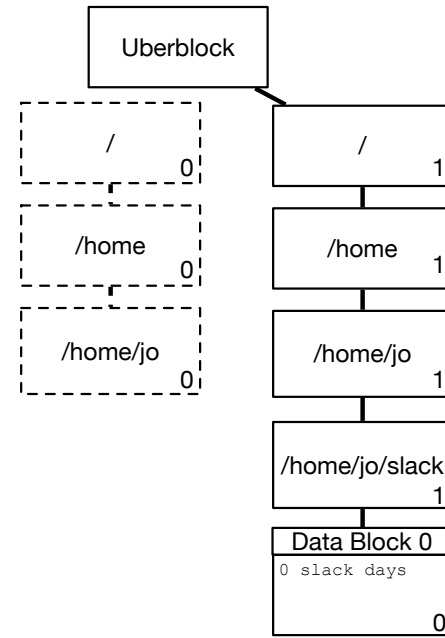
Figure 2: Copy-on-write filesystem: adding a data block



(a) Initial State



(b) System allocates and creates new versions of all modified blocks.

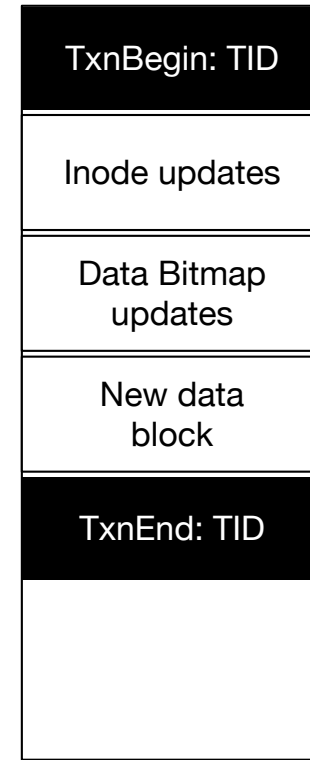


(c) System updates Uberblock to point to new version of blocks.

Figure 3: Copy-on-write filesystem: creating a file



ext3 disk layout



ext3 journal layout

Figure 4: Redo logging in a filesystem