


PERSISTENCE: FAST FILE SYSTEM

Shivaram Venkataraman

CS 537, Spring 2023

ADMINISTRIVIA

Midterm grades  Solutions with explanation
Next 1 or 2 days

P6 progress?

↳ Test cases? → Update by end of day

→ Two equal keys: does order matter? Not?

AGENDA / LEARNING OUTCOMES

How does file system represent files, directories?

What steps must reads/writes take?

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
- hierarchical
- traverse once
- offsets precisely defined

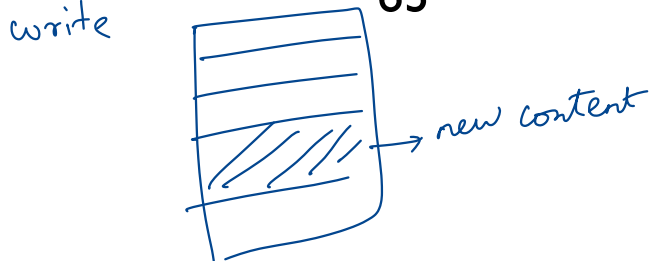
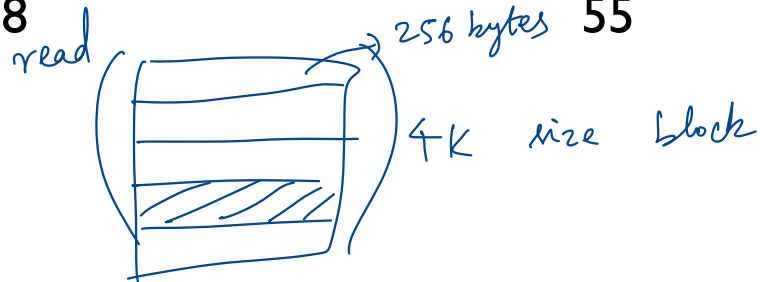
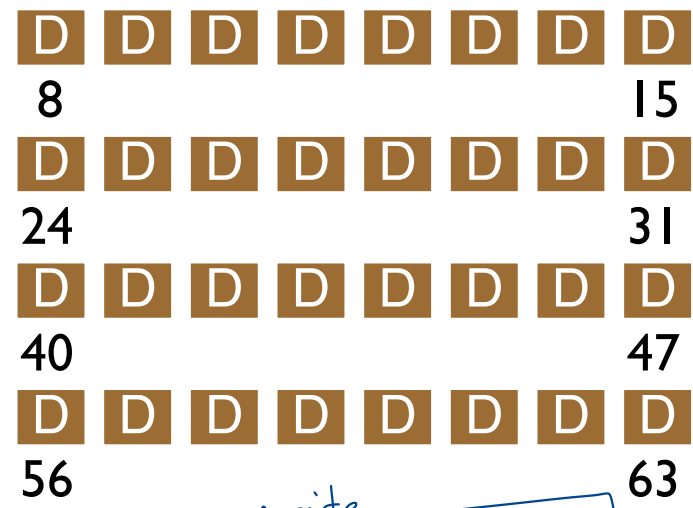
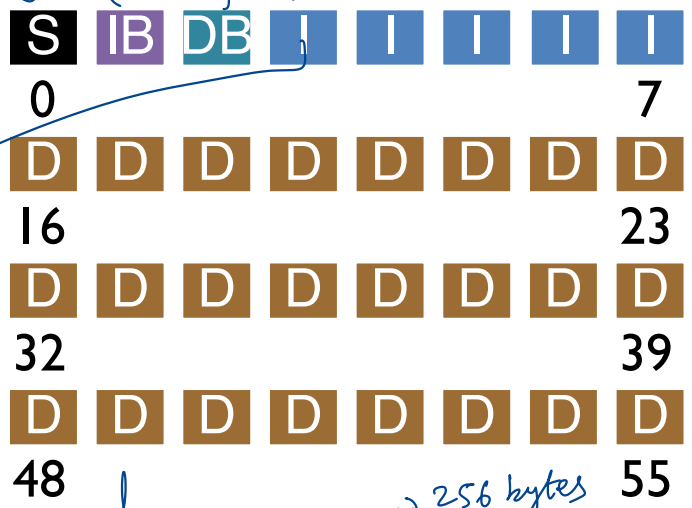
FILE SYSTEM LAYOUT

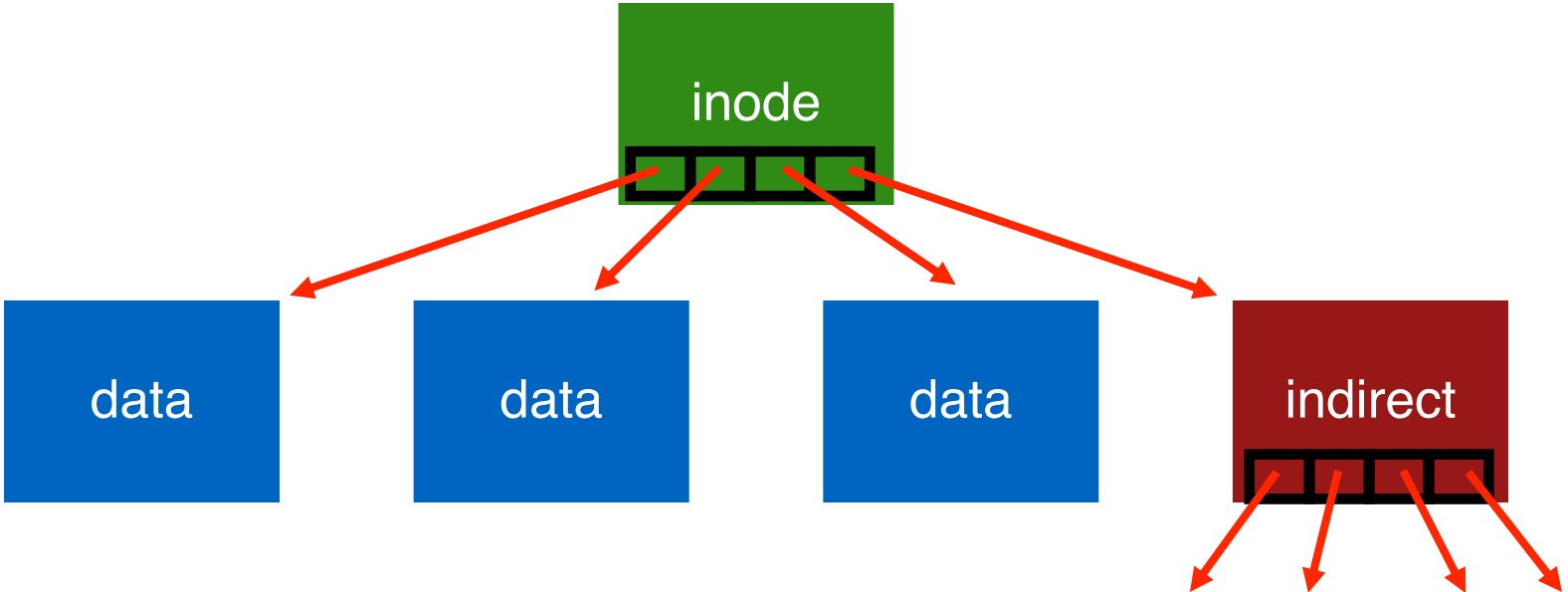
FS level metadata
block size
number of inodes

Bitmaps used for allocation

inodes → metadata, access time, permissions
ptrs to data blocks

store file contents
store directories



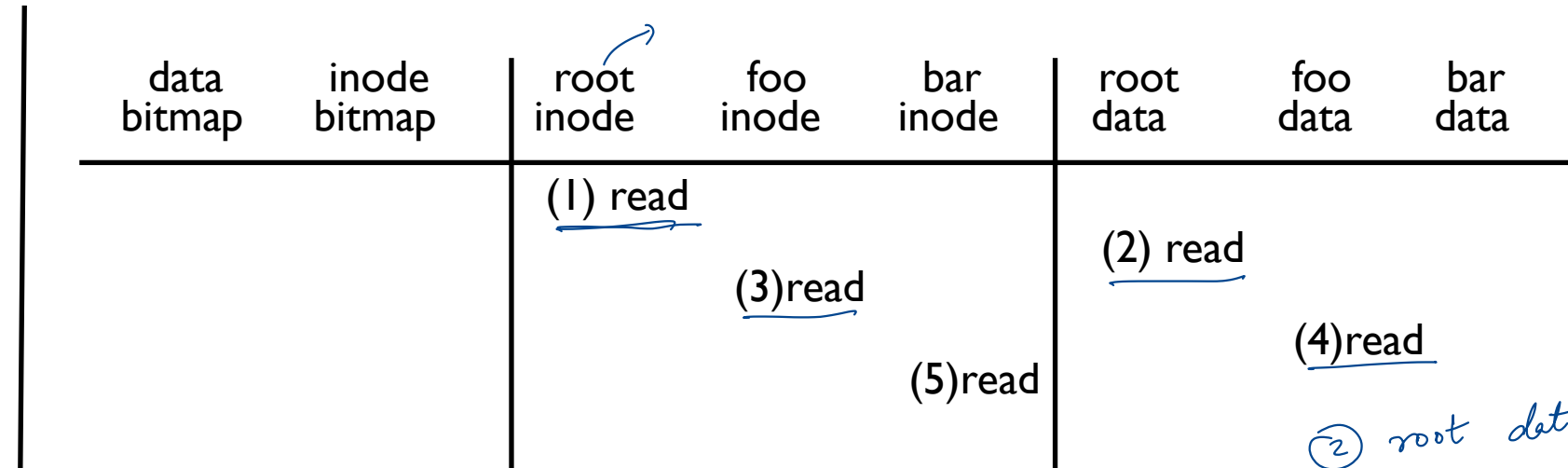


FS OPERATIONS

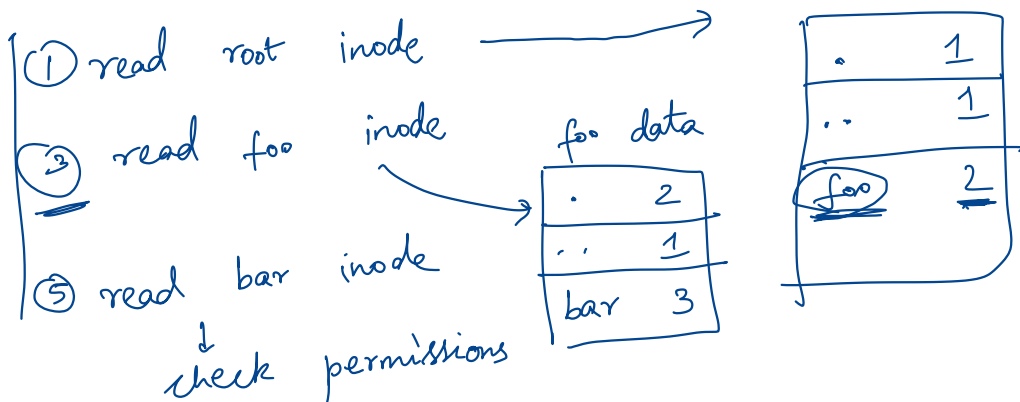
- open
- read
- close
- create file
- write

What does the FS need to do on disk? open /foo/bar

TIME

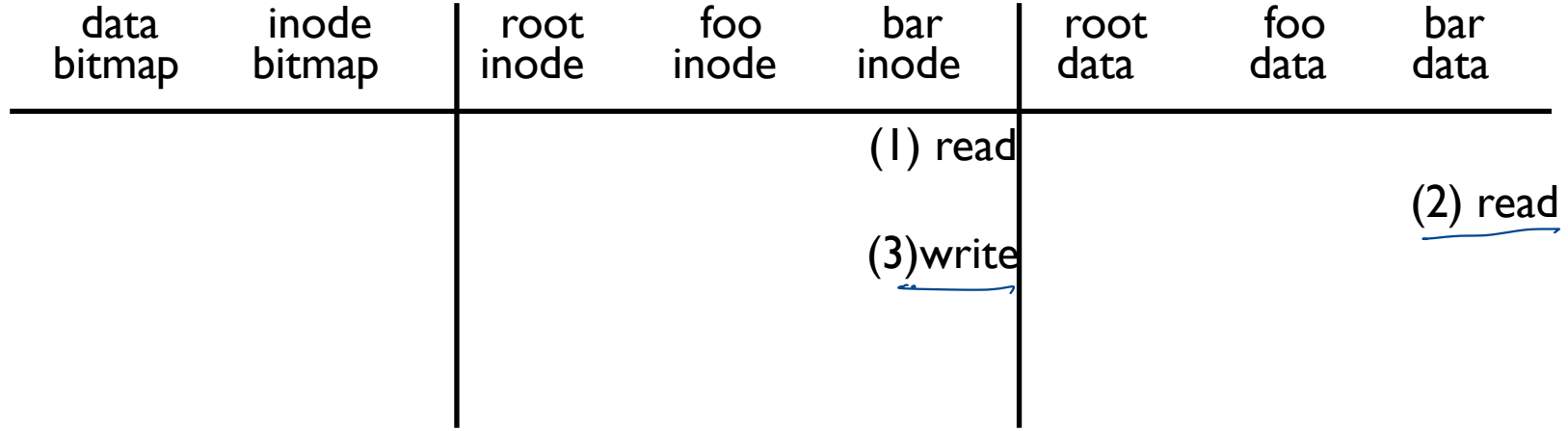


- ① Check if file exists
- ② Check if user has permissions



read /foo/bar – assume opened → OS has a FD (offset) for this file

TIME



read : start offset , read N bytes

- ① read inode : figure out which data blocks to read
- ② read data blocks calculated in prev step
- ③ Update the last accessed time for this file

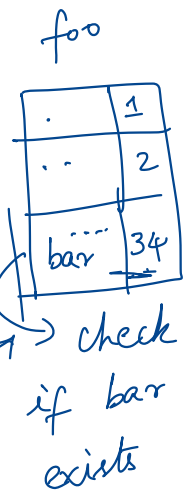
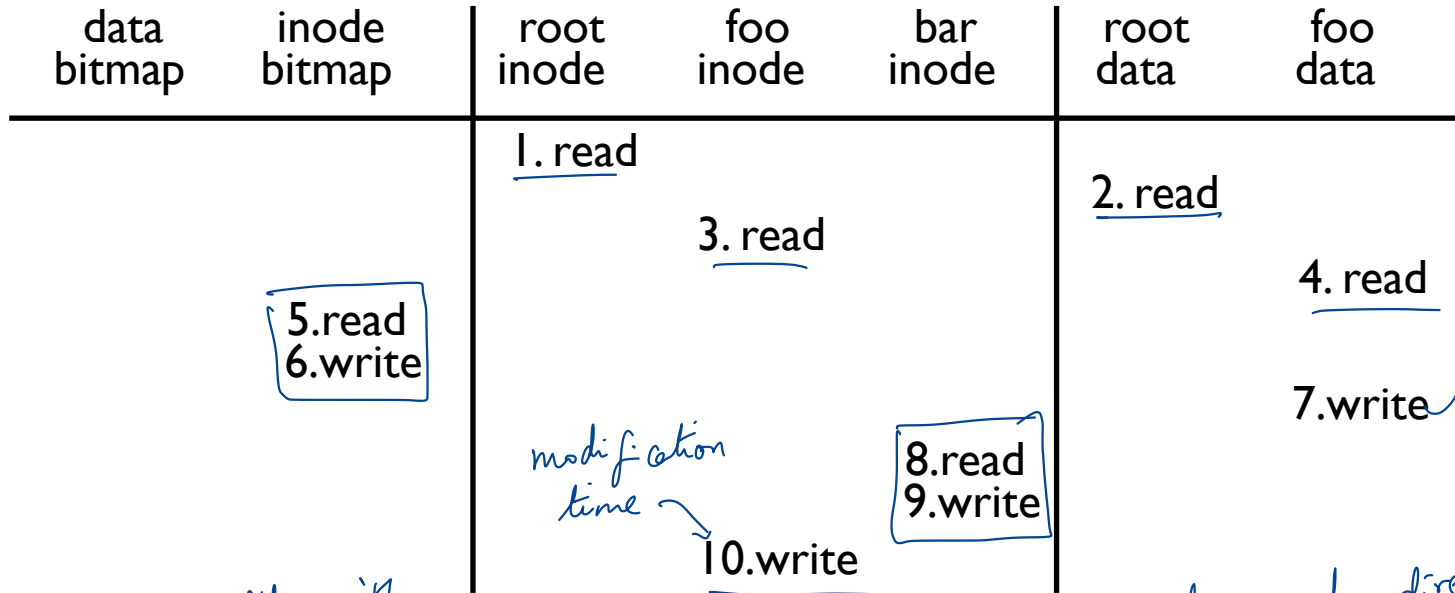
close /foo/bar

data bitmap	inode bitmap	root inode	foo inode	bar inode	root data	foo data	bar data
<i>No I/O operations triggered</i>							

nothing to do on disk!

TIME

create /foo/bar



① Check if a file with same name exists

→ Perform traversal and read directory entry for foo

② Allocate inode for file: read IB & write IB

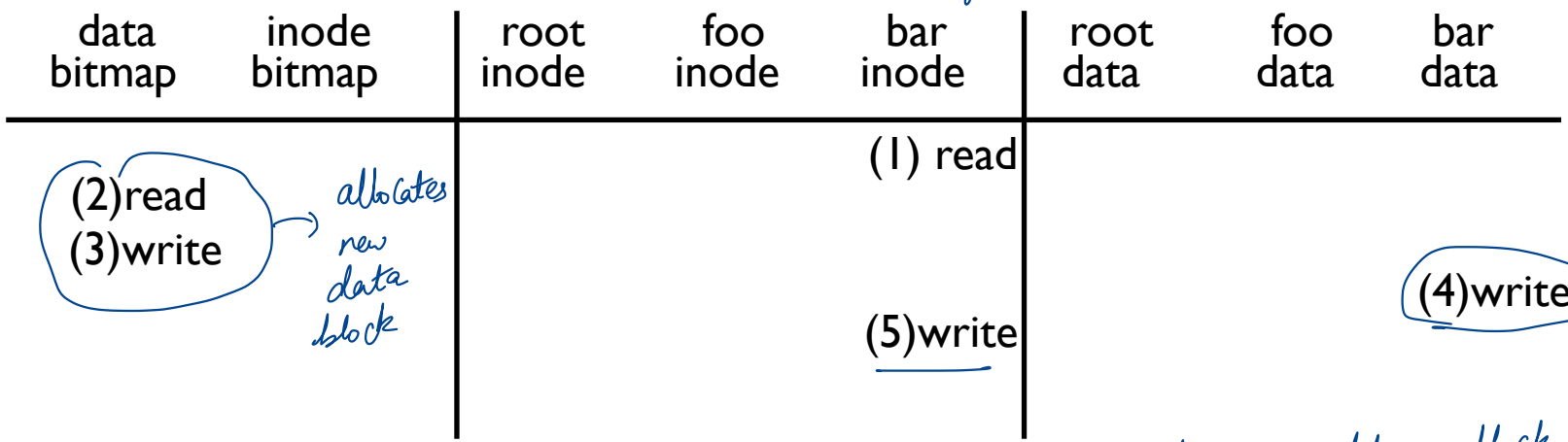
③ Initialize inode → set permissions who created, what time

read inode block
write inode block

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

TIME

offset based → (start, length)



allocates new data block

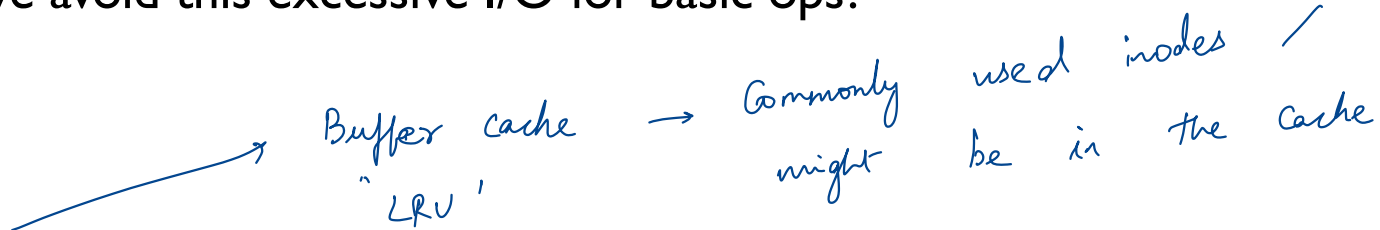
- ① Read inode to figure out whether I need new data block
- ② Allocate DB
- ③ Write to data block
- ④ Update inode to contain ptr to new data block and metadata

EFFICIENCY

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

Cache for:

- reads
- write buffering



WRITE BUFFERING

temporary files

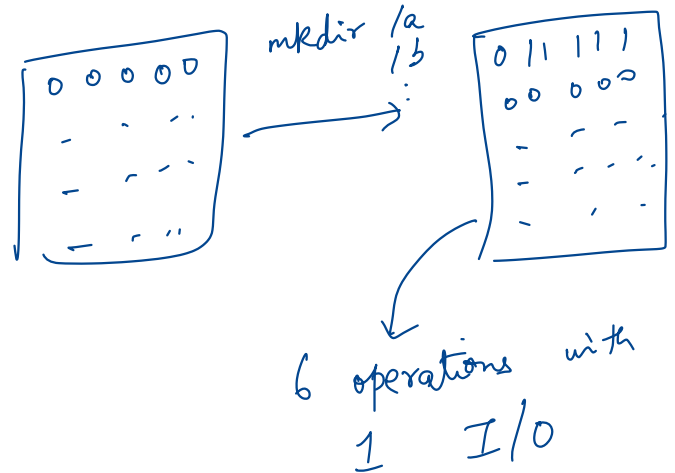
Overwrites, deletes, scheduling

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

Tradeoffs: how much to buffer, how long to buffer

Crash consistency

Inode Bitmap → Every create only
Data Bitmap → changes 1
bit



QUIZ 26

<https://tinyurl.com/cs537-sp23-quiz26>



```
inode bitmap  ??????????
inodes        [d a:0 r:3] [d a:1 r:2] [] [] [] [] [] []
data bitmap   11000000
data          [(.,0) (.,0) (n,1)] [(.,1) (.,0)] [] [] [] [] [] []
```

Handwritten annotations:
- "directory" with arrows pointing to the first two inodes.
- "ref count" with an arrow pointing to the first two inodes.
- "data block" with an arrow pointing to the first two data blocks.

11 000000

```
inode bitmap  11000000
inodes        [d a:0 r:2] [f a:- r:1] [] [] [] [] [] []
data bitmap   10000000
data          [CORRUPT!] [] [] [] [] [] [] []
```

Handwritten annotations:
- "second entry" with an arrow pointing to the second inode.
- A box around the first inode in the inodes row.
- An arrow pointing from the first data block to the "CORRUPT!" text.

↑

.	0
..	0
z	1

creat ("1/2")

QUIZ 26

<https://tinyurl.com/cs537-sp23-quiz26>

```
inode bitmap 11100000
inodes [d a:0 r:4] [d a:1 r:2] [d a:2 r:2] [] [] [] [] []
data bitmap 11100000
data [(.,0) (.,0) (d,1) (w,2)] [????] [(.,2) (.,0)] [] [] [] [] []
```

/
/d
/w

-	1
..	0

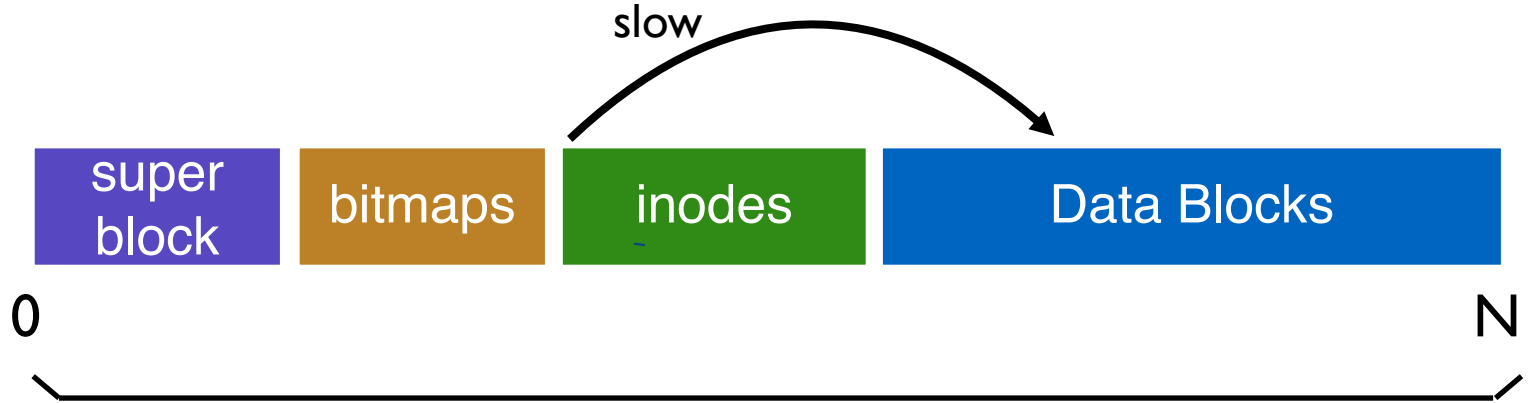
mkdir (/d)
mkdir (/w)

```
inode bitmap 11000000
inodes [d a:0 r:2] [f a:1 r:2] [] [] [] [] [] []
data bitmap ????????
data [(.,0) (.,0) (c,1) (m,1)] [foofoofoo] [] [] [] [] [] []
```

- ① Both /c and /m are links to same file
- ③ Proper root directory
- ④ First block has foofoofoo

FAST FILE SYSTEM

FILE LAYOUT IMPORTANCE



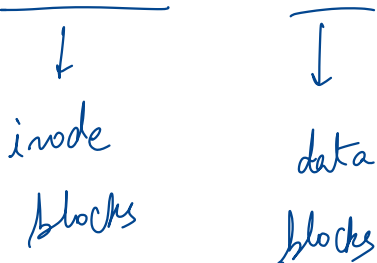
Layout is not disk-aware!

DISK-AWARE FILE SYSTEM

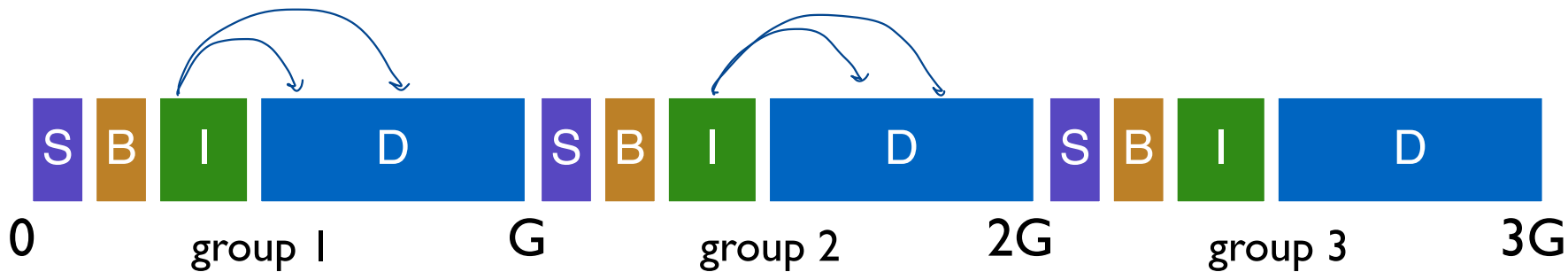
How to make the disk use more efficient?

FFS → Fast File
System

Where to place meta-data and data on disk?



PLACEMENT TECHNIQUE: GROUPS



Key idea: Keep inode close to data

Use groups across disks;

Strategy: allocate inodes and data blocks in same group.

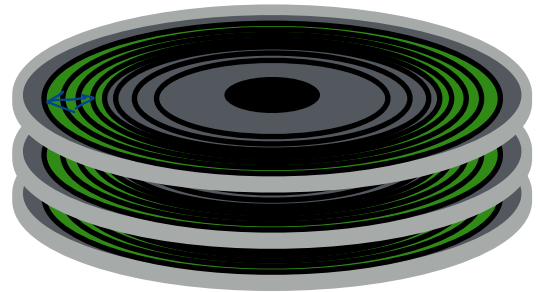
for files/directories

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

*inodes and data blocks
would be in nearby cylinders*



REPLICATED SUPER BLOCKS

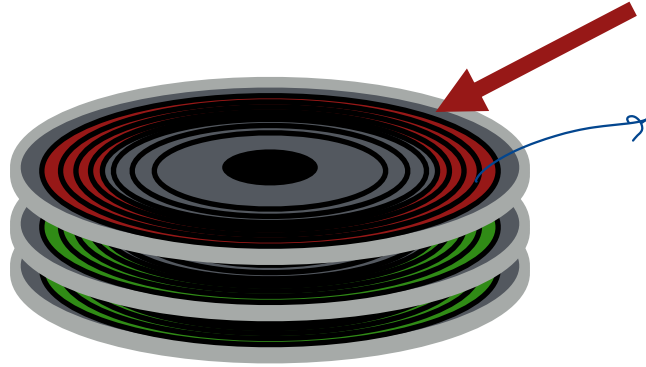
if SB is corrupted → no way to mount the FS



0 →

G
replicated super blocks

2G
Fault tolerance



top platter damage?

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

to overcome all
replicated SBs
becoming unavailable

SMART POLICY



Policy:

Where should new inodes and data blocks go?

for files and directories

PLACEMENT STRATEGY

Put related pieces of data near each other.

Rules:

1. Put directory entries near directory inodes.
2. Put inodes near directory entries.
3. Put data blocks near inodes.



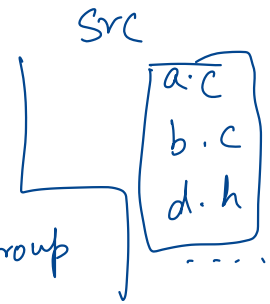
Problem: File system is one big tree

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



Put "more-related" pieces of data near each other

Put less-related pieces of data **far**

ls -R /
 /a/b
 /a/c
 /a/d
 /b/f

mkdir (/) 0
 mkdir (/a) 1
 creat (/a/c) 1

new directory
 new group

Files inside a
 dir are in same
 group

disk
 which inodes
 are in which group

group	inodes	data
0	/-----	/-----
1	<u>acde</u> -----	<u>accddee</u> ---
2	<u>bf</u> -----	<u>bff</u> -----
3	-----	-----
4	-----	-----
5	-----	-----
6	-----	-----
7	-----	-----
...		

POLICY SUMMARY

File inodes: allocate in same group with dir → parent

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

First data block: allocate near inode

data in same
grp as inode

↪ load balancing

Other data blocks: allocate near previous block

Most files are small

Most space is used
by large files

PROBLEM: LARGE FILES

Single large file can fill nearly all of a group

Displaces data for many small files

Big file → large number of
data blocks

group	inodes	data
0	/a-----	/aaaaaaaaa aaaaaaaaaa aaaaaaaaaa a-----
1	-----	-----
2	-----	-----
...		

Most files are small!

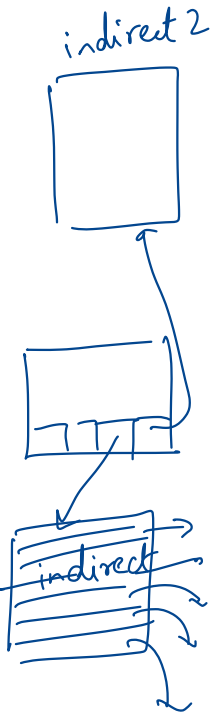
→ Better to do one seek for large file than
one seek for each of many small files →

how compilers
might access
lots of small
files

SPLITTING LARGE FILES

group	inodes	data
0	/a-----	<u>/aaaaa</u> -----
1	-----	<u>aaaaa</u> -----
2	-----	<u>aaaaa</u> -----
3	-----	aaaaa-----
4	-----	aaaaa-----
5	-----	aaaaa-----
6	-----	-----

save space which can be used by small files



...

Define “large” as requiring an indirect block

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

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 1MB.

OTHER FFS FEATURES

FFS also introduced several new features:

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

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

Next class: Filesystem consistency