

PERSISTENCE: FILE SYSTEMS & FFS

Shivaram Venkataraman

CS 537, Spring 2019

ADMINISTRIVIA

Project 4b: Due next week 4/16

Project 5: One project 9%. Updated due dates on website

Discussion this week: Review worksheet, More Q&A for 4b

AGENDA / LEARNING OUTCOMES

How does file system represent files, directories?

What steps must reads/writes take?

How does FFS improve performance?

↓
Fast file myKcm

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 ← path that is passed to open
- hierarchical
- traverse once ← traverse
- offsets precisely defined

FILE, DIRECTORY API SUMMARY

Using multiple types of name provides convenience and efficiency

Mount and link features provide flexibility.

Special calls (fsync, rename) let developers communicate requirements to file system

atomic writes using fsync, rename (old-name, new-name)
↳ atomicity

cp file.txt file.txt.tmp
(operate on tmp) ||

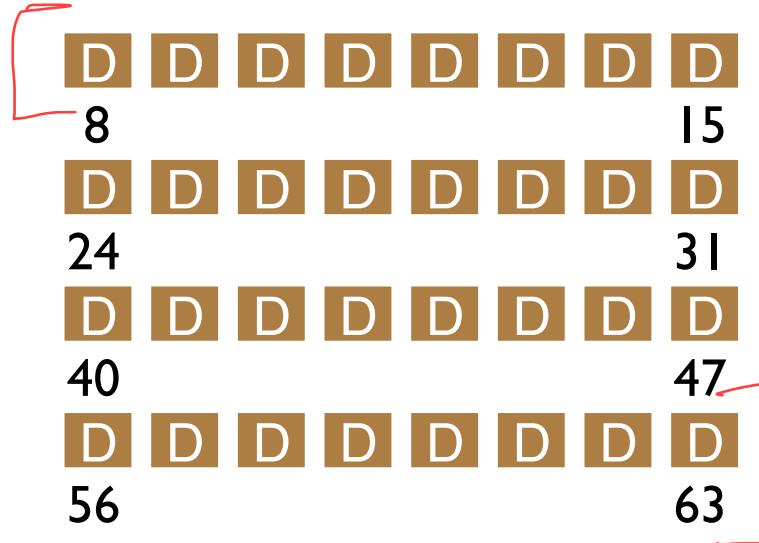
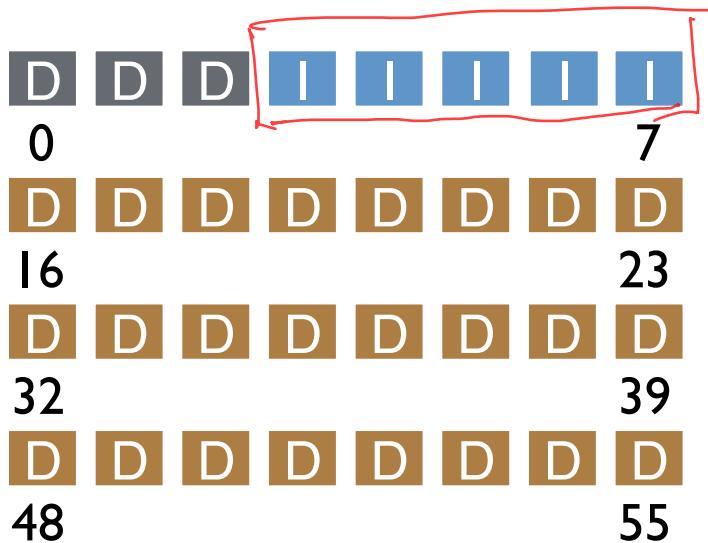
fsync file.txt.tmp
rename (file.txt.tmp, file.txt)

Metadata
which
has data for
which node

FS LAYOUT

→ Very Simple FS

inode blocks



Data
blocks

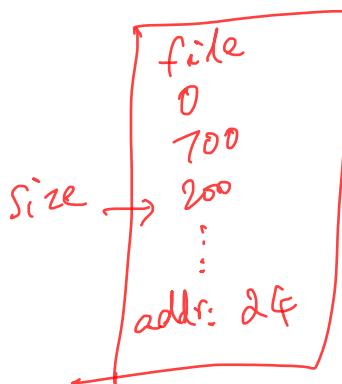
INODE

↳ Note

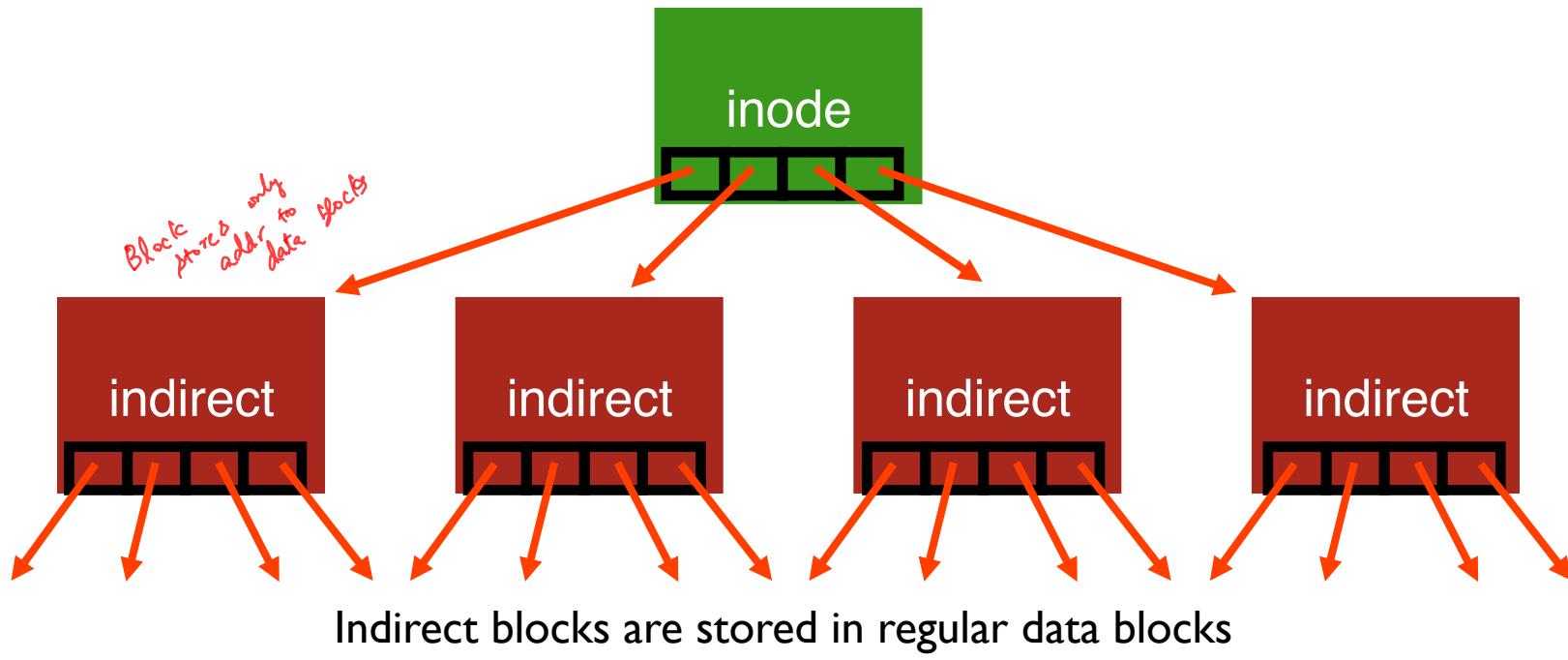
type (file or dir?)
uid (owner)
rwx (permissions)
size (in bytes)
Blocks
time (access)]]
ctime (create)]
links count (# paths)
→ addrs[N] (N data blocks)

What is max file size with single level?

Assume 256-byte inodes
(all can be used for pointers)
Assume 4-byte addrs



Direct pointers
↳ each block is 4 KB
Inode size is 256 byte
Each addr 4 byte
⇒ 64 addr in 1 inode
⇒ $64 \times 4 \text{ KB} = 256 \text{ KB}$



Largest file size with 64 indirect blocks?

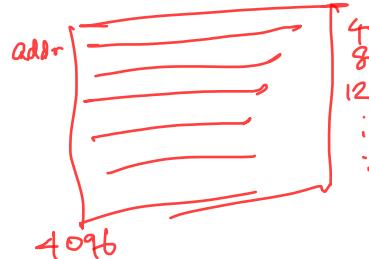
inode is 256 byte 64 ptrs to indirect block

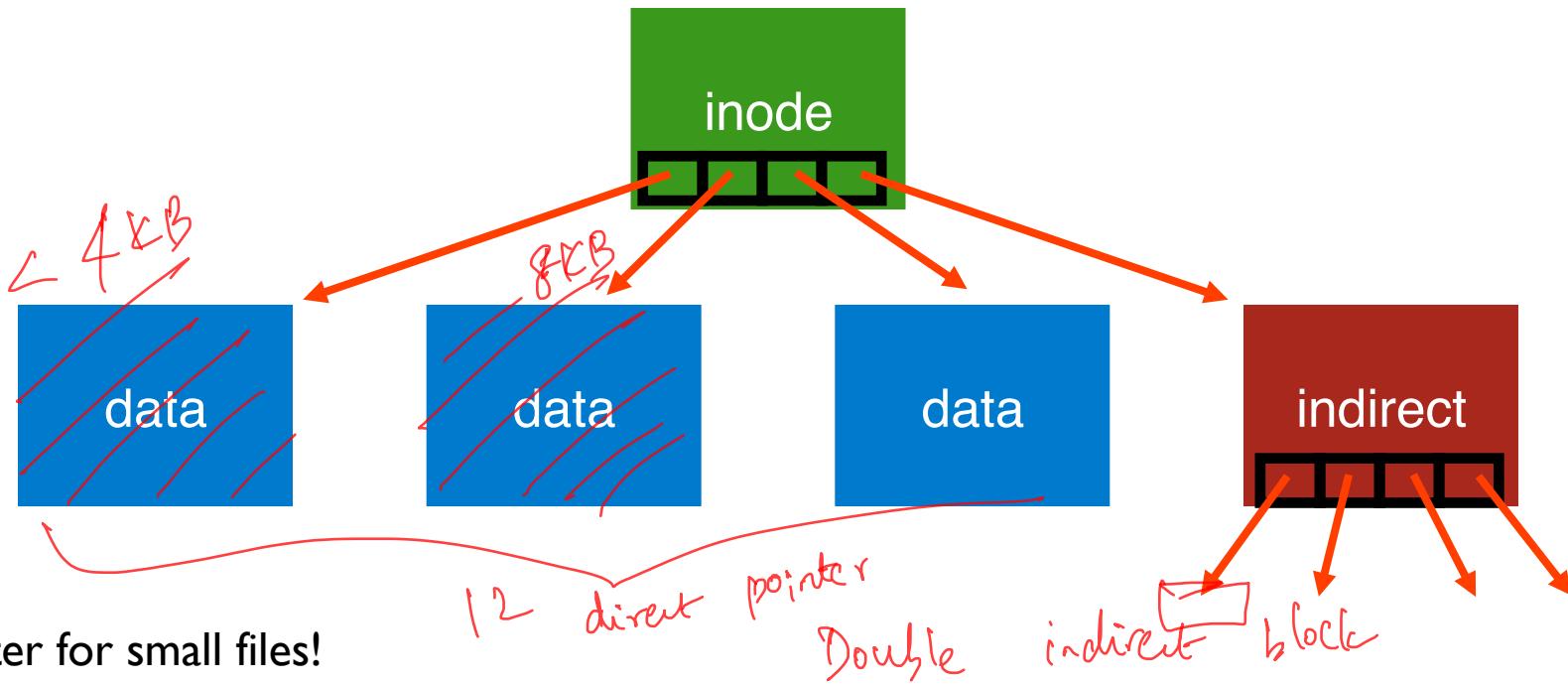
1 indirect block $\leq 4KB$, each addr 4 bytes
 ≤ 1024 addr in 1 indirect block

$$\leq 1024 \times 4KB = 4MB$$

$$64 \text{ indirect} = 64 \times 4MB = 256 MB$$

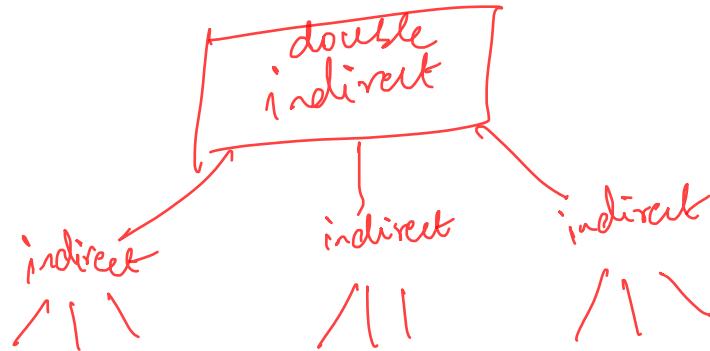
Any Cons?





Better for small files!

How to handle even larger files?



BUNNY 15



<https://tinyurl.com/cs537-sp19-bunny15>

BUNNY 15

Assume 256 byte inodes (16 inodes/block).
What is the offset for inode with number 0?

12 kB

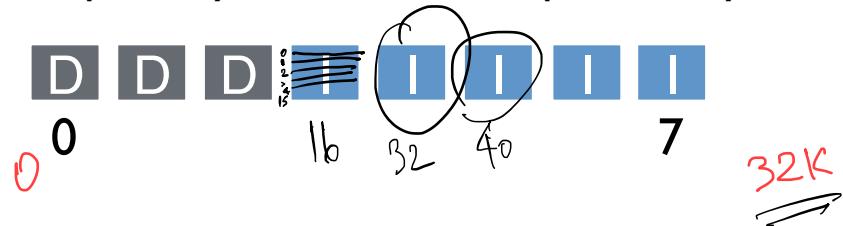
What is the offset for inode with number ~~0~~ 4?

$$12\text{ kB} + 4 \times 256 \equiv 13\text{ kB}$$

What is the offset for inode with number ~~0~~ 40?

$$12\text{ kB} + 40 \times 256 \equiv 22\text{ kB}$$

<https://tinyurl.com/cs537-sp19-bunny15>



DIRECTORIES

File systems vary

Common design:

Store directory entries in data blocks

Large directories just use multiple data blocks

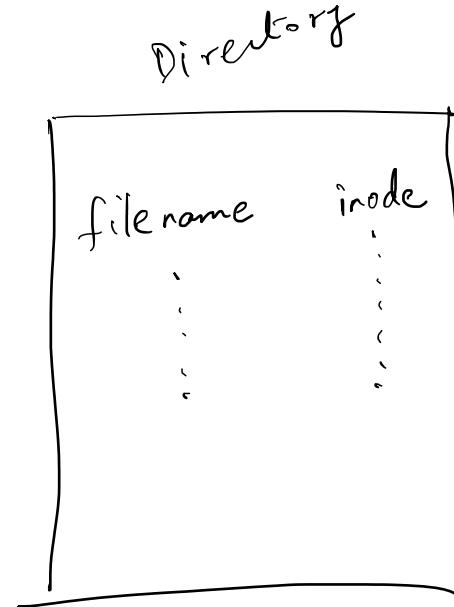
Use bit in inode to distinguish directories from files

type bit

Various formats could be used

- lists

- b-trees



SIMPLE DIRECTORY LIST EXAMPLE

valid	name	inode
-	.	134
-	..	35
✓ 0	foo	80
-	bar	23

special entries

could also be
directory

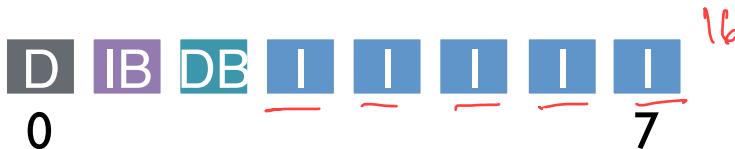
these

→ `unlink("foo")`

Create ("os")

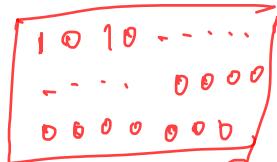
FS STRUCTS: BITMAPS

How do we find free data blocks or free inodes?



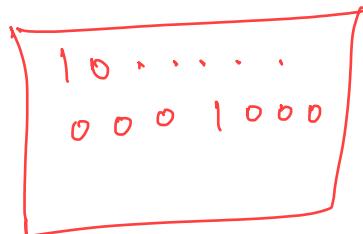
→ free list
list of addr which
are free

Data bitmap = data structure that has 1-bit to indicate data block is used or not



→ 56 bits for our layout

Inode bitmap



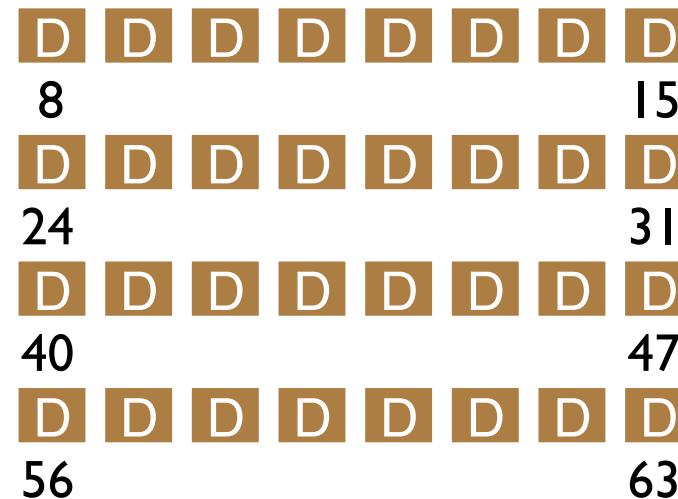
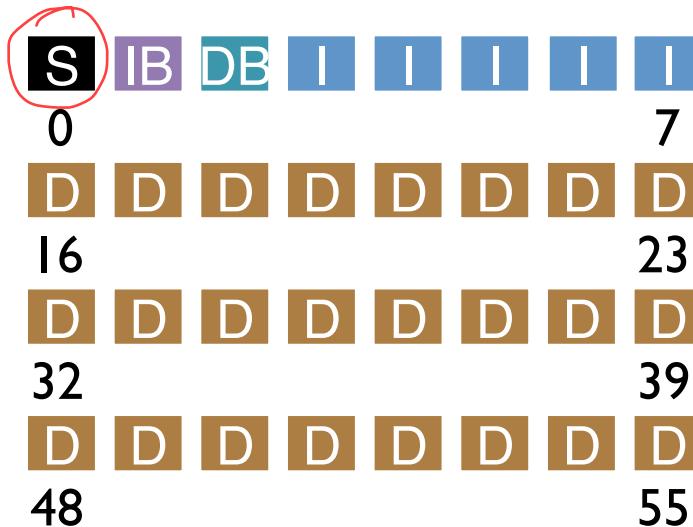
1-bit to indicate if an I-node is used or not

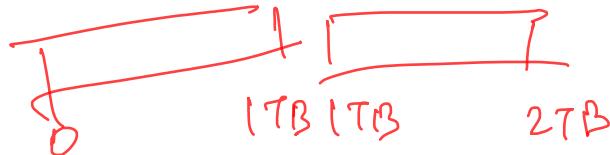
→ 80 bits for our layout

“magic NC-ing”

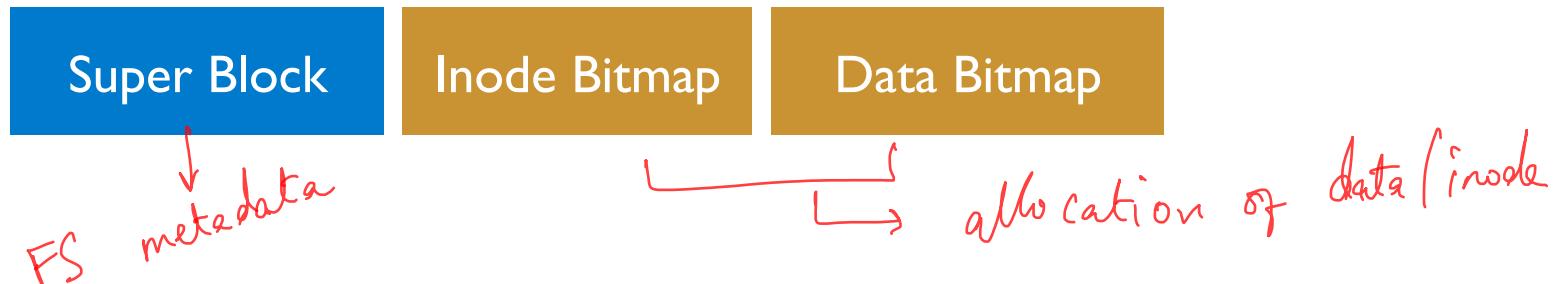
FS STRUCTS: SUPERBLOCK

Basic FS configuration metadata, like block size, # of inodes





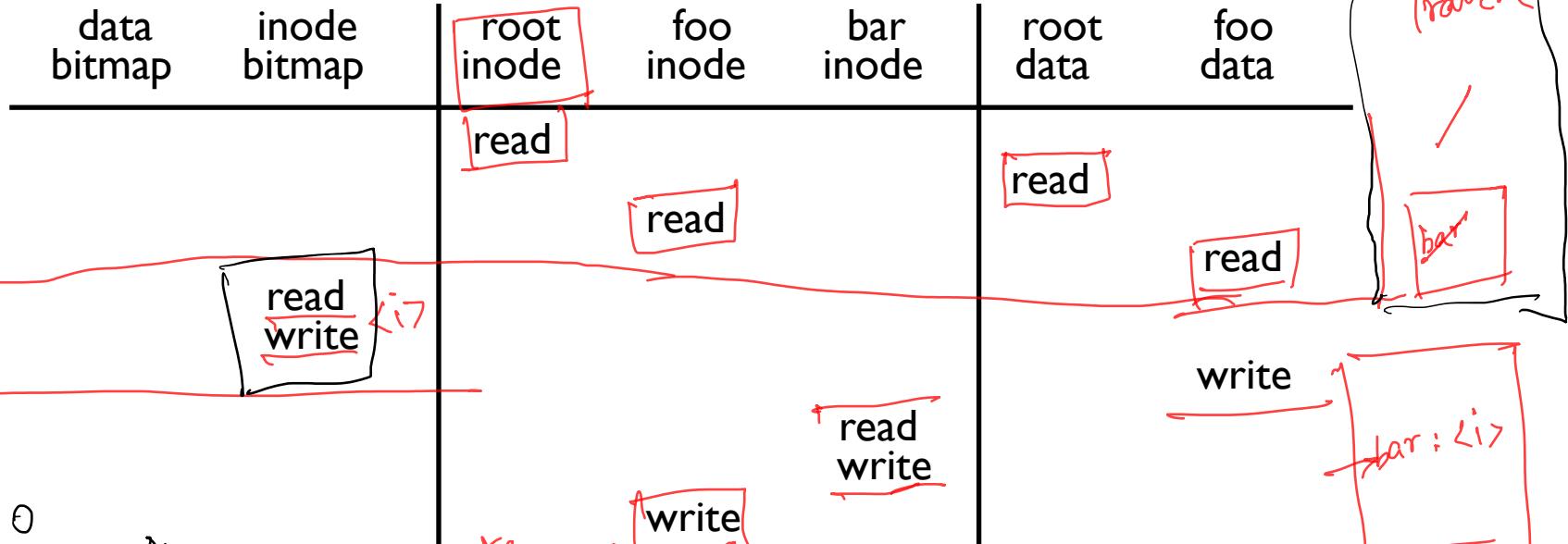
SUMMARY



FS OPERATIONS

- create file
- write
- open
- read
- close

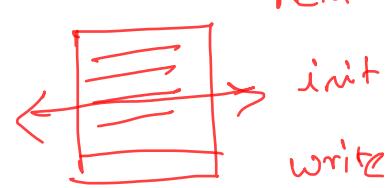
create /foo/bar



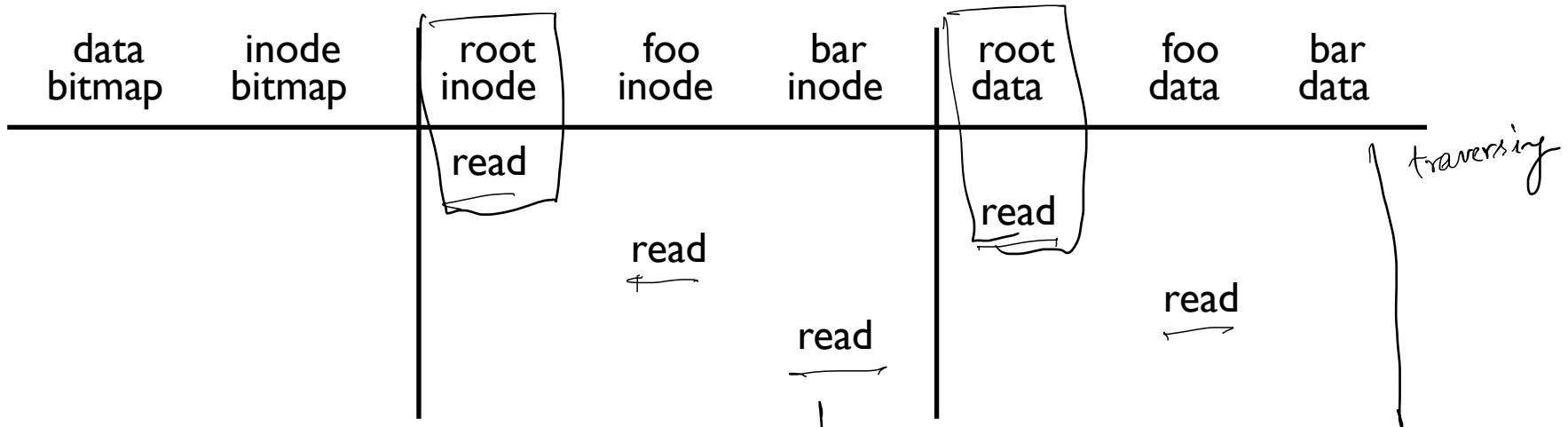
Why must **read** for bar inode?

Initialize

owner
type
permission



open /foo/bar

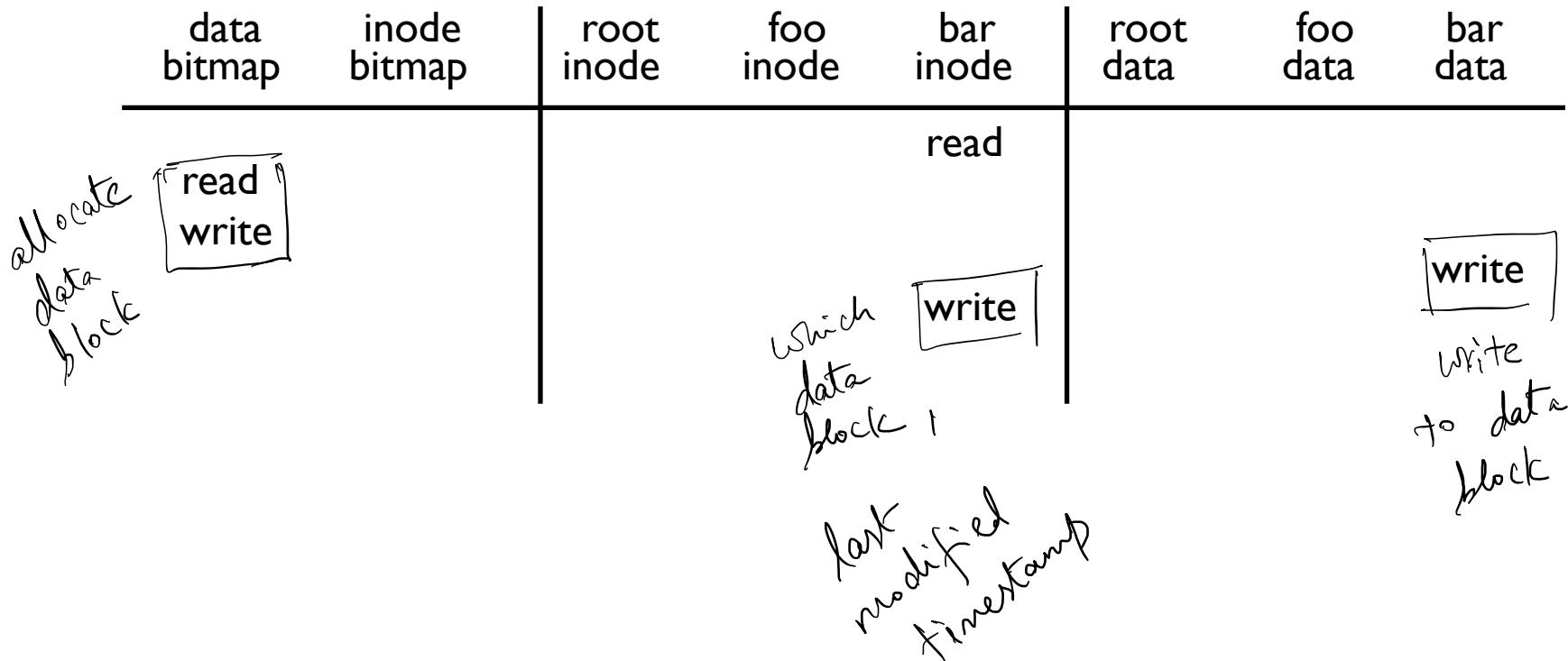


meal
for
a file

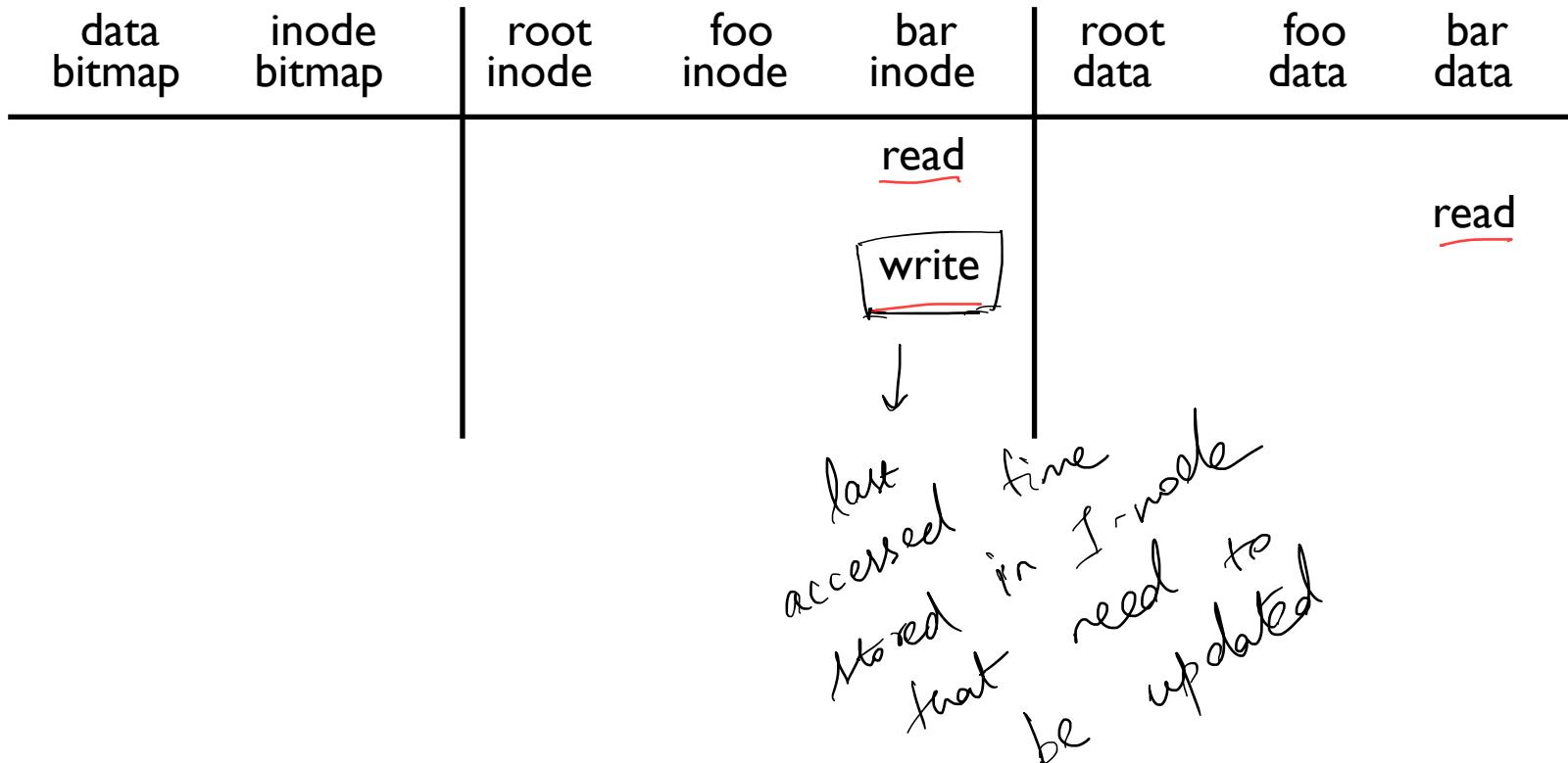
opening
a file

Check
permission
if user
can open file

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



read /foo/bar – assume opened



close /foo/bar

→ garbage collect
fd data str. in
memory

data bitmap	inode bitmap	root inode	foo inode	bar inode	root data	foo data	bar data

nothing to do on disk!

EFFICIENCY

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

Cache for:

- reads
- write buffering

commonly read
data or I-node in memory

" buffer cache" → pages of
 used to cache
 & I-nodes
 data

WRITE BUFFERING

Overwrites, deletes, scheduling

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

Tradeoffs: how much to buffer, how long to buffer

1. write to file || OS needs to
 delete file || do no I/O
2. Sequence of writes (fd, "hello") ; write (fd, "world") || Only needs to respond to write

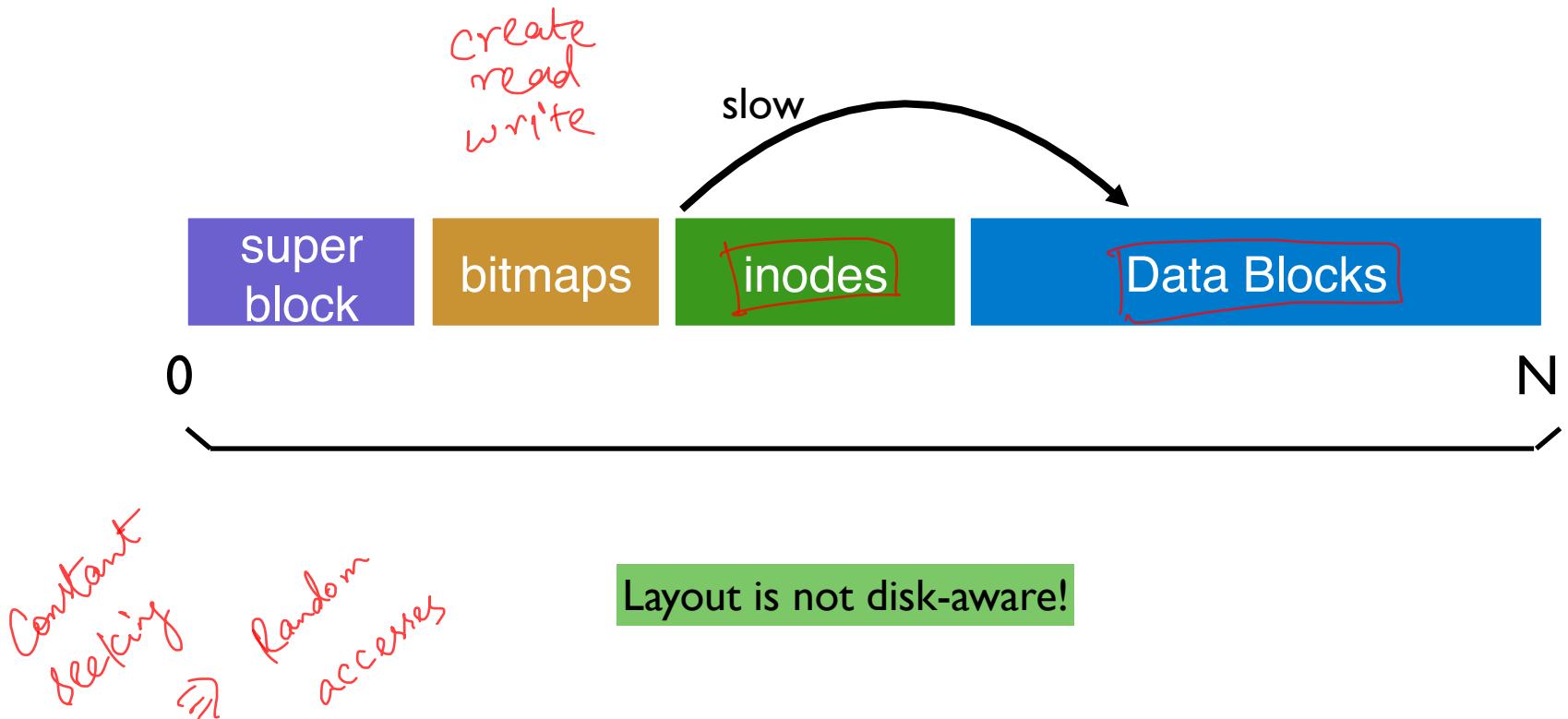
write (fd, "hello")

OS
fsync (fd) → buffer
write

FAST FILE SYSTEM

1980's Software Distribution
Berkeley BSD "based on UNIX"
Open File Net

FILE LAYOUT IMPORTANCE



DISK-AWARE FILE SYSTEM

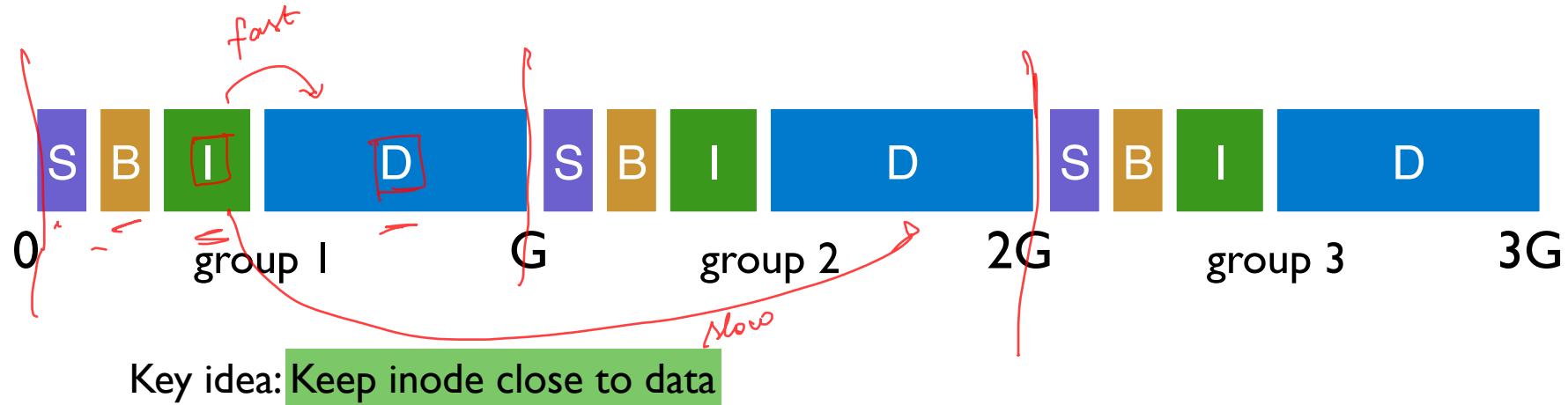
Given the same API

How to make the disk use more efficient?

Where to place meta-data and data on disk?



PLACEMENT TECHNIQUE: GROUPS



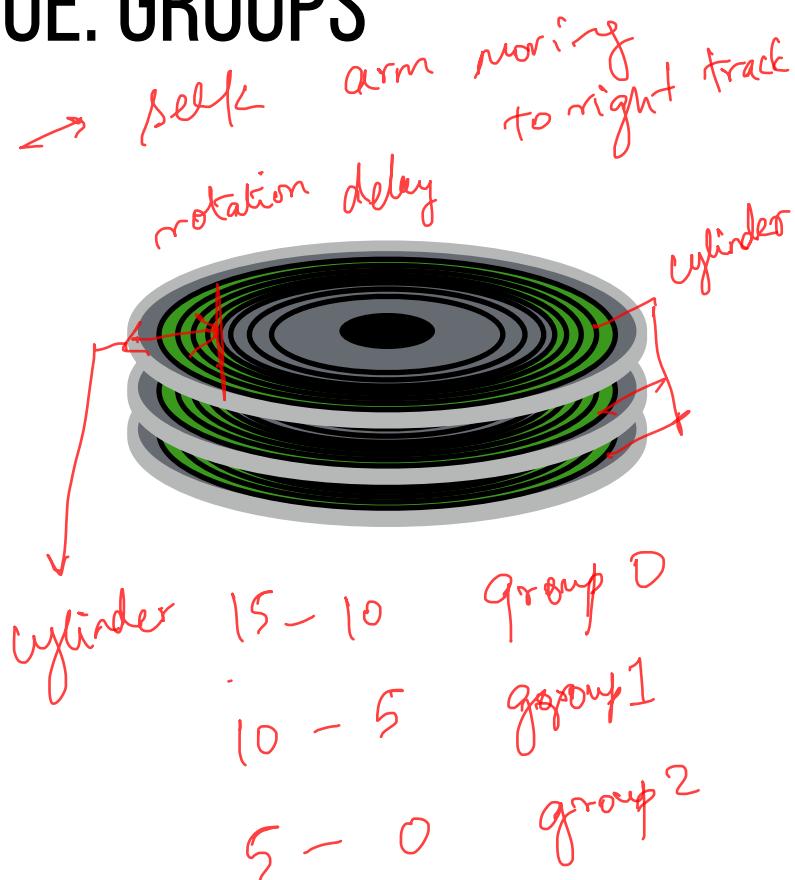
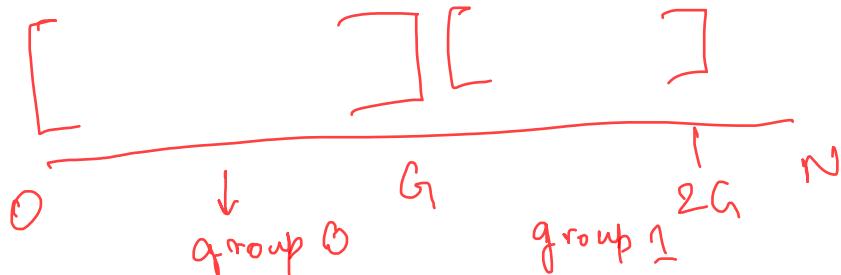
Use groups across disks;

Strategy: allocate inodes and data blocks in same group.

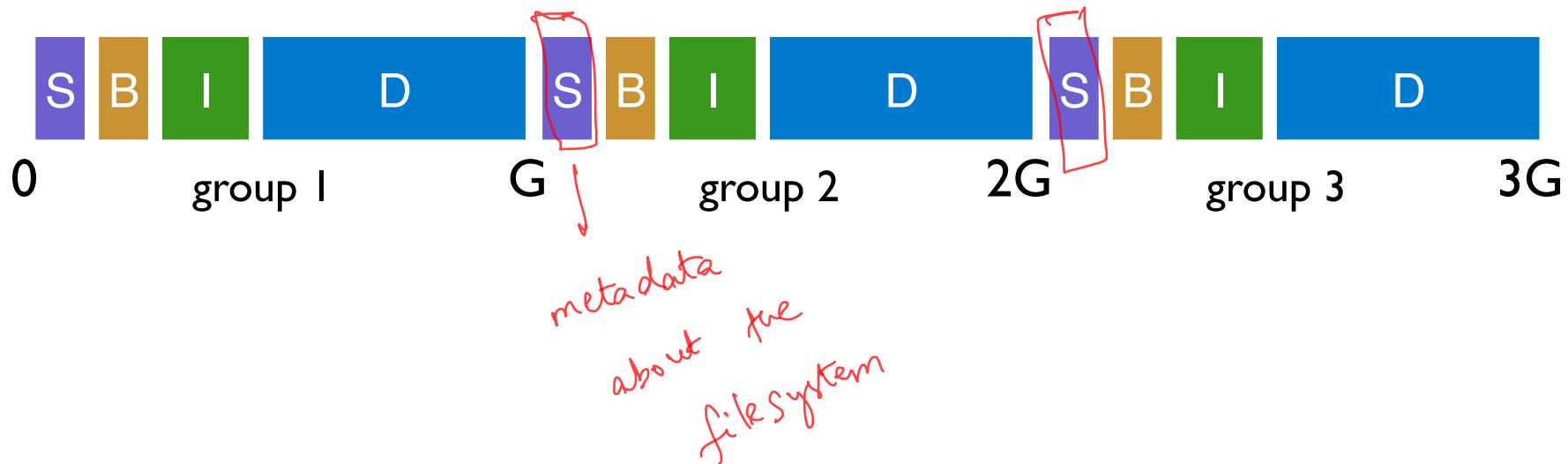
PLACEMENT TECHNIQUE: GROUPS

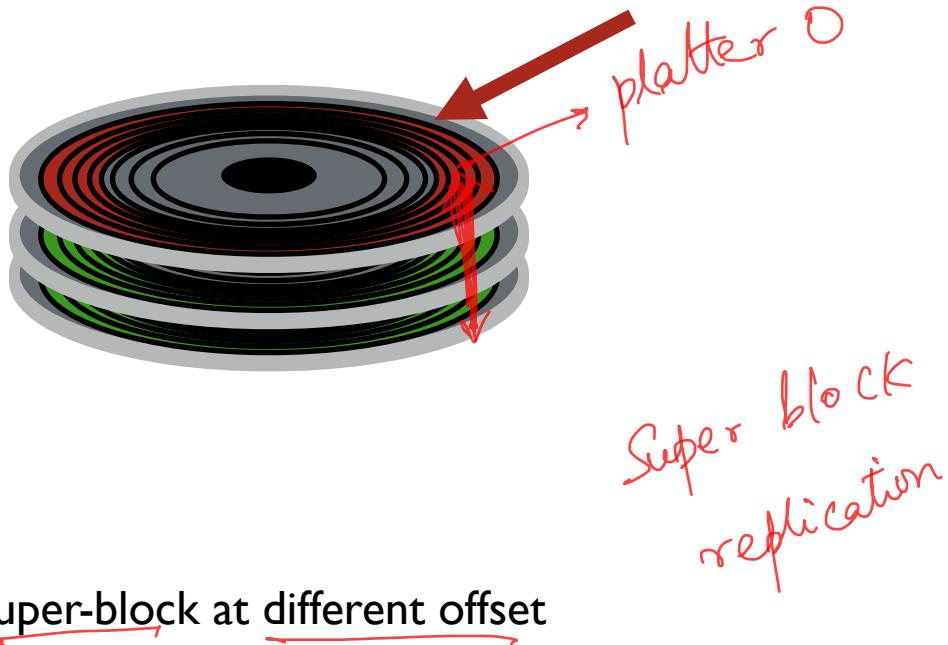
In FFS, groups were ranges of cylinders
called cylinder group

In ext2, ext3, ext4 groups are ranges of blocks
called block group



REPLICATED SUPER BLOCKS





SMART POLICY

↳ *create /a dir
create /a/b file*



Where should new inodes and data blocks go?

PLACEMENT STRATEGY

Put related pieces of data near each other.

Rules:

1. Put directory entries near directory inodes.
2. Put inodes near directory entries. file inodes
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**

Trace of operation

Create /a
/a/b/e
/a/c
/a/d
/b/f

Create /b/g

Compiling all .c
into a binary

mkdir /a

mkdir /b

Create /a/c

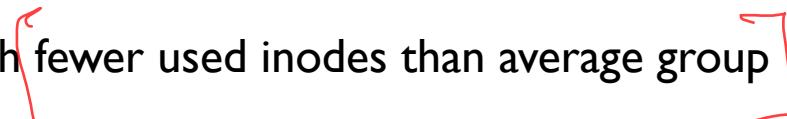
/a/d

/b/f

group	inodes	data
0	/	/
1	acde	accddee
2	bffg	bfffggg
3		
4		
5		
6		
7		
	...	

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

PROBLEM: LARGE FILES

Single large file can fill nearly all of a group

Displaces data for many small files

/a* used up all data
/b different group
/c

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

Most files are small!

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

SPLITTING LARGE FILES

group inodes data

group	inodes	data	Chunk	meta 1	group 2
0	/a-----	/aaaaaa ^b			
1	-----	aaaaaa ^b			
2	-----	aaaaaa			
3	-----	aaaaaa			
4	-----	aaaaaa			
5	-----	aaaaaa			
6	-----	-----			
...					

number
of
levels

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

$1024 \times 4\text{KB} = 4\text{MB}$ contiguous “chunk”

BUNNY 16



<https://tinyurl.com/cs537-sp19-bunny16>

BUNNY 16

Assume that the average positioning time (i.e., seek and rotation) = 10 ms.

Assume that disk transfers data at 100 MB/s.

If FFS large file chunk size is 4MB, what is the effective throughput we are getting?

What is the effective throughput with 8MB chunk size?

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.

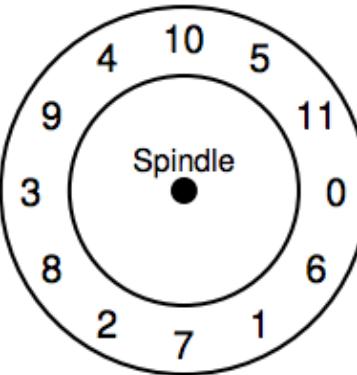
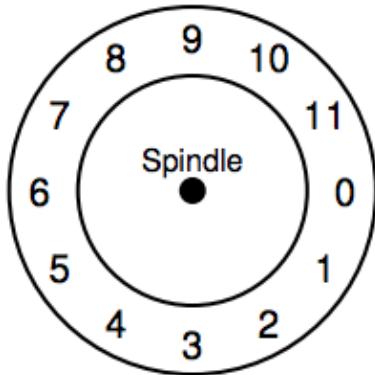
Move to another group (w/ **fewer than avg blocks**) every subsequent 1MB.

OTHER FFS FEATURES

FFS also introduced several new features:

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

FFS: SECTOR PLACEMENT



Similar to track skew in disks chapter

Modern disks:
Disk cache

FFS SUMMARY

First disk-aware file system

- Bitmaps
- Locality groups
- Rotated superblocks
- Smart allocation policy

Inspired modern files systems, including ext2 and ext3

OTHER TAKEAWAYS

All hardware is unique

Treat disk like disk!

Treat flash like flash!

Treat random-access memory like random-access memory!

NEXT STEPS

Next class: How to provide consistency despite failures?

Discussion today: Worksheet with problems, Q&A for project 4b