

Hello!

# PERSISTENCE: FILE SYSTEMS

Shivaram Venkataraman

CS 537, Fall 2024

# ADMINISTRIVIA

Midterm 2? — *This week*

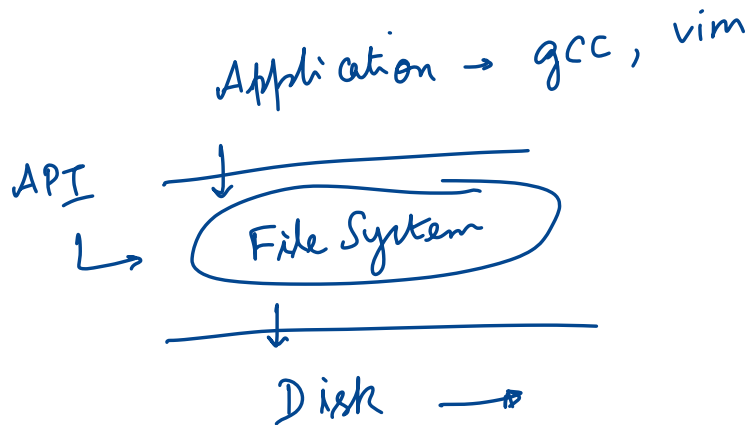
Project 5 — Due Nov 19<sup>th</sup> (one week from today)

Shivaram's OH / next week schedule → *Zoom office hours*  
*next week*  
↳ *Today 3-4 pm*

# AGENDA / LEARNING OUTCOMES

How does file system represent files, directories?

What steps must reads/writes take?



**RECAP**

# FILE API WITH FILE DESCRIPTORS

*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)
```

→ *Impl*

advantages:

- string names
- hierarchical
- traverse once
- offsets precisely defined

→ *track*

# STRACE

```
prompt> echo hello > foo
```

```
prompt> cat foo
```

```
hello
```

```
prompt>
```

```
prompt> strace cat foo -- prints system calls performed by program
```

```
...
```

```
open("foo", O_RDONLY|O_LARGEFILE) = 3
```

fd for foo

```
read(3, "hello\n", 4096) = 6
```

```
write(1, "hello\n", 6) = 6
```

6 bytes

```
hello
```

```
read(3, "", 4096) = 0
```

stdout

```
close(3) = 0
```

```
...
```

```
prompt>
```

# LINKS

metadata associated  
with a file

Hard links: Both path names use same inode number

File does not disappear until all hard links removed; cannot link directories

short cut

Soft or symbolic links: Point to second path name; can softlink to dirs

Can cause confusing behavior: "file does not exist"!

ls /usr/lib/libpthread.so  
→ libpthread-1.5.so

# FILE API SUMMARY

Using multiple types of name provides convenience and efficiency

Hard and soft link features provide flexibility.

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



# FILESYSTEM DISK STRUCTURES

# FS STRUCTS: EMPTY DISK

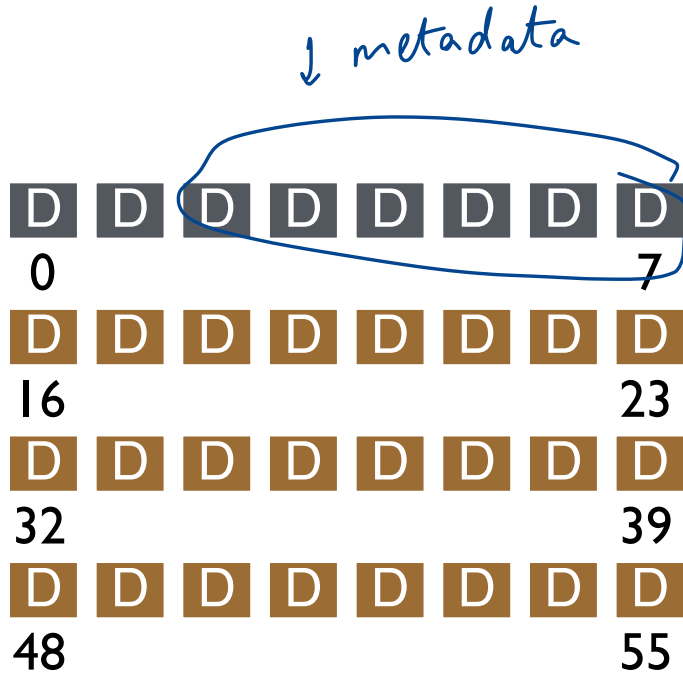
layout a  
filesystem on this disk

block = 4KB  
size



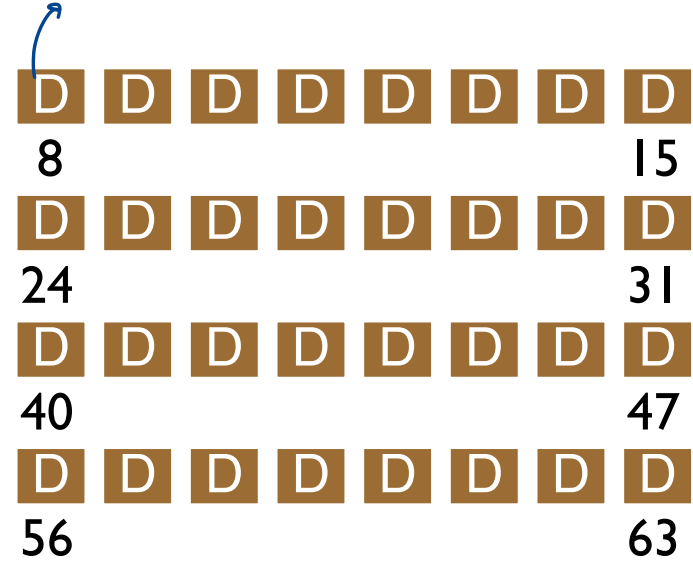
Assume each block is 4KB

# FS STRUCTS: DATA BLOCKS



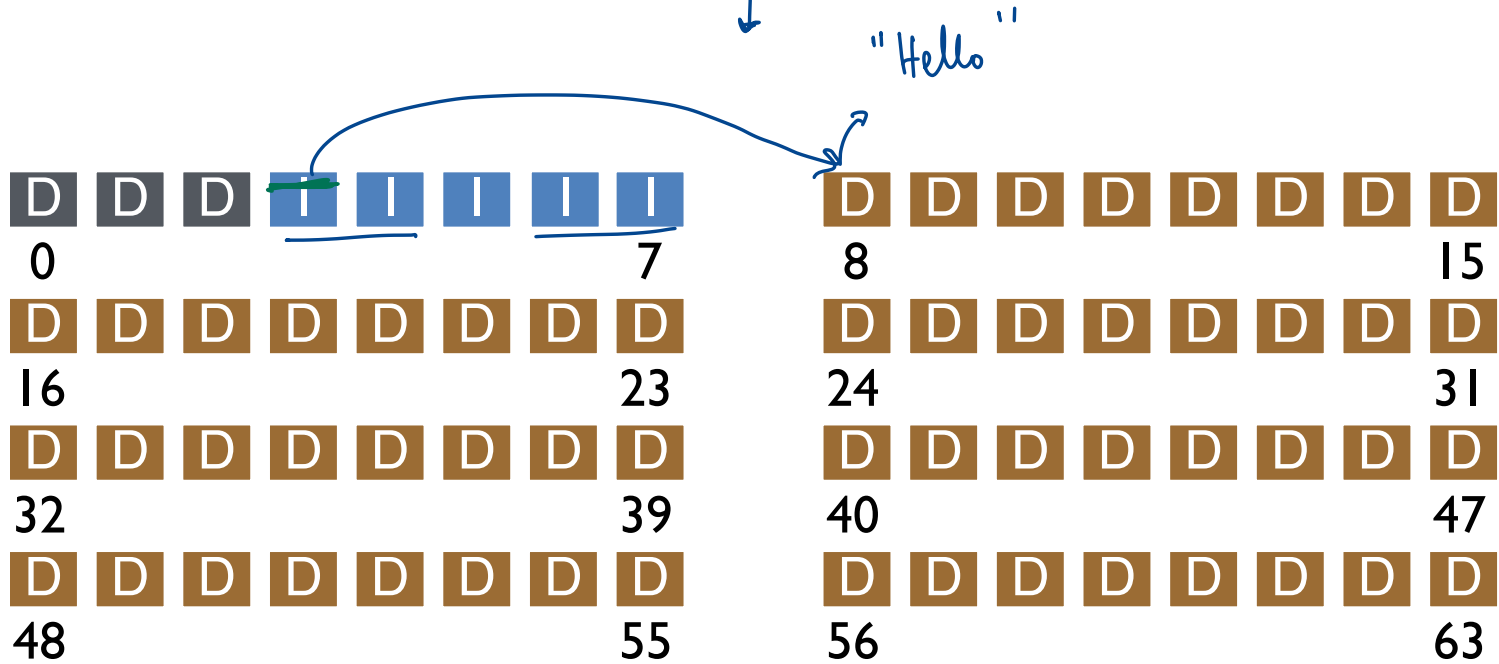
"Hello"

echo "Hello" > foo.txt



Simple layout → Very Simple File System

# FS STRUCTS: INODE DATA POINTERS



/foo.txt  
↳ Directory

# INODE

foo.txt

file

user 1

can read/write etc.

type (file or dir?)

uid (owner)

rxw (permissions)

size (in bytes)

Blocks

time (access)

ctime (create)

links\_count (# paths)

addrs[N] (N data blocks)

number of hard links

links to where this file's content are present

# ONE INODE BLOCK

Each inode is typically 256 bytes (depends on the FS, maybe 128 bytes)

4KB disk block

16 inodes per inode block. =  $\frac{4KB}{256} = 16 \text{ inodes}$

5 blocks for inodes  
=> 80 inodes in my simple FS

*inside 1 block*

inode 16	inode 17	inode 18	inode 19
inode 20	inode 21	inode 22	inode 23
inode 24	inode 25	inode 26	inode 27
inode 28	inode 29	inode 30	inode 31

# INODE

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

Assume single level (just pointers to data blocks)

What is max file size?

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

Assume 4-byte addr

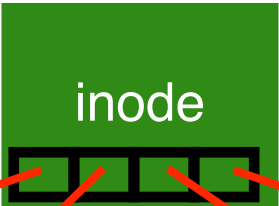
⇒ 64 addr / pointers

How to get larger files?

Max file size

$$= 64 \times 4 \text{ KB} = 256 \text{ KB}$$

*direct pointers*



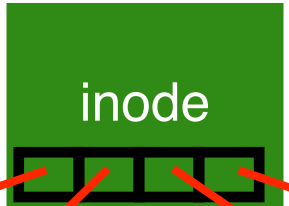
*single-level of pointers*



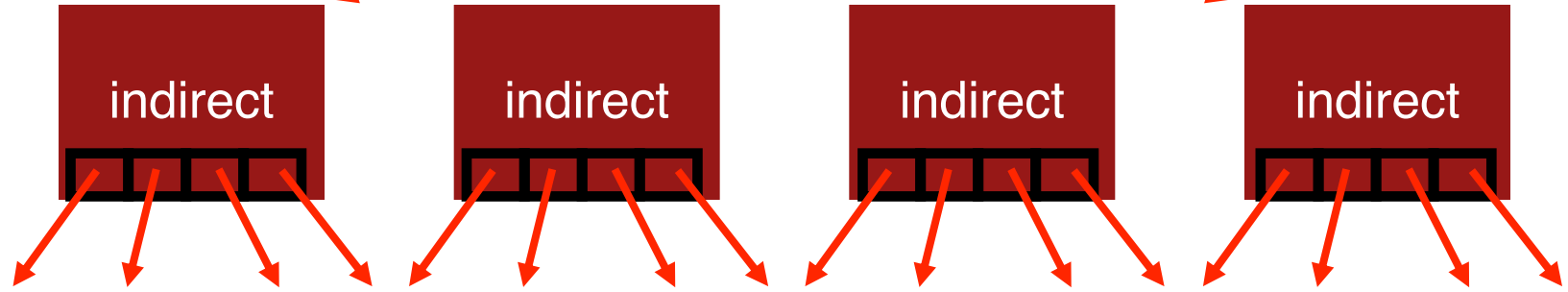


① Access will be slow  
→ more levels

② space wasted  
↳ especially small files



64 ptrs



Indirect blocks are stored in regular data blocks

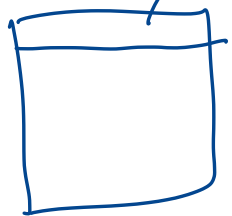
64 x 4MB  
= 256 MB

Largest file size with 64 indirect blocks?

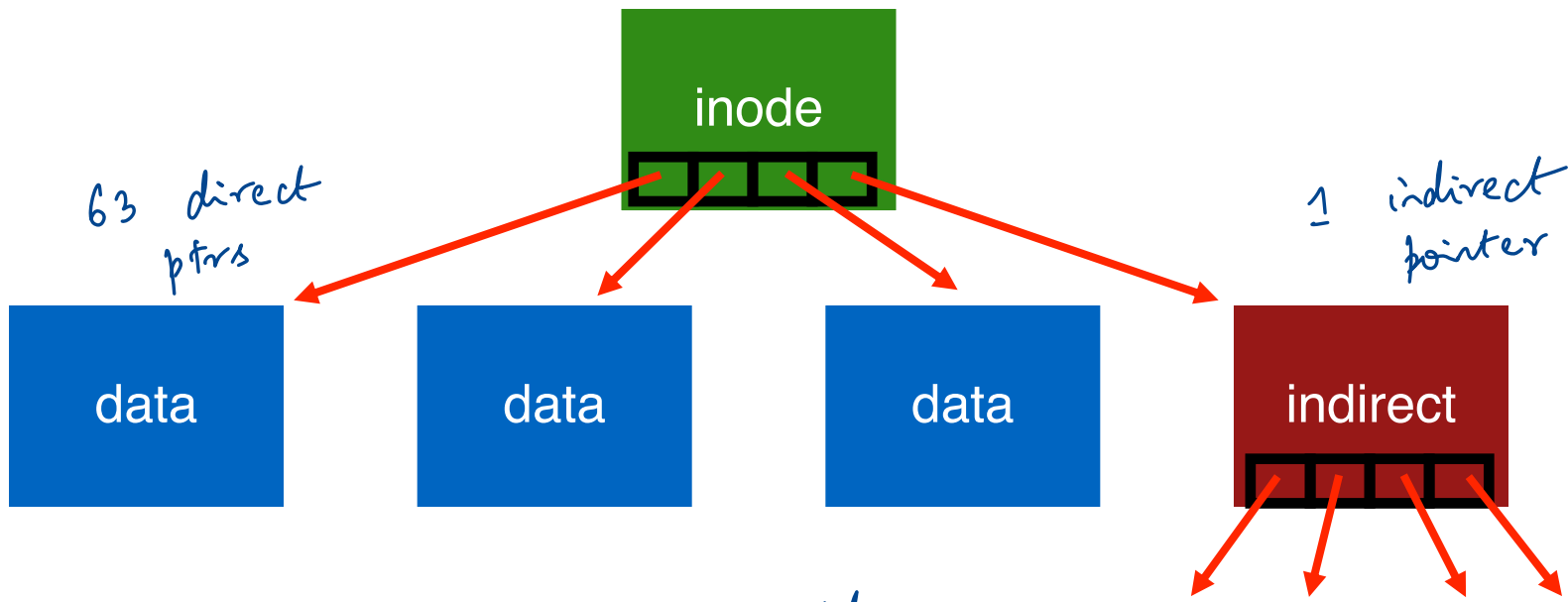
addr data block = 4 bytes  
4KB size

Any Cons?

$\frac{1 \text{ indirect block}}{4} = 1024 \text{ address}$



$1024 \times 4KB = 4MB$   
↳ 1 indirect block

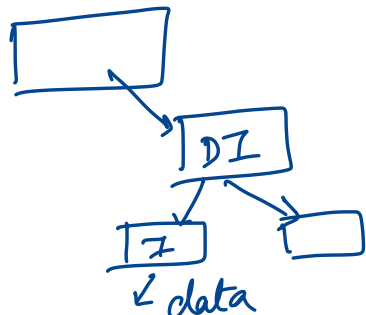


Better for small files!  
 How to handle even larger files?

directly store  
 on data blocks

& 1 traversal

→ double indirect block  
 → triple indirect block



# DIRECTORIES

File systems vary

Common design:

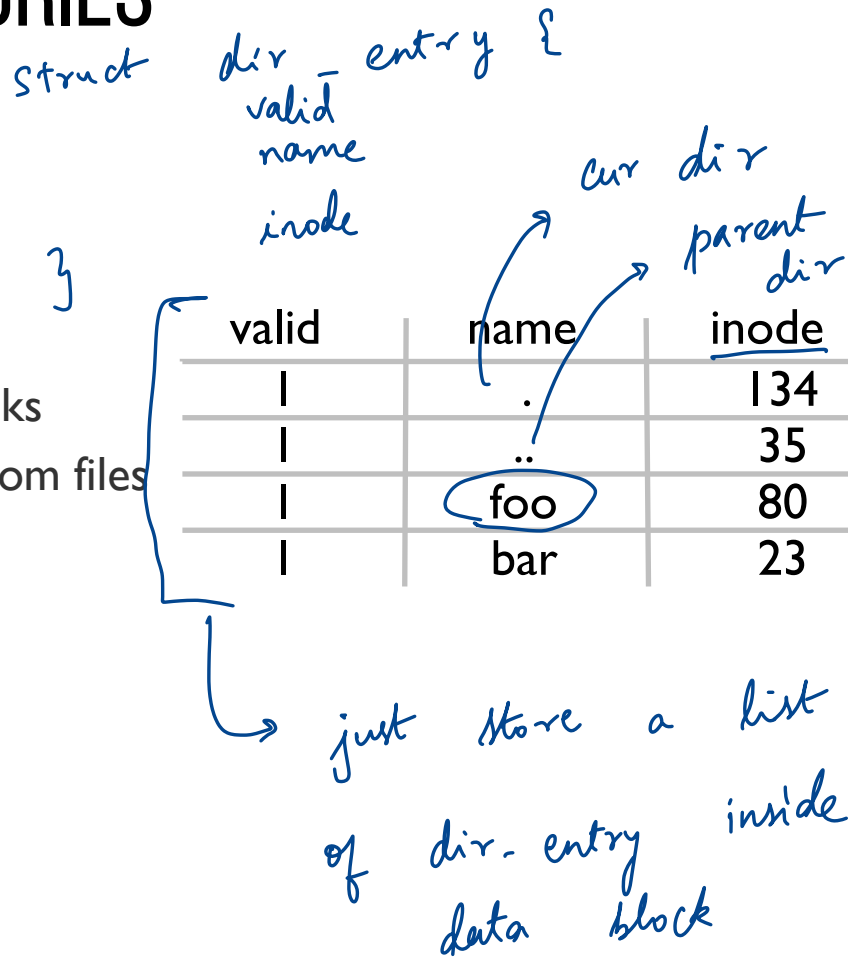
Store directory entries in data blocks

Large directories just use multiple data blocks

Use bit in inode to distinguish directories from files

Various formats could be used

- lists
- b-trees



# FS STRUCTS: BITMAPS

How do we find free data blocks or free inodes?

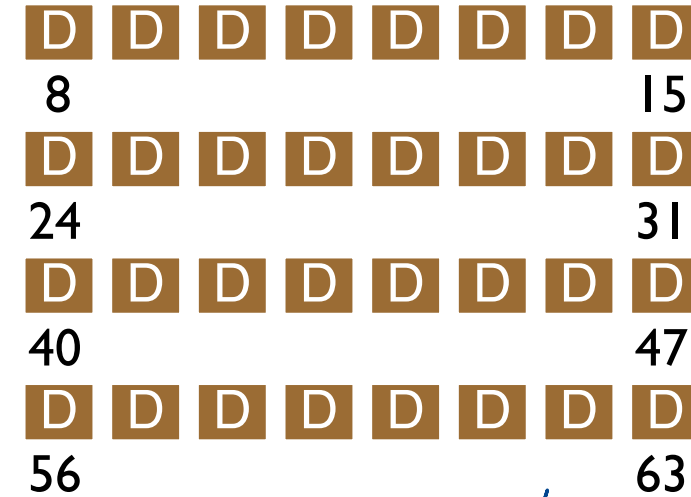
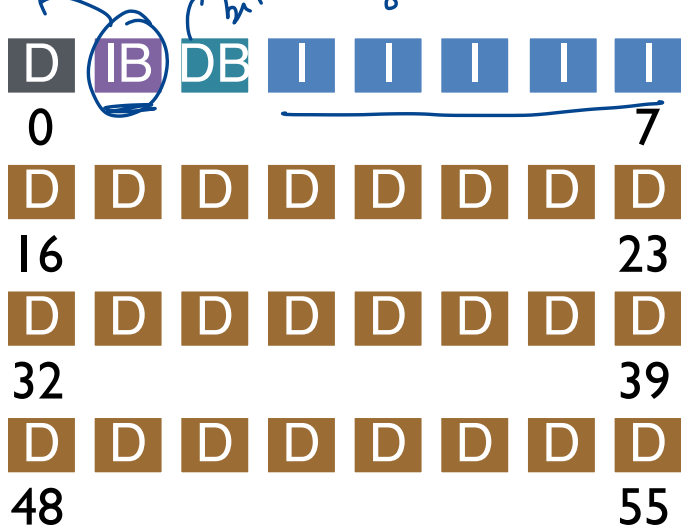
metadata

allocate a DB

80 bits

64 bits

80 inodes



Bitmap : 1 bit to indicate used or not  
=> inode 3 used? bit 3 in IB

# SUPERBLOCK

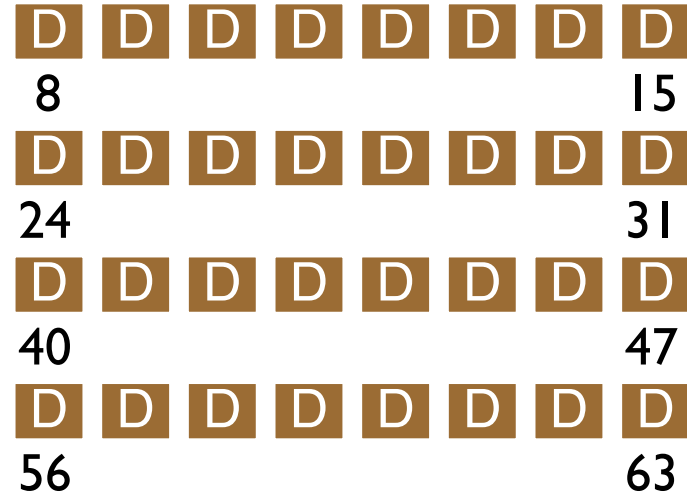
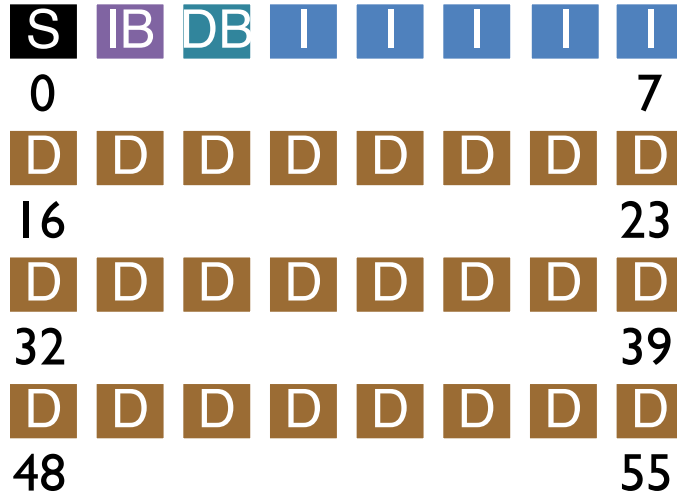
Need to know basic FS configuration metadata, like:

- block size
- # of inodes

Store this in superblock

↳ 1st block on disk  
read the SB when you mount a FS  
↳ block size = 4KB  
→ 80 inodes, 256 bytes inode

# FS STRUCTS: SUPERBLOCK



# QUIZ 16

The following command is executed:

```
echo "Hello World" > foo.txt
```

What produces

```
openat(AT_FDCWD, "foo.txt", O_RDONLY) = 3
read(3, "Hello World\n", 131072)    = 12
write(1, "Hello World\n", 12)      = 12
read(3, "", 131072)                = 0
close(3)                            = 0
close(1)                            = 0
```



*Cat foo.txt*

# QUIZ 16

```
openat(AT_FDCWD, "foo.txt", O_RDONLY|O_NOCTTY) = 3
read(3, "Hello World\n", 98304)           = 12
read(3, "", 98304)                         = 0
close(3)                                   = 0
close(1)                                   = 0
```

→ *grep "Goodbye" foo.txt*

```
openat(AT_FDCWD, "/proc/filesystems", O_RDONLY|O_CLOEXEC) = 3
close(3)                                                 = 0
statx(AT_FDCWD, "foo.txt", ...) = 0
write(1, "-rw-rw-r-- 1 oliphant oliphant 1"..., 55) = 55
close(1)
```

→ *output ls produces*



```
int fd1 = open("file.txt"); // returns 12
int fd2 = open("file.txt"); // returns 13
read(fd1, buf, 16);
int fd3 = dup(fd2);          // returns 14
read(fd2, buf, 16);
lseek(fd1, 100, SEEK_SET);
```

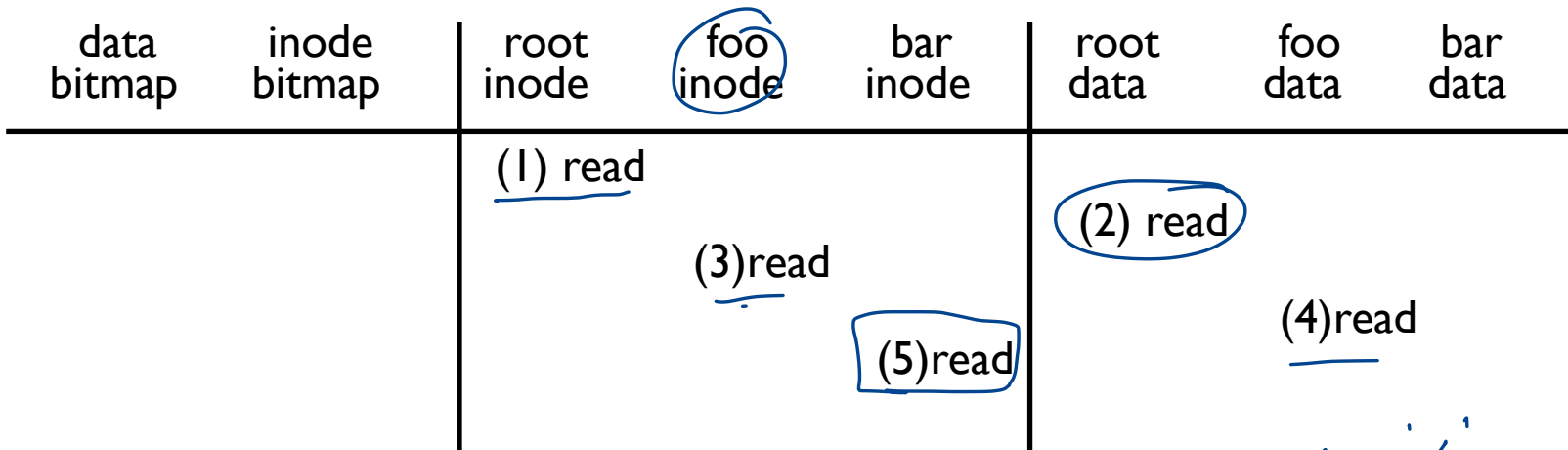
offset for fd! → 100

# FS OPERATIONS

- open
- read
- close
- create file
- write

TIME

open /foo/bar → fd to the calling process

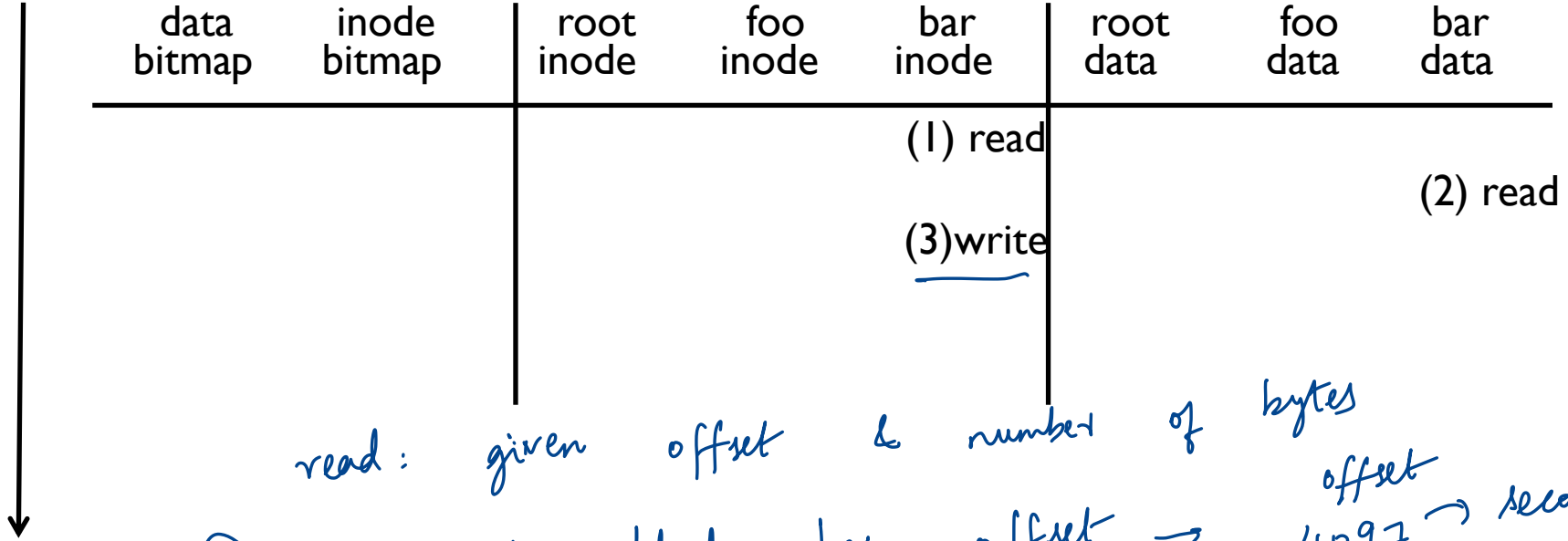


- ① file exists →
- ② permission to read this files
  - ↳ read inode for /foo/bar

- first read '/'
  - ↳ inode → data block
  - ↳ read data block
    - ↳ inode /foo

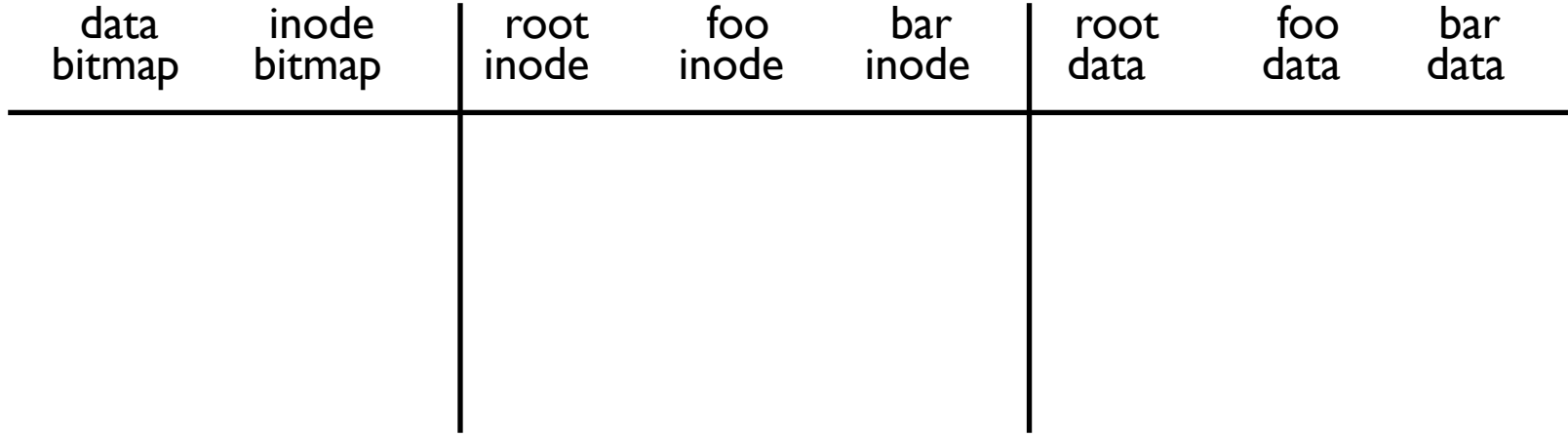
read /foo/bar – assume opened → fd

TIME



- read: given offset & number of bytes
- ① which data block has offset → offset 4097 → second data block
  - ② read block
  - ③ update last accessed ts.

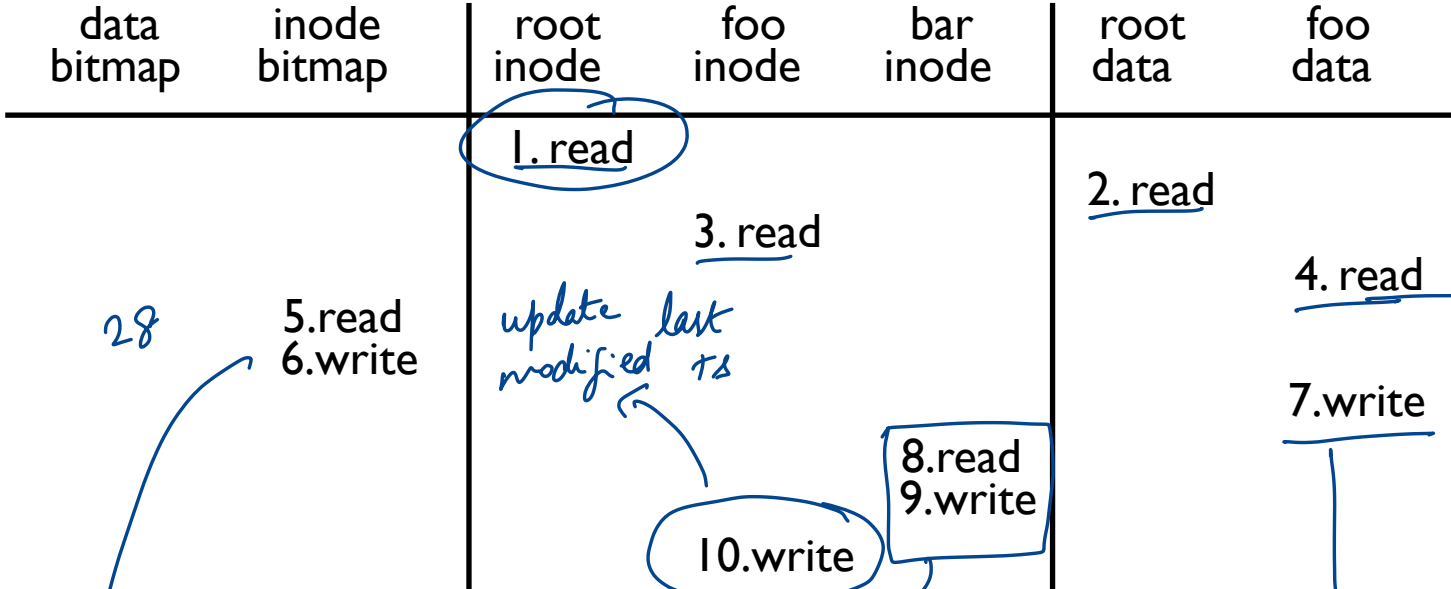
# close /foo/bar



nothing to do on disk!

TIME

# create /foo/bar



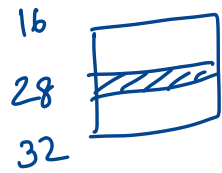
ensure /foo /bar doesn't already exist

28  
 find unused inode & write out bit map updated

update last modified ts

10. write

initializing bar inode

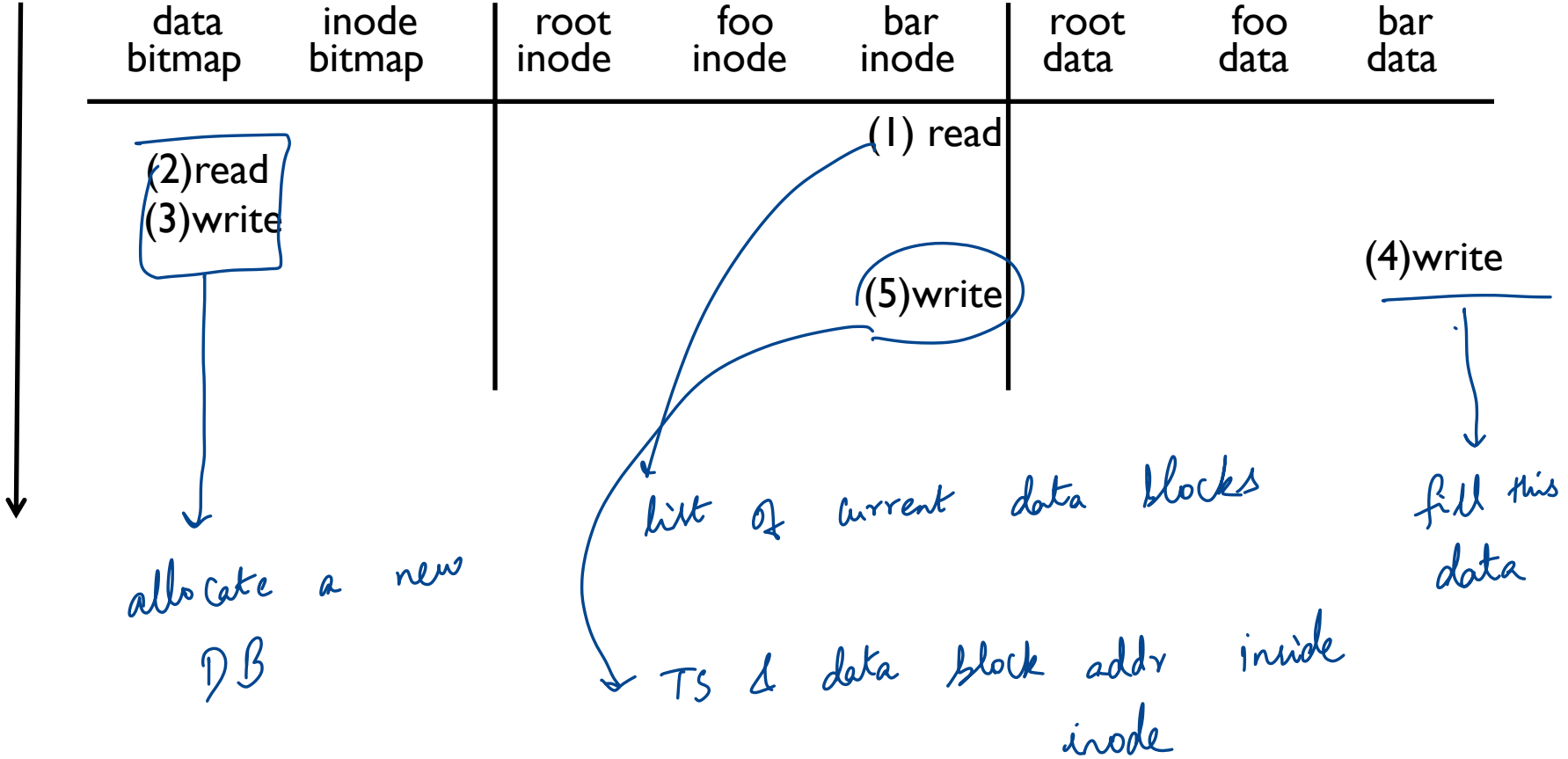


write it back

valid	name	inode
	.	
	..	
2	bar	28

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

TIME



# EFFICIENCY

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

Cache for:

- reads

- write buffering

→ delay when writes are made to a FS

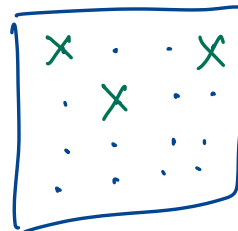
Crashes with writes buffered!?

Overwrites, deletes, scheduling

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

Tradeoffs: how much to buffer, how long to buffer

create a file ⇒ update inode bitmap



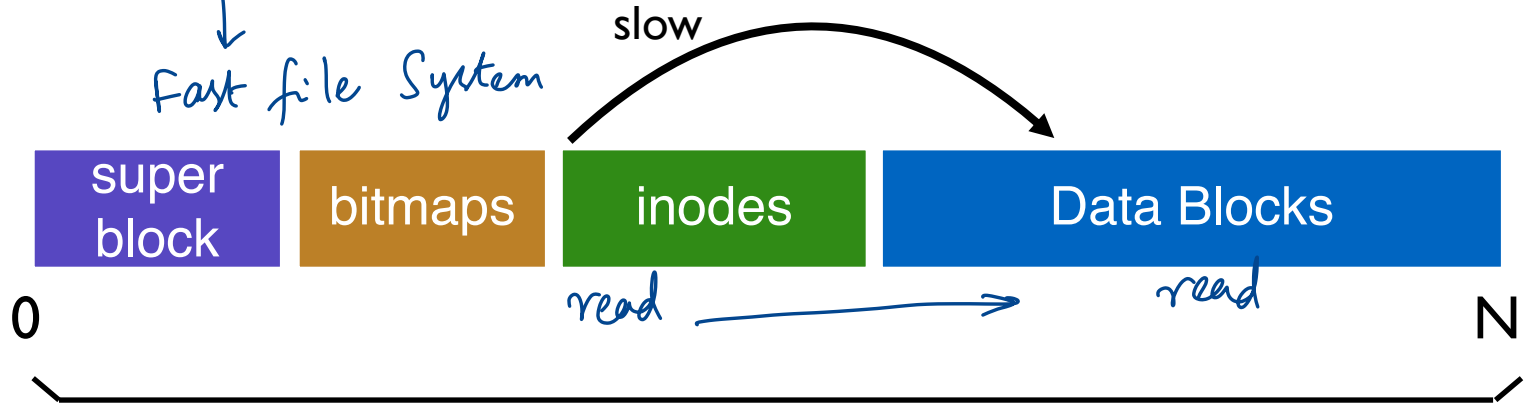
→ 1 disk write for 3 bits



# FFS: FILE LAYOUT IMPORTANCE

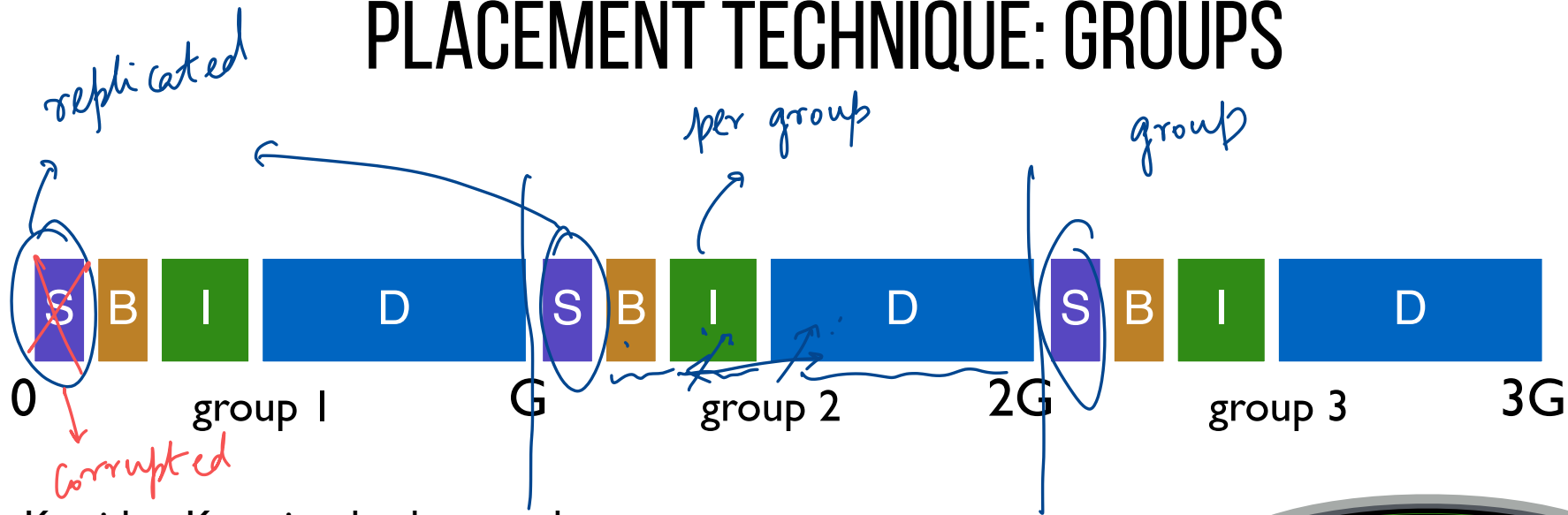
mid ~ 1980s

Fast file System



Layout is not disk-aware!

# PLACEMENT TECHNIQUE: GROUPS

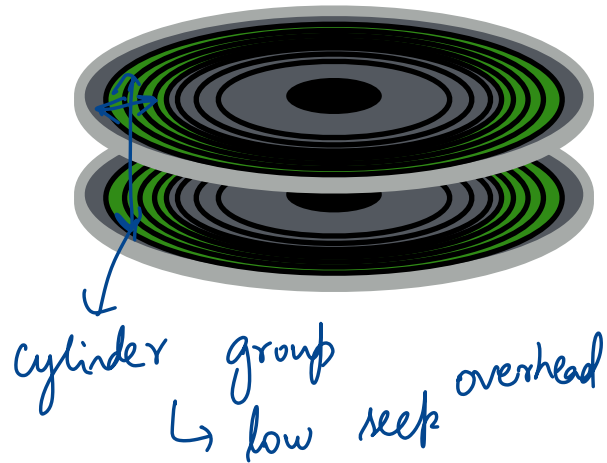


Key idea: Keep inode close to data

Use groups across disks

Strategy: allocate inodes and data blocks in same group.

Replicated superblocks



# REPLICATED SUPER BLOCKS



# 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**

/a/b  
/a/c  
/a/d  
/b/f

group	inodes	data
0	/-----	/-----
1	acde-----	accddee---
2	bf-----	bff-----
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

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

# SPLITTING LARGE FILES

group	inodes	data			
0	/a-----	/aaaaa-----	-----	-----	-----
1	-----	aaaaa-----	-----	-----	-----
2	-----	aaaaa-----	-----	-----	-----
3	-----	aaaaa-----	-----	-----	-----
4	-----	aaaaa-----	-----	-----	-----
5	-----	aaaaa-----	-----	-----	-----
6	-----	-----	-----	-----	-----
...					

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 4MB.

# NEXT STEPS

Next class: Journalling