4/24  welcome to
        CS 537!
should we hold
class outside?
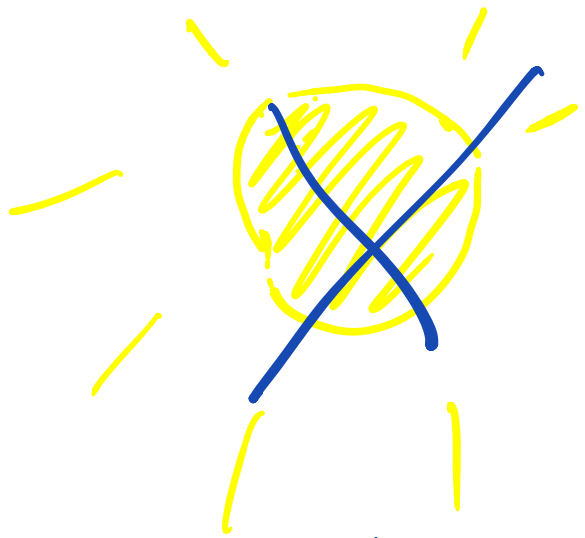
☐    Yes

▨    No

Do you believe
that this is
a  REAL

a **QUESTION?**

No

Today

Wrong Thing:
(Enjoy Life)  Later

Today: File Systems
(Implementation)
→ Locality }②
→ [Crash Consistency] }①

Review:

read/write

Disk
{ | | . . ⌇ . . | } }
0 1 ⌣ N-1

block
(~4KB)

File System: API ↙ open,
uses ↰ ⌐ close,
read
etc
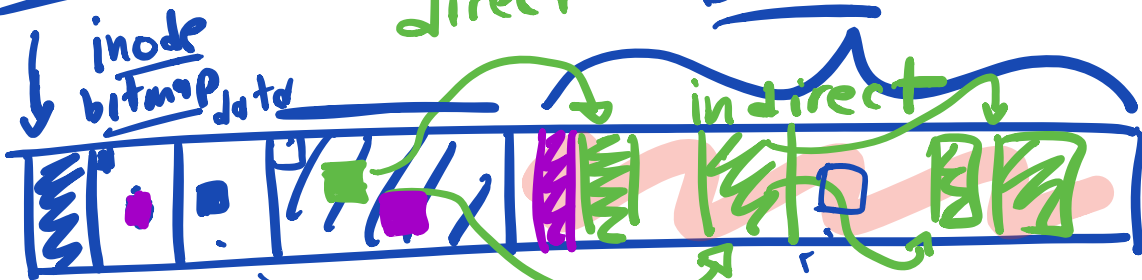
main
memory
cache/
write
buffer

ops on disk
=> low-level
reads/writes

FS:
{ => On-disk structures
  => Access Methods ↑

superblock
↓ inode
  bitmap data

direct          Data

indirect

alloc.     per-file
structures metadata
           => inode
              table

create:
"/foo"
(empty file)

Data:
-> user data
-> dir data
-> indirect
   blocks

Do

Sett on Thurs:

# Crash: what is a "crash"?

$\Rightarrow$ power loss $\leftarrow$ UPS

$\Rightarrow$ BSOD:
    kernel panic, $O \Leftarrow \neq O$
    bug

$\Rightarrow$ user restarts

# why important for FS?

$\Rightarrow$ in the middle of
    an update

# Examples: updates

=> File Creation:
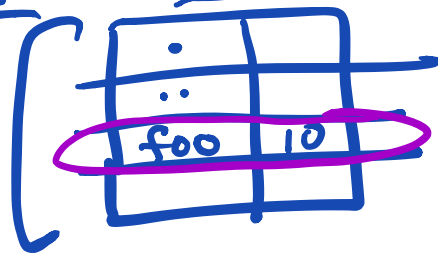
{ what data structures are updated? }

/ <- root    =>    creat("/foo");

e.g. =>  inode bitmap    ( find free, mark it used )

10 <-
   =>  inode : all its various contents

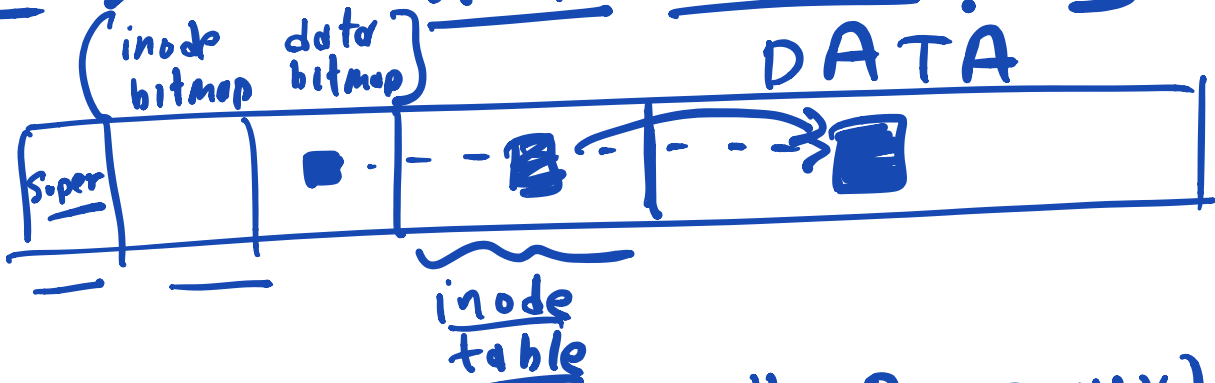=>  root directory: data



name =>
low-level
name
(inode #)

=> root inode:
   update times

[ could also (if dir grows)
   => data bitmap,

new data block ↵

→ existing file, goal ⇒ append
↓ new data block

⇒ how API? (system calls)
⇒ which blocks written?

(inode data
bitmap bitmap)

DATA

| Super | | ▪ - | - - ▣ - → | ▣ |

inode
table

int fd = open ("/foo", O_WRONLY);
↓                              ⇕

lseek (fd, 0, SEEK_END);

[ write (fd, buffer, size); ]
                    └─┬─┘
                     4KB

close (fd);

Data/Metadata:

Data Bitmap    Data    Memory
         inode

Disk

DB

inode
table

(D)
Data , Inode, Data Bitmap:
$\overline{(D)}$ $\overline{(I)}$ $\overline{(DB)}$

orders of writing:

| D | D | I | I | DB | DB |
|---|---|---|---|---|---|
| I | DB | D | DB | I | D |
| DB | I | DB | D | D | I |

crash here
or
here

[ what
happens ? ]

DB          inode  i̇        address:
                              100

Before

| ▣ | size ● IN USE | · | ← nothing here |

inode                    100

After

| 1 | ⊟ ⟋ | ▨ |

size     4KB
pointer  100    disk addr
         ▨

(D)

# Data , Inode, Data Bitmap:

$$\underline{(D)} \quad \underline{(I)} \quad \underline{(DB)}$$

orders of writing:

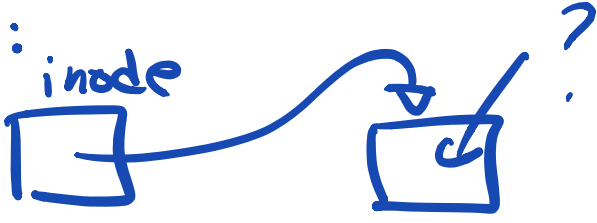| D₁ | D₃ | Ⓘ₅ | I | DB | DB |
|----|----|-----|---|----|----|
| I₂ | DB₄ | D₆ | DB | I | D |
| DB | I | DB | D | D | I |

① ③ looks good! file system metadata "consistent"

② inode ⟷ data bitmap don't agree! ⟹ inconsistent

④ inconsistent : space leak

⑤ inconsistent : inode
   could
   read
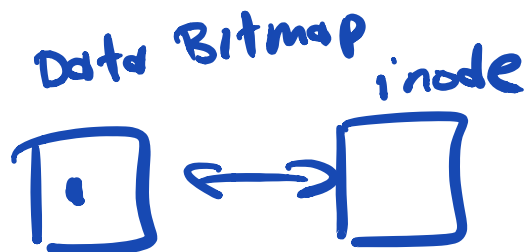     garbage
       (or zeroes)
       (or privacy problem)

Solutions?

Check: Is the FS consistent?

⇒ before usage ( before mount )

⇒ file system check
   (fsck)

⇒ [ focus of (P5) ( good news ) ]
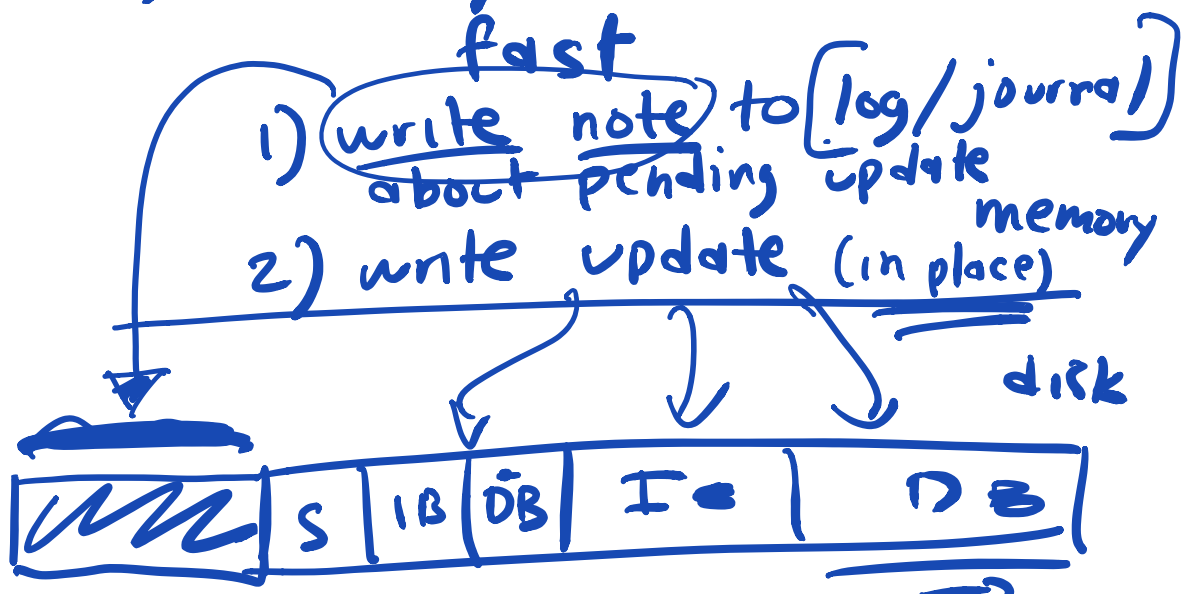
Data Bitmap  inode
[ ▪ ] ⟷ [ ]

But: Disks got large
(RAID)
too

=) too slow?

Something else: Journaling
( OR write-ahead logging)

=) Eager:
→ some work w/ every
update

→ recovery (after crash):
fast
1) write note to [log/journal]
about pending update
memory
2) write update (in place)

disk
| ⟍⟍⟍ | S | IB | DB | I● | DB |

Journal
(or
Log)

Details :  example
→ append to existing
file
info about update  ⇒ location

① $T_{Begin}$  | DB | I. | D |  $T_E$

memory
disk

Journal | DB | I | D |

All or none : Atomicity

② writes "in place"
⇒ DB, I, D

Crash:

[during step **2** (after step 1))
recover from journal

during step 1:
or (issue all writes)
to journal

Aside: disk guarantees
that any 512 byte write
is atomic

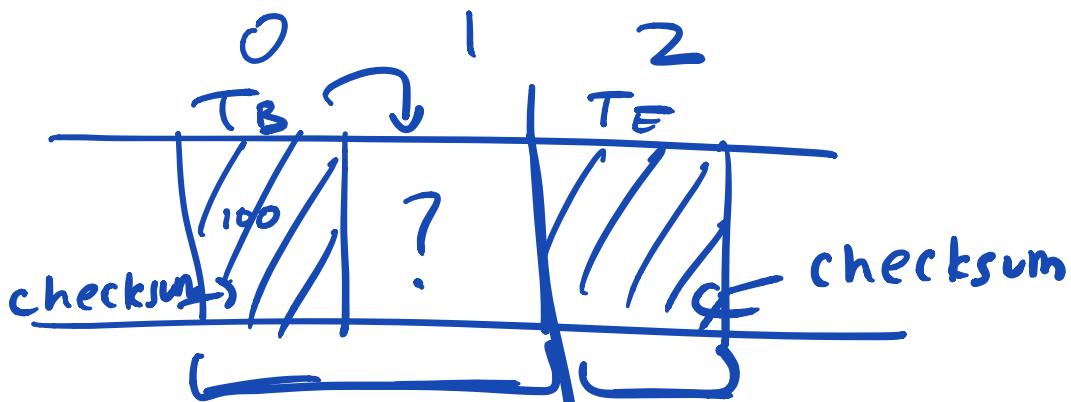but larger write may
partially complete

Step1 => 2 parts
1a) write Trans. Begin +
contents (not $T_{END}$)
to journal
=> wait
1b) write $T_{END}$

$\Rightarrow$ transaction
commit

wait

2) in place updates

```
     O        I        2
   T_B    ↱↓        T_E
   //    ?      ///
  /100/         ////
checksum            checksum
```

transaction
2005:  checksum    $\Rightarrow$ {Linux ext4}

Crash Consistency:
$\Rightarrow$ [PS] $\in$ initial deadline

xv6 file system image
↑
C/Linux    look for
         inconsistencies

|   | M | T | W | Th | F |
|---|---|---|---|---|---|
|   |   |   | X |   |   |
|   |   |   | ⊗ |   | F P5 |

~1/9    ~25%

Final

8th ?

A

100%

AB

B

0%

D  Og

Ba  (D)

## (Locality)

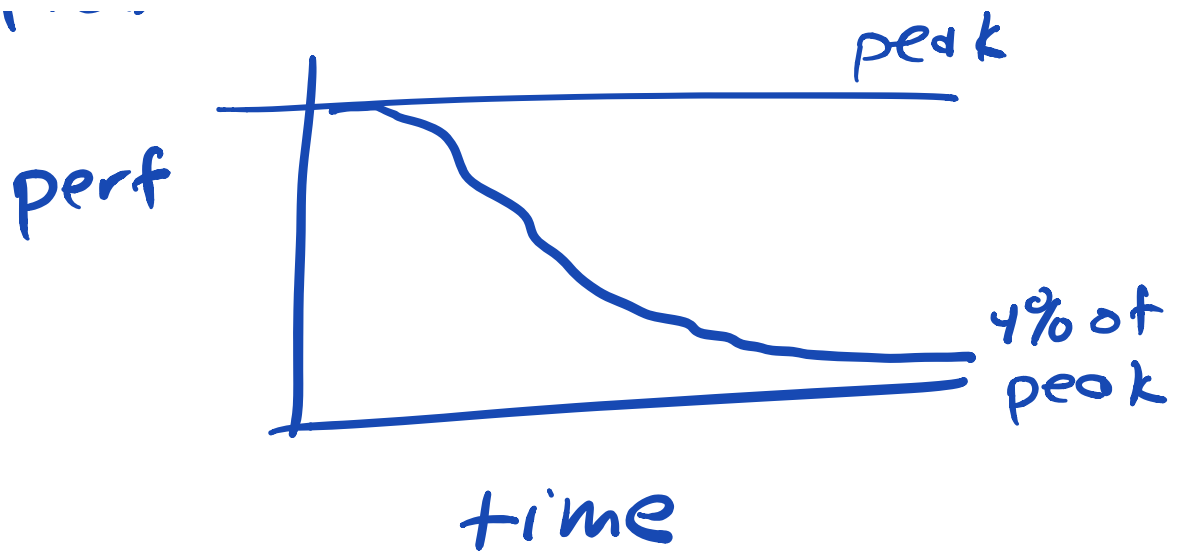### Hard Disks: long seeks costly



### SSDs: locality important too

goal:
put files/etc. on
disk ⇒ accessed together
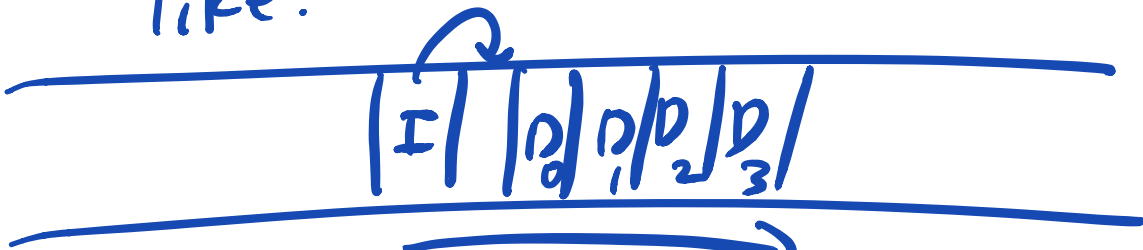should be near
one another

## Fast File System (FFS)
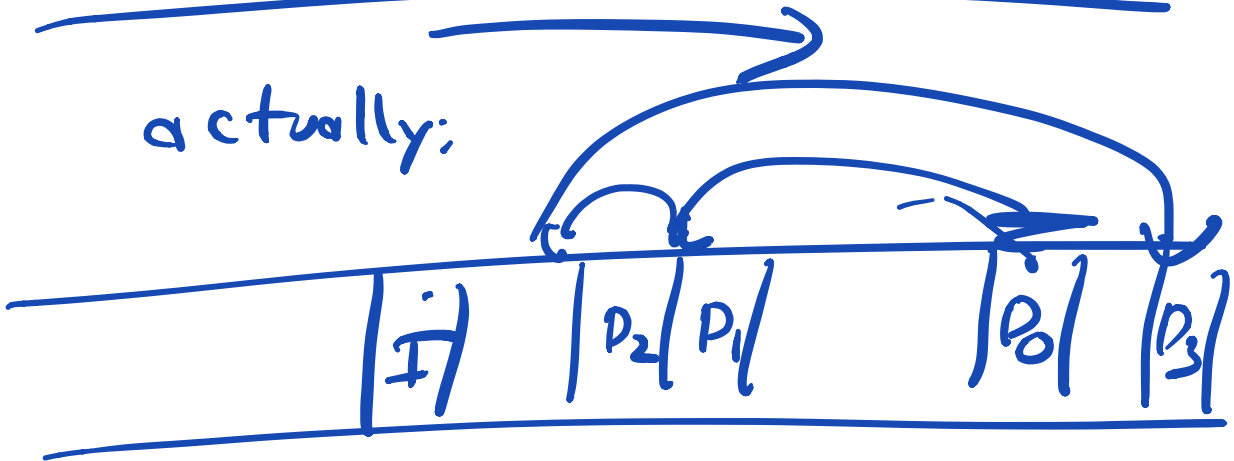
⇒ Performance

Problem: old unix was slow

peak

perf

4% of
peak

time

=> data blocks, inodes
scattered through
disk

like:

| I |  | D₀ | D₁ | D₂ | D₃ |

actually:

| I |  | D₂ | P₁ |  | P₀ |  | D₃ |

=>

# FFS: Treat Disk like a Disk!

## On-disk structure: change

old:

| S | FB | DB | Inodes | Data |
|---|----|----|--------|------|

somewhat large   large

## new FFS format:

| $G_0$ | $G_1$ | $G_2$ | $\lessgtr\lessgtr\ldots$ | $G_{n-1}$ |
|-------|-------|-------|------|-----------|

| S | FB | DB | Inodes | Data |
|---|----|----|--------|------|

somewhat    "large"
large

Goal: Put related stuff
     into same  group
     (unrelated -> different)

what  is "related"?

{
  -> data same file
  -> files in same dir.
  -> inodes/data
}

make new dir:
  pick some group X
     (how?) -> (free space,
                         inodes)
  all files in dir
     (inodes/data) => group
                          X

$a'$      $b$

$f_1$   $f_2$   $f_3$

## Large Files : exception

only put first N blocks
of file => desired
group
put $\overset{each}{\vee}$ next chunk of
file in other
groups

$G_0$   $G_1$   $G_2$   .  .  .  .  .  .  .  $G_{N-1}$

3       1       1                           2

Problem: Performance?

[ how big is each chunk? ]

⇒ need to know:
 ⇒ how big file is ✓
 ⇒ group size ✓
 ⇒ memory ✓

⇒ cost of seek
⇒ (transfer speed)

read:
( read chunk )
seek ⌉

most ⤵
transfer :

⤷ 10% ← tolerable

⎡ transfer: ( 100 MB/s )
⎣ seek (avg) : 10 ms

transfer                    seek

                          | 10ms |

$$100\text{ms} \longrightarrow \boxed{.1 \text{ sec}}$$

$$\boxed{10^{X} \text{ MB} \cdot \frac{\text{sec}}{100 \text{ MB}} =}$$

$$\Rightarrow \boxed{\underline{Final}} \Leftarrow \text{study this}$$

$\underline{FFS}$ : $\begin{bmatrix} \text{first time:} \\ \text{same API,} \\ \text{new } \underline{implementation} \end{bmatrix}$

$\Rightarrow$ symbolic link

$\Rightarrow$ long file names

$\Rightarrow$ other perf. opts
(not as relevant)

Treat
Disk like Disk
$\Rightarrow$ (peak, sus$\underline{tained}$

perf.