

Today (5/1)

---

Almost Done!

(Happy, or Sad?)

Correct Answer: Both

Now: It's warm ...

Today: Virt

Conc

→ LFS

(log-structured  
file system)

Persistence

→ Not Hard drives:

Flash-based

SSDs

Admin → (logging)  
→ End: Words of wisdom?

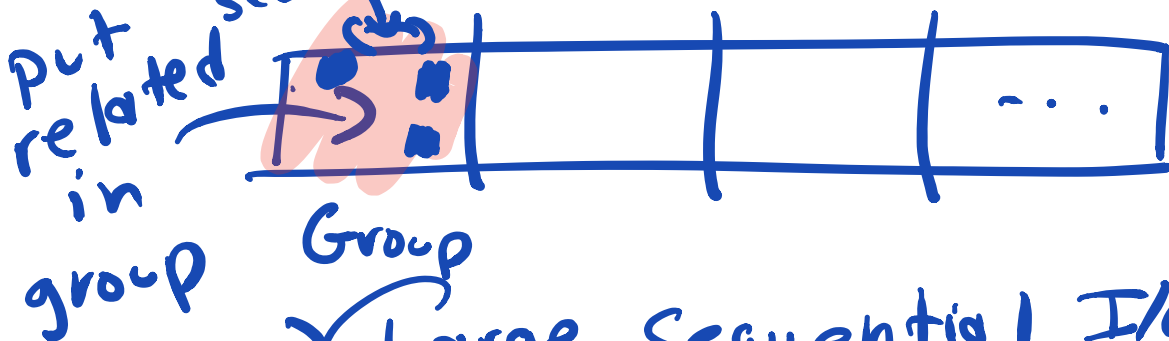
## Log-Structured File Systems (LFS)

Motivation:

Hard Drive Performance

⇒ Random I/Os: Poorly  
(small) (seeks, rotations)

short seek, but what about FFS?

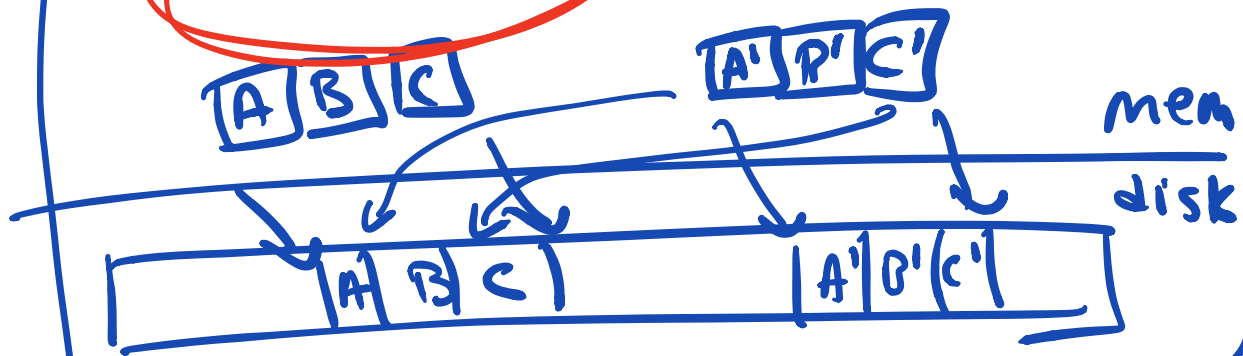


⇒ Large, Sequential I/Os:  
perform well

Problem:

Reads: hard to ensure large, contiguous

Writes: we have choice:



⇒ insight: can write data anywhere

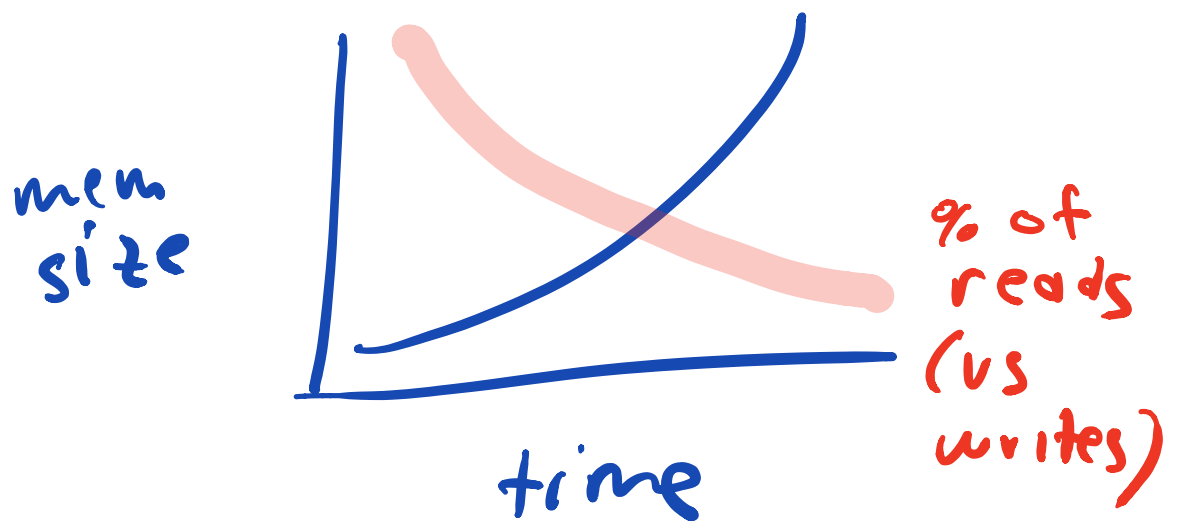
Read Problem: Memory has gotten larger

⇒ 1990 Main memory [1 MB]

⇒ now many GB:

Cache: main memory is cache for

file data



LFS: Structure of LFS

=> Data Structures  
=> mostly, keep "same"  
(inodes, data blocks, directories)

+ few new structures

=> Access Methods

=> writing, reading  
what's different?

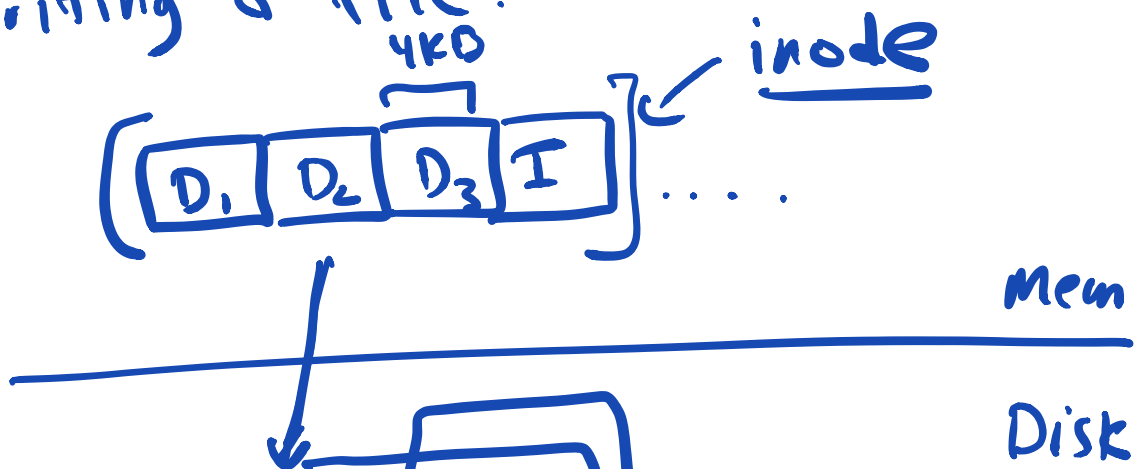
Basic idea: focus on writes

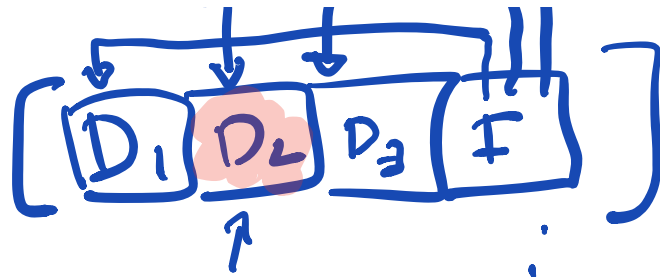
1) buffer many  
writes in memory  
(=> in-memory "segment")  
~ a few MB in size

2) write to disk  
key: what to write  
to enable later  
reads?

Example

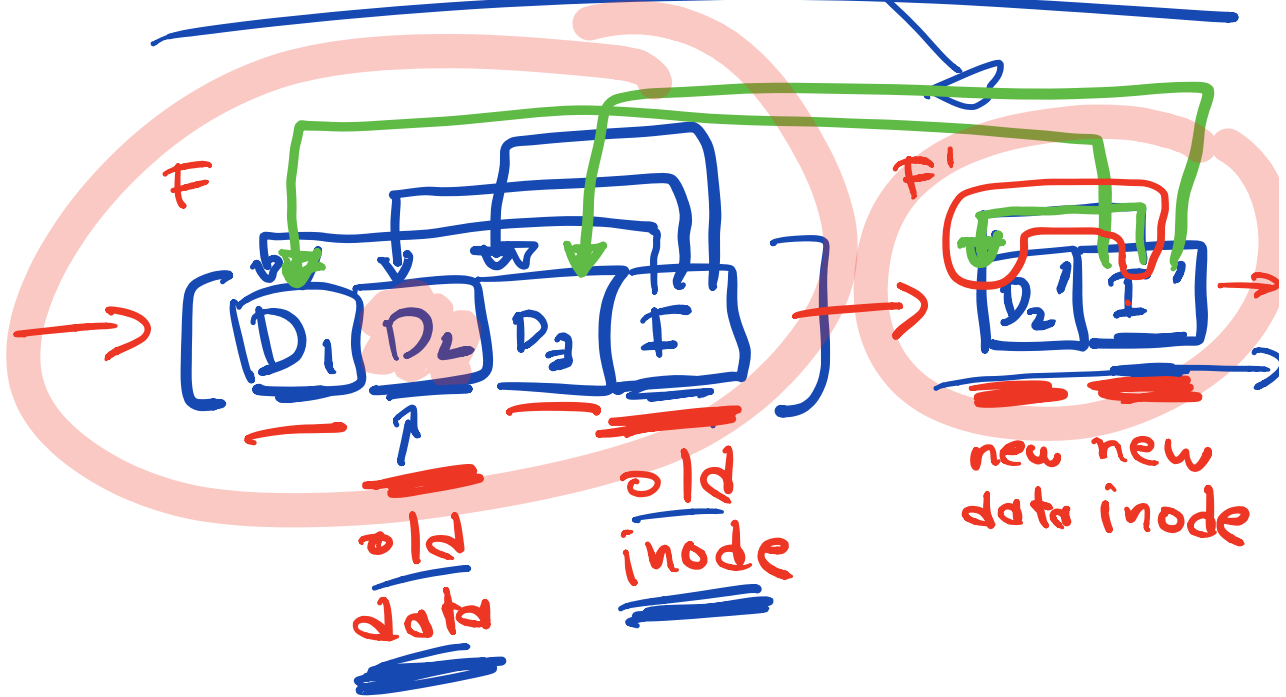
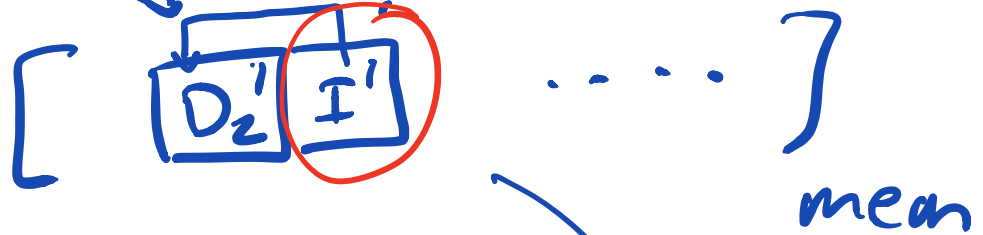
writing a file:





Later: user overwrites D2  
 [write (D2')] ...

in-memory  
segment



Addresses

0

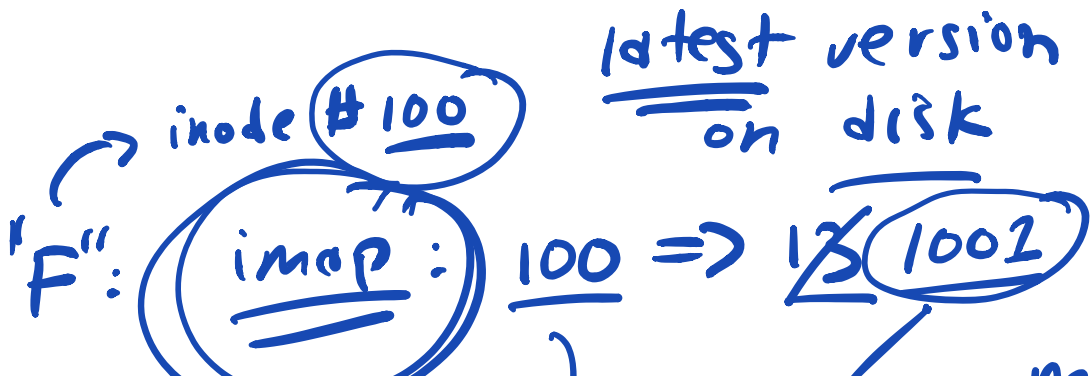
max

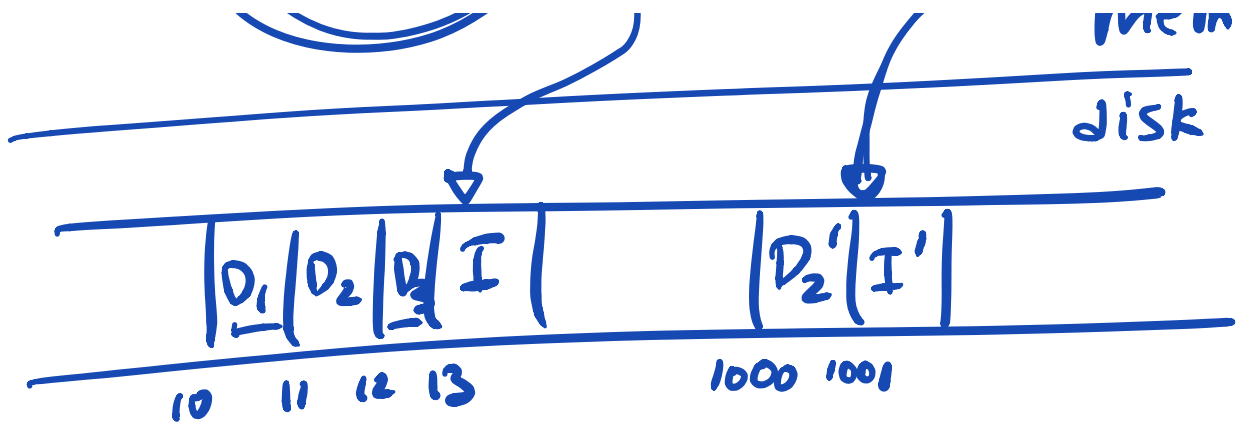
Problems:

- ⇒ how to find latest version of inode?  
(I' not I)
- ⇒ garbage: old versions of inodes, data lying around: what to do?

how to find inode (latest?)

⇒ in-memory "inode map" array: for each inode #, ⇒ location of





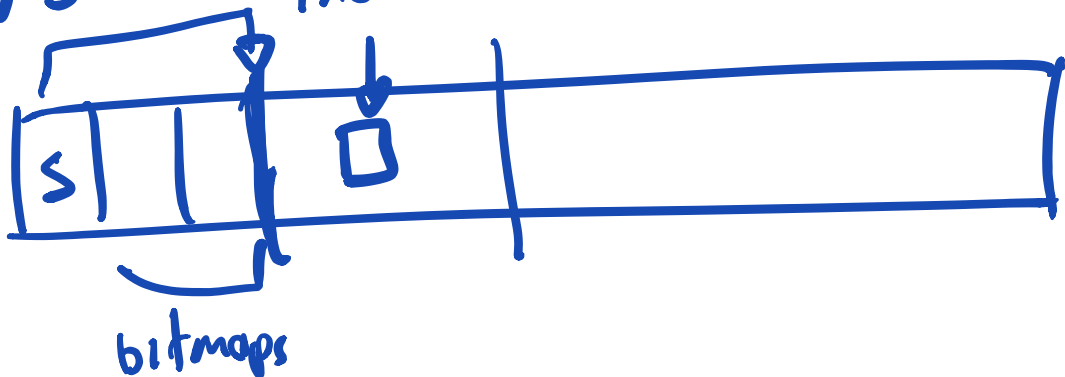
Reads:

directory contents:  
same as before

|            |            |
|------------|------------|
| <u>"F"</u> | <u>100</u> |
|            |            |

inode #

VFS: inode table





But: crash! (lose inode map)

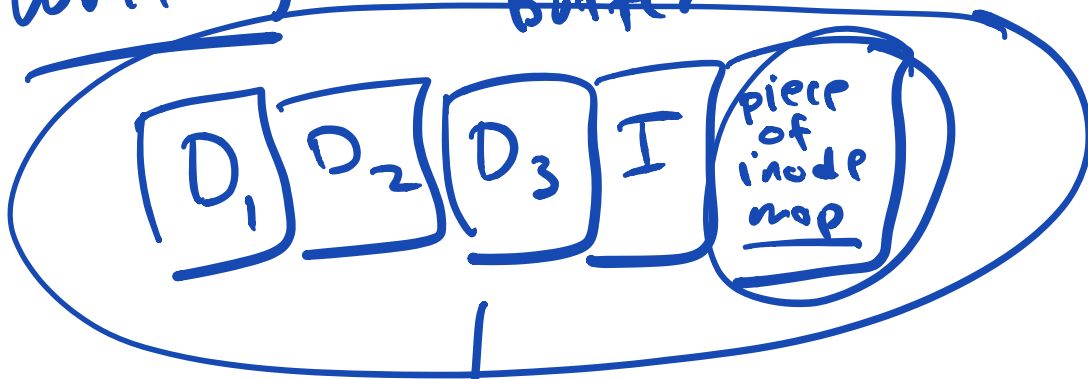
options:

1) write out pieces of inode map

2) reboot: scan entire disk, find "latest" of each inode  
=> rebuild inode map

writing:

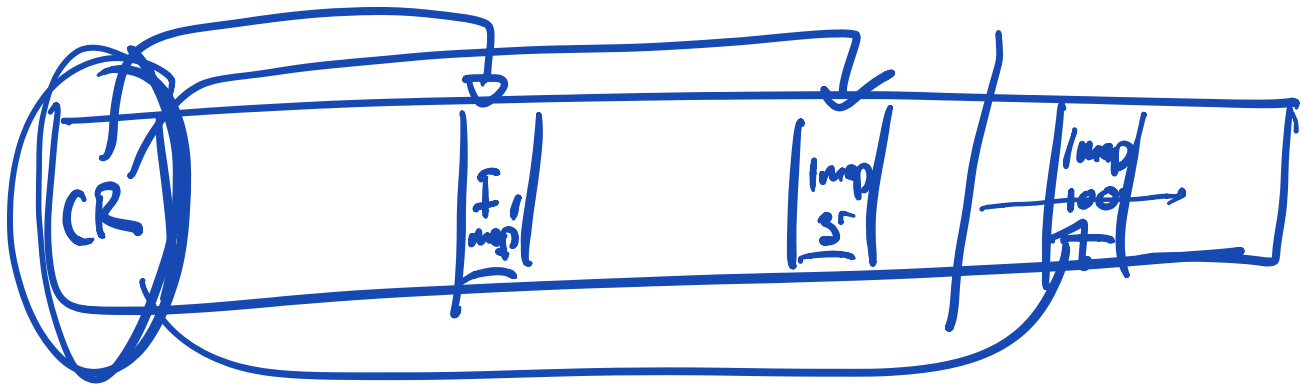
buffer in mem



write to disk



periodically: write out  
 one more data structure:  
 Checkpoint Region

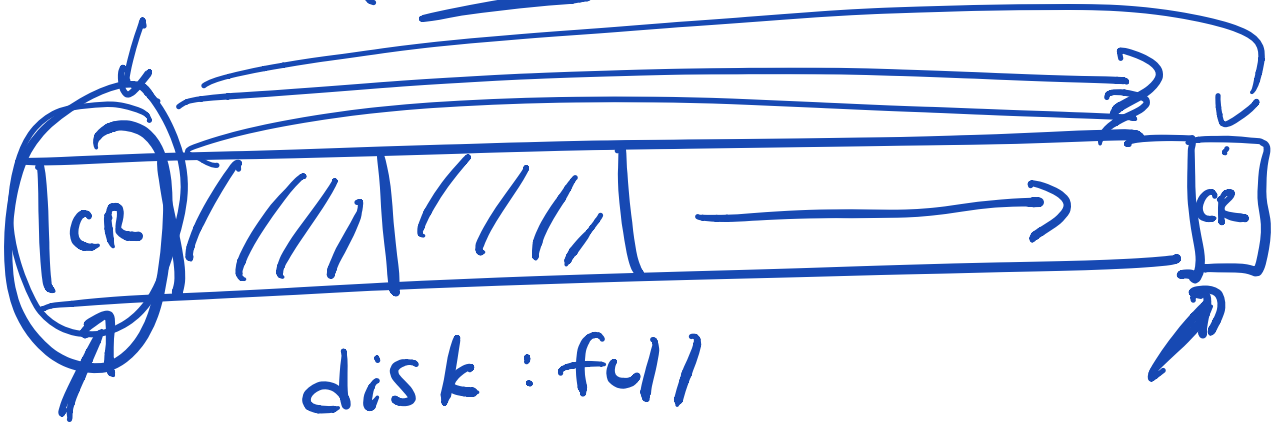


Old: "overwrite in place"  
 => fixed location for inodes  
 => once written,  
 data stays in place

New: "copy on write" (cow)  
 LFS: "cow"-based file  
 system

Problem: lots of old copies

=> (Garbage) ("dead")

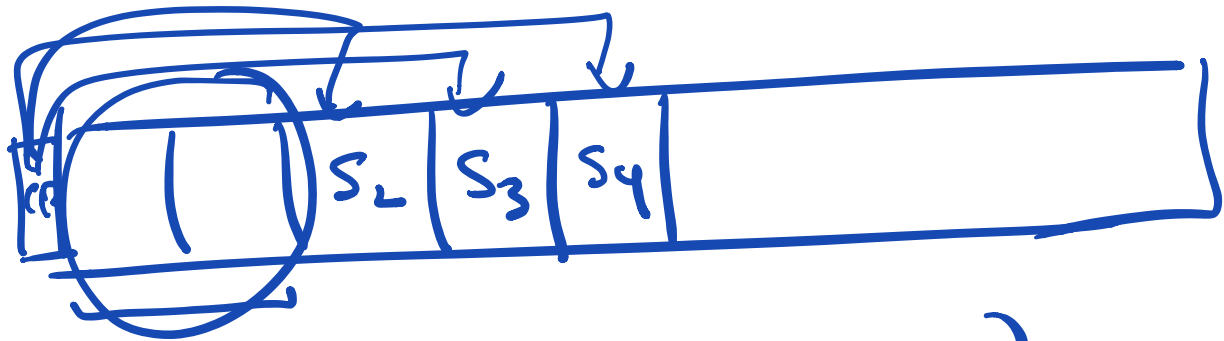


"Background Process":

- => cleaner
- read old segments
  - determines what is live / not
  - write out live data
  - free old segments



cleaner:  $\rightarrow$   
 reads  $S_0, S_1$   
 determine liveness  
 write  $S_4$  (w/ live info)  
 free  $S_0, S_1$

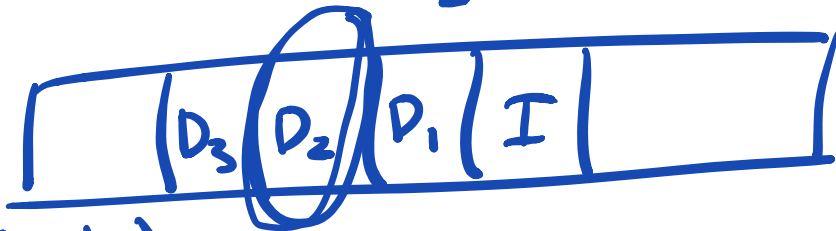


segment (up close):  
 how to tell if something  
 is "live"?

know: inode#,  
 offset

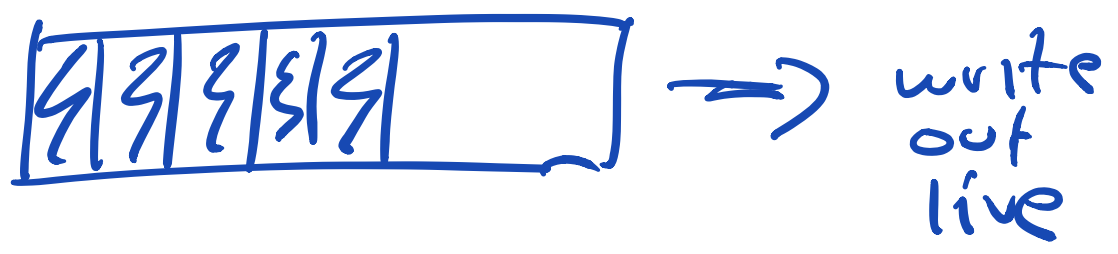
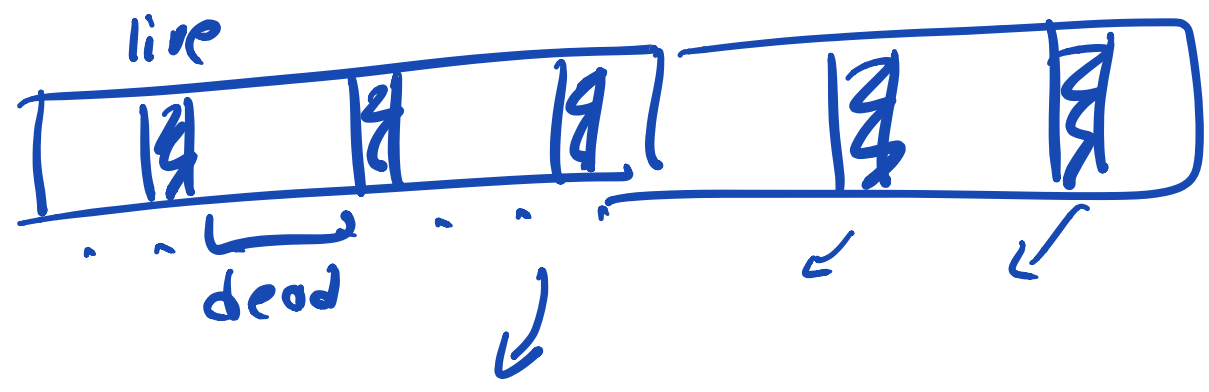
---

$D_2$ : live? dead?

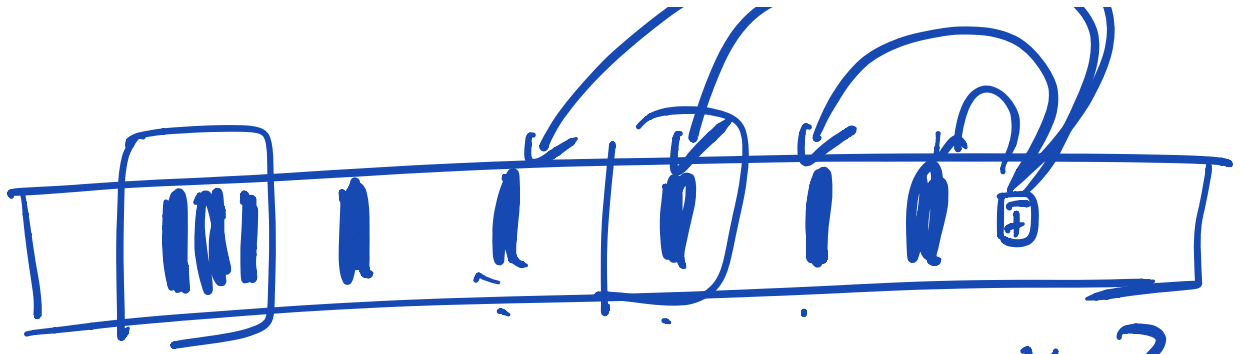


(data)  
 live:  $\leftarrow$  is this pointed to  
 by an inode that's

in inode map?  
extra info w/ segment:  
segment summary block



Performance : LFS  
→ sum of foreground  
+ background  
I/O traffic

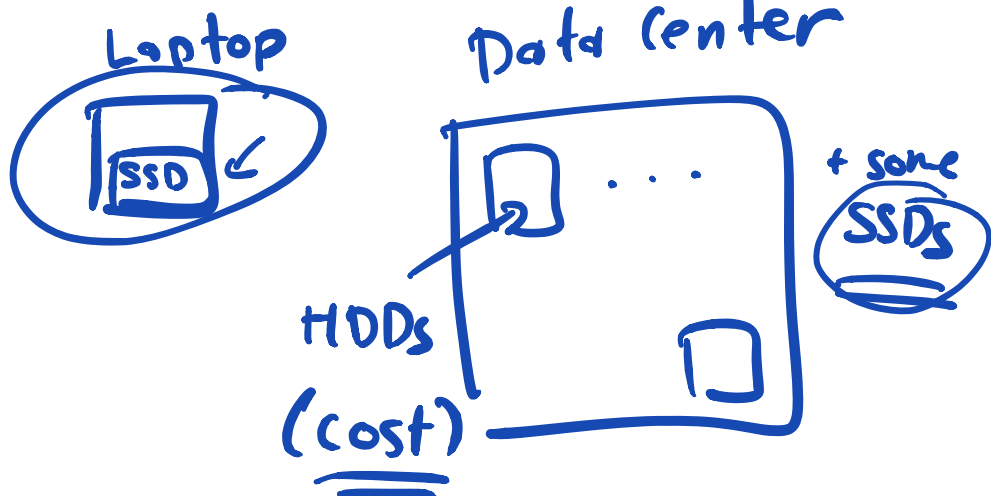


fragmentation: how to handle?

⇒ fragmented file:

read, write it out  
(whole thing)  
(user program)

SSDs:



Flash-based

⇒ Solid-State Devices

=> (computer chips) <sup>non-volatile</sup>  $\rightarrow$   
=> serves as persistent memory

(can survive power loss)

Internals:

=> Flash chips

Interface:

just same as hard drive

=> blocks (sectors)

=> read/write

SSD: has many chips

Chip: organized as follows

pages: ~2KB, 4KB



block => "erase" block  
(~256KB)

store bit:

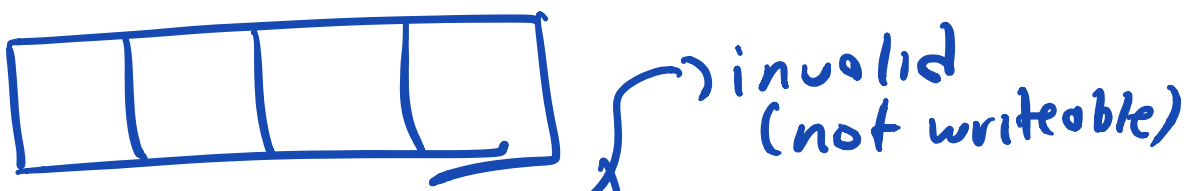
"trapping  
1,0 =<= charge"

operations: different than  
just read/write

read (page) => gives data  
update: 2 steps ~256KB

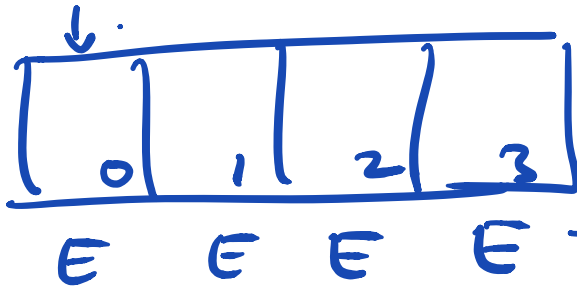
erase (block) =>  
makes block writeable

program (page) => set its  
contents  
(w/ in block)

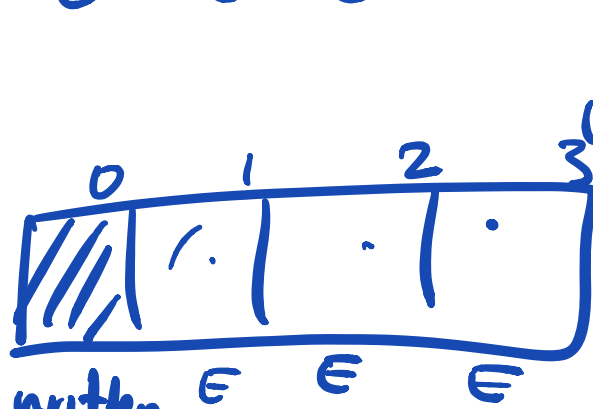




I I I I } erase (block)



erased  
(writable)



Program (0) =  
data

written  
(v)  
↓  
valid

once programmed  
=> to overwrite  
have to  
erase entire  
block

Performance / Reliability

micro

Perf:

↑ ... (no. of μs) ↓

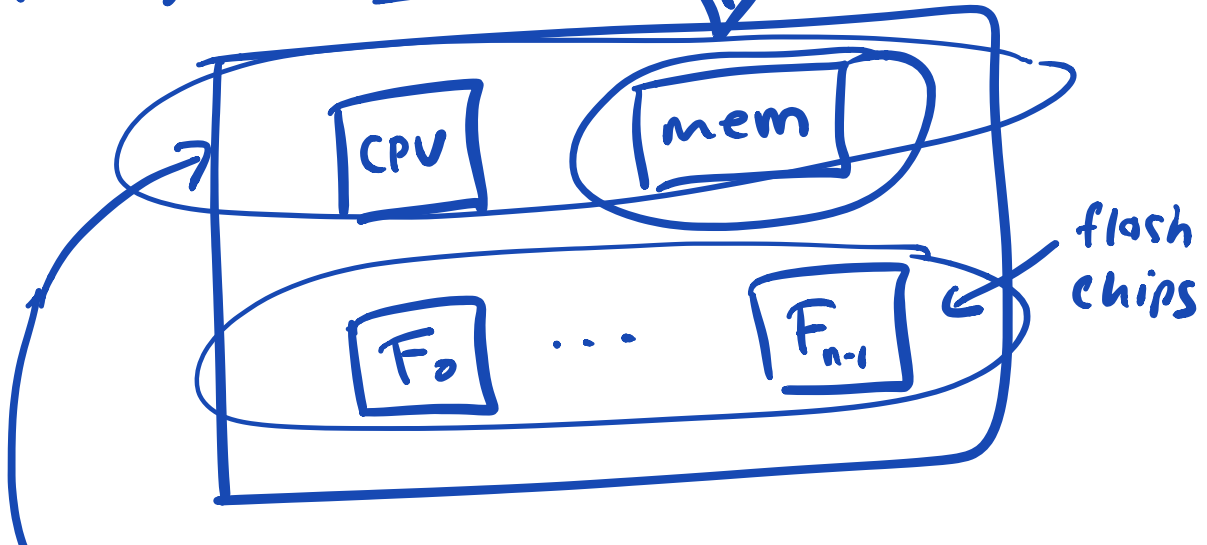
|                   |                      |                        |
|-------------------|----------------------|------------------------|
| reads:            | <u>fast</u>          | (100 ns)               |
| erase:<br>(block) | <u>slow</u>          | (~ a few ms)           |
| program<br>(page) | some<br>what<br>fast | ↑<br>milli<br>(100 μs) |

Reliability:

wear out:

erase / program block  
 "too many" times  
 => unusable

Design: SSD ↑ interface

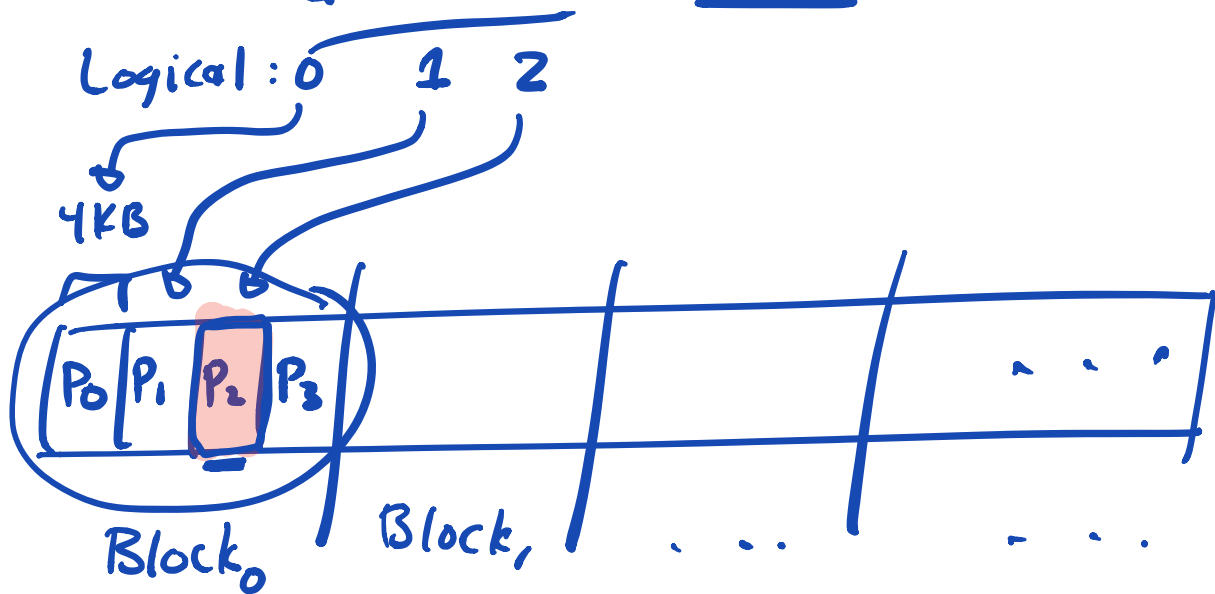


# Flash Translation Layer (FTL)

=> takes reads/writes  
to interface

=> map them into  
low-level flash  
ops (reads, <sup>erases</sup> programs)

Bad way to build FTL:  
"direct mapped"



reads: easy read (address)

writes: e.g. write (addr: 2)

read  $P_0, P_1, P_3$

(put it somewhere: where?)

erase block(0)

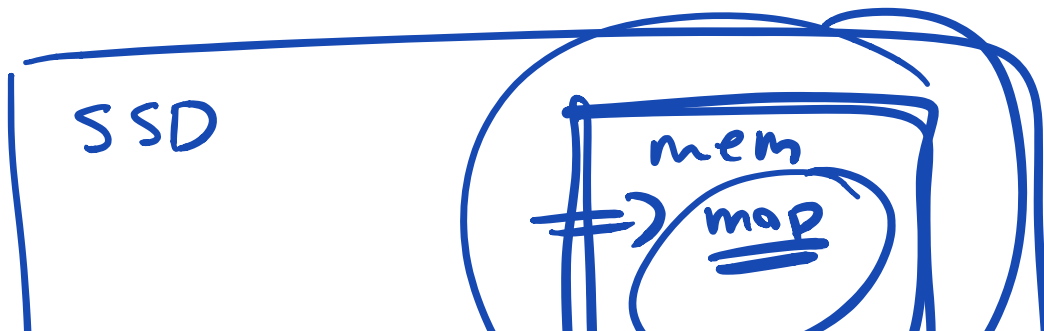
program  $P_0, P_1, P_2', P_3$

$\Rightarrow$  (slow)  $\Rightarrow$  wearout

What to do?

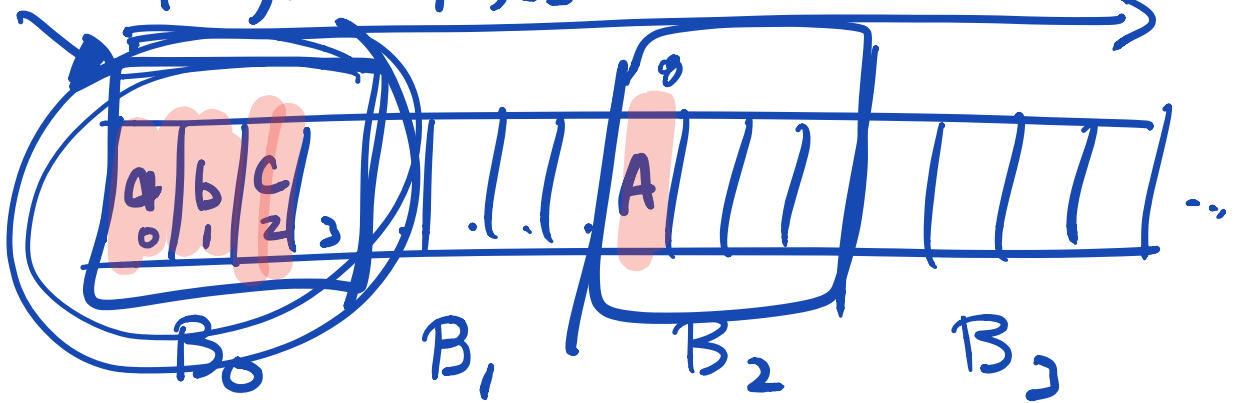
FTL : better w/  
logging

(copy-on-write)





logical write S (4KB in size) : log-style write  
 → Erase block once  
 → program pages one at a time



→ record (in memory)  
 → in FTL

logical block → physical page

mapping

→ write (100, "a")      write (200, "c")  
 write (101, "b")

FTL: 100 → 8    200 → 2  
102 → 1

write (100) "A")

Problems:

⇒ old: garbage

⇒ old: how to recover  
FTL translations  
after crash

⇒ new: wear leveling  
(in FTL)

⇒ new: how big is FTL?

⇒ 1 TB SSD  
per page mappings } ~

page: 1 KB

how many mappings?

⇒ goal: not need

all those mappings  
in device  
memory

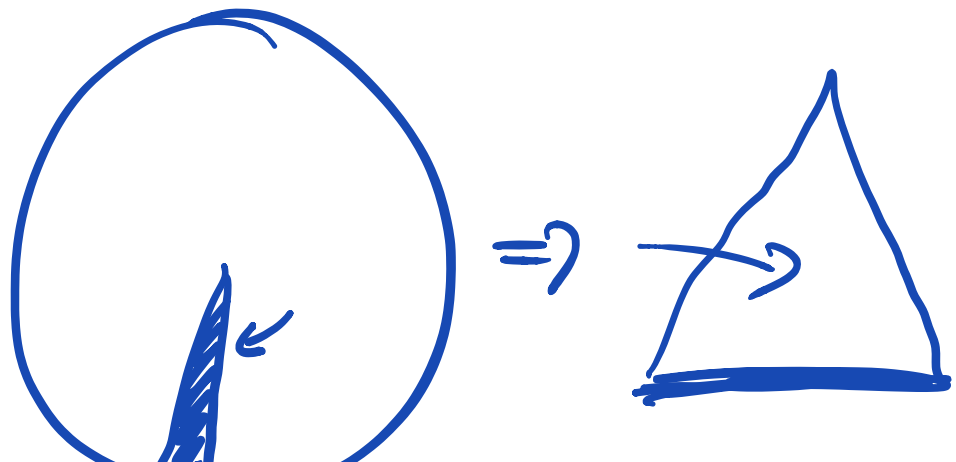
=> better structure

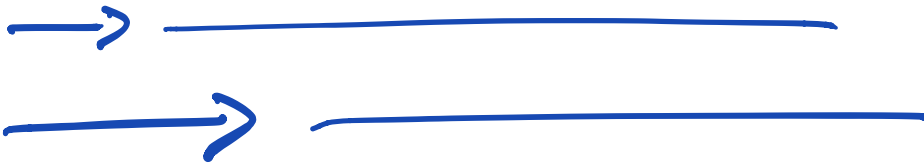
=> only keep "active"  
pieces in memory  
(caching, swapping)

=> block mappings:  
(refer to entire  
blocks, not pages)

---

Topics not covered:  
Most Stuff





## Maintenance:

=> P5

(Friday)  
→ send mail

=> Final

Tuesday (5/8)  
evening ~~5/8~~

Review Session:

Monday 4-6 room TBA

▷ ("Final" => emphasis on "since midterm")

Cheat Sheet: reclaim  
old: ok

=> 2 sep. pages (4 sides)



- => ~~Recommended Courses~~
- => ~~Intel etc. (Meltdown)~~  
~~Spectre~~
- => ~~Favorite Architecture:~~
- => ~~inode, imap => ipads, iPhones~~
- => which OS book?  
did I use
- => (Tanenbaum) => not an A

