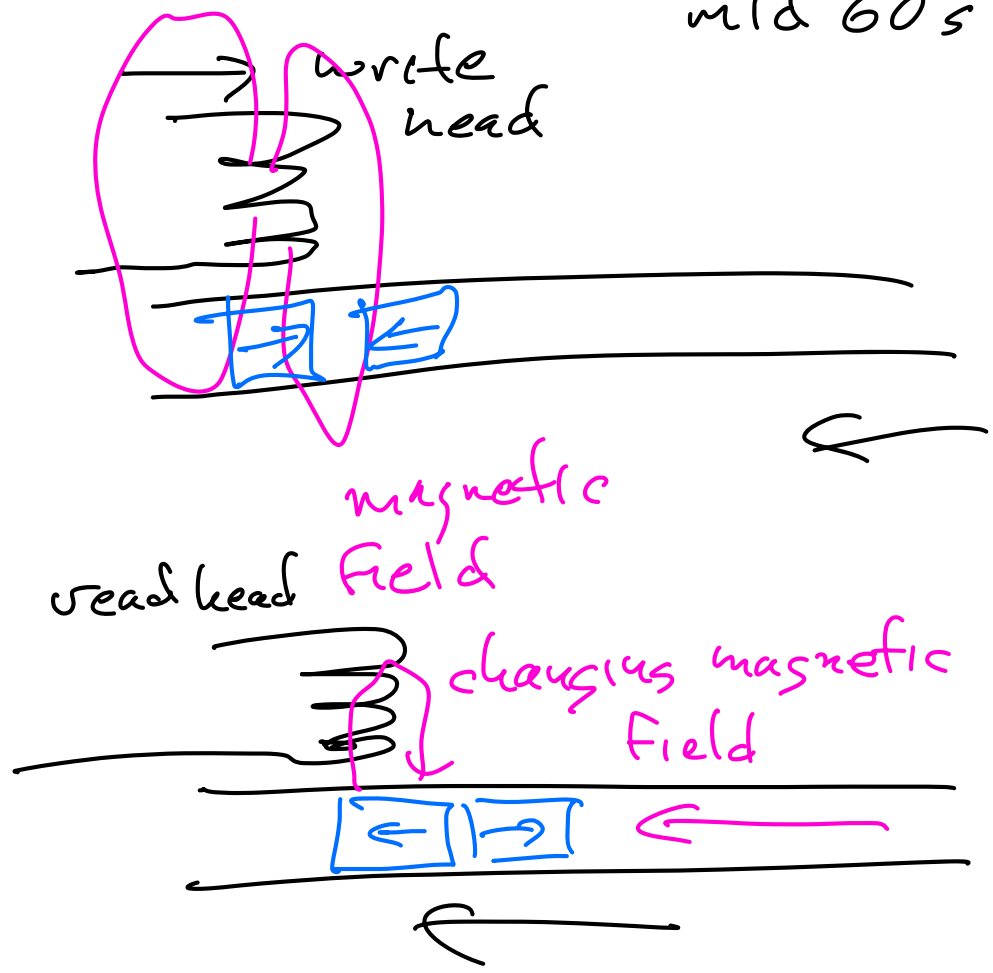


CS 3650 – Computer Systems
Spring 2024
Peter Desnoyers

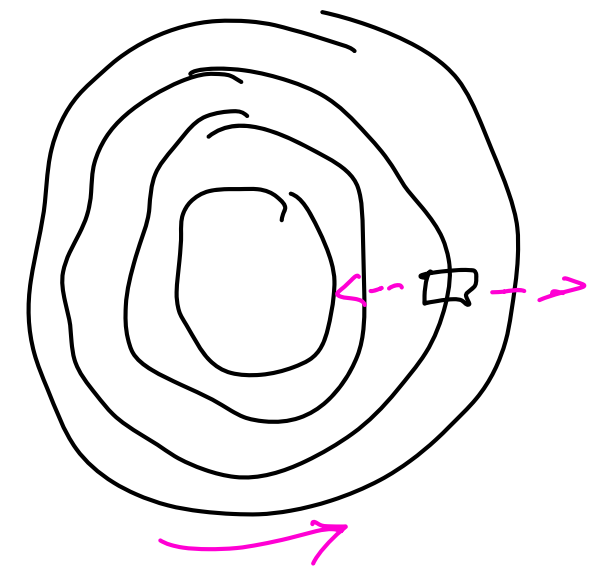
Lecture 19, Tue Mar 19, 2024

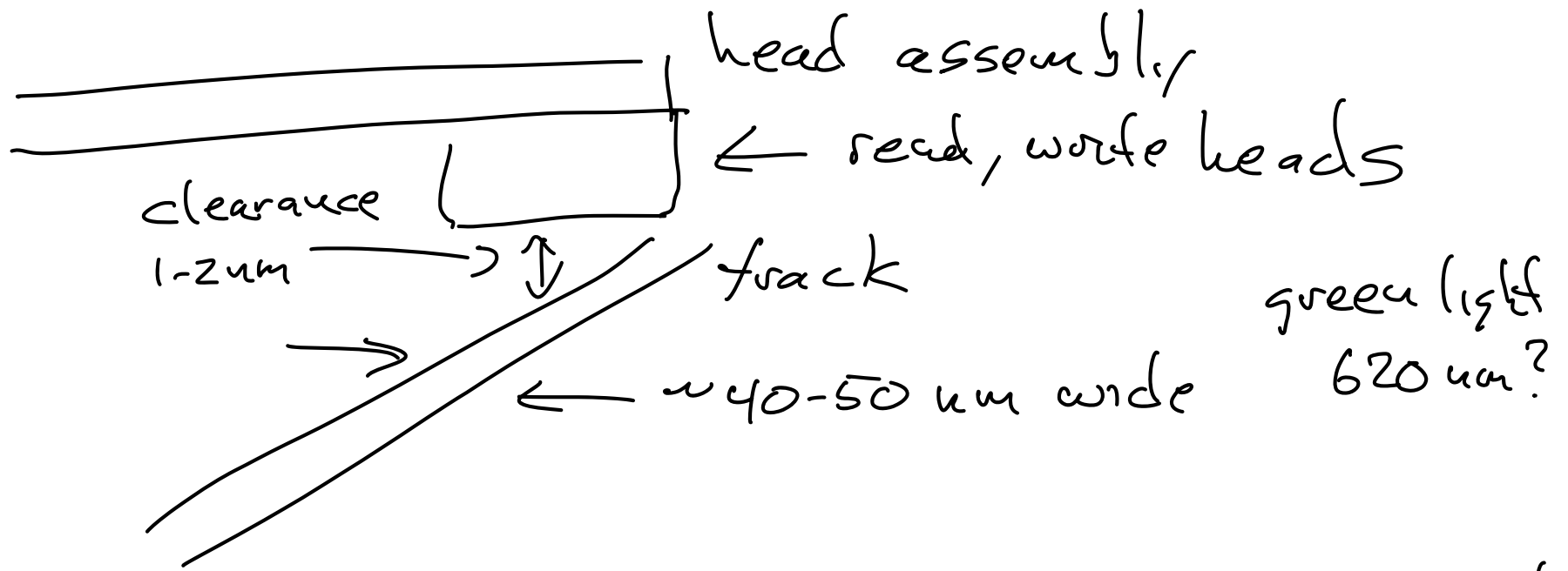
File systems, disks & SSDs, Lab 4

Hard drive - 1959 (58?) - 2015 primary storage
mid 60s



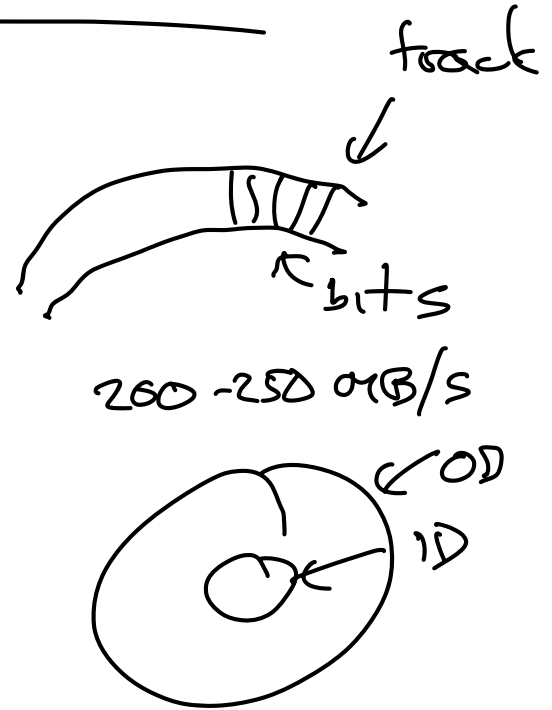
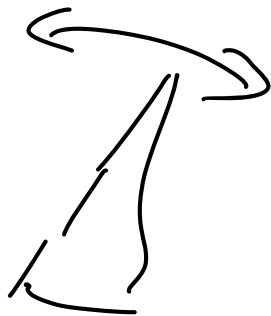
magnetic



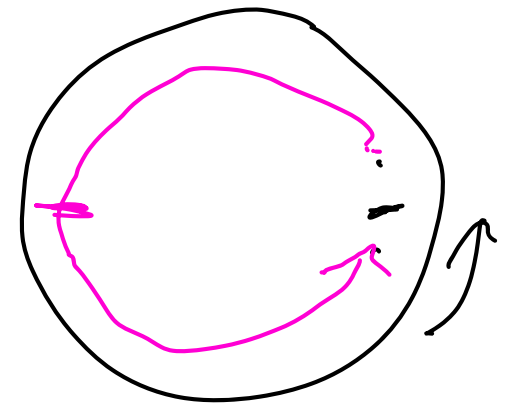


Hard drive characteristics:

- rotation speed 7200 RPM
- transfer rate - - - 200 MB/s
- seek time
15 μs max



transfer time for random I/O: 7200 RPM
 $= 8.3 \text{ ms}$
 seek time + ($\frac{1}{2}$ max = average)
 $\frac{1}{2}$ rotation (average) +
 transfer time



random 4KB block:

seek: 7.5 ms

rotation: 4.2 ms

$\frac{4 \text{ KB}}{200 \text{ MB/s}}$ transfer: $20 \mu\text{s} = 0.020 \text{ ms}$
 $\approx 11.7 \text{ ms}$

fixed $\approx 20 \mu\text{s}$
 variable (e.g. per-byte)

10 MB: 11.7 overhead

random 10 MB:

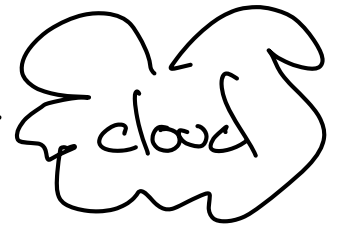
transfer: 5 ms

50 ms
 data
 transfer

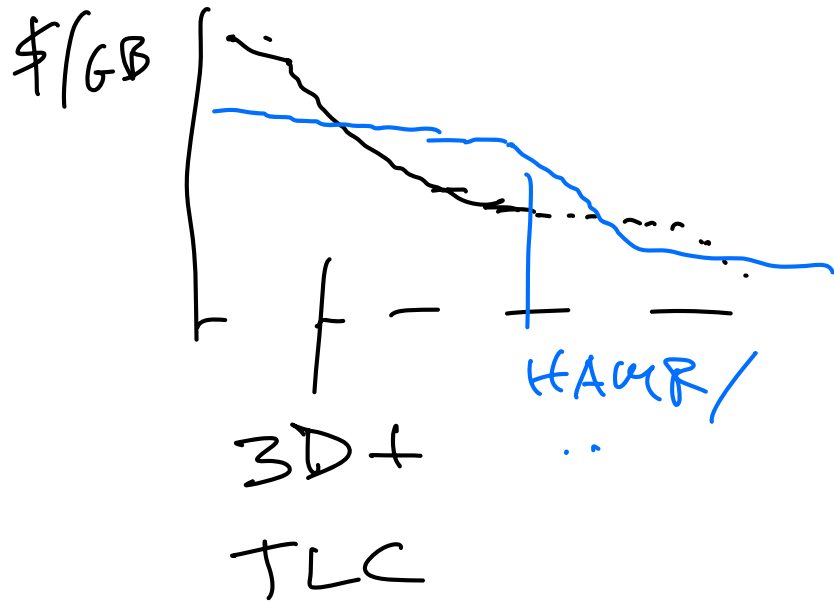
Are hard drives dead?

bulk storage

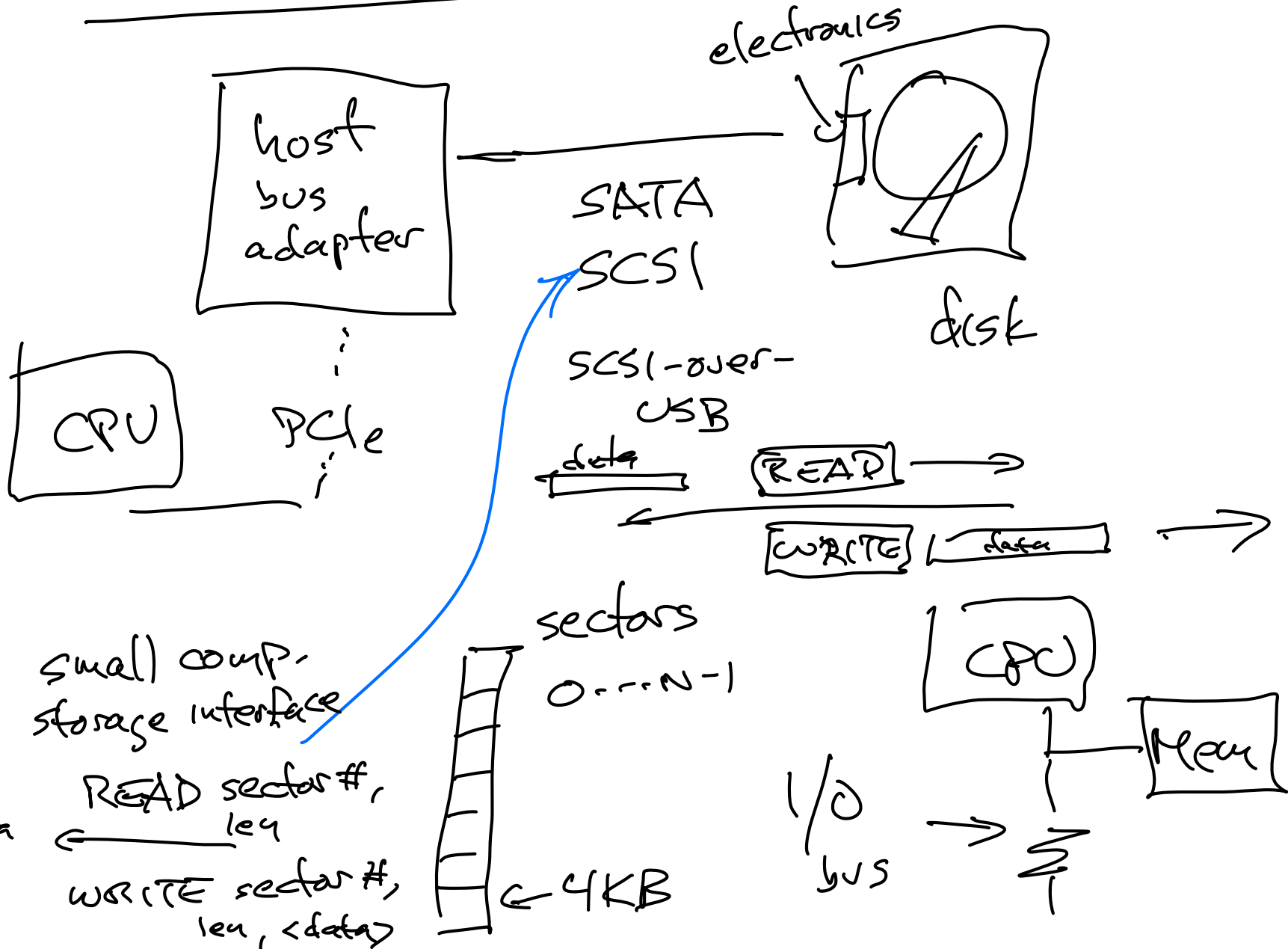
cheap, performance
not critical



backup
drives

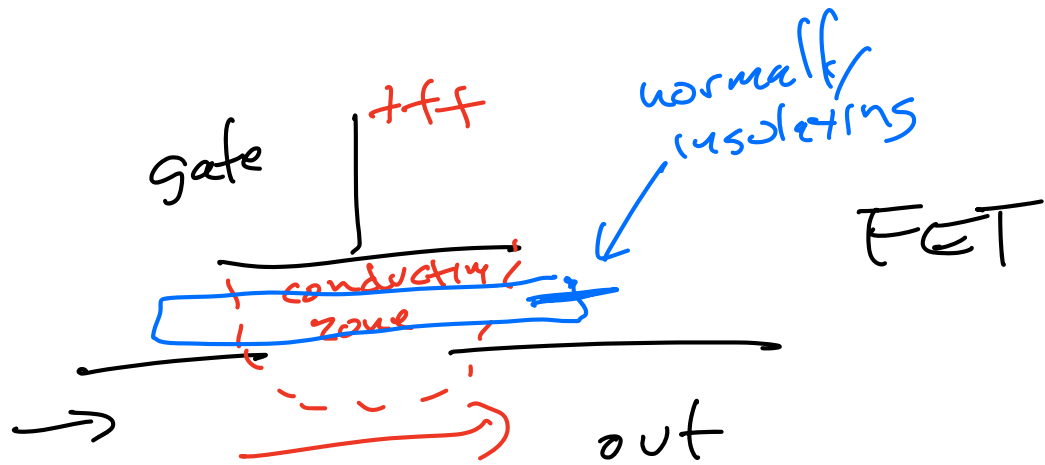


How your computer talks to a disk

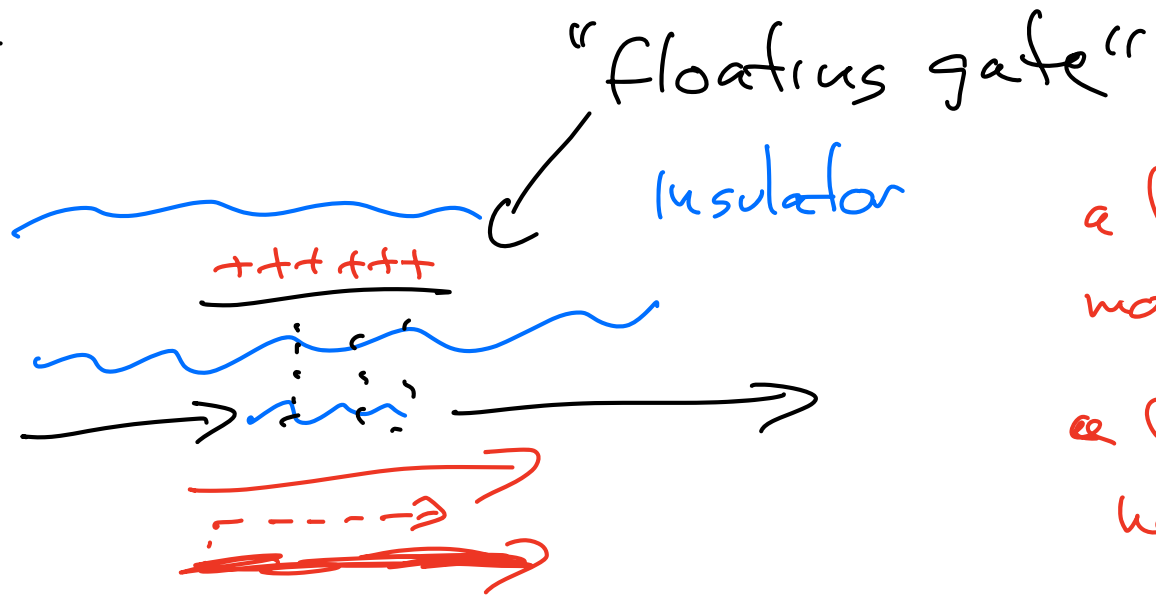
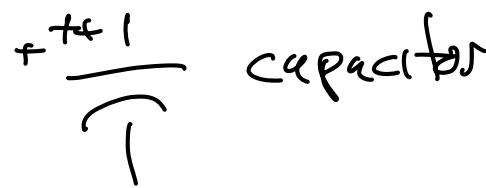


SSDs

NAND Flash

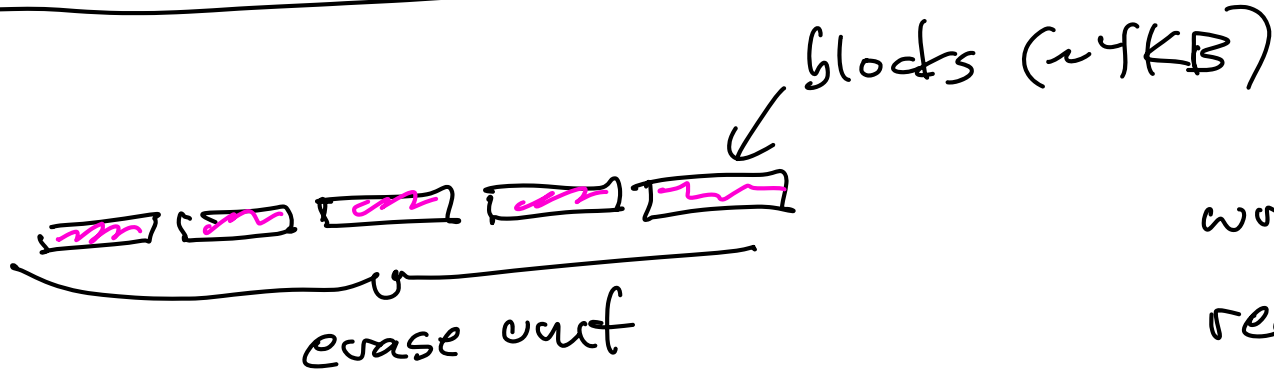


DRAM:



- a lot : 1 1
- more : 1 0
- a little : 0 1
- none : 0 0

Flash translation layer & garbage collection

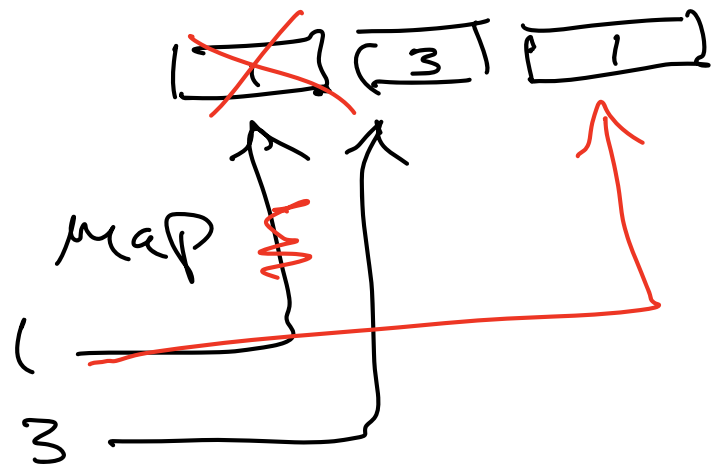


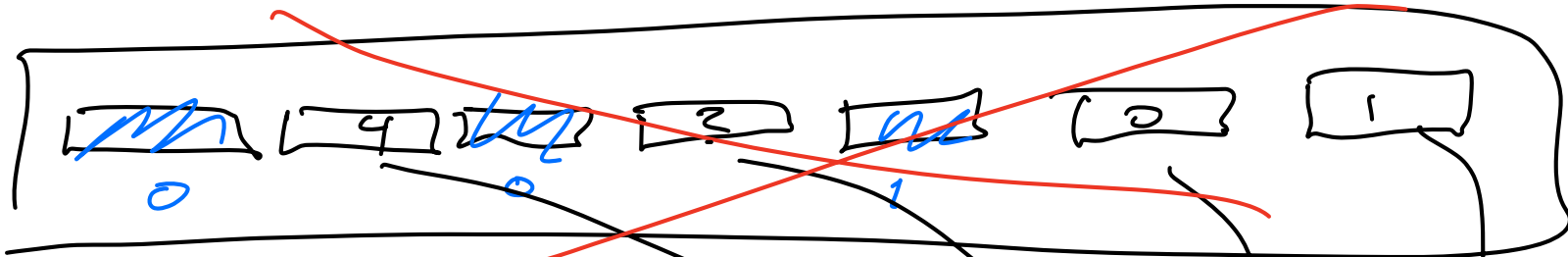
block write, bulk erase

write once (block)
read (block)
erase (erase unit)

- write 1
- write 3
- write 1.

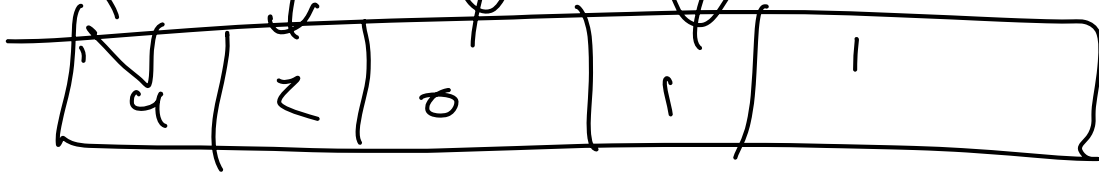
out-of-place write



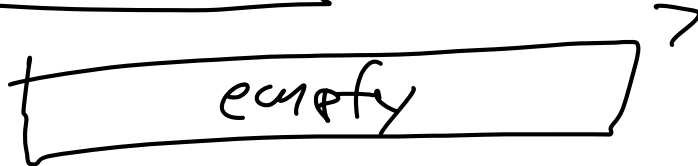


erase

4 copies



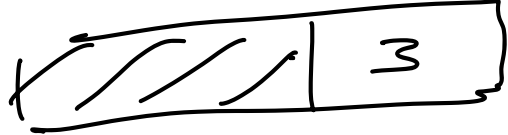
7 free



7 free



10 free



write amplification:

internal flash write ops =

"real" writes (from computer)

+ GC writes (moving data so can
erase a unit)

higher write amp. if:

- random small writes

- low free space %

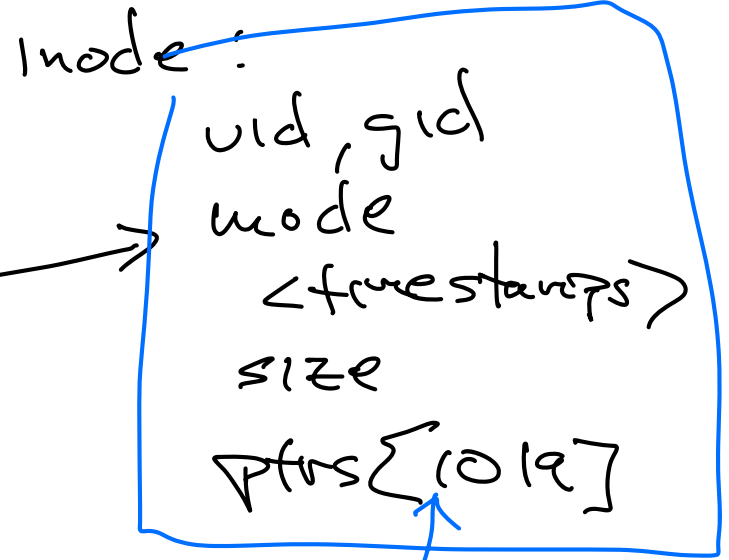
TRIM → tell SSD to forget blocks not
(discard) used by file system

Lab 4

python3 gen-disk1.py test.img
generate disk image

python3 print-disk.py test.img
views it

read-only
sort-of-unix
file system



type

reg. file = 100

directory = 040

S_ISREG(mode) S_ISDIR(mode)

permission

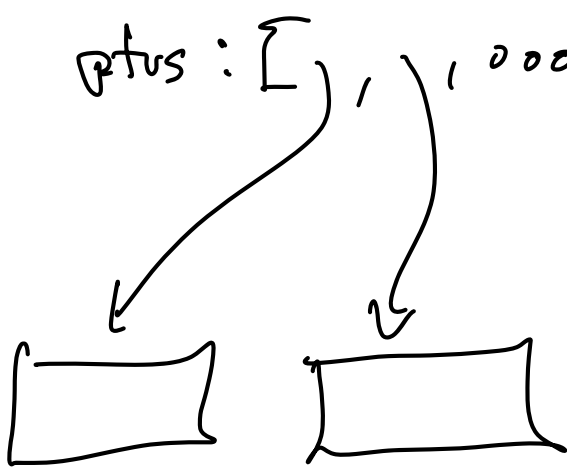
Rwx Rwx Rwx
~ ~ ~
user group "world"

4096B
block#5

block_read(buf, block#,
#blocks)

inode
size: 5000

ptrs: [, , 0000...]



$\text{DIV_ROUND_UP}(\text{size}, 4096)$

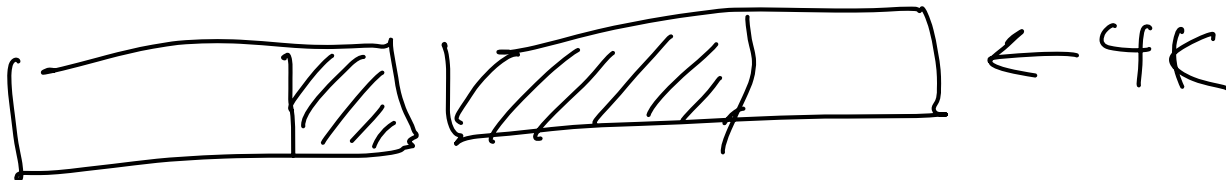
= 2

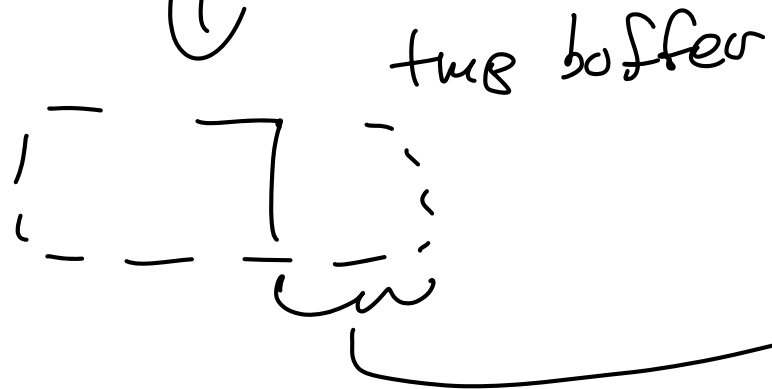
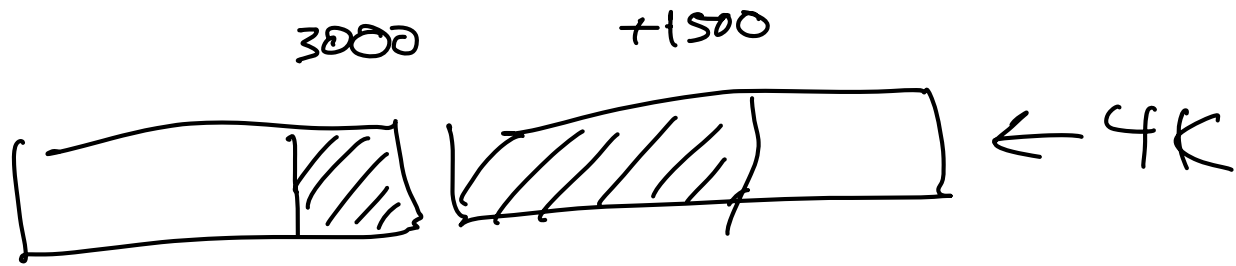
$$= \left\lceil \frac{\text{size} + N - 1}{N} \right\rceil$$

read(offset, len)

byte
granularity

3000 → +1500





read (buf, len, ...)



copy

