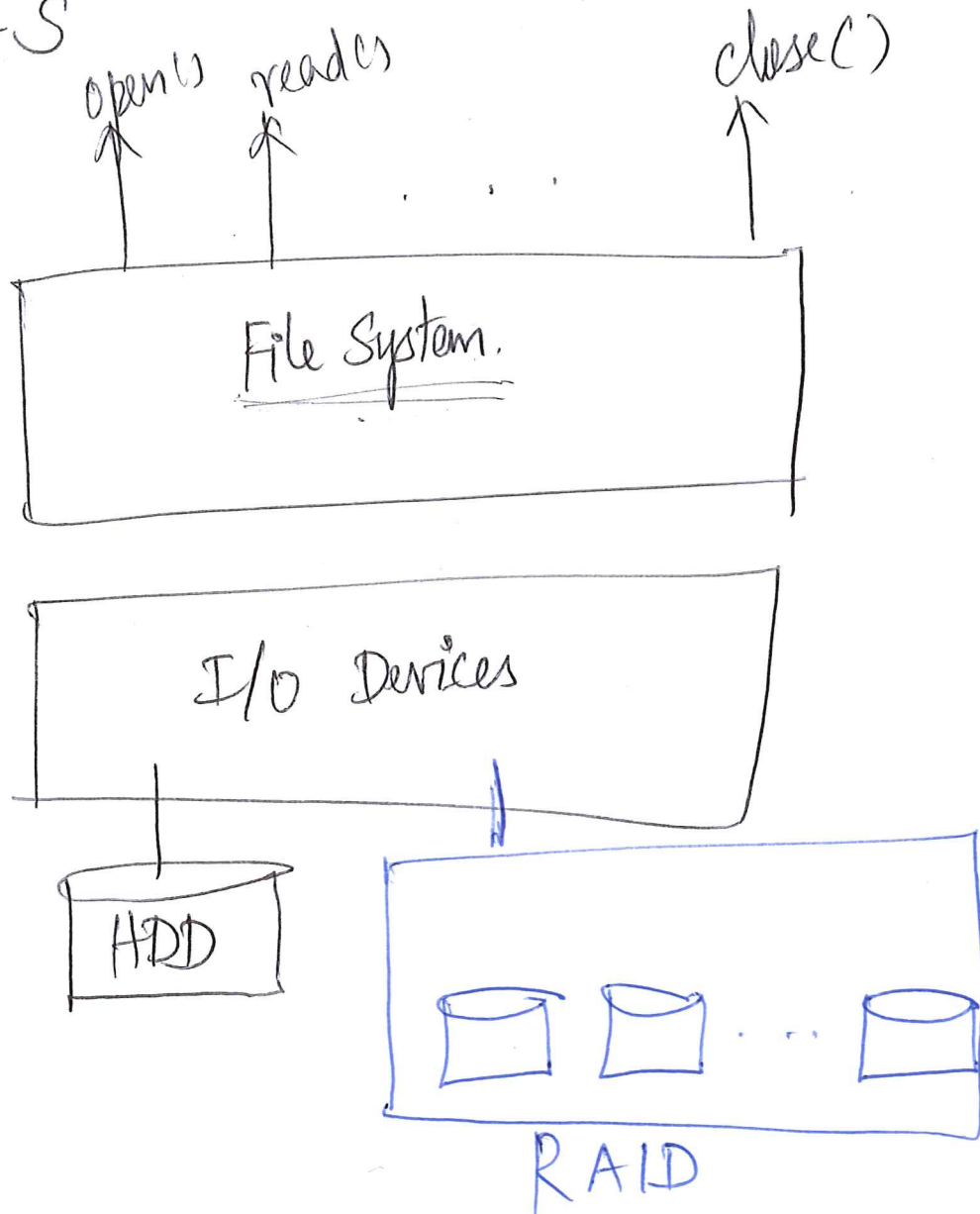


CS 537- Lecture 14

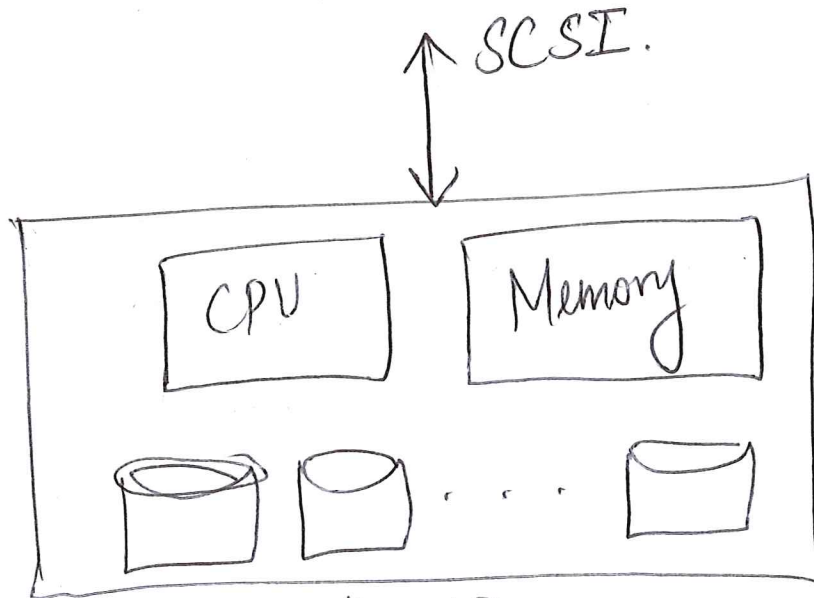
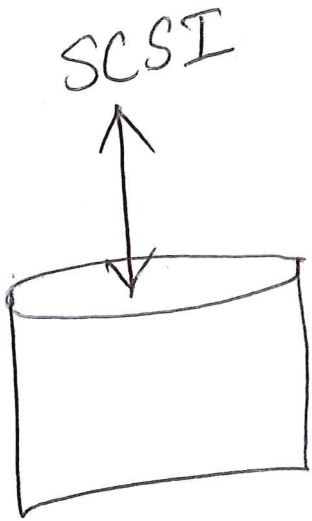
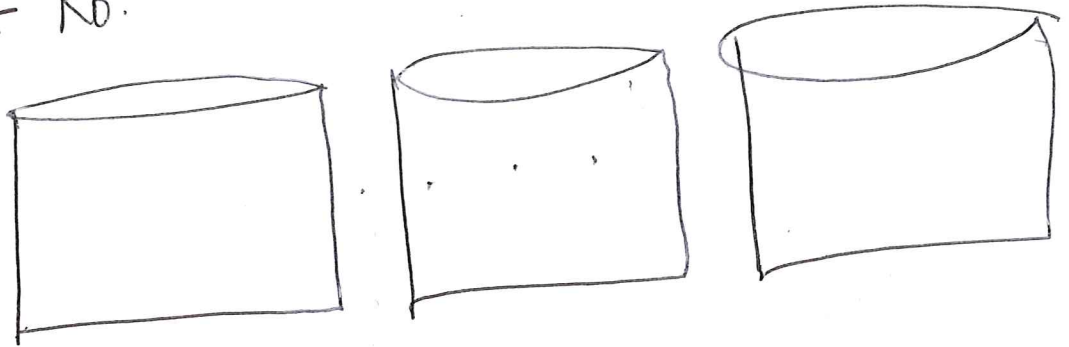
1. RAID
2. FFS
3. LFS



RAID

Redundant Array of Inexpensive Disks
(Independent)

- 1) Reliability - R_0, R_1, R_5
- 2) Performance:
- 3) Capacity - R_0 .

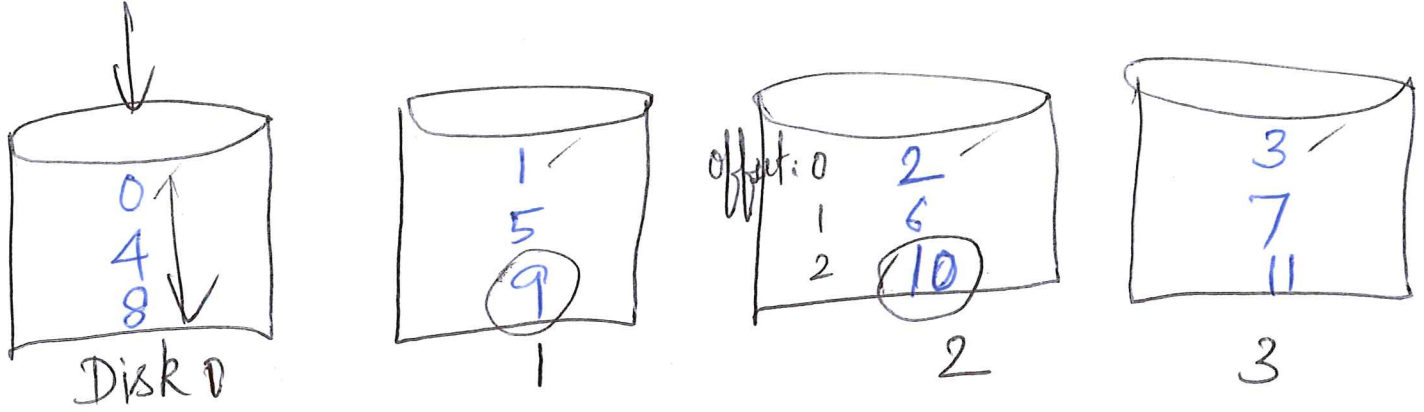


RAID.

Fault Model



RAID-0 (striping)



$$\text{disk \#} = \text{Addr} \% \text{NUMDISKS}$$

$$\text{offset} = \text{Addr} / \text{NUMDISKS}$$

chunk size

1 chunk = 4KB (1 block) ✓

1 chunk = 8KB (2 blocks)

D ₀	D ₁	D ₂	D ₃
0	2	4	6
1	3	5	7
8	10	12	14
9	11	13	15

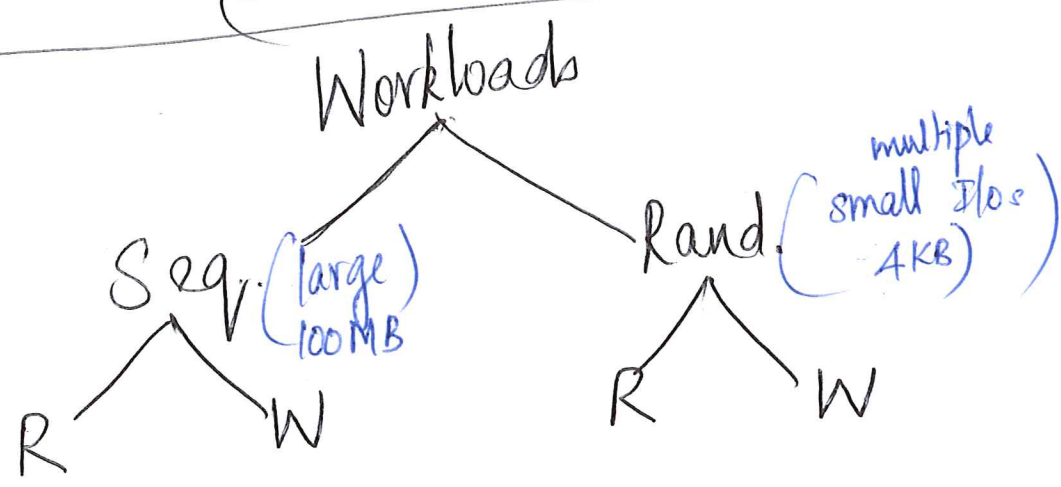
	small chunk	big chunk
positioning time of disk head	more	less
intra-file parallelism	more	less

1) Capacity = $N \times B$
disks # blocks/disk

2) Reliability = 0

3) Performance Analysis

1. ~~Latency~~ Latency of a single request
2. Throughput of many requests
(stead-state throughput)

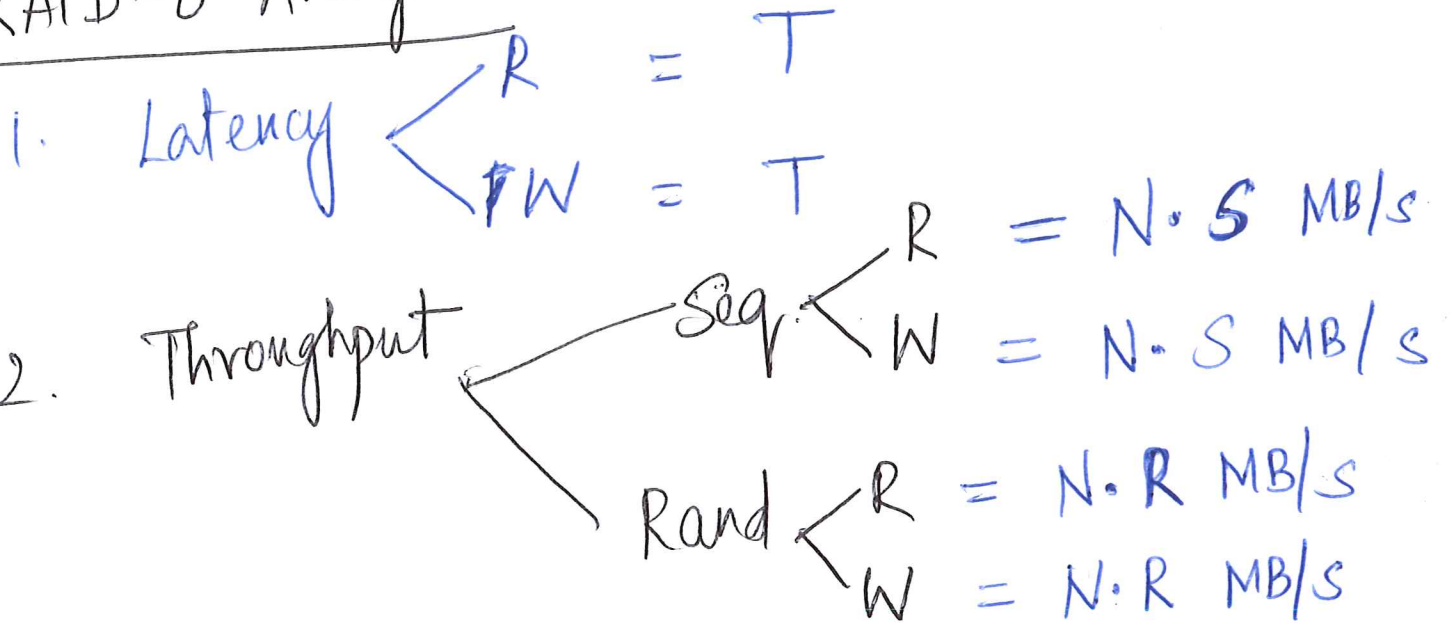


Data transfer Rate

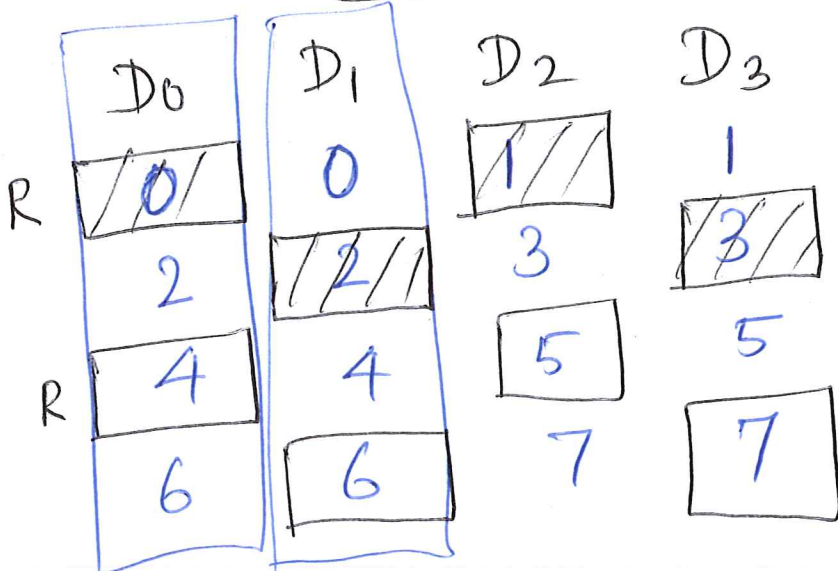


T = Time to read/write a block on a single disk.

RAID-0 Analysis



RAID-1 (Mirroring)



$$\text{capacity} = \frac{N}{2} \times B.$$

$$\text{Reliability} = \frac{N}{2} \text{ (if lucky)}$$

$$= 1$$

RAID-1 Analysis

1. Latency $\begin{cases} R = T \\ W = T \end{cases}$

2. Throughput

1. Sequential

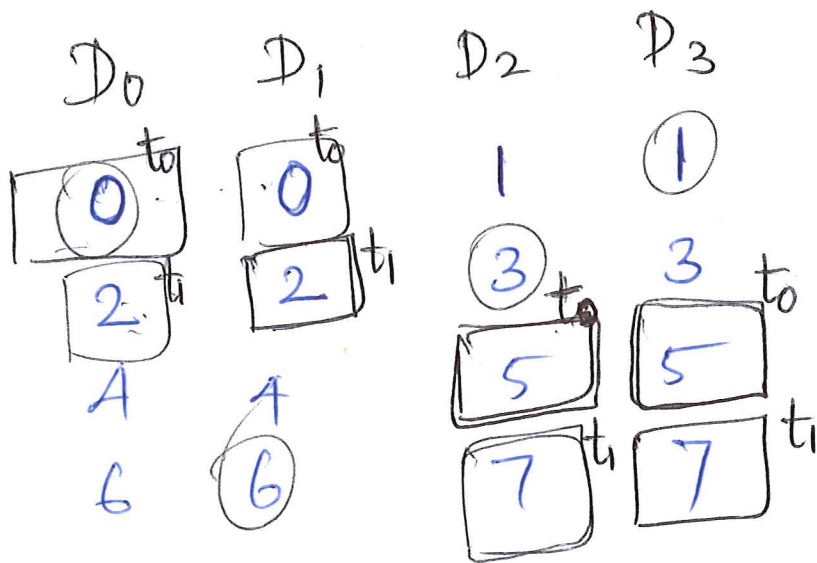
Reads = $\frac{N \times S}{2}$ MB/s.

Writes = $\frac{N \times S}{2}$ MB/s

2. Random

Reads = $N \times R$ MB/s

Writes = $\frac{N \times R}{2}$ MB/s

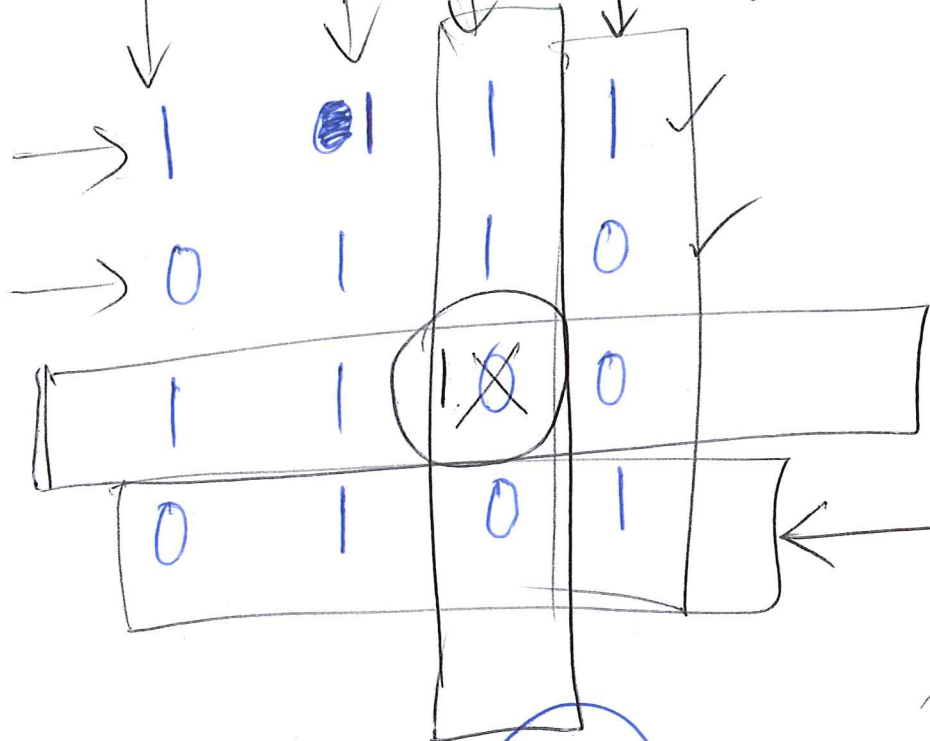


Write
0-3.

Random write: 0, 2, 5, 7

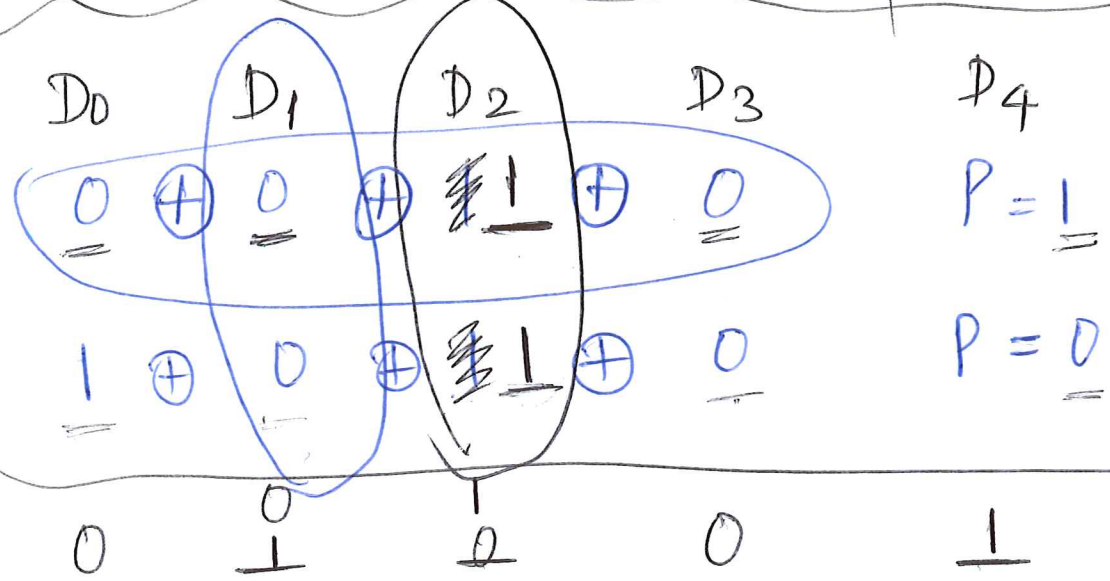
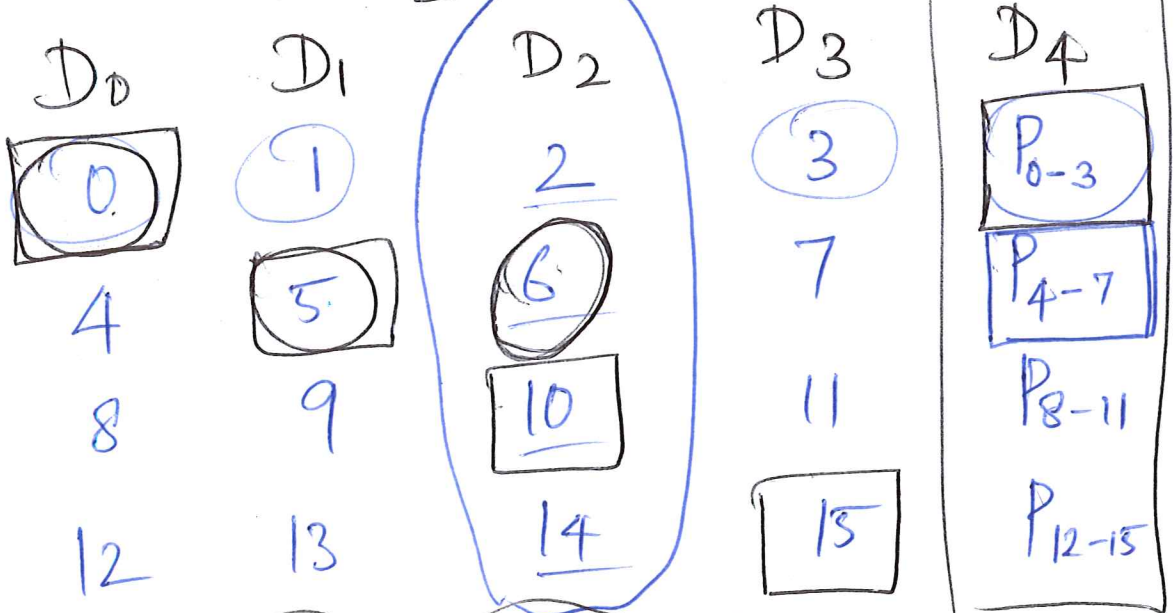
RAID - 4 (parity)

parity bits



parity bits

Parity Disk



RAID-4 Analysis

$$\text{capacity} = (N-1) \times B$$

$$\text{reliability} = 1$$

$$\text{Latency} \begin{cases} \text{Read} = T \\ \text{Write} = 2T \end{cases}$$

How writes happen?

D ₀	D ₁	D ₂	D ₃	D ₄
0	1	2 ^{Cold}	3	P ^{Pold}

$$0 \oplus 1 \oplus 2^{\text{New}} \oplus 3 = P^{\text{New}}$$

How to compute parity?

Write to block 2

1. Read blocks 0, 1, 3.
2. ~~Write~~ compute new parity
3. Write blocks 2' and P'

Additive parity

Subtractive Parity

$$P_{new} = (C_{old} \oplus C_{new}) \oplus P_{old}$$

$$= (0 \oplus 1) \oplus 1 \Rightarrow 0$$

1. Read C_{old} } Parallel
2. Read P_{old} } Parallel
3. $P_{new} = (C_{old} \oplus C_{new}) \oplus P_{old}$ // compute new parity
4. Write C_{new} } Parallel
5. Write P_{new} }

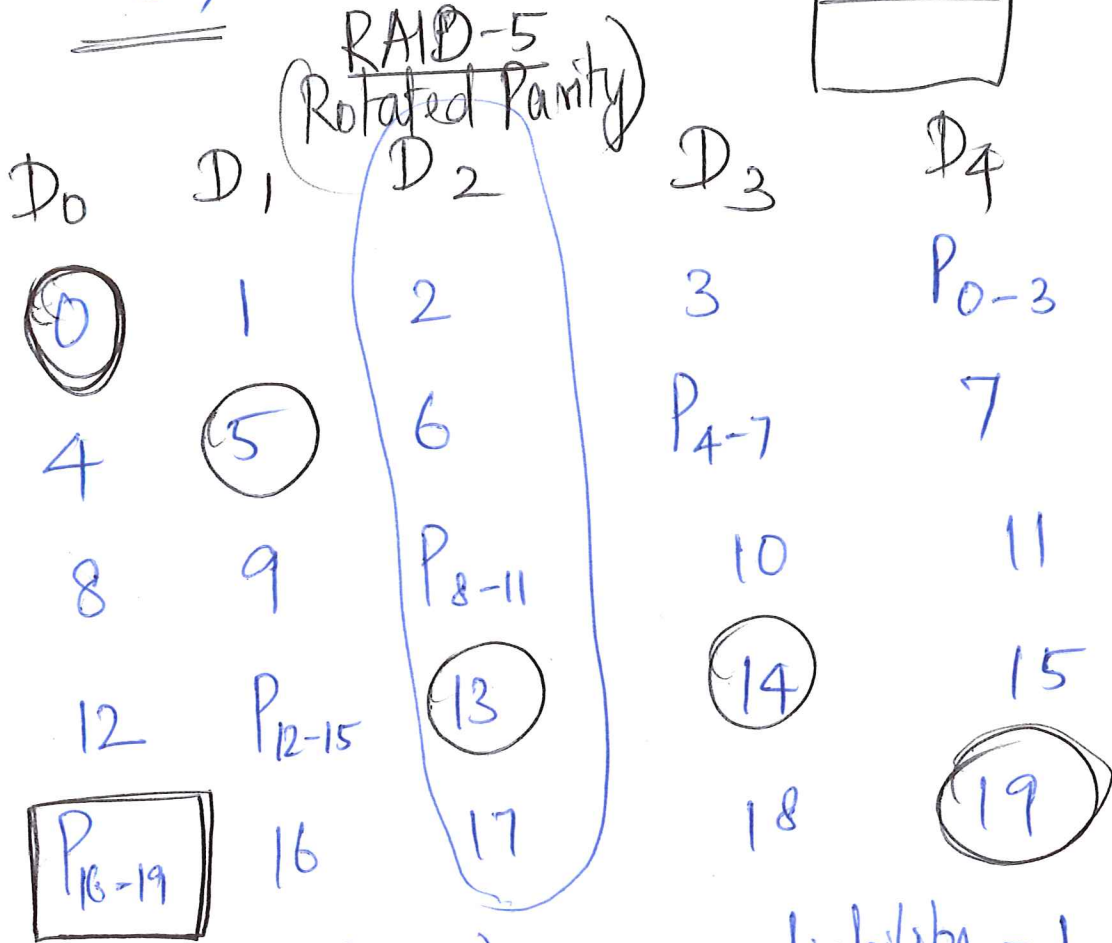
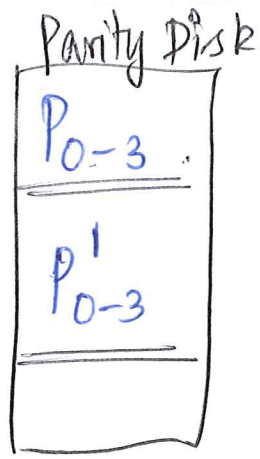
Throughput

1. Seq. $\left\{ \begin{array}{l} \text{Read} = (N-1) S \text{ MB/s} \\ \text{Write} = (N-1) S \text{ MB/s} \\ C_0' \oplus C_1' \oplus C_2' \oplus C_3' = P_{0-3}' \\ \downarrow \\ \text{new content of } C_0 \end{array} \right.$

2. Rand. $\left\{ \begin{array}{l} \text{Reads} = (N-1) \cdot R \text{ MB/s} \\ \text{Writes} = \frac{R}{2} \text{ MB/s} \end{array} \right.$

t₀: Read 0, 5, 10, 15,

t₁: Writes 1, 0,



capacity = (N-1)B

reliability = 1

latency $\left\{ \begin{array}{l} R = T \\ W = 2T \end{array} \right.$

throughput

Seq

R = (N-1) · S

W = (N-1) · S

Rand

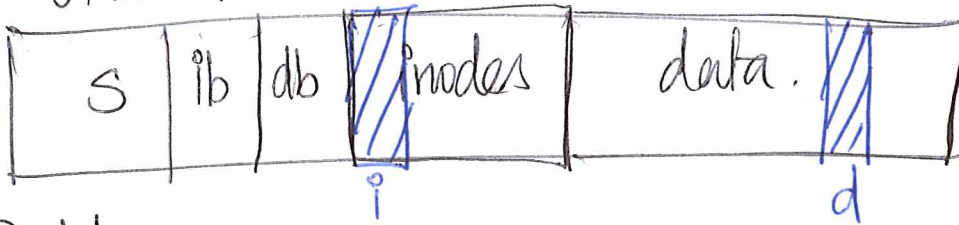
R = N · R

W = $\frac{N \cdot R}{4}$ MB/s.

FFS

Fast File System.

UNIX FS



Problems

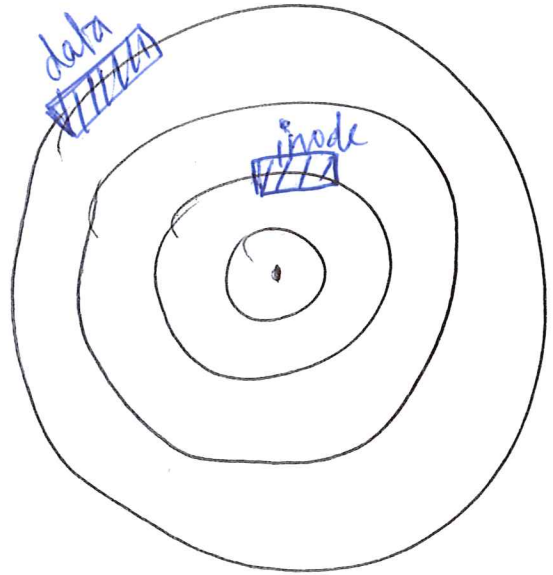
read a file

- 1) read inode
- 2) read data blocks

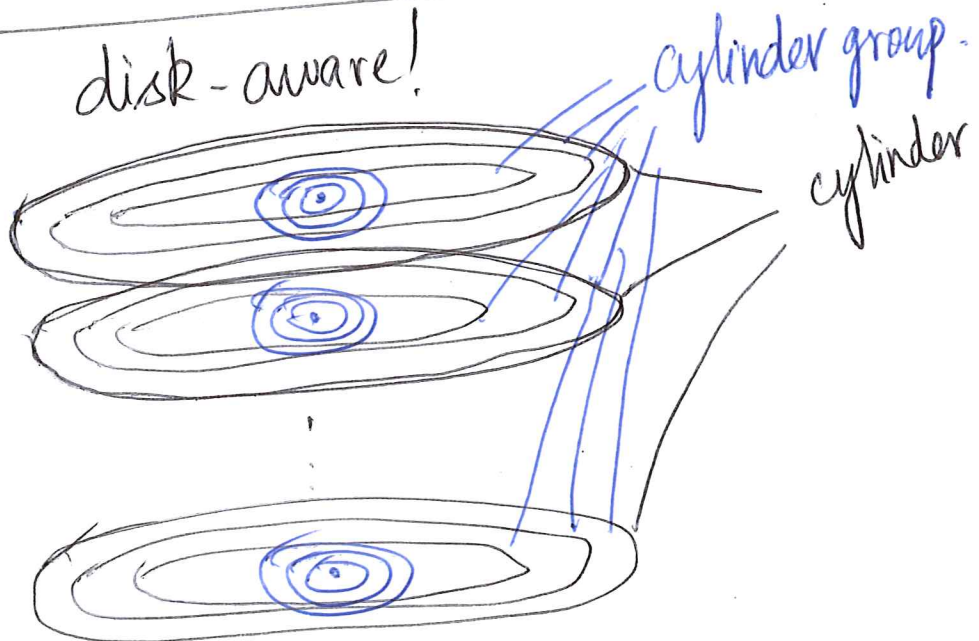
1. Performance

2. UNIX FS - Small blocks

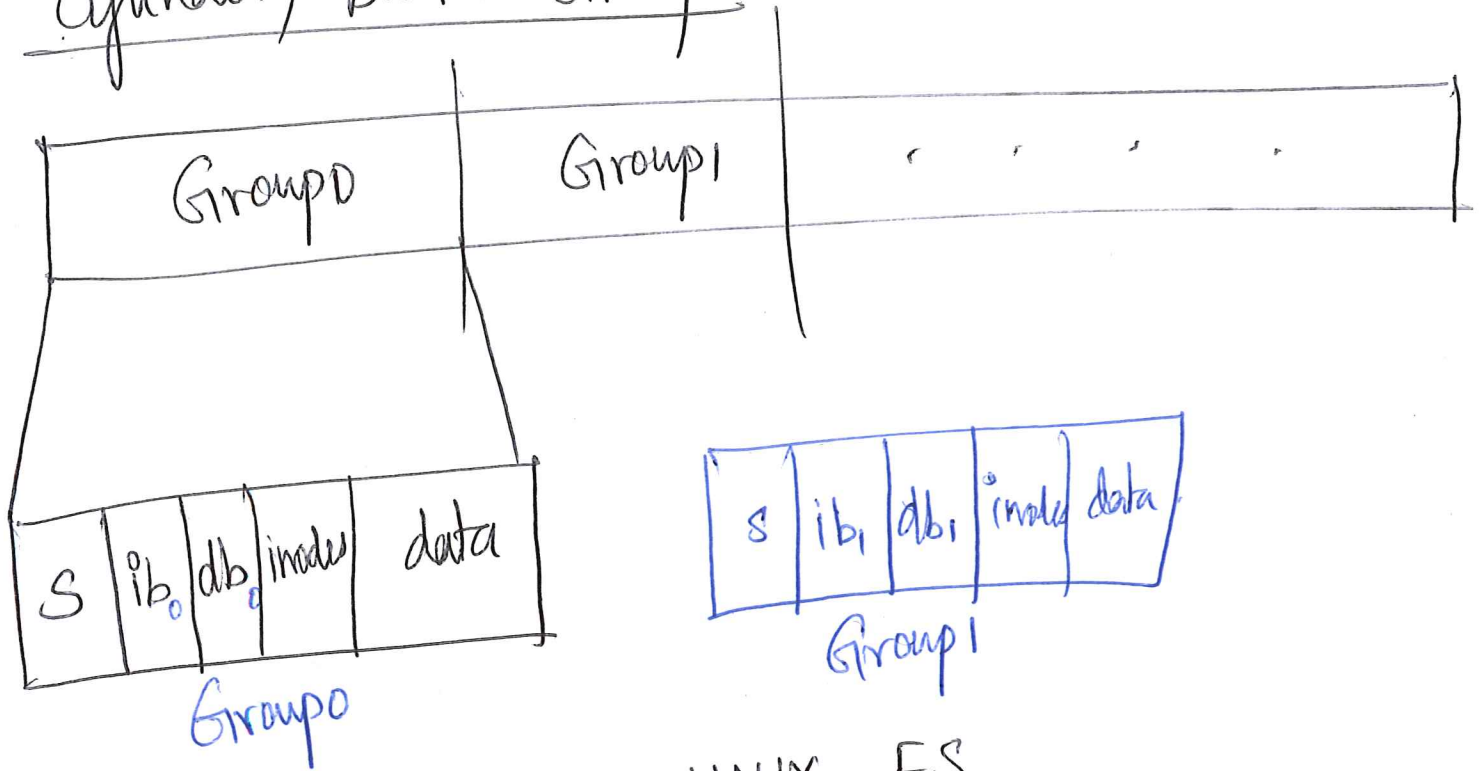
- reduce internal frag. (512 bytes)
- increases positioning costs.



Soln: FS disk-aware!



Cylinder/Block Groups



APIs - same as UNIX FS
open(), read(), write(), close().

Allocation Policies

Rule: Keep related stuff together.
" unrelated " apart.

What is related?

- 1) a file's inode + data \Rightarrow same group.
- 2) files within the same directory.

What is unrelated?

diff. directories.

$/d_1/f_1$

$/d_1/f_2$

$/d_2/d/f_3$

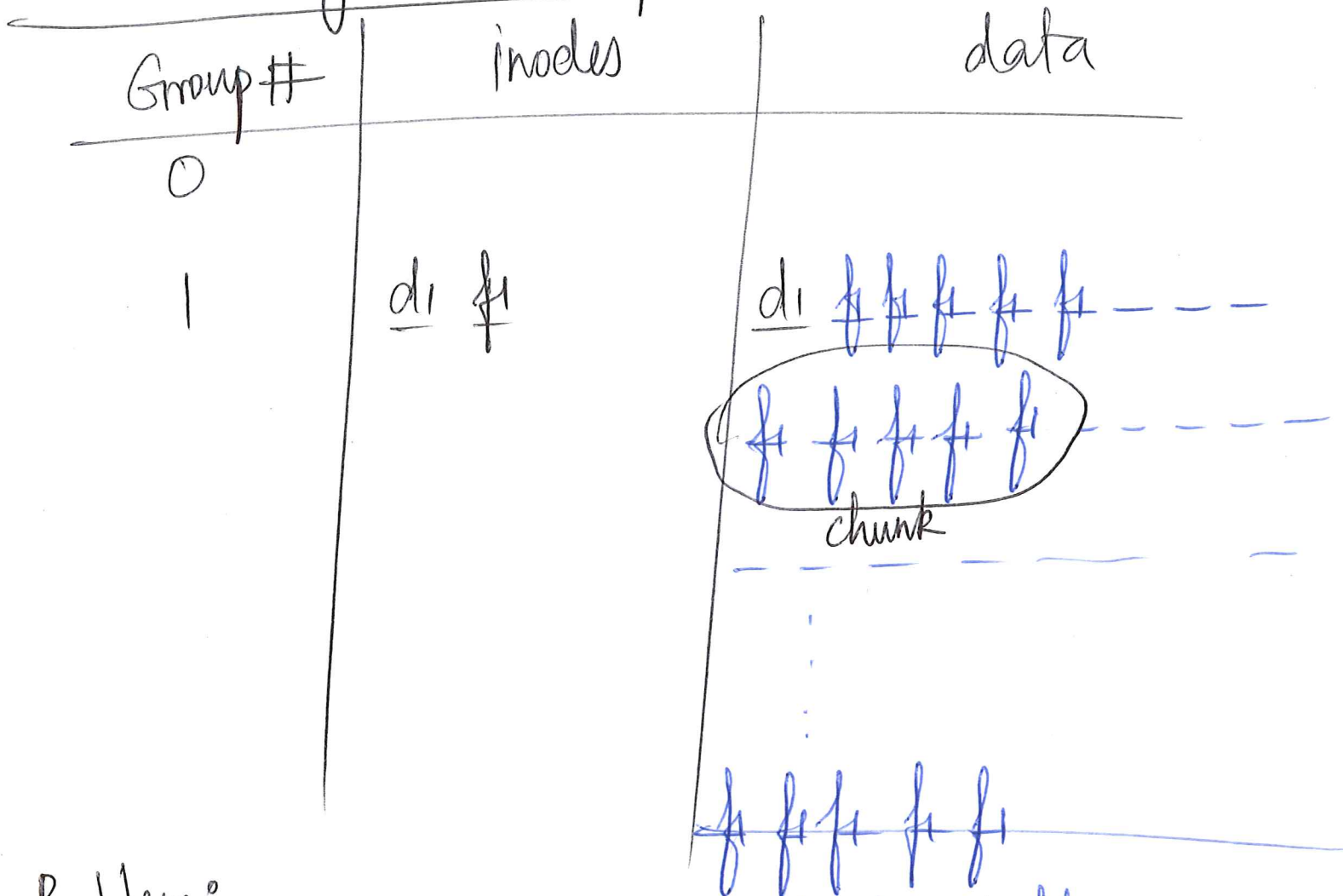
Group #	inodes		data	
	/	-----	/	-----
0				
1	d_1	$f_1 f_2$	d_1	$f_1 f_2 f_2$
2	d_2	f_3	d_2	$f_3 f_3$
3				
4				
5				

To allocate directories

- 1) which groups have low # of directories
 - 2) " " " a high # of free inodes
- How to write large files?

Group #	inodes		data	
	/	-----	/	-----
0				
1	d_1	f_1	d_1	$f_1 f_1 f_1 \dots f_1$
2			f_1	f_1
			f_1	f_1

FFS - Big File Exception



Problem:

Perf. of seq. file access of large files
 Positioning Time + Data transfer time
 (Seek + Rotation)

10 ms

10 ms

Data transfer rate = 40 MB/s

1 s → 40 MB

$$10 \text{ ms} \rightarrow \frac{40 \text{ MB} \times 10 \text{ ms}}{1 \text{ s}} = 40 \times 10^6 \times 10^{-3}$$

$$= 400 \times 10^3 = \underline{\underline{400 \text{ KB}}}$$

Positioning
10% time
10ms

Data transfer
90% time
90ms

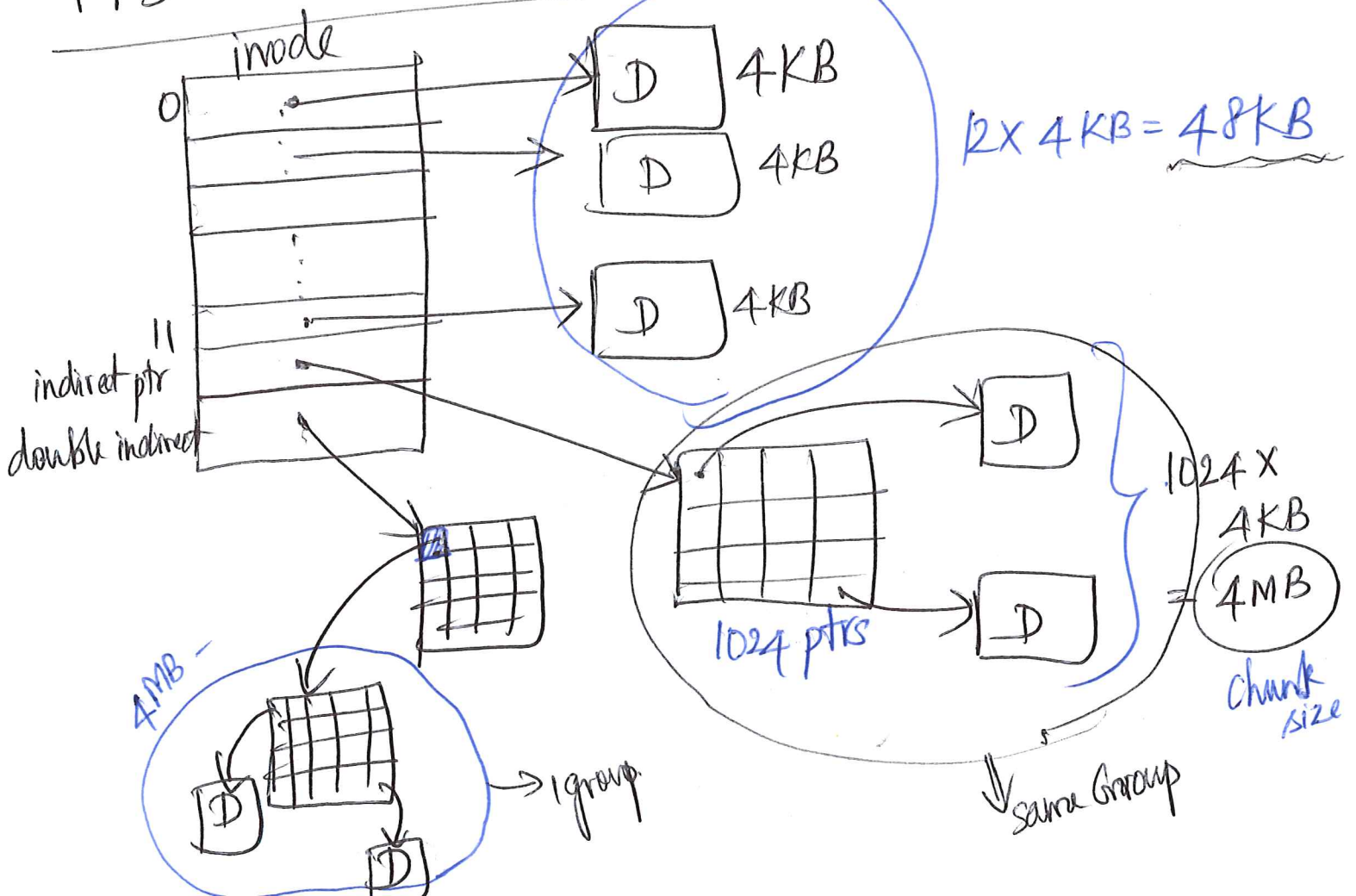
$$1s \rightarrow 40MB$$

$$90ms \rightarrow \frac{40MB \times 90ms}{1s}$$

$$= 3600 \times 10^3 \approx \boxed{3.6MB}$$

(chunk size)

FFS - did NOT use amortization.

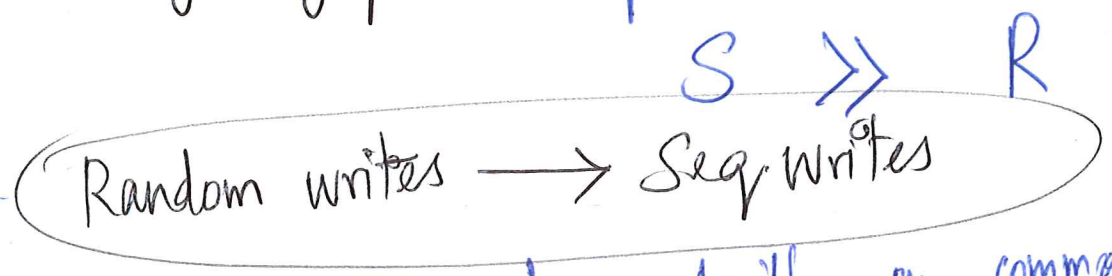


LFS Log-structured File System.

Observations

1. Systems have larger memories
 - Reads are serviced from memory.
 - Disk traffic - writes *

2. Large gap. Seq I/O Rand I/O



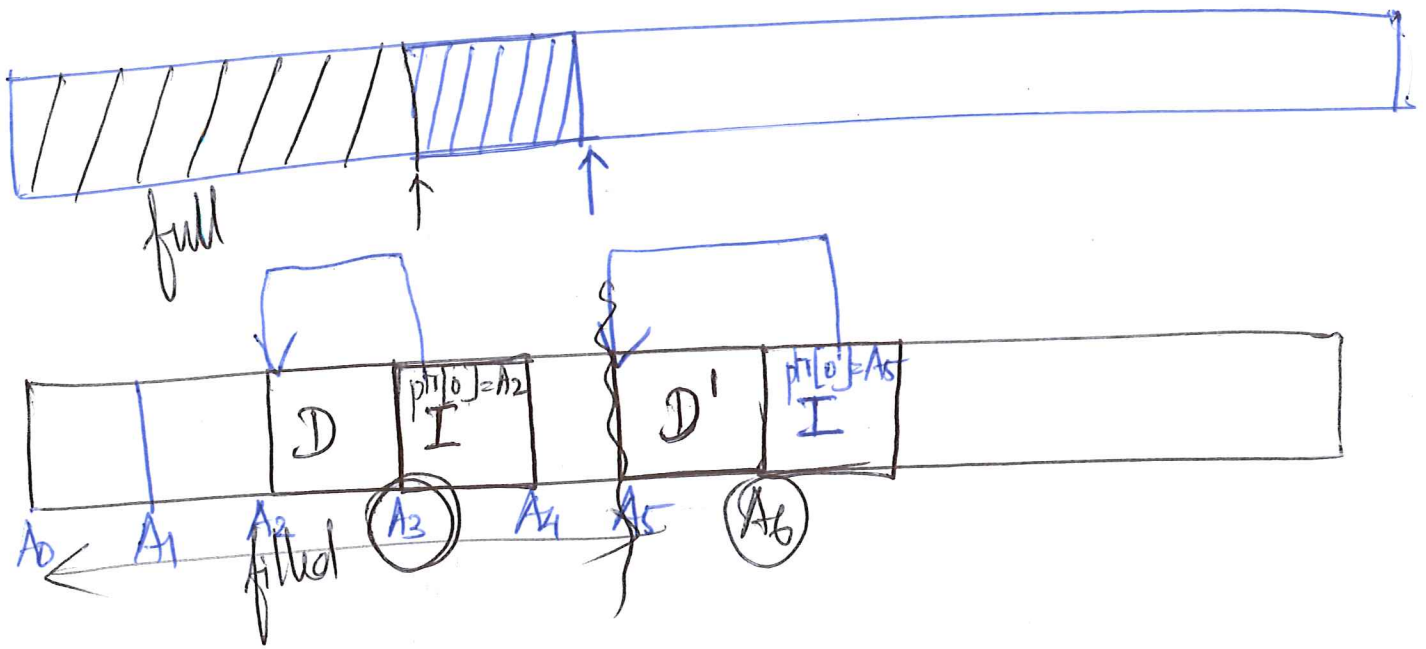
3. Existing FS perform badly on common workloads
e.g. file creation.

1. inodes \leftarrow dir
file
2. ' bitmap \leftarrow inode
data.
3. data \leftarrow dir
file.

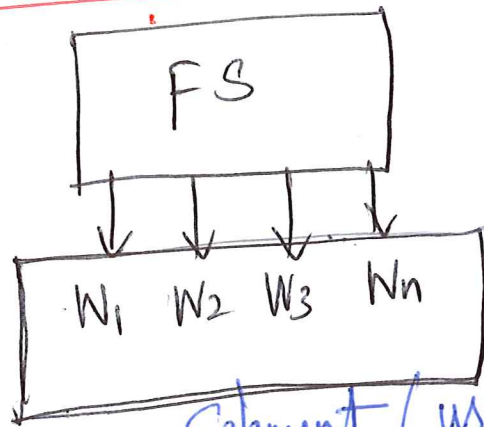
4. FS \rightarrow not RAID aware.

"small writes problems"

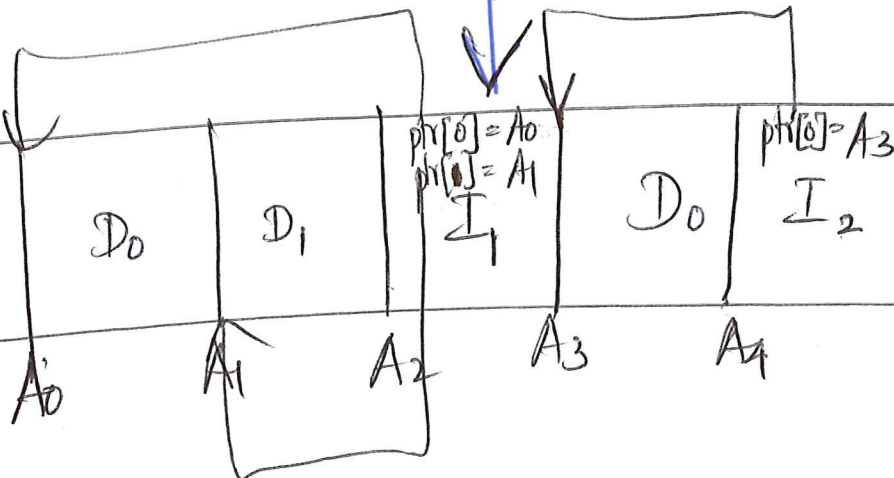
IDEA:

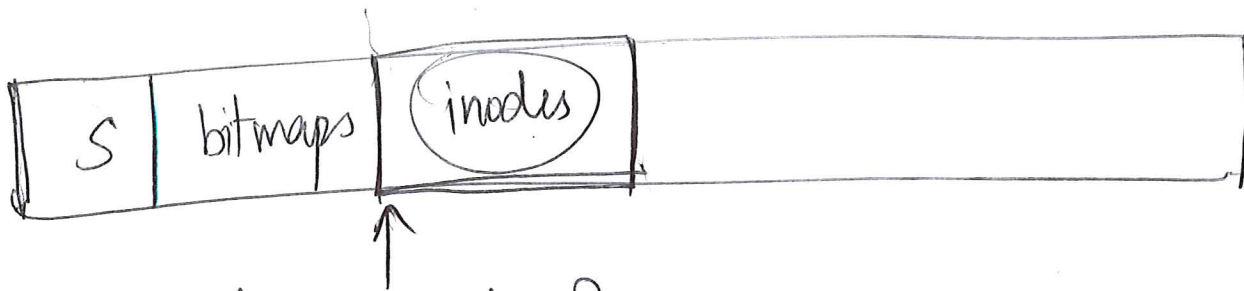


WRITE BUFFERING.



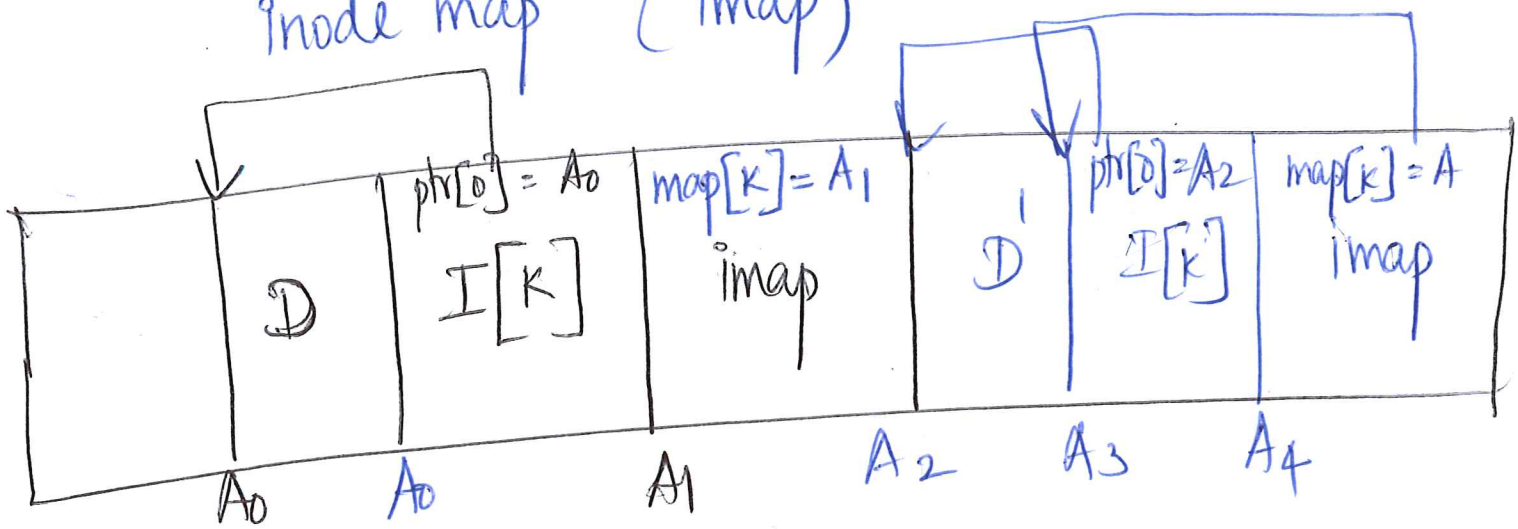
Segment (usually in MB)



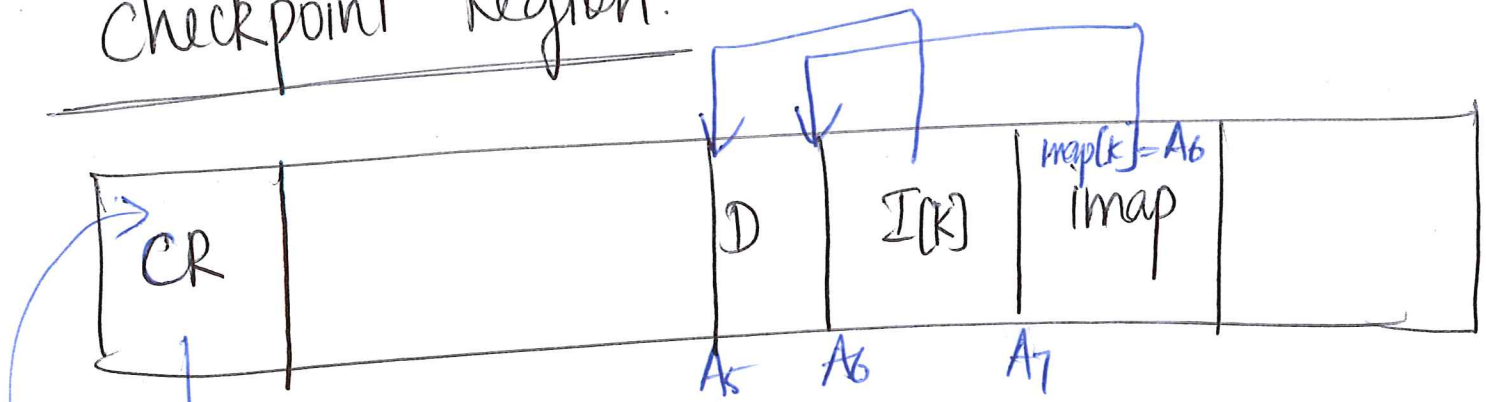


How to find inodes?

inode map (imap)



Checkpoint Region.



fixed location.

imap[k...k+n] = A₇

Append

