

Persistence: RAID

Ram Alagappan
CS 537, Spring 2019

Administrivia

P2b, midterm, P3 grades posted

Midterm solutions available
for regrades, see Piazza

P4 due on 4th , Thursday

Learning Outcomes

Why more than one disk?

What are the different RAID levels?

(striping, mirroring, parity)

How do we compare the RAID levels? What metrics?

Only One Disk?

Sometimes we may need many... Why?

- 1 Capacity
- 2 Performance
- 3 Reliability

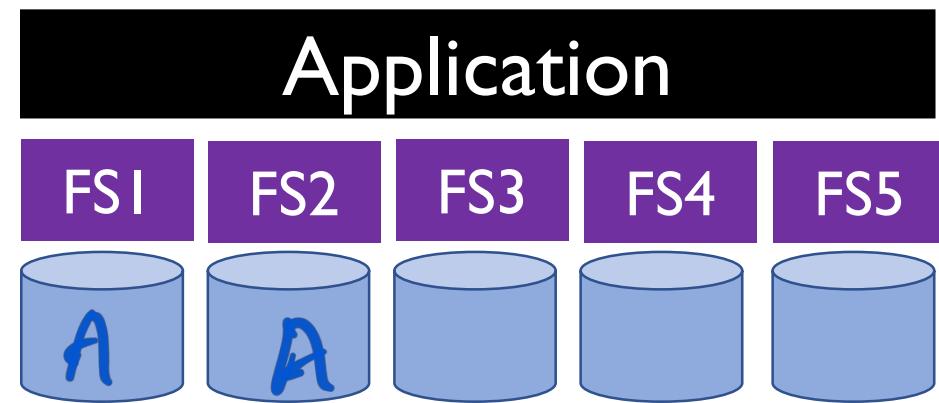


Challenge: most file systems work with one disk

Solution- I: JBOD

Just a Bunch Of Disks

Applications store data on different FS

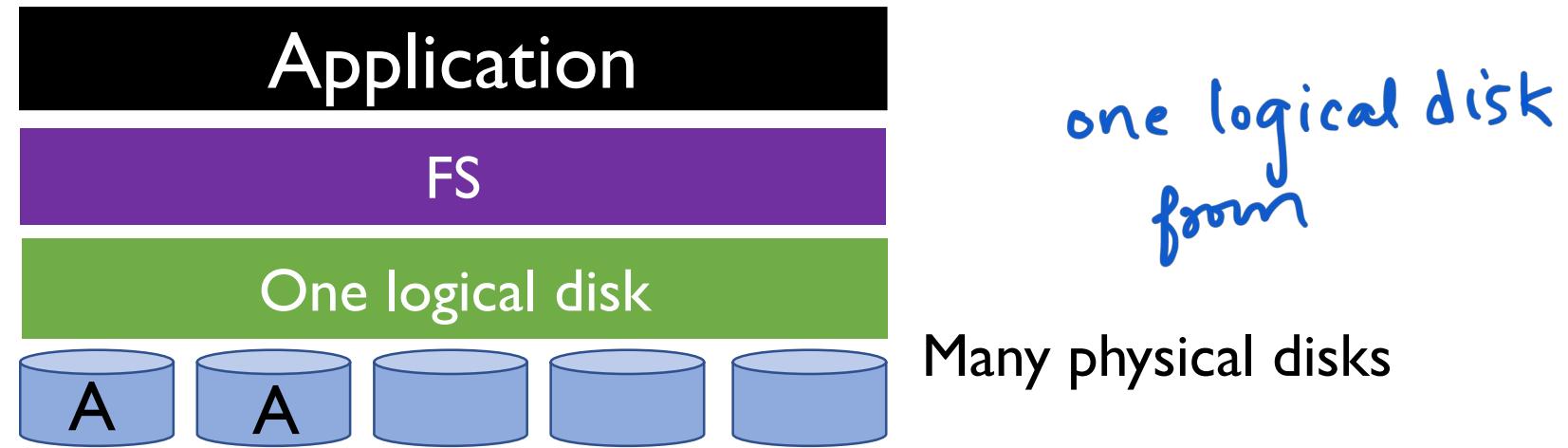


e.g.,
critical
data: app
decides to replicate

Downsides: need to know multiple devices,
need to be rewritten, not deployable

Solution-2: RAID

Redundant Array of Inexpensive (Independent) Disks



Advantages: transparent to apps, deployable

Improved capacity, performance, and reliability!

Fault Model

Simple: fail-stop model

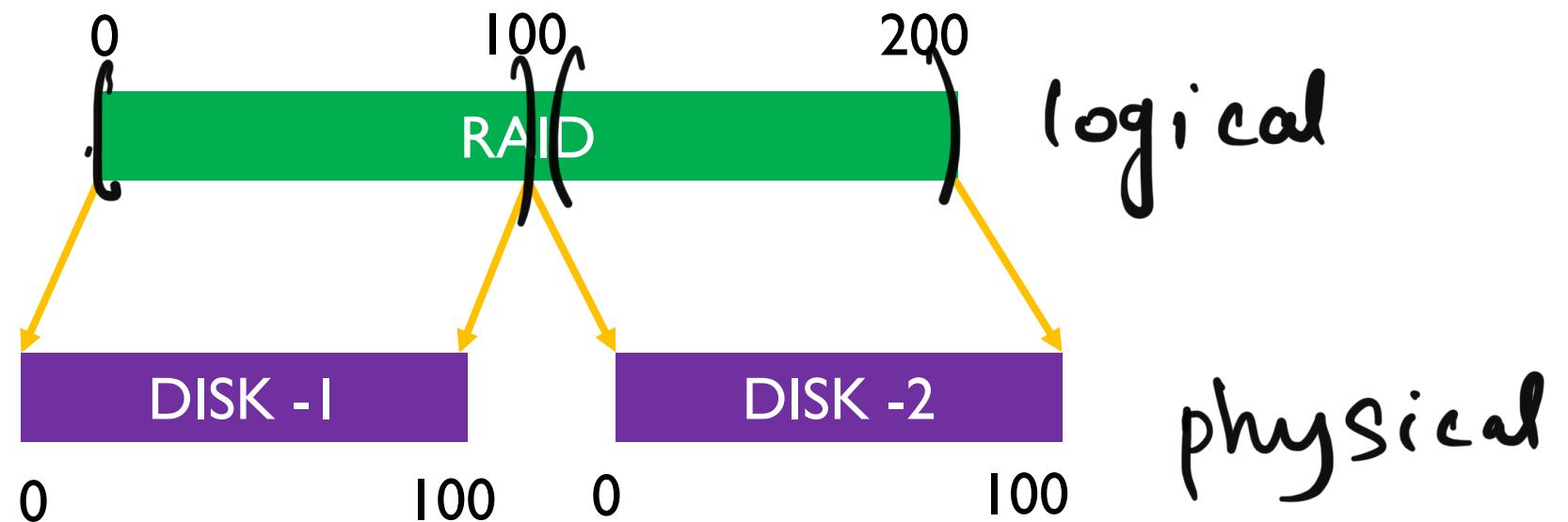
Works correctly or fails entirely

Easy to detect: system knows when not working

No silent failures: no data corruption etc.,

General Strategy: Mapping

Build fast, large disk from smaller ones



Mapping

How should we map logical block addresses to physical block addresses?

Is this similar to something?

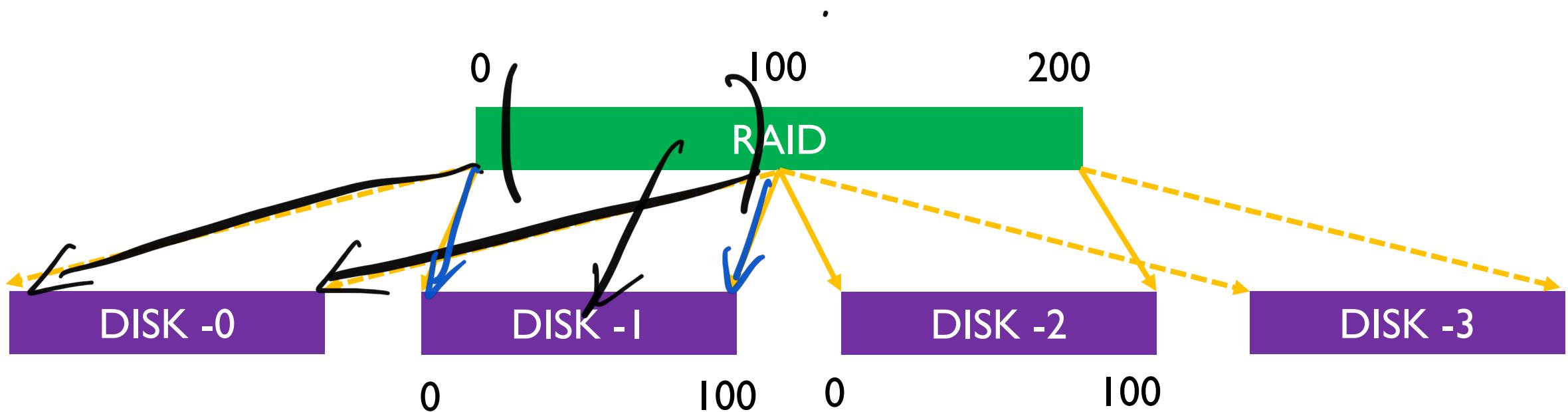
virtual memory

Virtual memory - dynamic mapping: page tables

RAID - Static mapping: simple calculation

General Strategy: Redundancy

Add even more disks for reliability



Redundancy

Trade-offs to amount of redundancy

Increase number of copies

improves reliability (and maybe performance)

Decrease number of copies (deduplication)

improves space efficiency

RAID Analysis

RAID: different levels

Workload: types of reads/writes issued by app

Metric: capacity, reliability, performance

Given **Workload**, Raid level, determine **Metric**

RAID Levels

Levels –

0 (striping)

1 (mirroring)

4 (parity)

5 (rotated parity)

2, 3, 6

We'll not discuss these

logical : 1 - 1000 blocks

Workload

Single request

read write

e.g., read(51)

Write (5)

latency - t

Many requests

random

reads

writes

read: 0, 51, 78, 767

...
500.

Throughput
· ut / BW
→ MBs / s

sequential

reads

writes

e.g., read(100 - 500)

Metrics

Capacity: how much space can apps use?

Reliability: how many disks can we safely lose?

Performance: latency? throughput?

Metrics

Normalize each to characteristics of one disk

N := number of disks

C := capacity of 1 disk

S := sequential throughput of 1 disk

R := random throughput of 1 disk

D := latency of one small I/O operation

$$C \sim 500 \text{ GB}$$

$$S \sim 100 \text{ MB/s}$$

$$R \sim 1-7 \text{ MB/s}$$

$$S \gg R$$

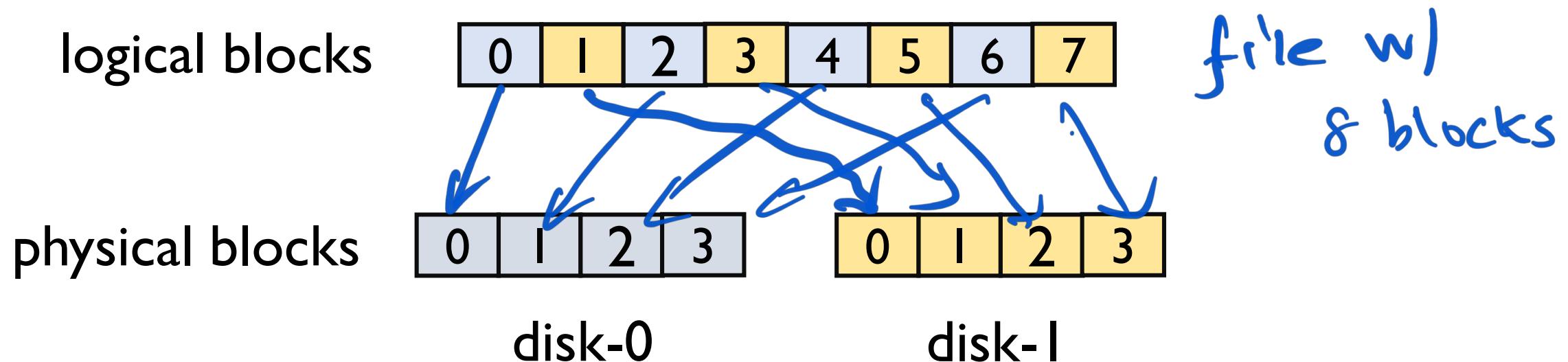
Reasonable S ? Reasonable R ? How does S compare to R ?

Assumptions.

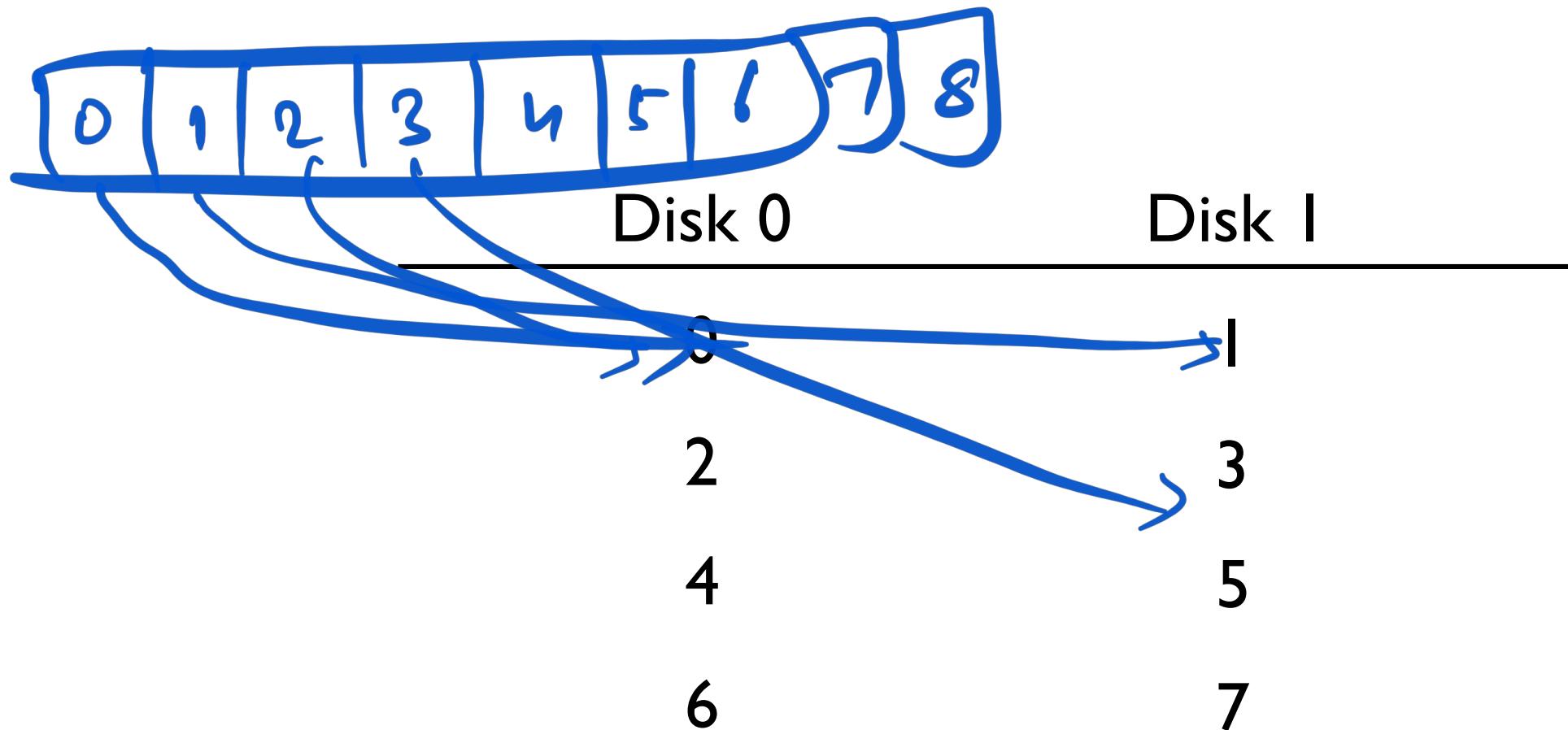
Symmetric perf. $\gamma = w$

RAID-0: Striping

Optimize for capacity; no redundancy!



RAID-0: Another View



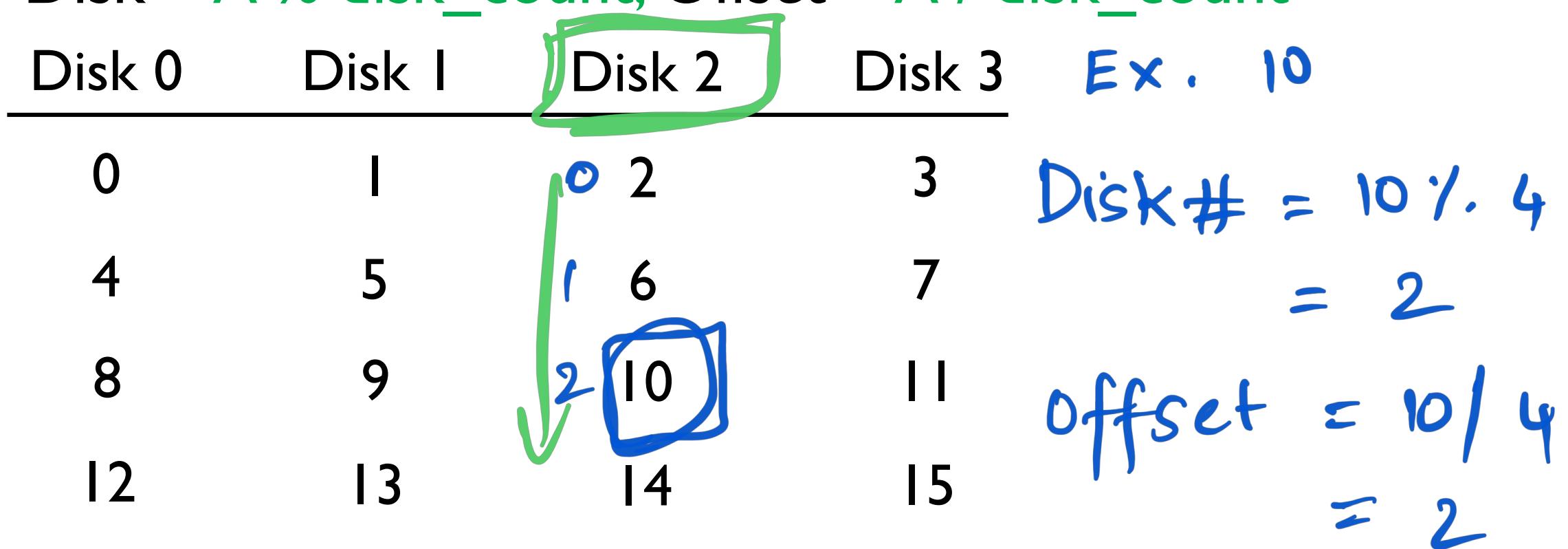
4 Disks

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

How to Map?

Given logical address A, find:

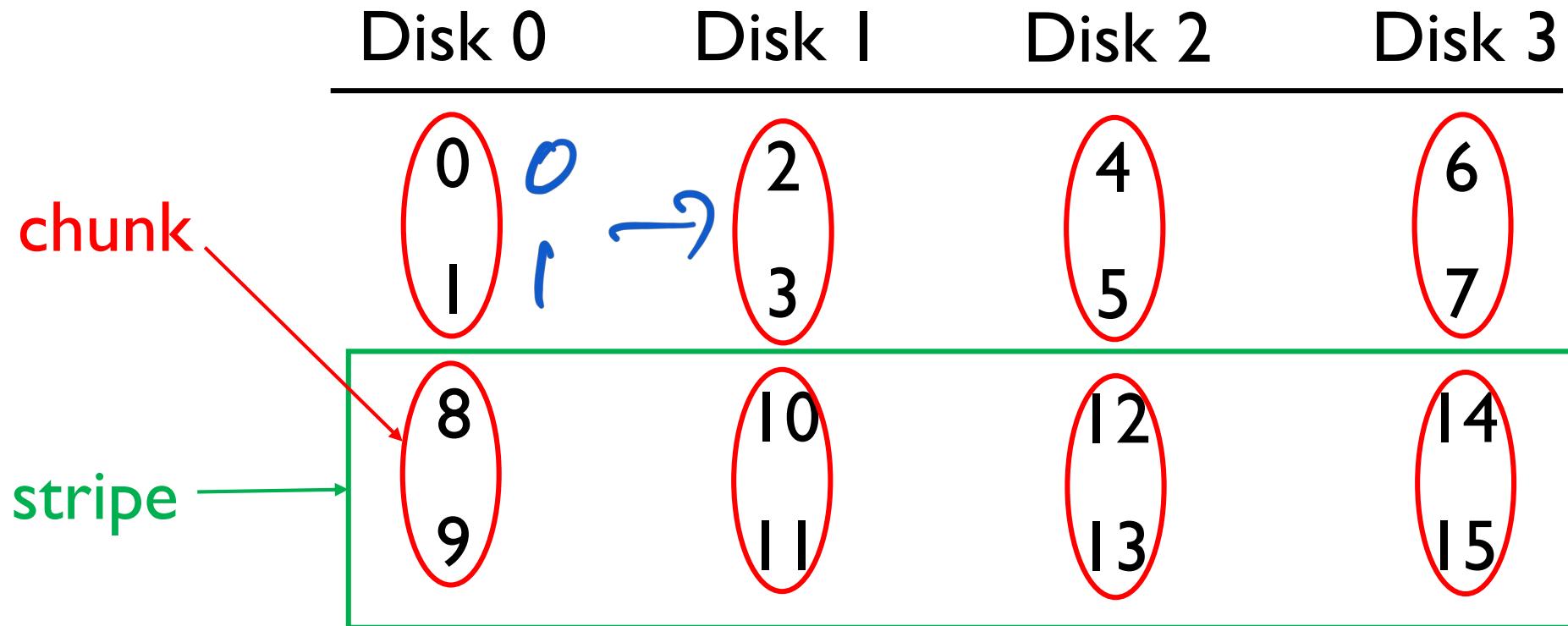
Disk = $A \% \text{disk_count}$, Offset = $A / \text{disk_count}$



Chunk Size = 1

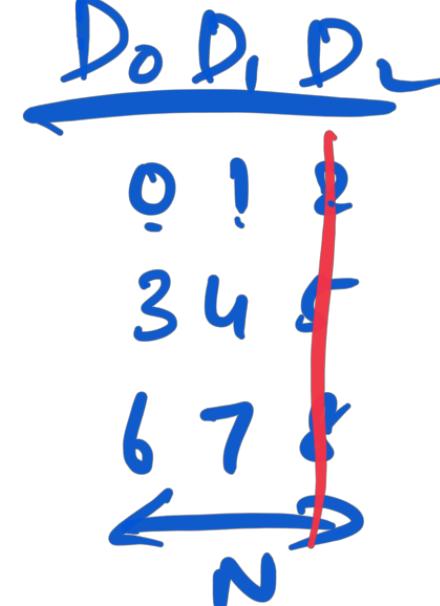
Disk 0	Disk 1	Disk 2	Disk 3
10	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

Chunk Size = 2



We'll assume chunk size of 1 for today. Sizes of 64KB are typical in deployment.

RAID-0: Analysis



What is the capacity?

N · C

How many disks can fail?

• 0 - Worse than single-disk in some cases?

Throughput (sequential, random)?

Seq / read — N · S

— write — N · S

Latency — D

rand < r — N · R

r : D, w : D

w — N · R

RAID-0: Analysis Results

What is the capacity?

$N * C$ *upper bound*

How many disks can fail?

0 *bad*

Throughput (sequential, random)?

$N*S, N*R$ *upper bound*

Latency

D

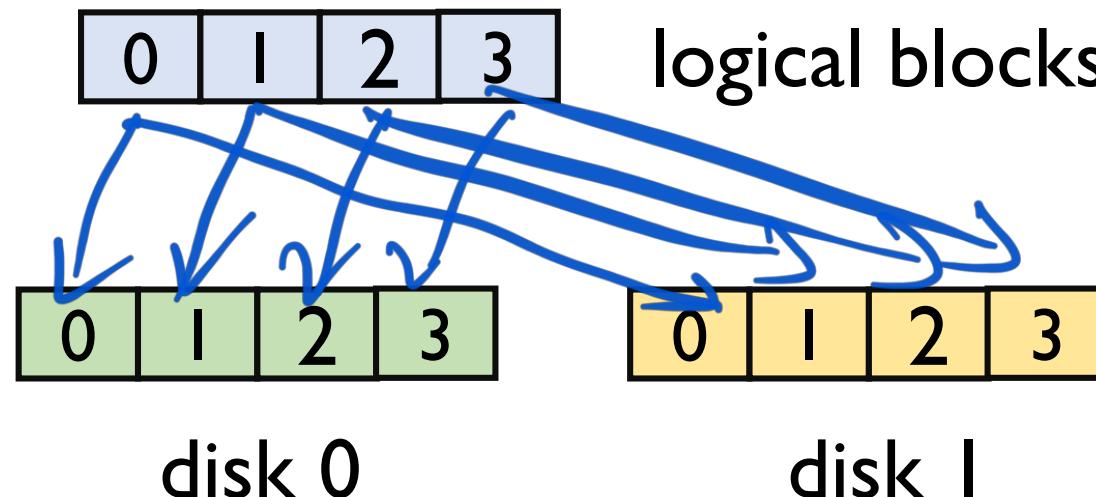
buying more disks improves throughput, but not latency!

RAID-1: Mirroring

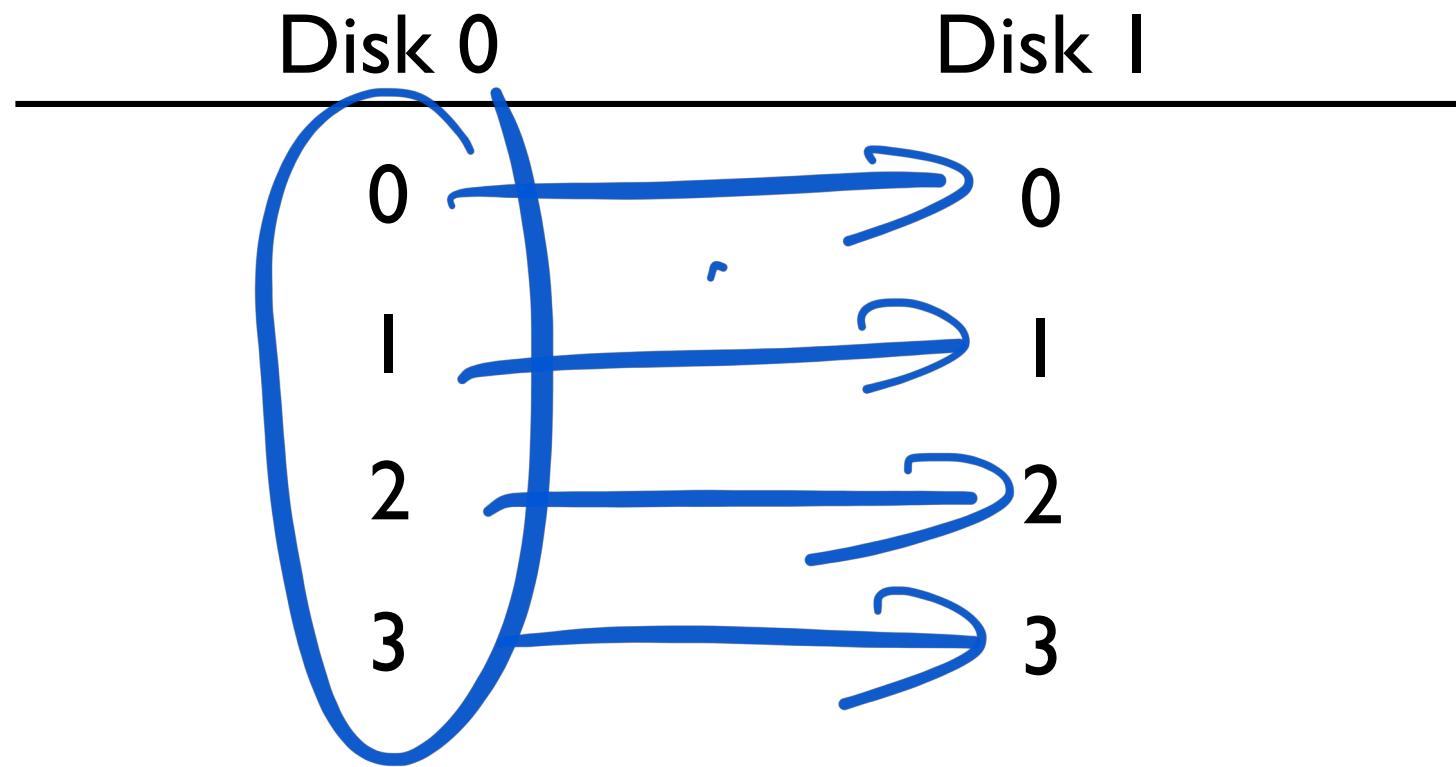
Problem w/
Raid - 0 :
Poor reliabi
lity.
Solution?

Keep two copies of all data

file w/ 4
blocks



2 Disks



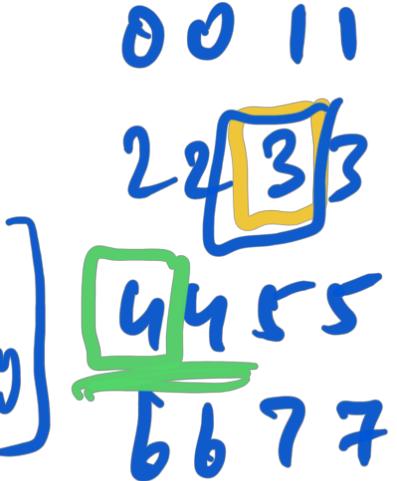
4 Disks

Disk 0	Disk 1	Disk 2	Disk 3
0	0	1	1
2	2	3	3
4	4	5	5
6	6	7	7

RAID-1: Analysis

What is the capacity?

$$\frac{N}{2} \cdot c \quad \left[\begin{array}{l} \text{only half} \\ \text{effective capacity} \end{array} \right]$$



How many disks can fail?

$\frac{N}{2}$, if lucky
1 for sure

Latency (read, write)?

are writes same as one disk? slower? faster?

r: D, w: D (2 wait but still
only D, because we can
do them in parallel)

RAID-1: Analysis Results

What is the capacity?

$N/2 * C$

How many disks can fail?

one (maybe $N/2$ if lucky)

Latency (read, write)?

read = D, write = D

RAID-1 Throughput

What is steady-state throughput for random reads?

$$N \cdot R \text{ MB/s}$$

random writes?

$$\frac{N \cdot R}{2}$$

Each logical write
become 2 physical writes

$$w(7)$$

7, 2, 0, 1, 3, 4, 5, 6

$$w(\text{mirror of } 7)$$

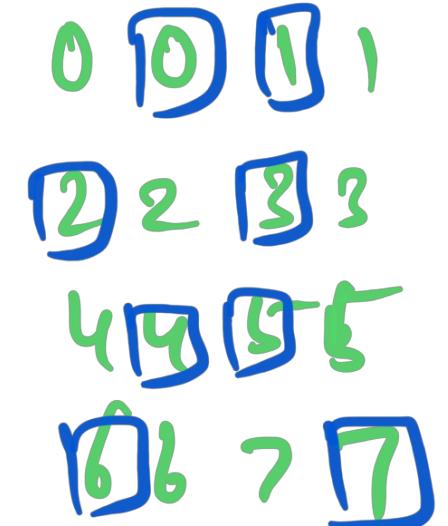
w(mirror
of 7)

0011

2233

4455

6677

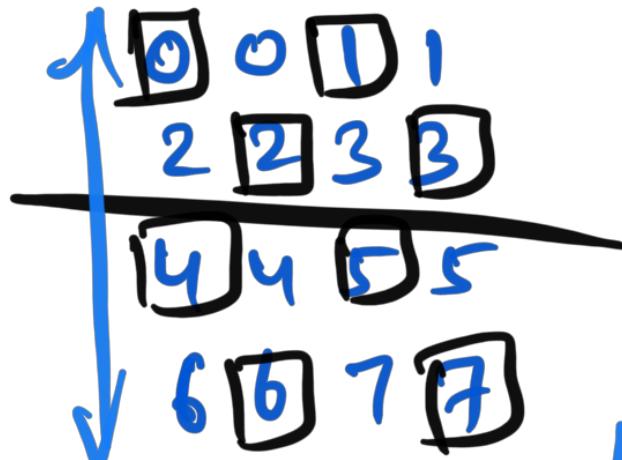


RAID-I Throughput

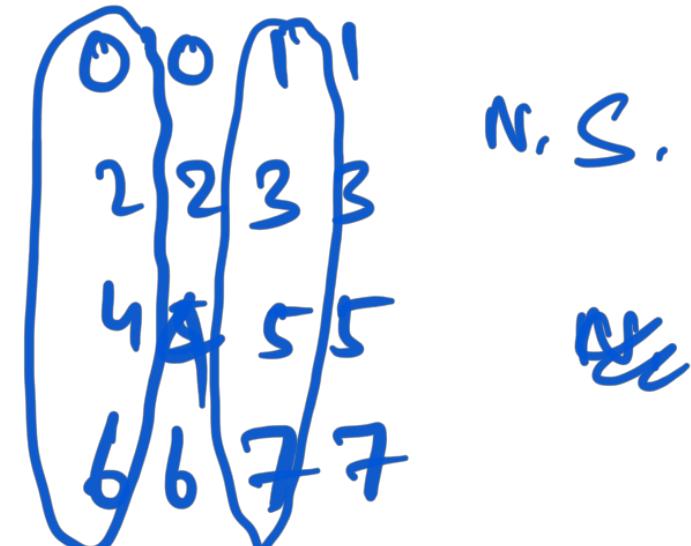
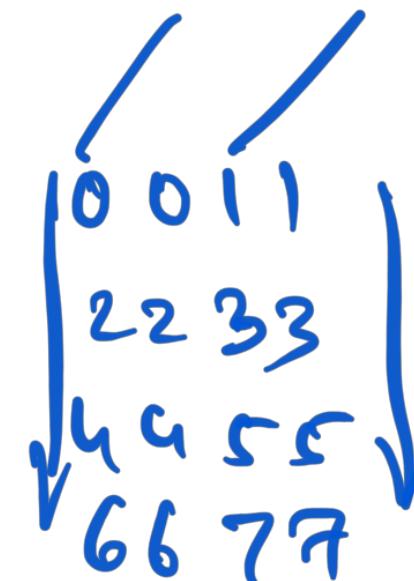
What is steady-state throughput for sequential writes?

$$\frac{N \cdot S}{2}$$

sequential reads?



$$\frac{N \cdot S}{2}$$



$$\frac{N}{2} \cdot S$$

RAID-1 Throughput Results

What is steady-state throughput for random reads?

$$N * R$$

same as raid-0

random writes?

$$N/2 * R$$

sequential writes?

$$N/2 * S$$

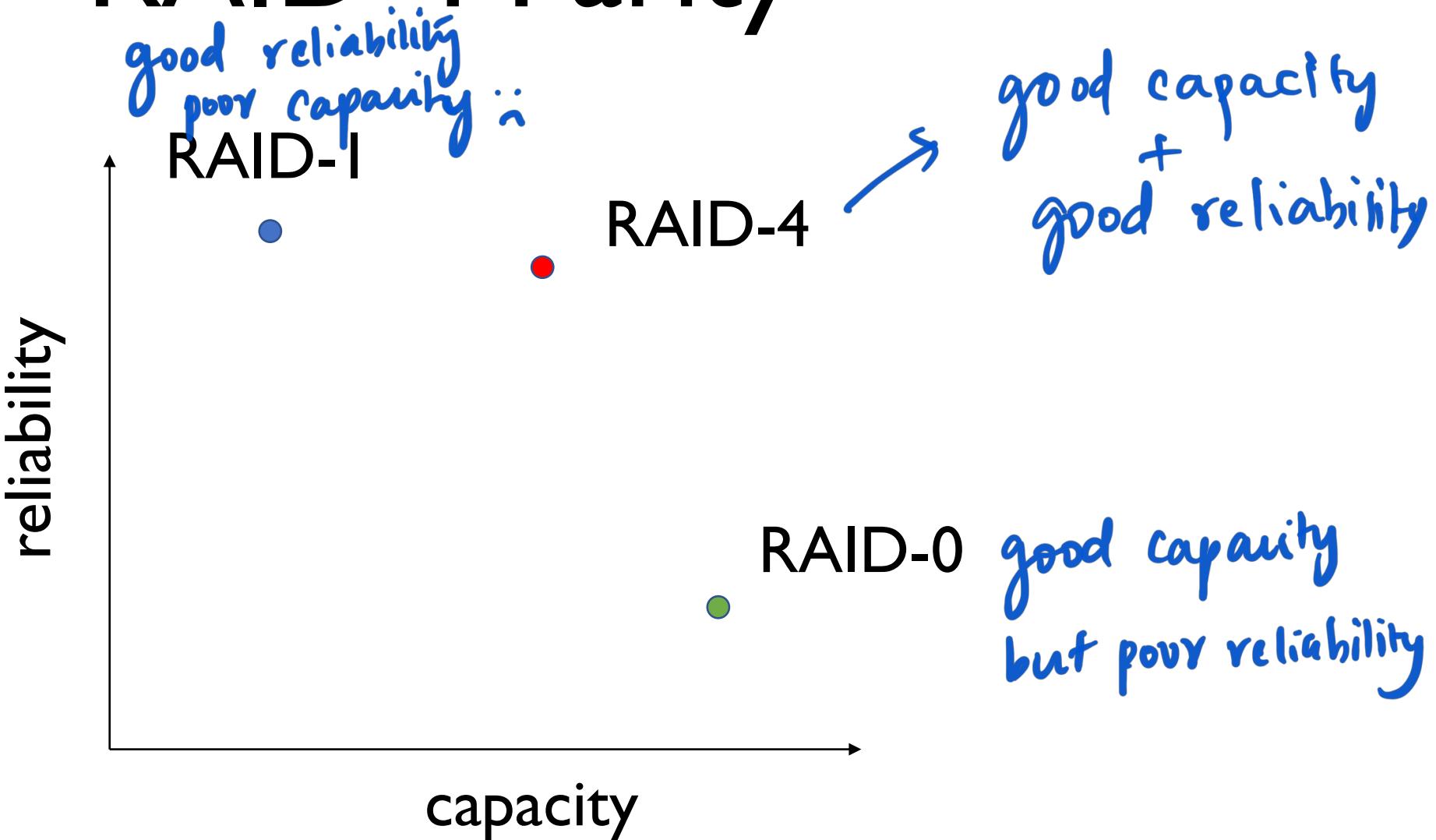
sequential reads?

$$N/2 * S$$



only half of raid-0

RAID-4 Parity



RAID-4 Strategy

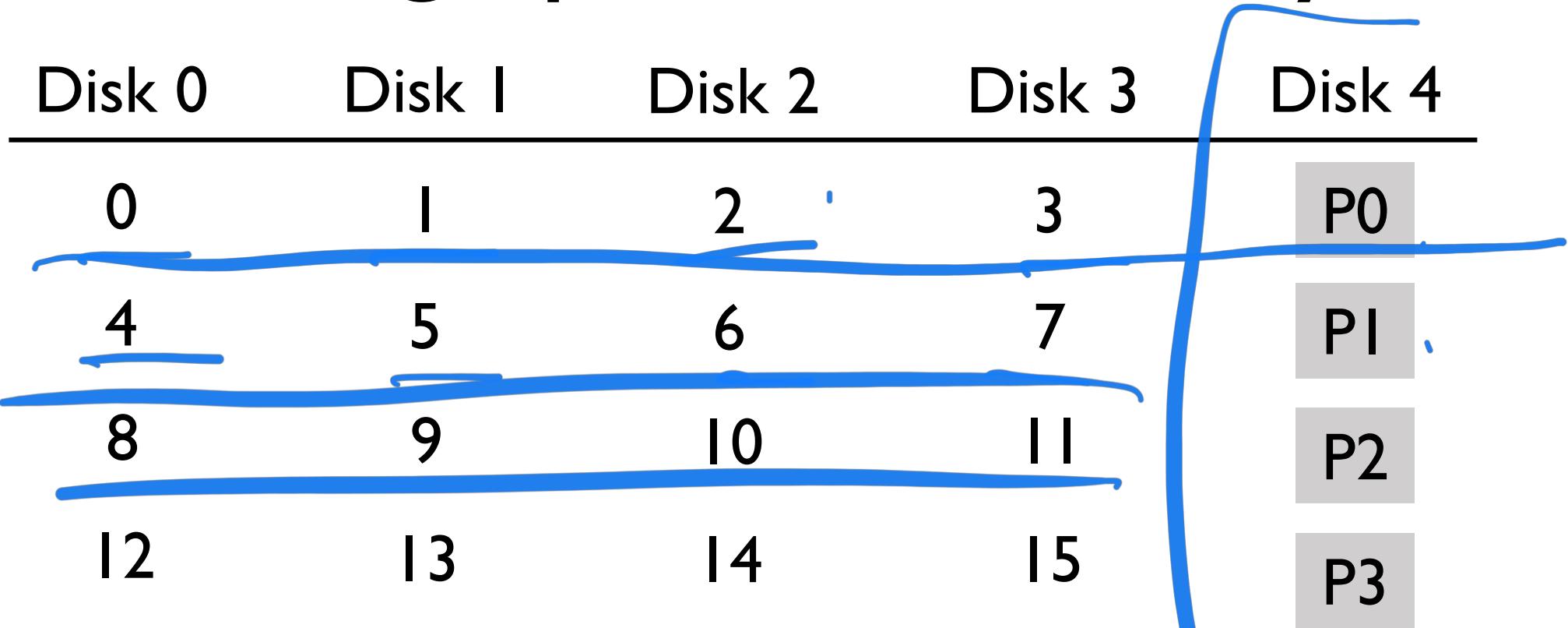
Stripe with parity

In algebra, if an equation has N variables, and $N-1$ are known, you can often solve for the unknown

Treat sectors across disks in a stripe as an equation

Data on bad disk is like an unknown in the equation

Saving Space with Parity



$$P_0 = 0 \oplus 1 \oplus 2 \oplus 3 \quad \text{XOR.}$$

XOR Parity

Disk 0	Disk 1	Disk 2	Disk 3	Disk 4
0	1	2	3	P0
4	5	6	7	P1
8	9	10	11	P2
12	13	14	15	P3

P0 = contents of 0 xor contents of 1 xor contents of 2 xor
contents of 3

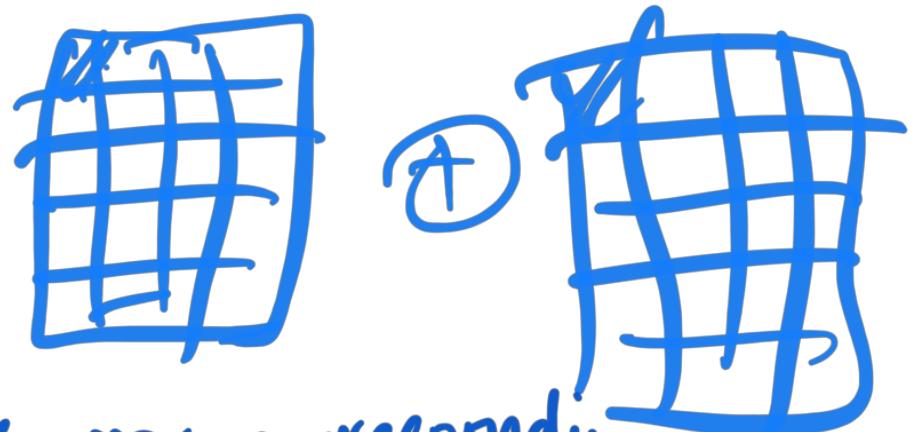
Parity Invariant

D_0	D_1	D_2	P
0	1	0	1
1	1	1	1
0	1	1	0
0	0	0	0

Pretend bits instead of

blocks

- for blocks, xor corresponding
bits of the blocks



1's to be even

Data Reconstruction in Raid-4



Access 2nd block on D₂
but D₂ has failed

Reconstruct from other
disks using parity invari-
-ant. We have three
ones - so missing must
be a $\frac{1}{=}$ (not 0)

RAID-4 Analysis

What is the capacity? $(N - 1) \cdot c$

How many disks can fail? 1

Latency (read, write)? need to read old data, old parity

D

→

2D

0	1	2	P
3	4	5	P
6	?	8	P
9	10	11	P

Updating Parity on Writes

Ex: 0 1 2 P

3 4 5 P
6 7 8 P
9 10 11 P

Write (8) $\delta_{\text{old}} \rightarrow \delta_{\text{new}}$

Approach 1: read 6, 7, we have δ_{new}

$$P_{\text{new}} = 6 \oplus 7 \oplus \delta_{\text{new}}$$

Approach 2: read $\delta_{\text{old}}, P_{\text{old}}$; calc P_{new}

if $\delta_{\text{new}} = P_{\text{old}}$:

$$P_{\text{new}} = P_{\text{old}}$$

else

$$P_{\text{new}} = (P_{\text{old}})'$$

Latency (write)?

RAID-4 Analysis Results

What is the capacity? $(N-I) * C$

How many disks can fail? one

Latency (read, write)? D, 2D

RAID-4 Throughput

What is steady-state throughput for –

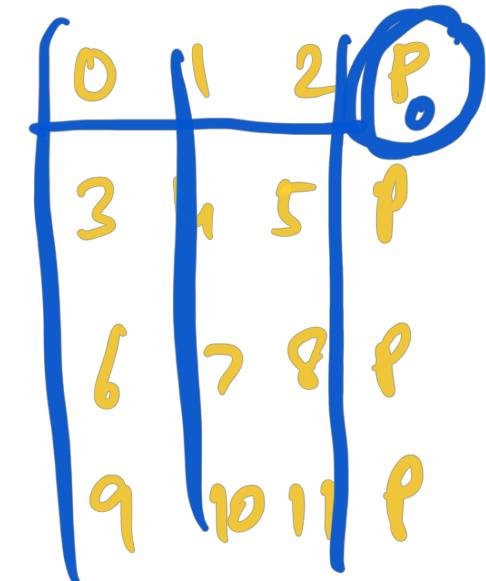
sequential reads? $(N - 1) \cdot s$

sequential writes (full-stripe write)?

$$P_0 = 0 \oplus 1 \oplus 2$$

$$(N - 1) \cdot s$$

$0, 1, 2, P_0 \Rightarrow$ calc. from memory



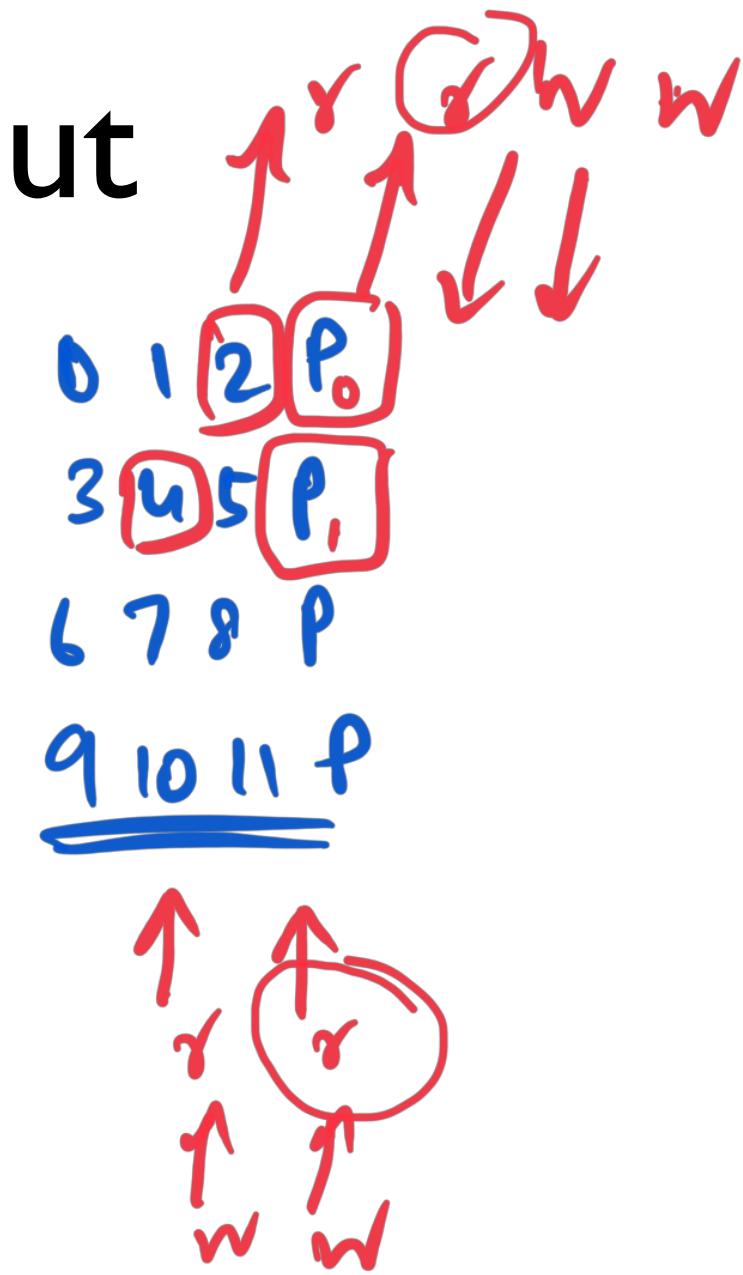
RAID-4 Throughput

What is steady-state throughput for –

random reads? $(N - 1) \cdot R$

random writes? $\frac{R}{2}$ (this is very
bad!)

random write perf. in RAID-4 is
bottlenecked by the parity disk



RAID-4 Small-write Problem

Performance of random writes?

$$R/2 \binom{\text{notice: no}}{N!}$$

RAID-4 Throughput Results

What is steady-state throughput for –

sequential reads? $(N-1) * S$

sequential writes? $(N-1) * S$

random reads? $(N-1) * R$

random writes? $R/2$

RAID-5: Rotated Parity

Disk 0	Disk 1	Disk 2	Disk 3	Disk 4
-	-	-	-	P0
-	-	-	P1	-
-	-	P2	-	-
-	P3	-	-	-
P4	-	-	-	-

RAID-5 Analysis

What is the capacity?

How many disks can fail?

Latency (read, write)?

RAID-5 Analysis Results

What is the capacity? $(N-I) * C$

How many disks can fail? one

Latency (read, write)? D, 2D

RAID-5 Throughput

What is steady-state throughput for –

sequential reads?

sequential writes?

random reads?

random writes?

RAID-5 Throughput Results

What is steady-state throughput for –

sequential reads? $(N-1) * S$

sequential writes? $(N-1) * S$

random reads? $(N) * R$

random writes? $N/4 * R$

RAID Summary

	Reliability	Capacity
RAID-0	0	$C*N$
RAID-1		$C*N/2$
RAID-4		$C*(N-1)$
RAID-5		$C*(N-1)$

RAID Summary

	Read latency	Write latency
RAID-0	D	D
RAID-1	D	D
RAID-4	D	2D
RAID-5	D	2D

RAID Summary

	Seq read	Seq write	Rand read	Rand write
RAID-0	$N*S$	$N*S$	$N*R$	$N*R$
RAID-1	$N/2 * S$	$N/2 * S$	$N*R$	$N/2 * R$
RAID-4	$(N-1)*S$	$(N-1)*S$	$(N-1)*R$	$R/2$
RAID-5	$(N-1)*S$	$(N-1)*S$	$N*R$	$N/4*R$

RAID-0 is always fastest and has best capacity (but at cost of reliability)

RAID-5 is strictly better than RAID-4

RAID-5 better than RAID-1 for sequential

RAID-1 better than RAID-4 for random write

Summary

RAID: a faster, larger, more reliable disk system

One logical disk built from many physical disk

Different mapping and redundancy schemes

Present different trade-offs