# 537 : Persistence

=> **Hard Drive** (last time)

=> cheap

=> (slow)    ops : milliseconds
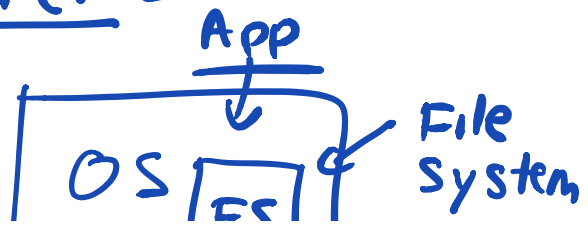
=> large ops, } fastest mode
   sequential    (~ 100 MB/s)

=> small ops, } slow mode
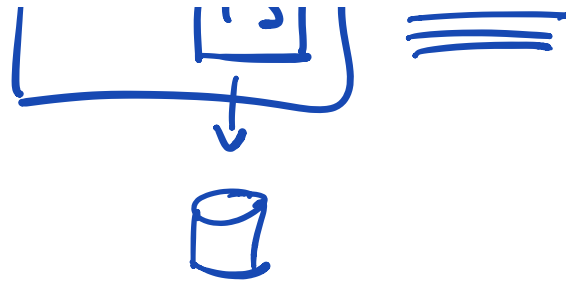   random       ~ 1 MB/s

**Today :**

=> (RAID) { Redundant
            Arrays of
            Independent
            Drives }

=> Intro to
   File Systems

App
↓
OS [FS]    File System

## RAID : many drives, not one

=> Why?

=> (Performance)
(Bandwidth,
IOPS )
-> I/Os per second

=> Reliability
many copies :
tolerate disk failure

=> Capacity

Interface:

RAID: "looks like" a disk

Hardware
RAID

PC → ← → CPU | Mem

interface:
array of blocks,
read/write

Fault model: Simple

How do hard drives fail?

⟹ #drives
→ working
→ entirely not (broken)
easy to detect

RAID: {Levels}

Levels 0, 1, 4, 5 → parity-
based

striping      mirroring
(no           (n copies)
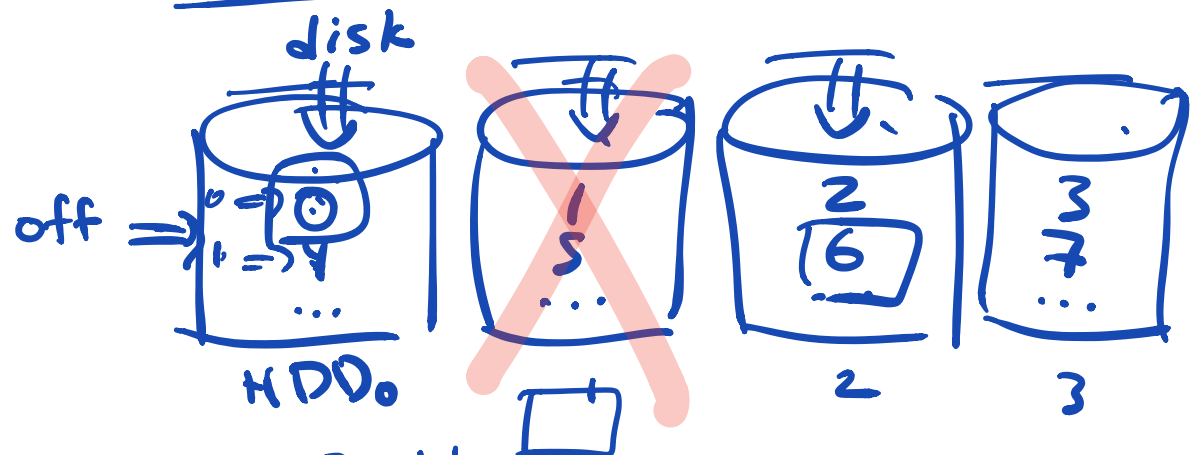redundancy

Metrics:
→ Perf
→ Reliability

# RAID - 0 : Striping

0 .... N-1          array of
read/write⇒blocks

disk
#

off ⇒

HDD$_0$          HDD$_0$          2          3
                                  6          7

## Mapping Problem:

interface address ⇒
(disk, offset)

Calculation:
disk: $\dfrac{address}{} \% \dfrac{num\ disks}{}$

offset: $address\ /\ num\ disks$

## Striping:

Reliability:     no redundancy
⇒ tolerate
0 failures
⇒ { "bad" } worse

Capacity: $\Rightarrow$ "good"

N disks,
each have D bytes

$$\boxed{\Rightarrow N \cdot D}$$

Performance : RAID-0

(Large)
Seq. Read     $\underline{N \cdot S}$ MB/s $\Big\}$ parallel
Seq. Write    $\underline{N \cdot S}$ MB/s

N disks :          Single Disk:

$\Rightarrow \underline{S}$ $\underline{MB/s}$
(seq I/o
bandwidth)

(Small)
Random Read: $\underline{N \cdot R}$
Random Writes: $N \cdot R$
$\Rightarrow R$ MB/s
(rand I/o
Bandwidth)

$R \ll S$

(e.g. $\underline{1 \, MB/s} \ll \underline{100 \, MBs}$)

RAID:1   (Mirroring)

Read: **2**   write: **3**   in parallel

Disk $_0$   1   2   3



## 2 copies (physical)

⇒ Capacity:  N Drives,
                D bytes/drive

$$\Rightarrow \boxed{\frac{N \cdot D}{2}}$$

⇒ Reliability:     Fault Model:
                    whole drive
                    failure

lucky: $\frac{N}{2}$

paranoid: $\boxed{1}$ ( then replace)

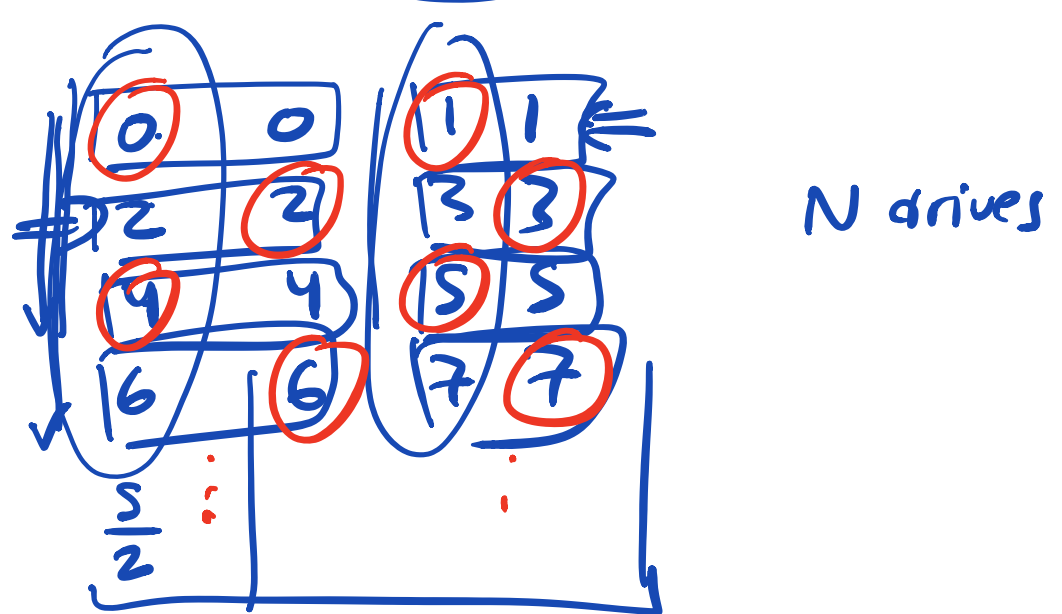⇒ Performance:

Random I/O: Bandwidth    Latency:
                          writes ② a little
                                  longer
Reads ( N · R MB/s )  than

Writes: $\frac{N}{2} \cdot R$ MB/s          write: to
                                   7

                                   reads ⇒ 1

single disk: R MB/s
             (~ 1 MB/s)

Seq R/w:     N · ~S   2·S

Reads: $\frac{N}{2} \cdot S$ MB/s ?        S MB/s

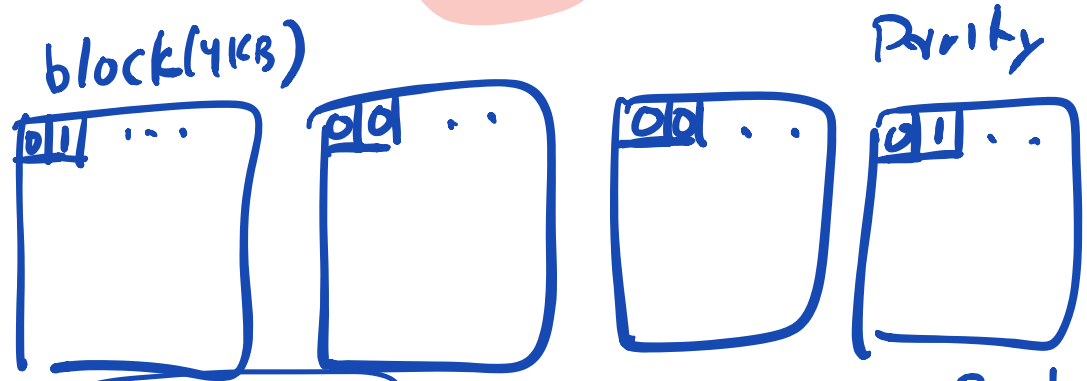Writes: $\frac{N}{2} \cdot S$ MB/s (~ 100 MB/s)



N drives

$\frac{S}{2}$

RAID: Parity - based (4/5)

Parity: XOR

## Data Disks          Parity Disk

$D_0$   $D_1$   $D_3$              $P$

block { [1]   [0]   [0]  ↙        [1]

pretend:
it's a **bit**

invariant:
for each **row**:
even # of 1's

(1)   (0)   (1.)  ↙   (0)
 |     |     |    ↙    |
 O     O     O   ↙     O

block(4KB)                              Parity

[0|1| ...]   [0|0| ..]   [0|0| ..]   [0|1| ..]

[0 , 1 , 2] ⇒ $P_0$        RAID-4:  Parity
                                     Disk

0        1        2        [$P_0$]

**Reliability:** 1 failure

**Capacity:** $(N-1) \cdot D$ MB

**Performance:**

|  | Seq $S$ MB/s | Random $R$ MB/s |
|---|---|---|

Random Reads: $(N-1) \cdot R$ MB/s

Random Writes: $\frac{R}{2}$ MB/s

Seq Reads: $(N-1) \cdot S$ MB/s

Seq writes: $(N-1) \cdot S$ MB/s

**Admin:**

=> 4b due ~~next~~ monday

=> 4a tomorrow      next next
   (concurrency)    Monday
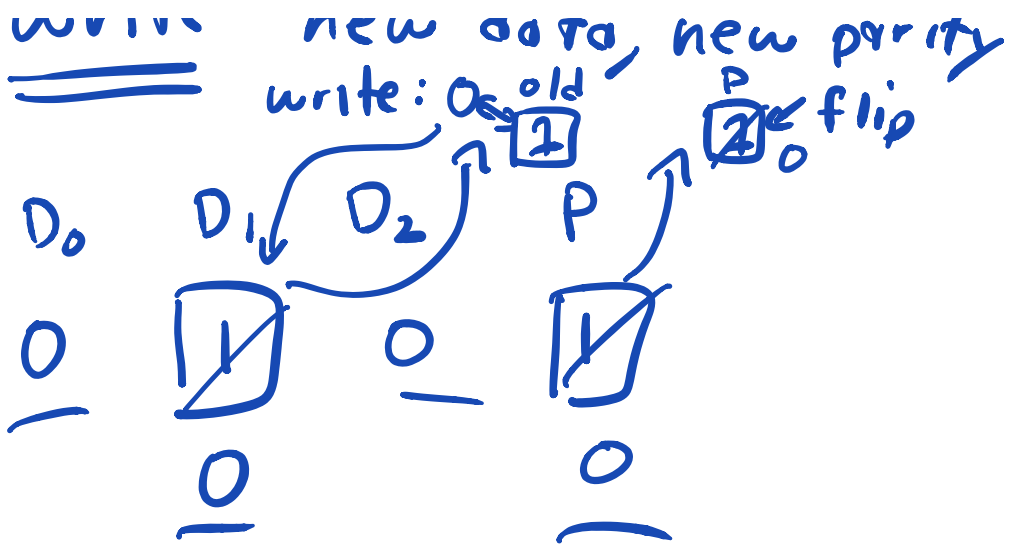
=> <u>5</u> ....   file systems  $\{$(Map Reduce$)\}$

=>"final" ⟸

---

RAID-4    write: $P_4'$!    write:

| 0 | 1 | 2 | $P_0$ |
|---|---|---|---|
| 3 | 4. | 5 | $P_1$ |
| 6 | 7 | 8 | $P_2$ |
| 9 | 10. | 11 | $P_3$ |

N-1

---

<u>Read</u> old data + <u>Read</u> old parity:

compare old data, new data  => when they differ, flip in bit old parity

write new d...

write    new data, new parity
write: $O$ old    $P$ flip

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

0        $\boxed{1}$    0    $\boxed{1}$

         0             0

**Latency:**    single RAID-4
                write ?
                time

        latency of
        read or write : $\boxed{T}$

⇒ how long RAID-4 write?

I/O
⇒ ( read old data , old parity )

[Compute]
⇒ compute new parity
        ( old , old
          data  parity )
              new,
              data
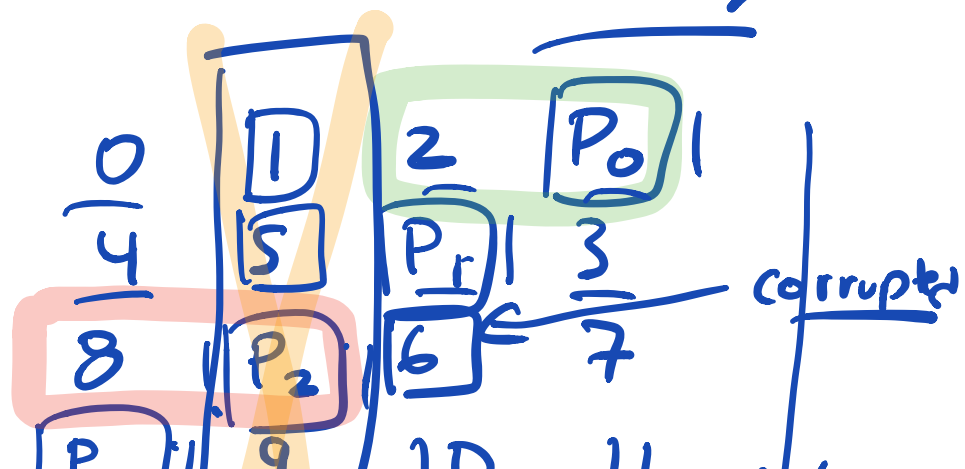
⇒ write new data, new parity

⇒ ~2T  RAID 4 single
write

RAID-4:  many small
writes

⇒ Parity disk
bottleneck

"small write problem"

R MB/s ⇒ random I/O
bandwidth

---

RAID-5:  Rotated
Parity

| | 0 | 1 | 2 | $P_0$ 1 |
|---|---|---|---|---|
| | 4 | 5 | $P_1$ 1 | 3 |
| | 8 | $P_2$ | 6 | 7 | — corrupted |
| | P | 9 | 10 | 11 |

$$\underbrace{\boxed{3 \mid\mid 1} \quad \cdots \quad \cdots}_{n}$$

## Capacity, Reliability : Same as RAID-4

## Performance:

Rand Read: $N \cdot R$ MB/s

Rand Writes: $\frac{N}{4} \cdot R$ MB/s

Seq Read : $\dfrac{(N-1) \cdot S}{\phantom{x}}$

Seq write : $\dfrac{(N-1) \cdot S}{\phantom{x}}$

RAID-1: mirroring

$\frac{N}{2} \cdot R$ MB/s

RAID-5: 4 I/Os / but Capacity, Seq I/O, logical write

vs

Mirroring : small write perf (high capacity cost)

## RAID:

faster, larger, more reliable

disk

$\Rightarrow$ 0,1,4,5 (6) $\Rightarrow$ 2 parity disks

$\Rightarrow$ checksums : detect/recover

from corruption

$\Rightarrow$