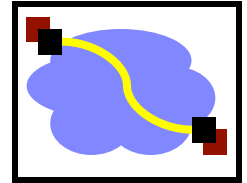# 416 Distributed Systems

## RAID, Feb 8 2017
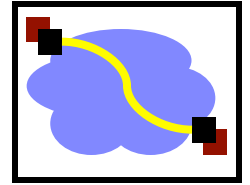
Thanks to Greg Ganger and Remzi Arapaci-Dusseau
for slides

# Replacement Rates
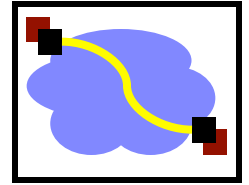
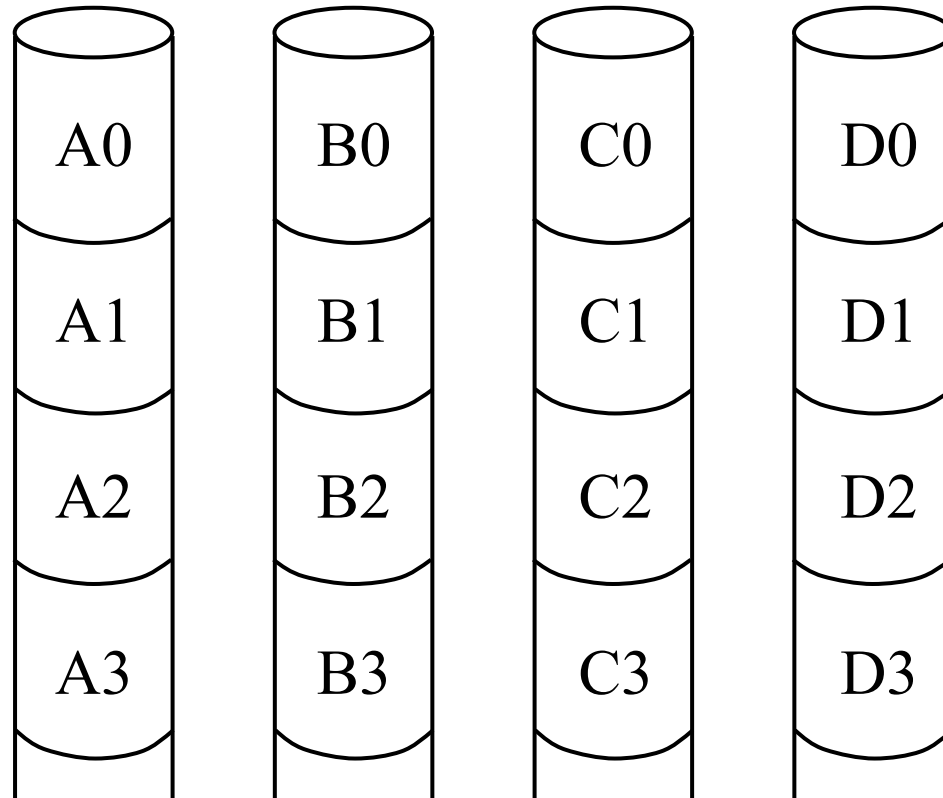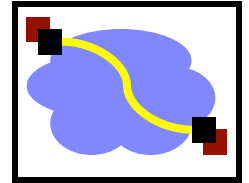| HPC1 | | COM1 | | COM2 | |
|---|---|---|---|---|---|
| **Component** | **%** | **Component** | **%** | **Component** | **%** |
| Hard drive | 30.6 | Power supply | 34.8 | Hard drive | 49.1 |
| Memory | 28.5 | Memory | 20.1 | Motherboard | 23.4 |
| Misc/Unk | 14.4 | Hard drive | 18.1 | Power supply | 10.1 |
| CPU | 12.4 | Case | 11.4 | RAID card | 4.1 |
| motherboard | 4.9 | Fan | 8 | Memory | 3.4 |
| Controller | 2.9 | CPU | 2 | SCSI cable | 2.2 |
| QSW | 1.7 | SCSI Board | 0.6 | Fan | 2.2 |
| Power supply | 1.6 | NIC Card | 1.2 | CPU | 2.2 |
| MLB | 1 | LV Pwr Board | 0.6 | CD-ROM | 0.6 |
| SCSI BP | 0.3 | CPU heatsink | 0.6 | Raid Controller | 0.6 |

# Outline

- Using multiple disks
  - Why have multiple disks?
  - problem and approaches

- RAID levels and performance
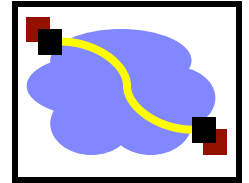
# Motivation: Why use multiple disks?

- Capacity
  - More disks allows us to store more data
- Performance
  - Access multiple disks in parallel
  - Each disk can be working on independent read or write
  - Overlap seek and rotational positioning time for all
- Reliability
  - Recover from disk (or single sector) failures
  - Will need to store multiple copies of data to recover

- So, what is the simplest arrangement?
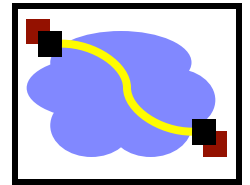
# Just a bunch of disks (JBOD)



| A0 | B0 | C0 | D0 |
| A1 | B1 | C1 | D1 |
| A2 | B2 | C2 | D2 |
| A3 | B3 | C3 | D3 |

- Yes, it's a goofy name
  - industry really does sell "JBOD enclosures"

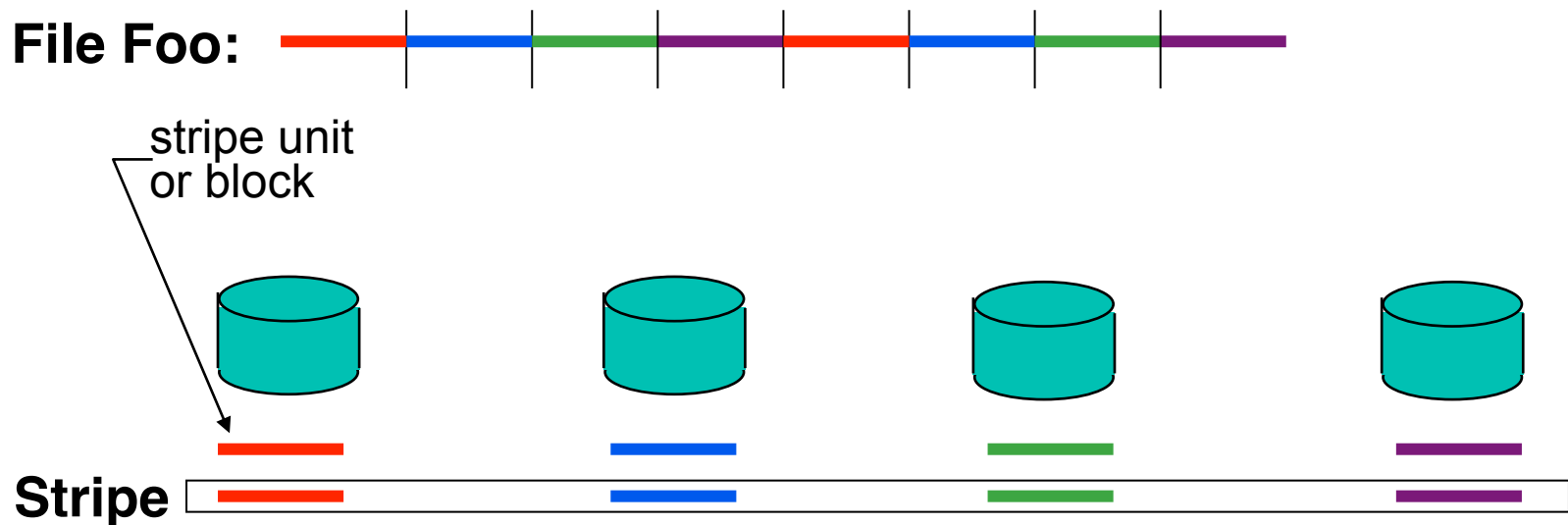# Disk Subsystem Load Balancing

- I/O requests are almost never evenly distributed
  - Some data is requested more than other data
  - Depends on the apps, usage, time, ...
- What is the right data-to-disk assignment policy?
  - Common approach: Fixed data placement
    - Your data is on disk X, period!
    - For good reasons too: you bought it or you're paying more...
  - Fancy: Dynamic data placement
    - If some of your files are accessed a lot, the admin(or even system) may separate the "hot" files across multiple disks
      - In this scenario, entire files systems (or even files) are manually moved by the system admin to specific disks
  - Alternative: Disk striping
    - Stripe all of the data across all of the disks
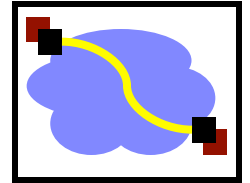
# Disk Striping

- Interleave data across multiple disks
  - Large file streaming can enjoy parallel transfers
  - High throughput requests can enjoy thorough load balancing
    - If blocks of hot files equally likely on all disks (really?)

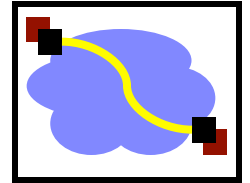**File Foo:**

stripe unit
or block

**Stripe**

# Disk striping details

- How disk striping works
  - Break up total space into fixed-size stripe units
  - Distribute the stripe units among disks in round-robin
  - Compute location of block #B as follows
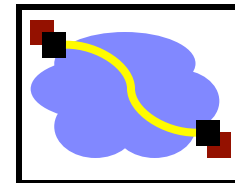    - disk# = B%N (%=modulo,N = #ofdisks)
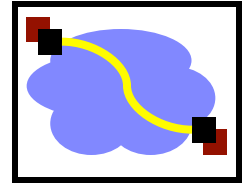
# Now, What If A Disk Fails?

- In a JBOD (independent disk) system
  - one or more file systems lost
- In a striped system
  - a part of each file system lost

- Backups can help, but
  - backing up takes time and effort
  - backup doesn't help recover data lost during that day
    - Any data loss is a big deal to a bank or stock exchange
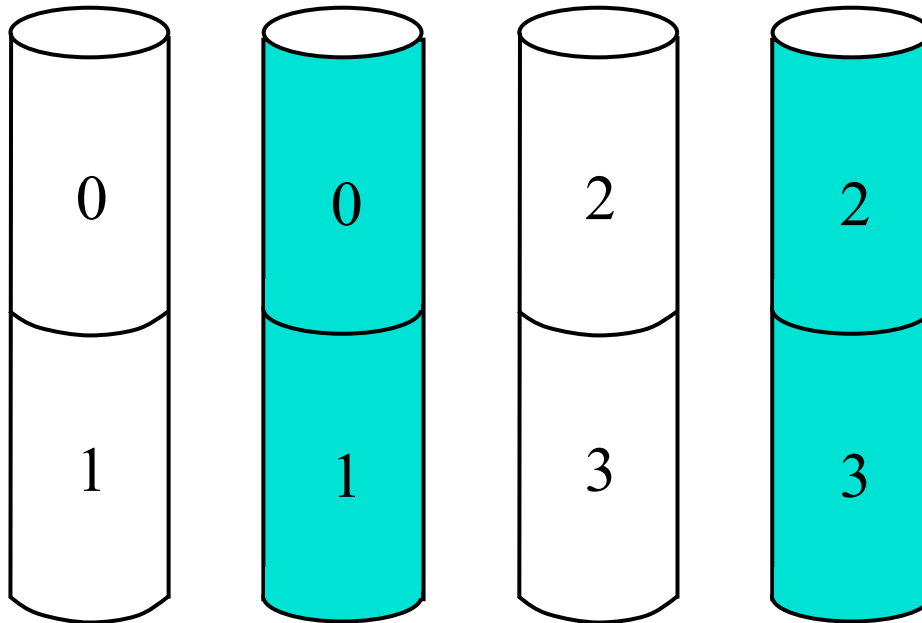
# Tolerating and masking disk failures

- If a disk fails, it's data is gone
  - may be recoverable, but may not be
- To keep operating in face of failure
  - must have some kind of data redundancy
- Common forms of data redundancy
  - replication
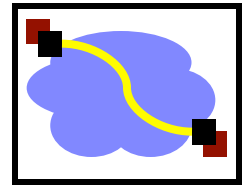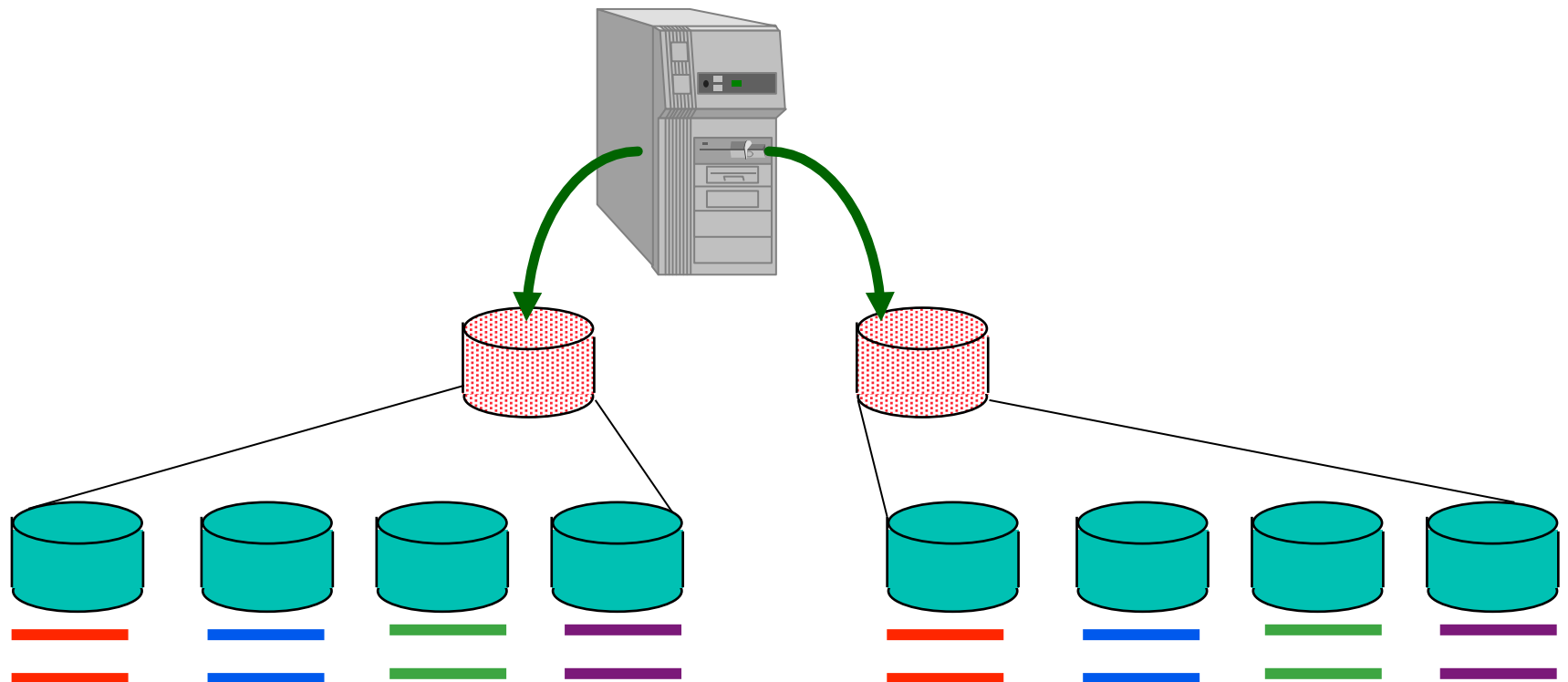  - error-correcting codes

# Redundancy via replicas

- Two (or more) copies
  - mirroring, shadowing, duplexing, etc.
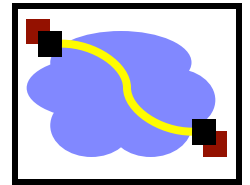- Write both, read either
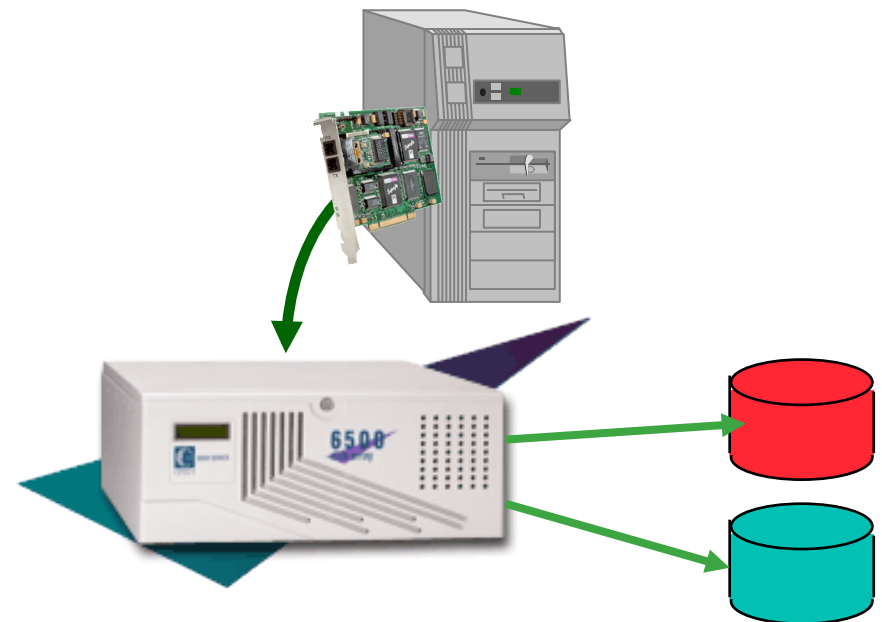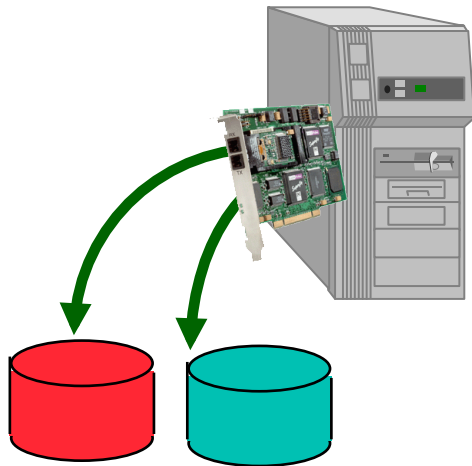
# Mirroring & Striping

- Mirror to 2 virtual drives, where each virtual drive is really a set of striped drives
  - Provides reliability of mirroring
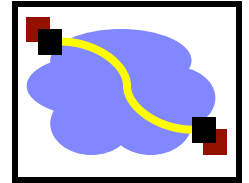  - Provides striping for performance (with write update costs)

# Implementing Disk Mirroring

- Mirroring can be done in either software or hardware
- Software solutions are available in most OS's
- Hardware solutions
  - Could be done in Host Bus Adaptor(s)
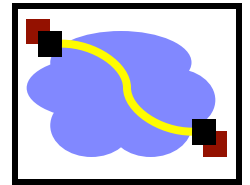  - Could be done in Disk Array Controller
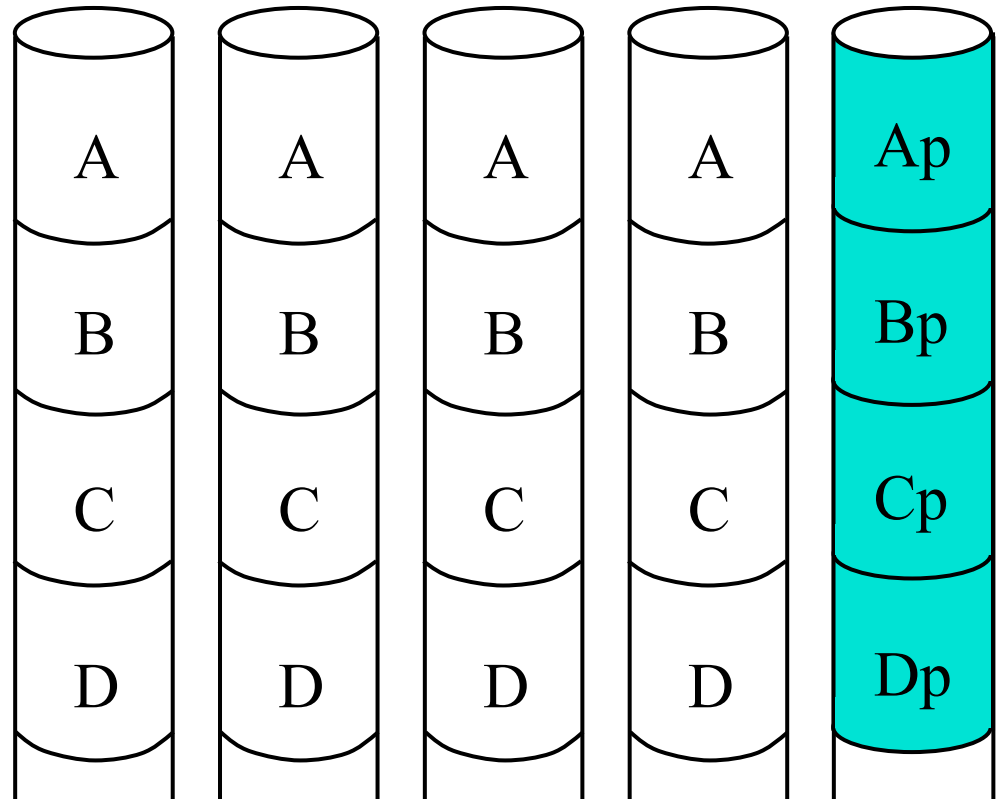
# Lower Cost Data Redundancy

- Single failure protecting codes
  - general single-error-correcting code is overkill
    - General code finds error and fixes it
- Disk failures are self-identifying (a.k.a. erasures)
  - Don't have to find the error
- Parity is single-disk-failure-correcting code
  - recall that parity is computed via XOR
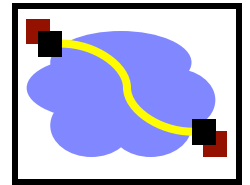  - it's like the low bit of the sum
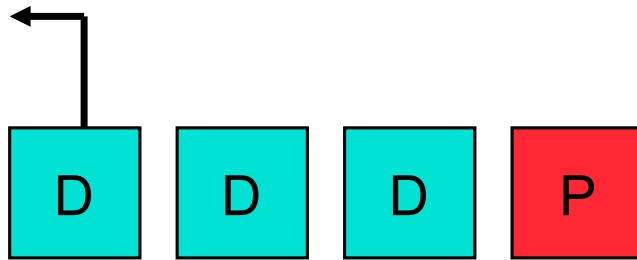
# Simplest approach: Parity Disk

- One extra disk
- All writes update parity disk
  - Potential bottleneck

  - (different data in different As, Bs, Cs, Ds)
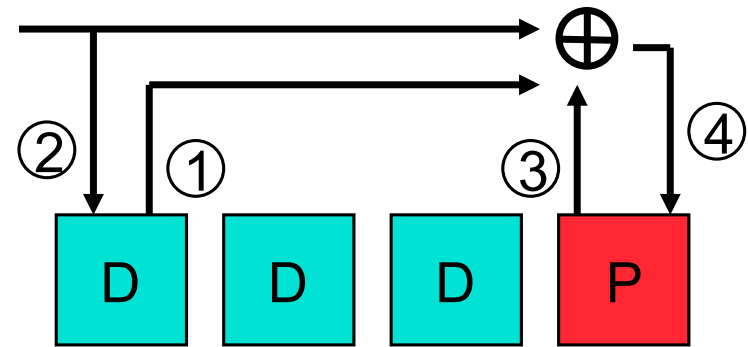  - (Ap contains parity for all As)

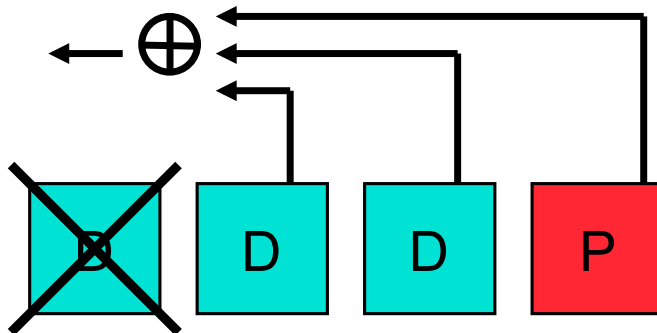| A | A | A | A | Ap |
| B | B | B | B | Bp |
| C | C | C | C | Cp |
| D | D | D | D | Dp |

# Updating and using the parity

**Fault-Free Read**

D  D  D  P

**Fault-Free Write**

② ① ③ ④

D  D  D  P

**Degraded Read**

⊕

D  D  D  P

**Degraded Write**

⊕
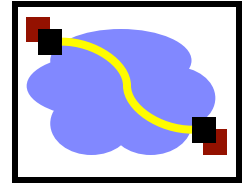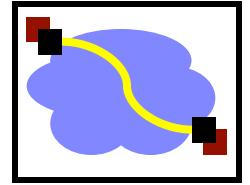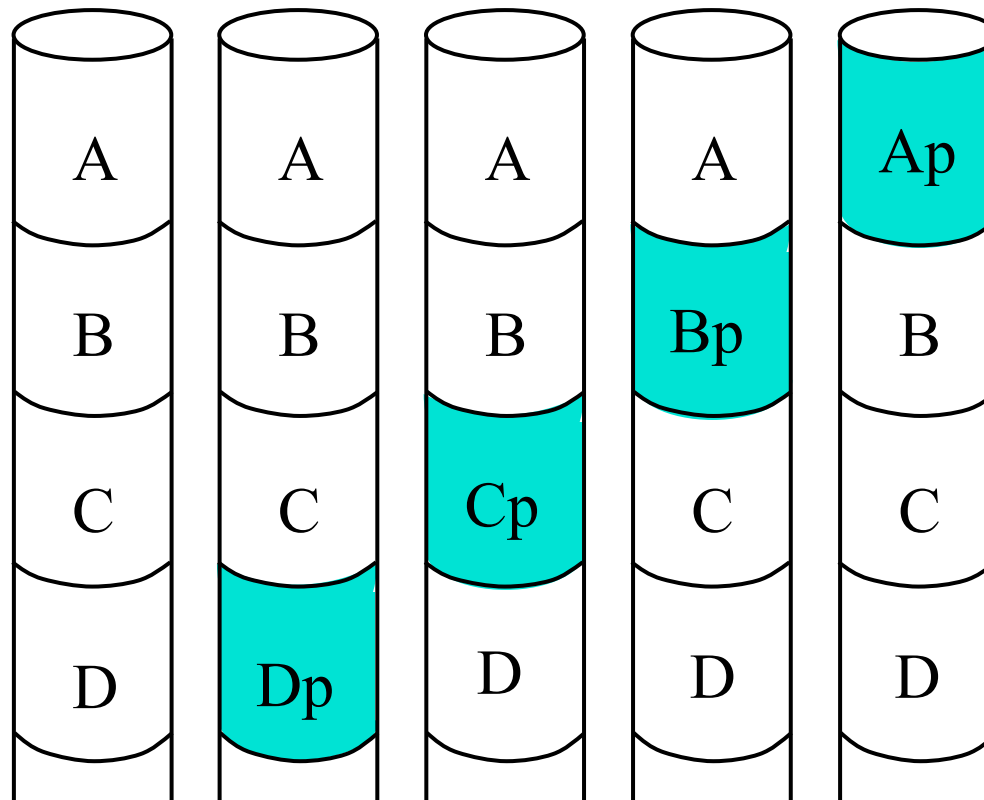
D  D  D  P

# The parity disk bottleneck

- Reads go only to the data disks
  - But, hopefully load balanced across the disks


- All writes go to the parity disk
  - And, worse, usually result in Read-Modify-Write sequence
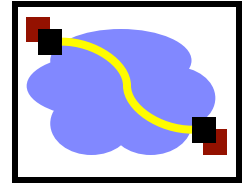  - So, parity disk can easily be a bottleneck

# Solution: Striping the Parity

- Removes parity disk bottleneck

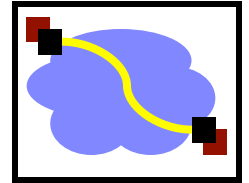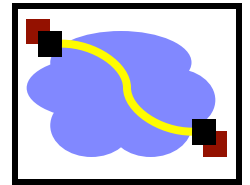| | | | | |
|---|---|---|---|---|
| A | A | A | A | Ap |
| B | B | B | Bp | B |
| C | C | Cp | C | C |
| D | Dp | D | D | D |

# Outline

- Using multiple disks
  - Why have multiple disks?
  - problem and approaches

- RAID levels and performance

# RAID Taxonomy
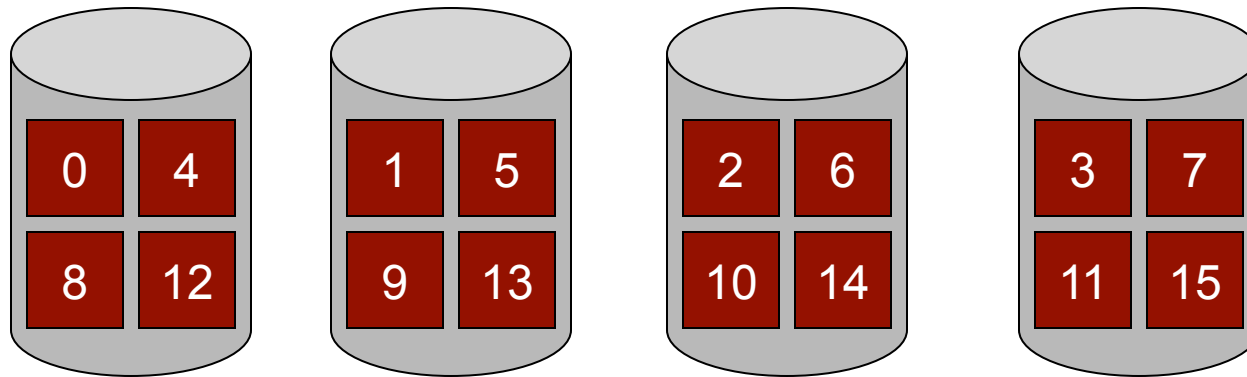
- Redundant Array of Inexpensive Independent Disks
  - Constructed by UC-Berkeley researchers in late 80s (Garth)
- RAID 0 – Coarse-grained Striping with no redundancy
- RAID 1 – Mirroring of independent disks
- RAID 2 – Fine-grained data striping plus Hamming code disks
  - Uses Hamming codes to detect and correct multiple errors
  - Originally implemented when drives didn't always detect errors
  - Not used in real systems
- RAID 3 – Fine-grained data striping plus parity disk
- RAID 4  – Coarse-grained data striping plus parity disk
- RAID 5  – Coarse-grained data striping plus striped parity
- RAID 6 – Coarse-grained data striping plus 2 striped codes

# RAID-0: Striping

- ## Stripe blocks across disks in a "chunk" size
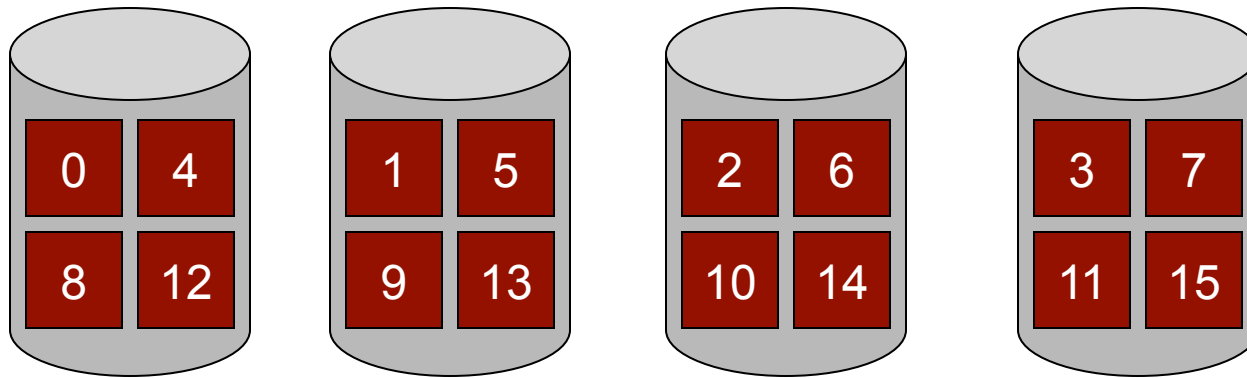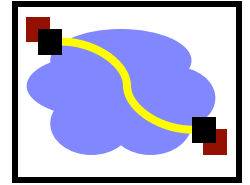  - ### How to pick a reasonable chunk size?
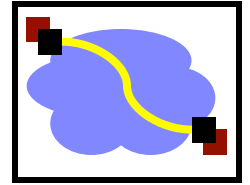


How to calculate where chunk # lives?

Disk #:

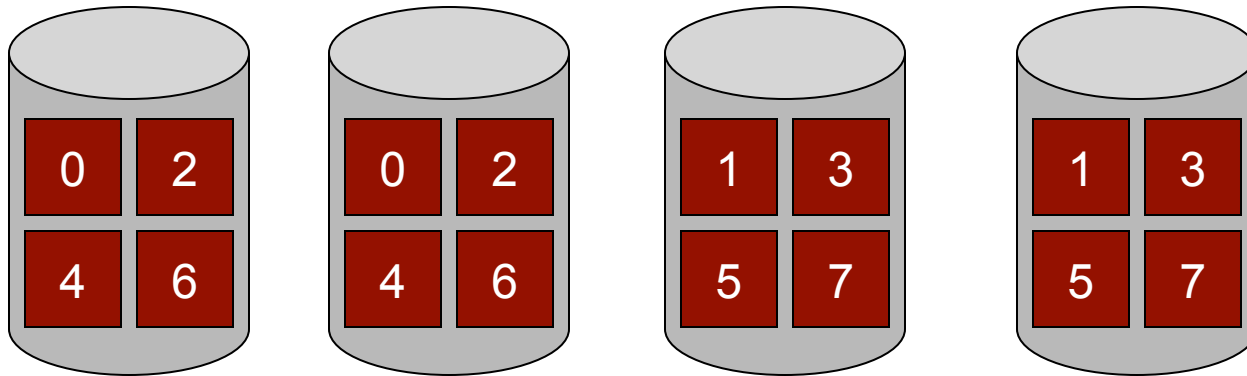Offset within disk:

# RAID-0: Striping



- Evaluate for D disks

- Performance: How much faster than 1 disk?

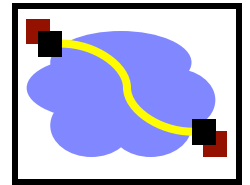- Reliability: More or less reliable than 1 disk?

# RAID-1: Mirroring

- Motivation: Handle disk failures
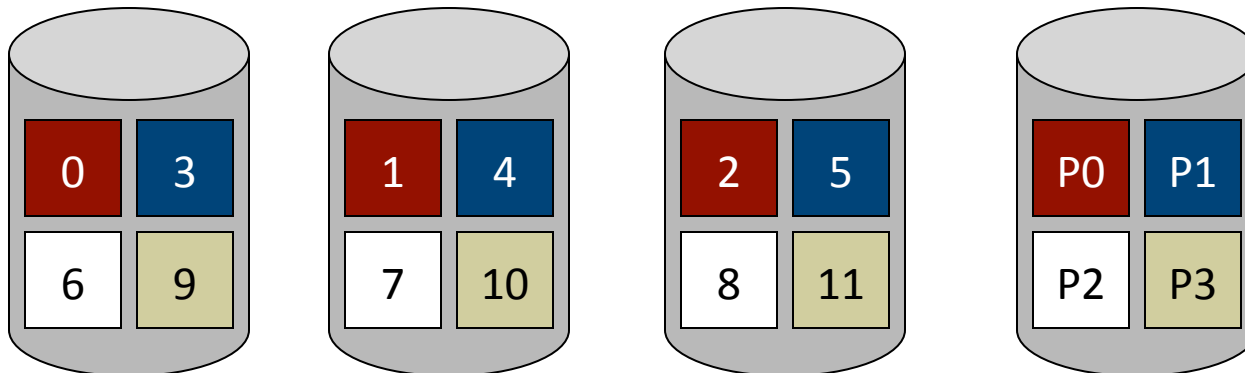- Put copy (mirror or replica) of each chunk on another disk



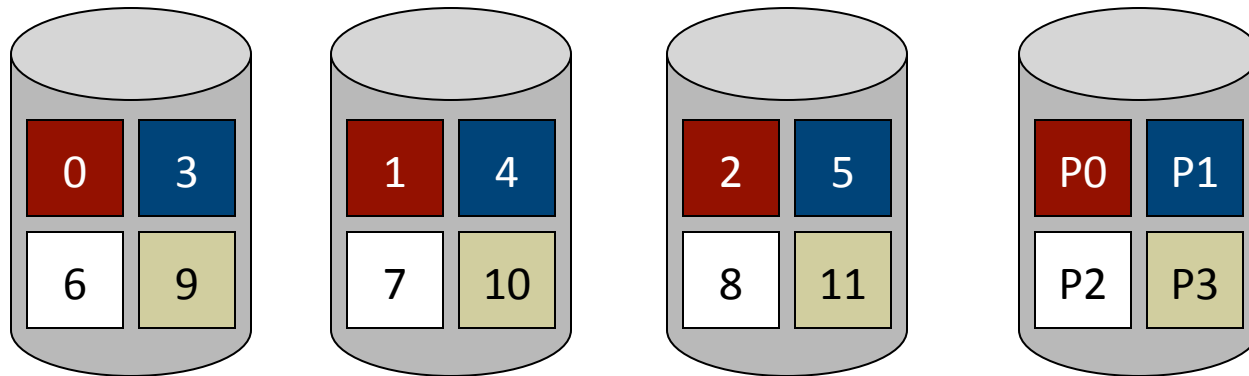- Capacity
- Reliability
- Performance

# RAID-4: Parity

- Motivation: Improve capacity
- Idea: Allocate parity block to encode info about blocks
  - Parity checks all other blocks in stripe across other disks
- Parity block = XOR over others (gives "even" parity)
  - Example: 0 1 0 → Parity value?
- How do you recover from a failed disk?
  - Example: x 0 0 and parity of 1
  - What is the failed value?

| $A$ | $B$ | XOR |
|-----|-----|-----|
| 0   | 0   | 0   |
| 0   | 1   | 1   |
| 1   | 0   | 1   |
| 1   | 1   | 0   |

| 0 | 3 | | 1 | 4 | | 2 | 5 | | P0 | P1 |
|---|---|---|---|---|---|---|---|----|----|
| 6 | 9 | | 7 | 10 | | 8 | 11 | | P2 | P3 |

# RAID-4: Parity



- Capacity:
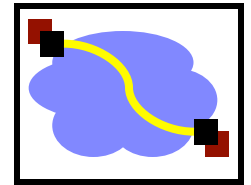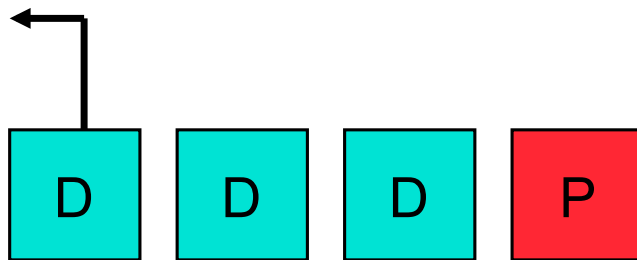- Reliability:
- Performance:
  - Reads
  - Writes: How to update parity block?
    - Two ways:
      - Use parity disk
      - Re-compute parity from non-parity disks
    - (Parity disk is the bottleneck)

# Updating and using the parity

## Fault-Free Read

D D D P

## Fault-Free Write

②  ①  ③  ④

D D D P

## Degraded Read

⊕

D D D P

## Degraded Write

⊕

D D D P

# RAID-5: Rotated Parity

Rotate location of parity across all disks



- Capacity:
- Reliability:
- Performance:
  - Reads:
  - Writes:
  - Still requires 4 I/Os per write, but not always to same parity disk

# Comparison

N: number of disks
S: throughput for sequential read
R: throughput for random read
D: delay to read/write

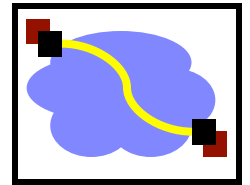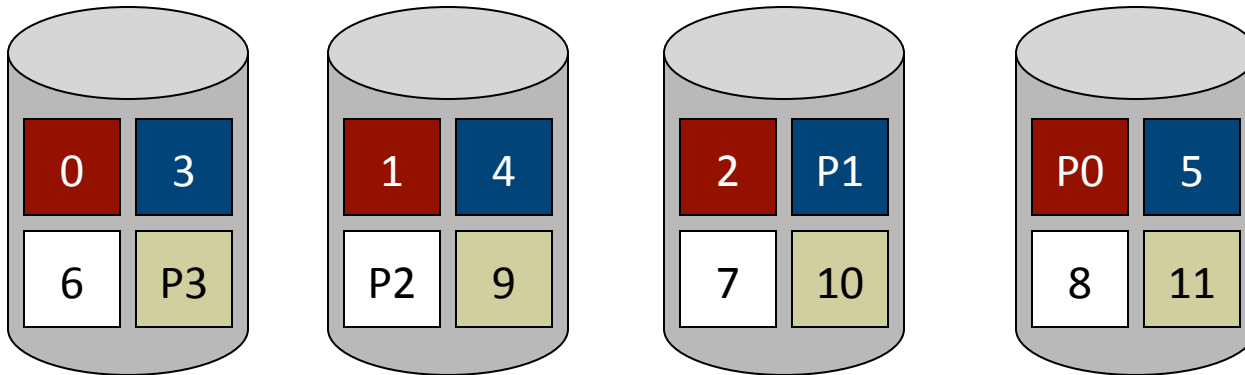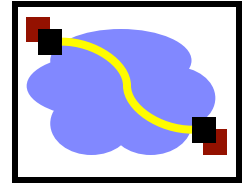|  | RAID-0 | RAID-1 | RAID-4 | RAID-5 |
|---|---|---|---|---|
| Capacity | $N$ | $N/2$ | $N-1$ | $N-1$ |
| Reliability | 0 | 1 (for sure) $\frac{N}{2}$ (if lucky) | 1 | 1 |
| Throughput |  |  |  |  |
| Sequential Read | $N \cdot S$ | $(N/2) \cdot S$ | $(N-1) \cdot S$ | $(N-1) \cdot S$ |
| Sequential Write | $N \cdot S$ | $(N/2) \cdot S$ | $(N-1) \cdot S$ | $(N-1) \cdot S$ |
| Random Read | $N \cdot R$ | $N \cdot R$ | $(N-1) \cdot R$ | $N \cdot R$ |
| Random Write | $N \cdot R$ | $(N/2) \cdot R$ | $\frac{1}{2} \cdot R$ | $\frac{N}{4} R$ |
| Latency |  |  |  |  |
| Read | $D$ | $D$ | $D$ | $D$ |
| Write | $D$ | $D$ | $2D$ | $2D$ |

Table 38.7: **RAID Capacity, Reliability, and Performance**

# Advanced Issues

- What happens if more than one fault?
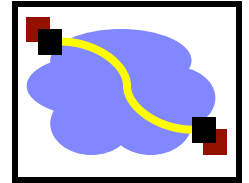  - Example: One disk fails plus "latent sector error" on another
  - RAID-5 cannot handle two faults
  - Solution: RAID-6: add multiple parity blocks
- Why is NVRAM useful?
  - Example: What if update 2, don't update P0 before power failure (or crash), and then disk 1 fails?
  - NVRAM solution: Use to store blocks updated in same stripe
    - If power failure, can replay all writes in NVRAM
  - Software RAID solution: Perform parity scrub over entire disk

| 0 | 3 |
|---|---|
| 6 | 9 |

| 1 | 4 |
|---|---|
| 7 | 10 |

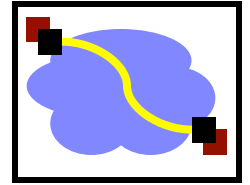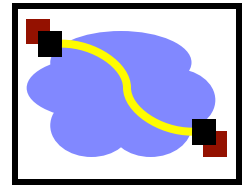| 2' | 5 |
|---|---|
| 8 | 11 |

| P0 | P1 |
|---|---|
| P2 | P3 |

# Conclusions

- RAID turns multiple disks into a larger, faster, more reliable disk

- RAID-0: Striping
Good when performance and capacity really matter, but reliability doesn't

- RAID-1: Mirroring
Good when reliability and write performance matter, but capacity (cost) doesn't

- RAID-4: Parity disk

- RAID-5: Rotating parity
Good when capacity and cost matter or workload is read-mostly
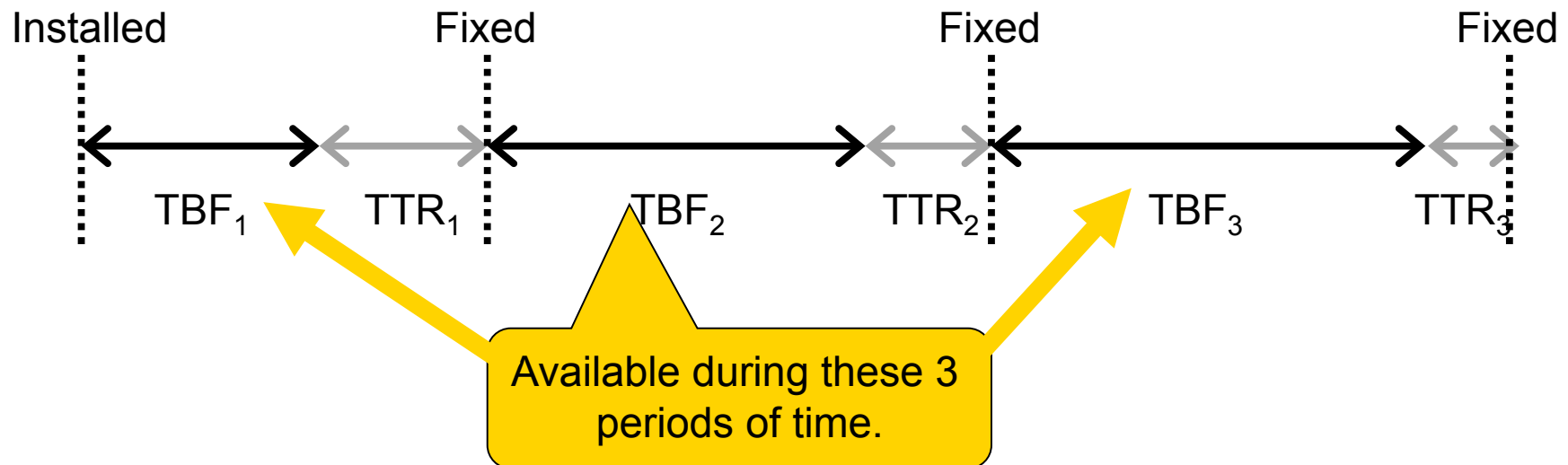  - Good compromise choice

# Outline

- Using multiple disks
  - Why have multiple disks?
  - problem and approaches

- RAID levels and performance

- Estimating availability
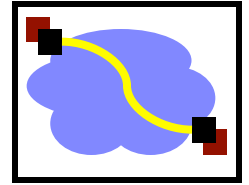
# Sidebar: Availability metric

- Fraction of time that server is able to handle requests
  - Computed from MTBF and MTTR (Mean Time To Repair)

$$\textbf{Availability} = \frac{\textbf{MTBF}}{\textbf{MTBF} + \textbf{MTTR}}$$

Installed        Fixed        Fixed        Fixed

$TBF_1$    $TTR_1$    $TBF_2$    $TTR_2$    $TBF_3$    $TTR_3$

Available during these 3 periods of time.

# How often are failures?

- MTBF (Mean Time Between Failures)
  - $MTBF_{disk}$ ~ 1,200,00 hours (~136 years, <1% per year)
- $MTBF_{mutli\text{-}disk\ system}$ = mean time to first disk failure
  - which is $MTBF_{disk}$ / (number of disks)
  - For a striped array of 200 drives
  - $MTBF_{array}$ = 136 years / 200 drives = 0.65 years