# Betriebssysteme

## 12. Tutorium - Files and Directories

Peter Bohner

30. Januar 2025

ITEC - Operating Systems Group

# Solid-State Drives

## NAND based flash memory

### Rejoice, TI might be useful once

How long do writes/reads take normally?

## Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page: 25µs

## Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page: 25μs
- Writing a page: 250μs

# NAND based flash memory

## Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page: 25µs
- Writing a page: 250µs
- Erasing a block:

### Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page:  25µs

- Writing a page:  250µs

- Erasing a block:  2ms

What happens when you just write to a random page?

### Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page:  25μs
- Writing a page:  250μs
- Erasing a block:  2ms

What happens when you just write to a random page?

### Speeding things up

What could you change so writing pages is faster?

### Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page: 25μs
- Writing a page: 250μs
- Erasing a block: 2ms

What happens when you just write to a random page?

### Speeding things up

What could you change so writing pages is faster?

- Keep around spare *erased* pages
- ⇒ You do not pay the erase penalty!
- When do you create / reserve / erase those spare pages?

## NAND based flash memory

### Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page: 25µs
- Writing a page: 250µs
- Erasing a block: 2ms

What happens when you just write to a random page?

### Speeding things up

What could you change so writing pages is faster?

- Keep around spare *erased* pages
- ⇒ You do not pay the erase penalty!
- When do you create / reserve / erase those spare pages? Probably in the background. Any problems?

### Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page: 25μs
- Writing a page: 250μs
- Erasing a block: 2ms

What happens when you just write to a random page?

### Speeding things up

What could you change so writing pages is faster?

- Keep around spare *erased* pages
- ⇒ You do not pay the erase penalty!
- When do you create / reserve / erase those spare pages? Probably in the background. Any problems? Might get exhausted if you write too much data in a short timeframe or the disk is full!

### Deleting files

What happens when you delete a file? What effect does that have on the SSD performance?
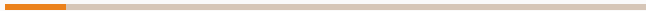
#### Deleting files

What happens when you delete a file? What effect does that have on the SSD performance? The block is not freed $\Rightarrow$ Can't be used as an erased empty page

What can the OS do to combat that?

#### Deleting files

What happens when you delete a file? What effect does that have on the SSD performance? The block is not freed $\Rightarrow$ Can't be used as an erased empty page

What can the OS do to combat that?

#### The `trim` command

Can be issued by the OS to tell the SSD firmware what pages can be safely erased.

# RAID

## What is that?

A **R**edundant **A**rray of **I**ndependent/**I**nexpensive **D**isks

What is that?

A **R**edundant **A**rray of **I**ndependent/**I**nexpensive **D**isks

Why would you use that?

## What is that?

A **R**edundant **A**rray of **I**ndependent/**I**nexpensive **D**isks

## Why would you use that?

- Probably cheaper than a SLED (Single Large Expensive Disk)

### What is that?

A **R**edundant **A**rray of **I**ndependent/**I**nexpensive **D**isks

### Why would you use that?

- Probably cheaper than a SLED (Single Large Expensive Disk)
- Might be more resilient
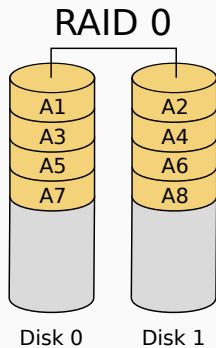
#### What is that?

A **R**edundant **A**rray of **I**ndependent/**I**nexpensive **D**isks
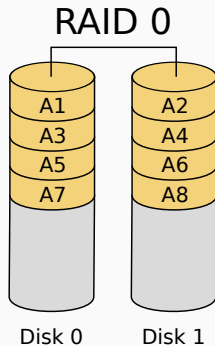
#### Why would you use that?

- Probably cheaper than a SLED (Single Large Expensive Disk)
- Might be more resilient
- Might be faster

Great, you now have multiple disks. How do you store your files on them?

- „I like to live dangerously" - RAID Level 0
- Mirroring: RAID Level 1
- Historic variants: RAID Level 2 and 3
- Block striping and parity: RAID Level 4
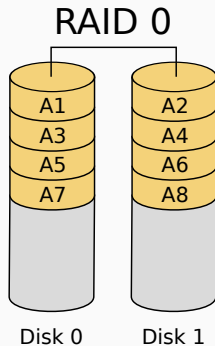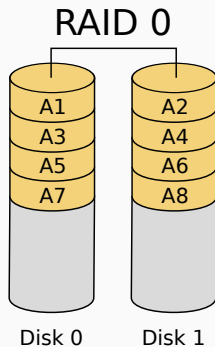- Block striping and *distributed* parity: RAID Level 5

# RAID 0

A1 A2
A3 A4
A5 A6
A7 A8

Disk 0    Disk 1

## Benefits / Drawbacks?

## RAID 0



Disk 0    Disk 1

### Benefits / Drawbacks?

+ Extremely fast (parallel reads and writes)

# RAID 0



```
        A1      A2
        A3      A4
        A5      A6
        A7      A8

      Disk 0   Disk 1
```

### Benefits / Drawbacks?

+ Extremely fast (parallel reads and writes)
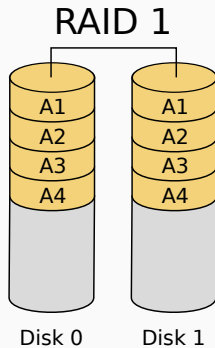+ Can use full capacity

# RAID 0



Disk 0      Disk 1

**Benefits / Drawbacks?**

+ Extremely fast (parallel reads and writes)
+ Can use full capacity
- If a single disk fails your files are toast

# RAID 1

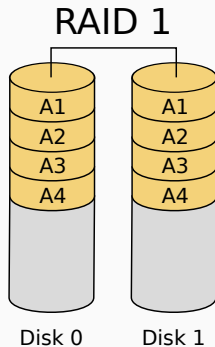A1
A2
A3
A4
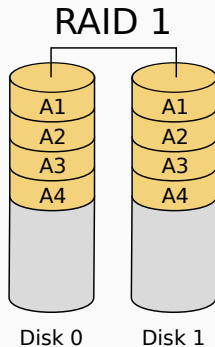
A1
A2
A3
A4

Disk 0    Disk 1

Benefits / Drawbacks?

# RAID 1

A1
A2
A3
A4

A1
A2
A3
A4

Disk 0    Disk 1

### Benefits / Drawbacks?

+ You can lose all but one disk without losing data

# RAID 1



Disk 0   Disk 1

### Benefits / Drawbacks?

+ You can lose all but one disk without losing data
+ Parallel reads possible

# RAID 1



Disk 0    Disk 1

### Benefits / Drawbacks?

+ You can lose all but one disk without losing data
+ Parallel reads possible
- Writes slower as they need to write to all disks

# RAID 1



Disk 0     Disk 1

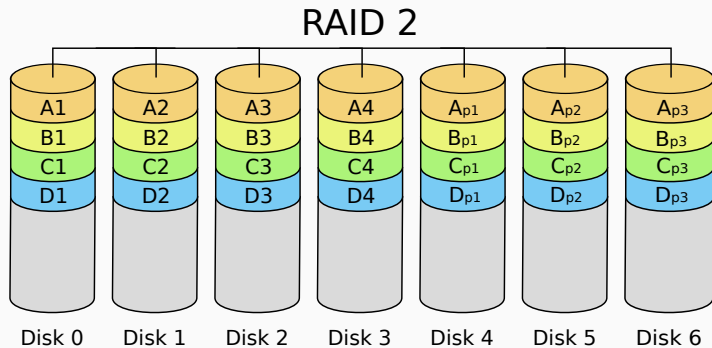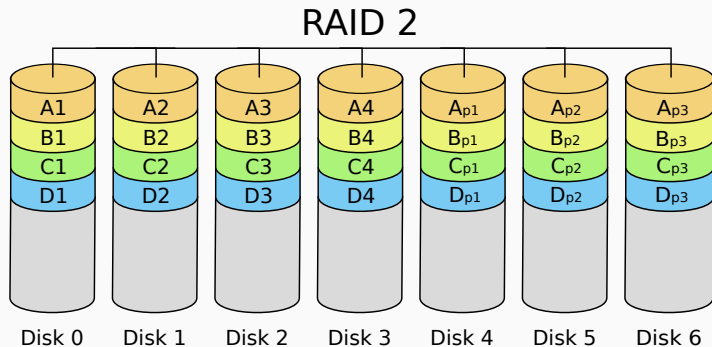### Benefits / Drawbacks?

+ You can lose all but one disk without losing data
+ Parallel reads possible
- Writes slower as they need to write to all disks
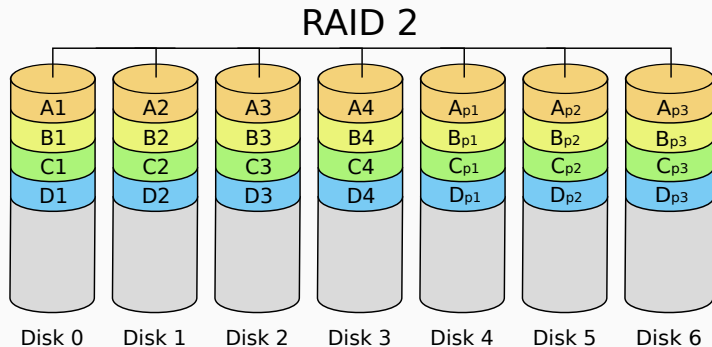- Size equals the size of a single disk

# RAID 2



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

## What is that?

- Have $\log_2(N)$ parity disk

# RAID 2



| Disk 0 | Disk 1 | Disk 2 | Disk 3 | Disk 4 | Disk 5 | Disk 6 |
|--------|--------|--------|--------|--------|--------|--------|
| A1 | A2 | A3 | A4 | $A_{p1}$ | $A_{p2}$ | $A_{p3}$ |
| B1 | B2 | B3 | B4 | $B_{p1}$ | $B_{p2}$ | $B_{p3}$ |
| C1 | C2 | C3 | C4 | $C_{p1}$ | $C_{p2}$ | $C_{p3}$ |
| D1 | D2 | D3 | D4 | $D_{p1}$ | $D_{p2}$ | $D_{p3}$ |

## What is that?

- Have $\log_2(N)$ parity disk
- Stripe data at the *bit* level

# RAID 2



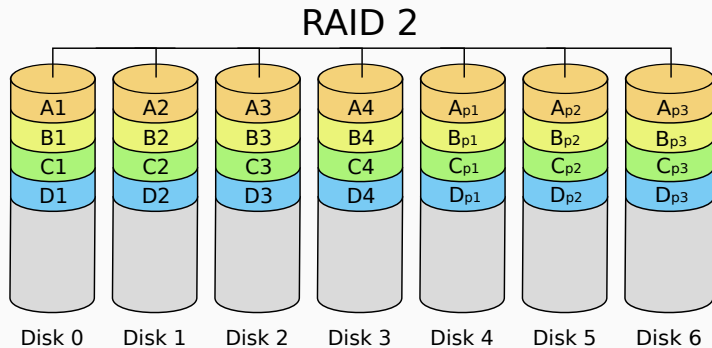| A1 | A2 | A3 | A4 | $A_{p1}$ | $A_{p2}$ | $A_{p3}$ |
| B1 | B2 | B3 | B4 | $B_{p1}$ | $B_{p2}$ | $B_{p3}$ |
| C1 | C2 | C3 | C4 | $C_{p1}$ | $C_{p2}$ | $C_{p3}$ |
| D1 | D2 | D3 | D4 | $D_{p1}$ | $D_{p2}$ | $D_{p3}$ |

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

## What is that?

- Have $\log_2(N)$ parity disk
- Stripe data at the *bit* level
- Use a hamming code of proper size

## RAID 2



| A1 | A2 | A3 | A4 | $A_{p1}$ | $A_{p2}$ | $A_{p3}$ |
| B1 | B2 | B3 | B4 | $B_{p1}$ | $B_{p2}$ | $B_{p3}$ |
| C1 | C2 | C3 | C4 | $C_{p1}$ | $C_{p2}$ | $C_{p3}$ |
| D1 | D2 | D3 | D4 | $D_{p1}$ | $D_{p2}$ | $D_{p3}$ |

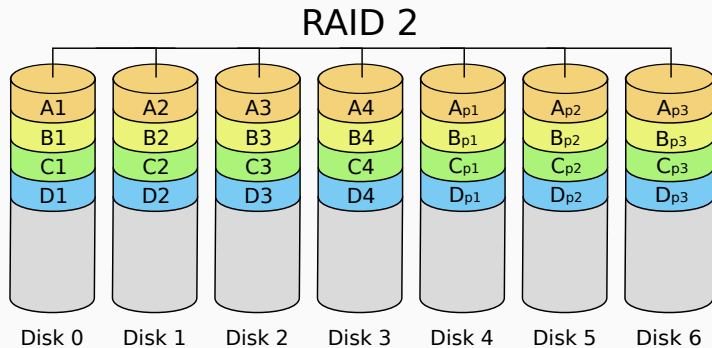Disk 0   Disk 1   Disk 2   Disk 3   Disk 4   Disk 5   Disk 6

### What is that?

- Have $\log_2(N)$ parity disk
- Stripe data at the *bit* level
- Use a hamming code of proper size
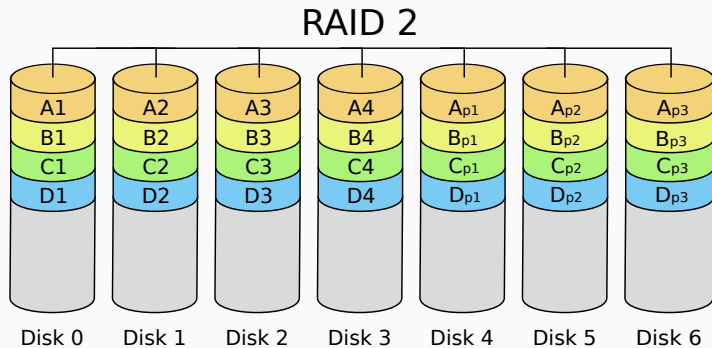- Spin the disks in lockstep (so you read all bits of your word at once)

# RAID 2



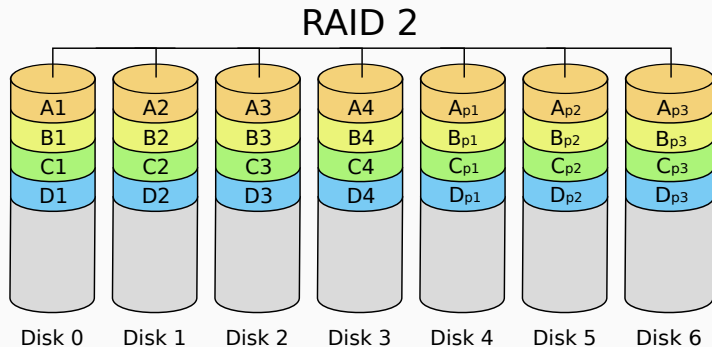| A1 | A2 | A3 | A4 | $A_{p1}$ | $A_{p2}$ | $A_{p3}$ |
| B1 | B2 | B3 | B4 | $B_{p1}$ | $B_{p2}$ | $B_{p3}$ |
| C1 | C2 | C3 | C4 | $C_{p1}$ | $C_{p2}$ | $C_{p3}$ |
| D1 | D2 | D3 | D4 | $D_{p1}$ | $D_{p2}$ | $D_{p3}$ |

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

## Benefits / Drawbacks

# RAID 2



| A1 | A2 | A3 | A4 | $A_{p1}$ | $A_{p2}$ | $A_{p3}$ |
| B1 | B2 | B3 | B4 | $B_{p1}$ | $B_{p2}$ | $B_{p3}$ |
| C1 | C2 | C3 | C4 | $C_{p1}$ | $C_{p2}$ | $C_{p3}$ |
| D1 | D2 | D3 | D4 | $D_{p1}$ | $D_{p2}$ | $D_{p3}$ |

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

## Benefits / Drawbacks

+ External error checking

# RAID 2

| Disk 0 | Disk 1 | Disk 2 | Disk 3 | Disk 4 | Disk 5 | Disk 6 |
|--------|--------|--------|--------|--------|--------|--------|
| A1 | A2 | A3 | A4 | $A_{p1}$ | $A_{p2}$ | $A_{p3}$ |
| B1 | B2 | B3 | B4 | $B_{p1}$ | $B_{p2}$ | $B_{p3}$ |
| C1 | C2 | C3 | C4 | $C_{p1}$ | $C_{p2}$ | $C_{p3}$ |
| D1 | D2 | D3 | D4 | $D_{p1}$ | $D_{p2}$ | $D_{p3}$ |

## Benefits / Drawbacks

+ External error checking

- Really slow

# RAID 2



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

## Benefits / Drawbacks

+ External error checking

- Really slow

- Not that useful as disks have internal error checking by now

# RAID 2



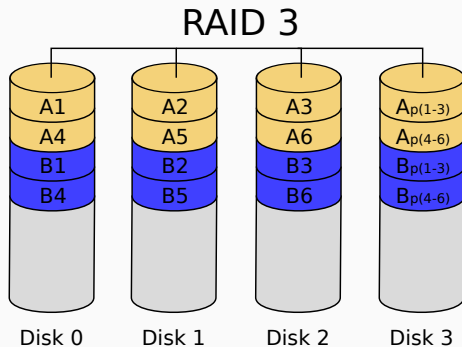Disk 0   Disk 1   Disk 2   Disk 3   Disk 4   Disk 5   Disk 6

## Benefits / Drawbacks

+ External error checking

- Really slow

- Not that useful as disks have internal error checking by now

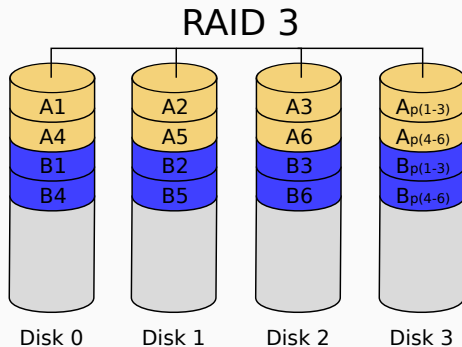- Spins in lockstep $\Rightarrow$ Can only service one request at a time

# RAID 3



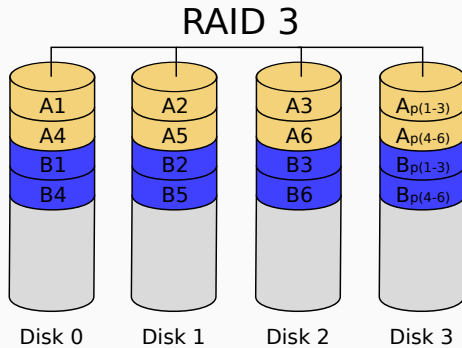### What is that?

- Have a dedicated parity disk

# RAID 3



Disk 0   Disk 1   Disk 2   Disk 3

### What is that?

- Have a dedicated parity disk
- Stripe data at the *byte* level

# RAID 3



| Disk 0 | Disk 1 | Disk 2 | Disk 3 |

A1, A2, A3, $A_{p(1-3)}$
A4, A5, A6, $A_{p(4-6)}$
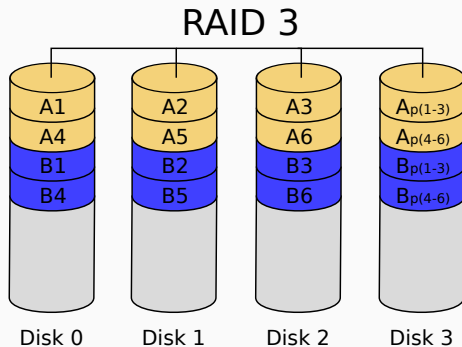B1, B2, B3, $B_{p(1-3)}$
B4, B5, B6, $B_{p(4-6)}$

## What is that?

- Have a dedicated parity disk
- Stripe data at the *byte* level
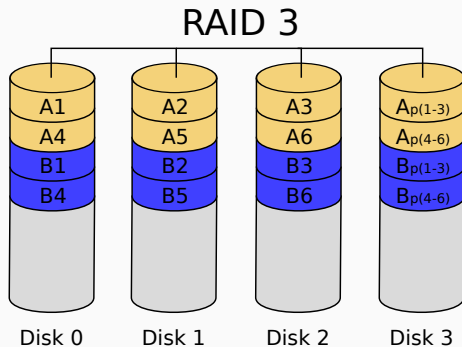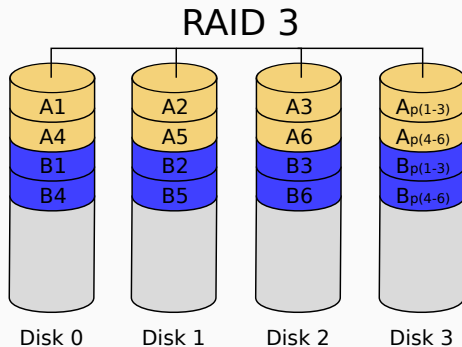- Spin the disks in lockstep (so you read all bytes of your word at once)

# RAID 3



| Disk 0 | Disk 1 | Disk 2 | Disk 3 |

A1, A2, A3, $A_{p(1-3)}$
A4, A5, A6, $A_{p(4-6)}$
B1, B2, B3, $B_{p(1-3)}$
B4, B5, B6, $B_{p(4-6)}$

## Benefits / Drawbacks

RAID 3

Disk 0    Disk 1    Disk 2    Disk 3

### Benefits / Drawbacks

+ You can lose a disk and restore it using the parity

## RAID 3



| Disk 0 | Disk 1 | Disk 2 | Disk 3 |

A1, A2, A3, $A_{p(1-3)}$
A4, A5, A6, $A_{p(4-6)}$
B1, B2, B3, $B_{p(1-3)}$
B4, B5, B6, $B_{p(4-6)}$

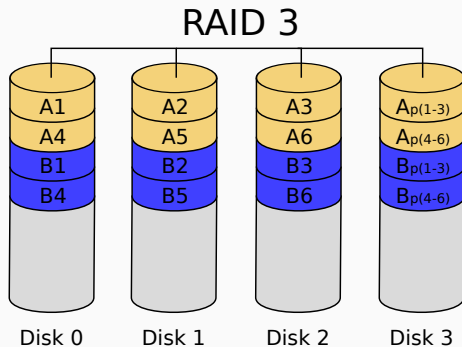### Benefits / Drawbacks

+ You can lose a disk and restore it using the parity
- Slow when reading/writing small files at random locations

8

RAID 3

| Disk 0 | Disk 1 | Disk 2 | Disk 3 |

Disk 0: A1, A4, B1, B4
Disk 1: A2, A5, B2, B5
Disk 2: A3, A6, B3, B6
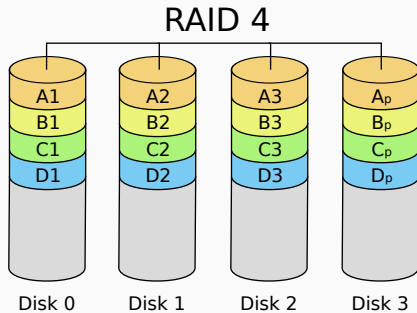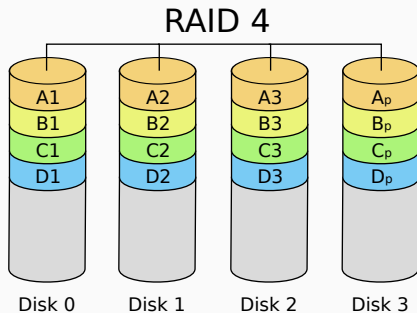Disk 3: $A_{p(1-3)}$, $A_{p(4-6)}$, $B_{p(1-3)}$, $B_{p(4-6)}$

## Benefits / Drawbacks

+ You can lose a disk and restore it using the parity

- Slow when reading/writing small files at random locations

- Spins in lockstep $\Rightarrow$ Can only service one request at a time

## RAID 3



Disk 0    Disk 1    Disk 2    Disk 3

### Benefits / Drawbacks

+ You can lose a disk and restore it using the parity

- Slow when reading/writing small files at random locations

- Spins in lockstep $\Rightarrow$ Can only service one request at a time
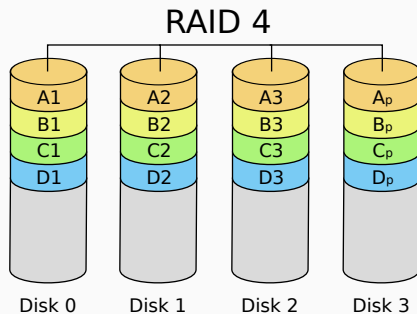
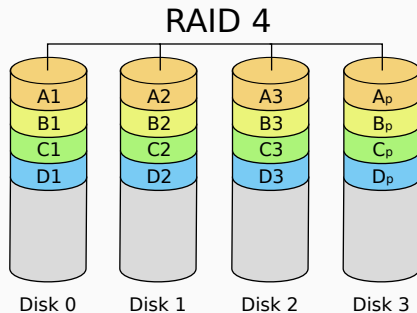- Every write and read hits the same single parity disk

# RAID 4



Disk 0     Disk 1     Disk 2     Disk 3

## What is that?

- Have a dedicated parity disk

RAID 4

## What is that?

- Have a dedicated parity disk
- Stripe data at the *block* level
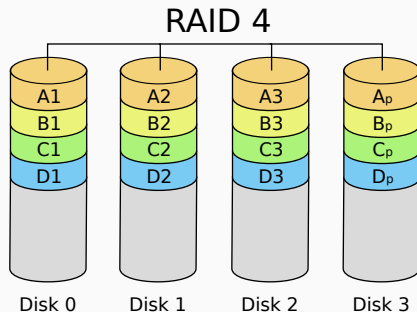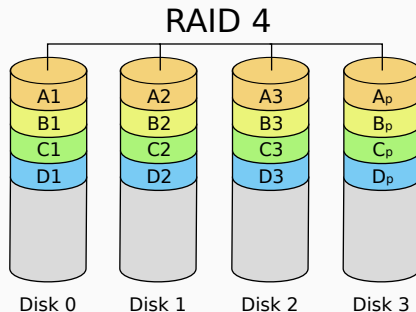
# RAID 4



Disk 0   Disk 1   Disk 2   Disk 3

## Benefits / Drawbacks

## RAID 4



| Disk 0 | Disk 1 | Disk 2 | Disk 3 |

### Benefits / Drawbacks

+ You can lose a disk and restore it using the parity

# RAID 4



Disk 0    Disk 1    Disk 2    Disk 3

### Benefits / Drawbacks

+ You can lose a disk and restore it using the parity

+ Good read performance

## RAID 4



| A1 | A2 | A3 | A_p |
| B1 | B2 | B3 | B_p |
| C1 | C2 | C3 | C_p |
| D1 | D2 | D3 | D_p |

Disk 0    Disk 1    Disk 2    Disk 3

### Benefits / Drawbacks

+ You can lose a disk and restore it using the parity
+ Good read performance
- Every write and read hits the same single parity disk $\Rightarrow$ Bottleneck, prone to failure
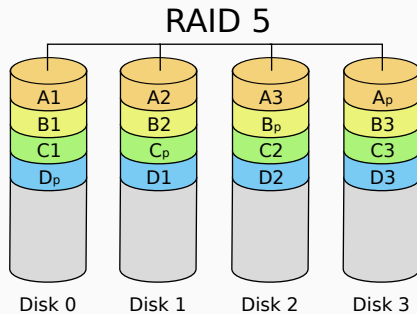
RAID 4

Disk 0    Disk 1    Disk 2    Disk 3

### Benefits / Drawbacks

+ You can lose a disk and restore it using the parity
+ Good read performance
- Every write and read hits the same single parity disk $\Rightarrow$ Bottleneck, prone to failure
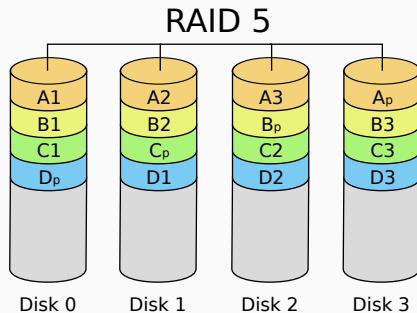- Slow writes (write to same parity disk)
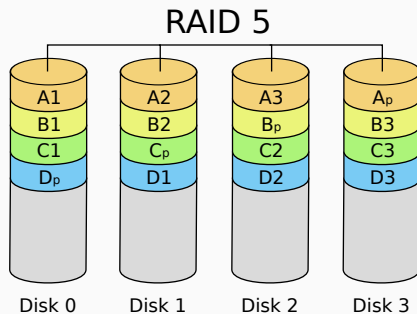
9

## What is that?

- Stripe data at the *block* level

# RAID 5



## What is that?

- Stripe data at the *block* level
- Distribute parity across your disks

# RAID 5



| | Disk 0 | Disk 1 | Disk 2 | Disk 3 |

Benefits / Drawbacks

# RAID 5



Disk 0    Disk 1    Disk 2    Disk 3

### Benefits / Drawbacks

+ You can lose a disk and restore it using the parity

## RAID 5



| Disk 0 | Disk 1 | Disk 2 | Disk 3 |

### Benefits / Drawbacks

+ You can lose a disk and restore it using the parity

+ Good read performance

# RAID 5



Disk 0    Disk 1    Disk 2    Disk 3

### Benefits / Drawbacks

+ You can lose a disk and restore it using the parity

+ Good read performance

+ Okay write performance

# RAID 5



Disk 0    Disk 1    Disk 2    Disk 3

### Benefits / Drawbacks

+ You can lose a disk and restore it using the parity

+ Good read performance

+ Okay write performance

- Still slower than RAID 0 or a SLED

### Compare SLED and RAID (Level 0, 1, 4, 5)

Each RAID uses 4 disks for actual data storage.

### How many disks do you need?

- SLED: **1**

### Compare SLED and RAID (Level 0, 1, 4, 5)

Each RAID uses 4 disks for actual data storage.

### How many disks do you need?

- SLED: 1
- RAID 0: 4

#### Compare SLED and RAID (Level 0, 1, 4, 5)

Each RAID uses 4 disks for actual data storage.

#### How many disks do you need?

- SLED: 1
- RAID 0: 4
- RAID 1: 8

#### Compare SLED and RAID (Level 0, 1, 4, 5)

Each RAID uses 4 disks for actual data storage.

#### How many disks do you need?

- SLED: **1**
- RAID 0: **4**
- RAID 1: **8**
- RAID 4: **5**

### Compare SLED and RAID (Level 0, 1, 4, 5)

Each RAID uses 4 disks for actual data storage.

### How many disks do you need?

- SLED: **1**
- RAID 0: **4**
- RAID 1: **8**
- RAID 4: **5**
- RAID 5: **5**

You want to modify one byte of data. How many blocks do you need to read/write?

- SLED: 1 read + 1 write

You want to modify one byte of data. How many blocks do you need to read/write?

- SLED: 1 read + 1 write
- RAID 0: 1 read + 1 write

You want to modify one byte of data. How many blocks do you need to read/write?

- SLED: 1 read + 1 write
- RAID 0: 1 read + 1 write
- RAID 1: 1 read + 2 write (1 data + 1 mirror)

You want to modify one byte of data. How many blocks do you need to read/write?

- SLED: 1 read + 1 write
- RAID 0: 1 read + 1 write
- RAID 1: 1 read + 2 write (1 data + 1 mirror)
- RAID 4: 2 read (data + old parity) + 2 write (data + new parity)

You want to modify one byte of data. How many blocks do you need to read/write?

- SLED: 1 read + 1 write
- RAID 0: 1 read + 1 write
- RAID 1: 1 read + 2 write (1 data + 1 mirror)
- RAID 4: 2 read (data + old parity) + 2 write (data + new parity)
- RAID 5: 2 read (data + old parity) + 2 write (data + new parity)

## You are using RAID

- You accidentally delete a file.

### You are using RAID

- You accidentally delete a file. *GONE*
- You accidentally overwrite a file.

## You are using RAID

- You accidentally delete a file. *GONE*
- You accidentally overwrite a file. *GONE*
- Some data gets corrupted on one disk.

### You are using RAID

- You accidentally delete a file. *GONE*
- You accidentally overwrite a file. *GONE*
- Some data gets corrupted on one disk. *GONE (probably)*
- The poor intern connects to the production database. (Or here)

### You are using RAID

- You accidentally delete a file. *GONE*
- You accidentally overwrite a file. *GONE*
- Some data gets corrupted on one disk. *GONE* *(probably)*
- The poor intern connects to the production database. (Or here) *GONE*
- A crypto-locker takes out your computer.

## You are using RAID

- You accidentally delete a file. *GONE*
- You accidentally overwrite a file. *GONE*
- Some data gets corrupted on one disk. *GONE* *(probably)*
- The poor intern connects to the production database. (Or here) *GONE*
- A crypto-locker takes out your computer. Believe it or not, ~~*JAIL*~~ *GONE*

## You are using RAID
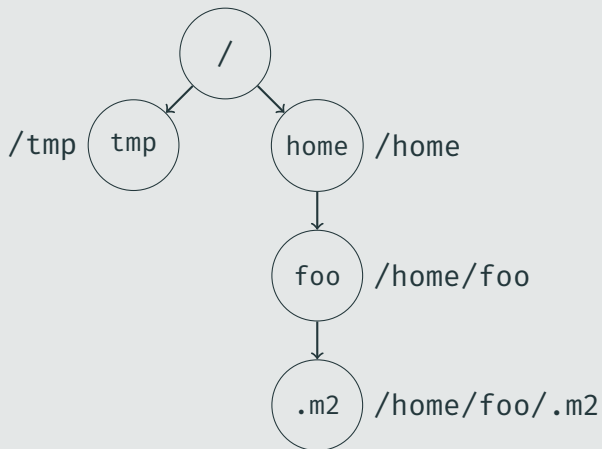
- You accidentally delete a file. *GONE*
- You accidentally overwrite a file. *GONE*
- Some data gets corrupted on one disk. *GONE* *(probably)*
- The poor intern connects to the production database. (Or here) *GONE*
- A crypto-locker takes out your computer. Believe it or not, ~~JAIL~~ *GONE*

## So what do we learn?

*RAID IS NO SUBSTITUTE FOR A BACKUP*

13

# Files And Directories

## Paths (Linux)

What are the two basic access methods (patterns) for reading a file?

- Sequential Access
  Accessed in order, reading sequential bytes. Writes append at the end.

- Random Access
  Reading and writing at arbitrary positions, programmer needs to specify
  where to write/read

Opening the file

Opening the file

fopen / open. Provide mode as argument (read/write, where, create, etc.)

Opening the file

fopen / open. Provide mode as argument (read/write, where, create, etc.)

Closing a file

Opening the file

`fopen` / `open`. Provide mode as argument (read/write, where, create, etc.)

Closing a file

`fclose` / `close`

## Which Library Functions and System Calls do we have to access files?

Opening the file

`fopen` / `open`. Provide mode as argument (read/write, where, create, etc.)

Closing a file

`fclose` / `close`

Reading from a file

Opening the file

`fopen` / `open`. Provide mode as argument (read/write, where, create, etc.)

Closing a file

`fclose` / `close`

Reading from a file

`fread` / `read`

## Which Library Functions and System Calls do we have to access files?

Opening the file

`fopen` / `open`. Provide mode as argument (read/write, where, create, etc.)

Closing a file

`fclose` / `close`

Reading from a file

`fread` / `read`

Writing to a file

Opening the file

`fopen` / `open`. Provide mode as argument (read/write, where, create, etc.)

Closing a file

`fclose` / `close`

Reading from a file

`fread` / `read`

Writing to a file

`fwrite` / `write`

## Which Library Functions and System Calls do we have to access files?

Opening the file

`fopen` / `open`. Provide mode as argument (read/write, where, create, etc.)

Closing a file

`fclose` / `close`

Reading from a file

`fread` / `read`

Writing to a file

`fwrite` / `write`

Flushing dirty buffers

## Which Library Functions and System Calls do we have to access files?

### Opening the file
`fopen` / `open`. Provide mode as argument (read/write, where, create, etc.)

### Closing a file
`fclose` / `close`

### Reading from a file
`fread` / `read`

### Writing to a file
`fwrite` / `write`

### Flushing dirty buffers
The library functions sometimes buffer to reduce the amount of syscalls.
`fflush` flushes those buffers.

Which system calls / library functions move the „current position pointer"?

Which system calls / library functions move the „current position pointer"?

- Linux: `fseek` / `lseek`
- Windows: `SetFilePointerEx` / `SetFilePointer`

Which system calls / library functions move the „current position pointer"?

- Linux: `fseek` / `lseek`
- Windows: `SetFilePointerEx` / `SetFilePointer`

What happens when you move the cursor behind the end of the file?

Which system calls / library functions move the „current position pointer"?

- Linux: `fseek` / `lseek`
- Windows: `SetFilePointerEx` / `SetFilePointer`

What happens when you move the cursor behind the end of the file?
It creates a hole filled with zeros!

Which system calls / library functions move the „current position pointer"?

- Linux: `fseek` / `lseek`
- Windows: `SetFilePointerEx` / `SetFilePointer`

What happens when you move the cursor behind the end of the file?
It creates a hole filled with zeros! Such a file is called a *sparse* file and some file systems might not store empty regions.

How could you implement random access without a dedicated file pointer?

Add an offset to the `read` and `write` system call.

How could you implement random access without a dedicated file pointer?

Add an offset to the `read` and `write` system call.

+ Save a system call

How could you implement random access without a dedicated file pointer?

Add an offset to the `read` and `write` system call.

+ Save a system call
- Caller need to keep track of the current offset for sequential reads

How could you implement random access without a dedicated file pointer?

Add an offset to the `read` and `write` system call.

+ Save a system call
- Caller need to keep track of the current offset for sequential reads

`mmap` the file (works quite well, might cause page faults. Probably faster for large files than read/write calls)

What system calls do you need to list files in a Linux directory?

What system calls do you need to list files in a Linux directory?

1. `opendir`

What system calls do you need to list files in a Linux directory?

1. `opendir`
2. `readdir`:

What system calls do you need to list files in a Linux directory?

1. `opendir`
2. `readdir`: Returns a `dirent` with a *relative path*

What system calls do you need to list files in a Linux directory?

1. `opendir`
2. `readdir`: Returns a `dirent` with a *relative path*
3. `closedir`

# Open Files

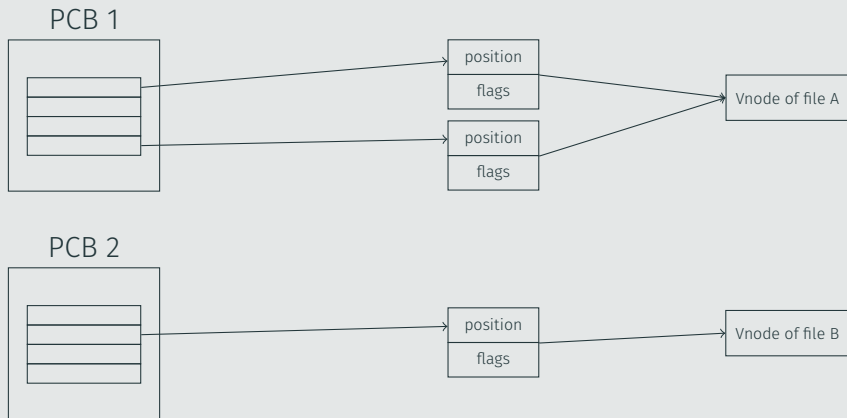## Kernel Data Structures For Open Files

What structures does the kernel use for open files?

## What structures does the kernel use for open files?

PCBs with local open file tables      Global open file table      vnode table

PCB 1

| position |
| flags |

| position |
| flags |

Vnode of file A

PCB 2

| position |
| flags |

Vnode of file B

XKCD 1084 - Server Problem

# F R A G E N ?



https://forms.gle/9CwJSKidKibubran9

Bis nächste Woche