

Betriebssysteme

11. Tutorium - Storage

Peter Bohner

23. Januar 2024

ITEC - Operating Systems Group

- Lief gut - von meiner Seite aus kein Besprechungsbedarf
- Fragen euerseits?

- Was ist ein Long und ein Short-Term scheduler?

- Was ist ein Long und ein Short-Term scheduler? LTS: Was kommt in die Run-Queue, STS: Was läuft auf der CPU?

- Was ist ein Long und ein Short-Term scheduler? LTS: Was kommt in die Run-Queue, STS: Was läuft auf der CPU?
- Unterschied mechanism - policy? Wieso wichtig?

- Was ist ein Long und ein Short-Term scheduler? LTS: Was kommt in die Run-Queue, STS: Was läuft auf der CPU?
- Unterschied mechanism - policy? Wieso wichtig?
- Was ist thrashing?

- Was ist ein Long und ein Short-Term scheduler? LTS: Was kommt in die Run-Queue, STS: Was läuft auf der CPU?
- Unterschied mechanism - policy? Wieso wichtig?
- Was ist thrashing?
- Welche threading Modelle kennt ihr? Vor- / Nachteile?

- Was ist ein Long und ein Short-Term scheduler? LTS: Was kommt in die Run-Queue, STS: Was läuft auf der CPU?
- Unterschied mechanism - policy? Wieso wichtig?
- Was ist thrashing?
- Welche threading Modelle kennt ihr? Vor- / Nachteile? One-to-One (Langsam, Kernel Support, einfach), Many-to-one (kein Kernel support, kein SMP), Many-to-Many (upcalls, komplexität)
- Unterschied KLT, KMT?

- Speicher - Segmentation: WS21 final1 T3
- Caching: WS1617 final1 T4 c,d

Device Management

What kind of I/O Devices do you find in a typical system?

What kind of I/O Devices do you find in a typical system?

- Block devices

What kind of I/O Devices do you find in a typical system?

- Block devices
- Character devices

What kind of I/O Devices do you find in a typical system?

- Block devices
- Character devices
- Network devices

What kind of I/O Devices do you find in a typical system?

- Block devices
- Character devices
- Network devices

Block Devices

- Offer random access to fixed-size blocks
- Applications typically deal with a file system on top of the device
- Examples?

What kind of I/O Devices do you find in a typical system?

- Block devices
- Character devices
- Network devices

Block Devices

- Offer random access to fixed-size blocks
- Applications typically deal with a file system on top of the device
- Examples? SSD, HDD, ...

What kind of I/O Devices do you find in a typical system?

- Block devices
- Character devices
- Network devices

Block Devices

- Offer random access to fixed-size blocks
- Applications typically deal with a file system on top of the device
- Examples? SSD, HDD, ...

Character Devices

- Provide a stream of characters
- Examples?

What kind of I/O Devices do you find in a typical system?

- Block devices
- Character devices
- Network devices

Block Devices

- Offer random access to fixed-size blocks
- Applications typically deal with a file system on top of the device
- Examples? SSD, HDD, ...

Character Devices

- Provide a stream of characters
- Examples? Mice, Keyboard, (classic) text terminals

Port Based I/O

Port Based I/O

Separate address space with dedicated instructions for reading/writing

- + Clear distinction in code \Rightarrow Optimizing easier (reordering, caching, ...)
- Less flexible, often lower performance

Port Based I/O

Separate address space with dedicated instructions for reading/writing

- + Clear distinction in code \Rightarrow Optimizing easier (reordering, caching, ...)
- Less flexible, often lower performance

Memory-mapped I/O

Device registers are mapped into the physical address space. How do you access that?

Port Based I/O

Separate address space with dedicated instructions for reading/writing

- + Clear distinction in code \Rightarrow Optimizing easier (reordering, caching, ...)
- Less flexible, often lower performance

Memory-mapped I/O

Device registers are mapped into the physical address space. How do you access that? Normal instructions!

- + Higher flexibility: Virtual memory, larger instruction set, mostly transparent
- Some special rules apply to I/O regions software needs to be aware of

Similar Names - Compare DMA, memory-mapped I/O and memory-mapped files

DMA

Similar Names - Compare DMA, memory-mapped I/O and memory-mapped files

DMA

Direct Memory Access

DMA

Direct Memory Access

- Devices can access the physical memory *without* involving the CPU

Similar Names - Compare DMA, memory-mapped I/O and memory-mapped files

DMA

Direct Memory Access

- Devices can access the physical memory *without* involving the CPU
- Needs special setup from the OS to know how and what to read/write

Memory Mapped Files

Similar Names - Compare DMA, memory-mapped I/O and memory-mapped files

DMA

Direct Memory Access

- Devices can access the physical memory *without* involving the CPU
- Needs special setup from the OS to know how and what to read/write

Memory Mapped Files

- OS abstraction: Treat a file like a normal range of virtual memory

Similar Names - Compare DMA, memory-mapped I/O and memory-mapped files

DMA

Direct Memory Access

- Devices can access the physical memory *without* involving the CPU
- Needs special setup from the OS to know how and what to read/write

Memory Mapped Files

- OS abstraction: Treat a file like a normal range of virtual memory
- No real relation to DMA, though the OS might use it to synchronize Memory Mapped Files with the underlying device

How does the OS know an I/O operation is finished?

How does the OS know an I/O operation is finished?

- Polling

How does the OS know an I/O operation is finished?

- Polling \Rightarrow Periodically check device registers

How does the OS know an I/O operation is finished?

- Polling \Rightarrow Periodically check device registers
- Interrupts

How does the OS know an I/O operation is finished?

- Polling \Rightarrow Periodically check device registers
- Interrupts \Rightarrow I/O devices send an interrupt signal

Hard Disks



What parts can you find in a hard disk?

What parts can you find in a hard disk?

- Heads

What parts can you find in a hard disk?

- Heads
- Arms

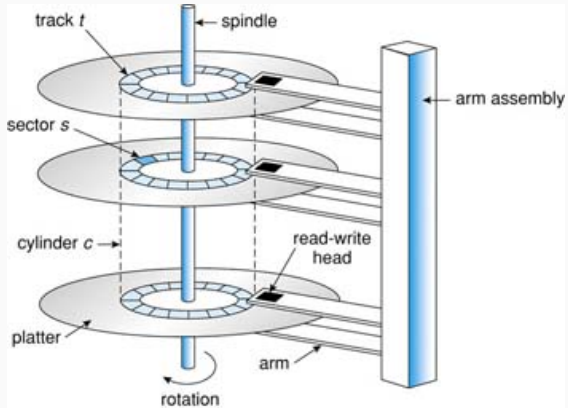
What parts can you find in a hard disk?

- Heads
- Arms
- Platters

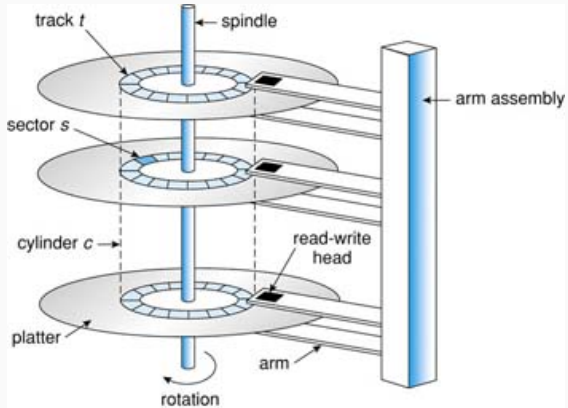
What parts can you find in a hard disk?

- Heads
- Arms
- Platters
- Spindle

Hard Disk Layout

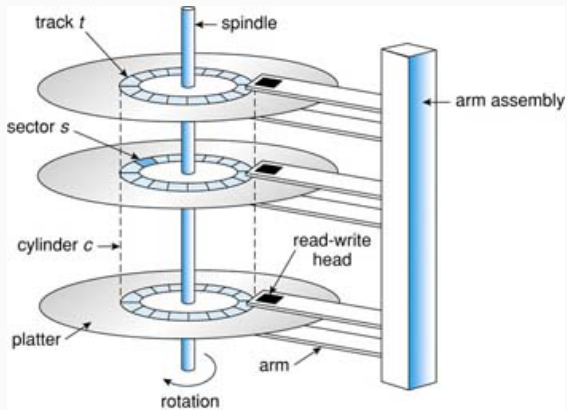


Hard Disk Layout



What do they do?

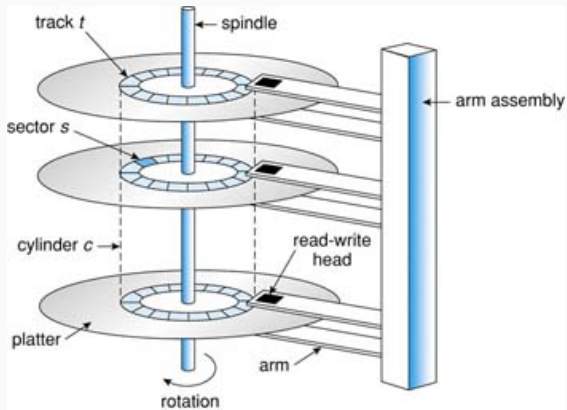
Hard Disk Layout



What do they do?

- Spindle: Spin connected platters!

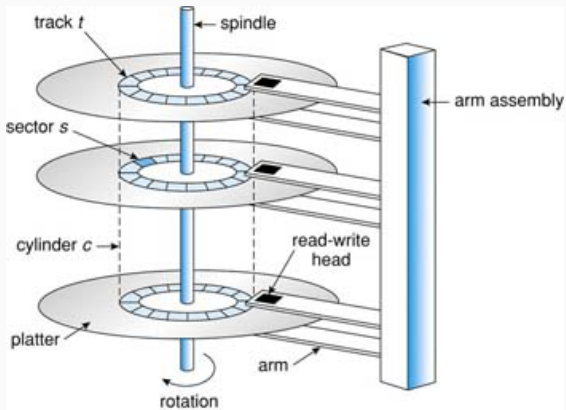
Hard Disk Layout



What do they do?

- Spindle: Spin connected platters!
- Head: Read/Write

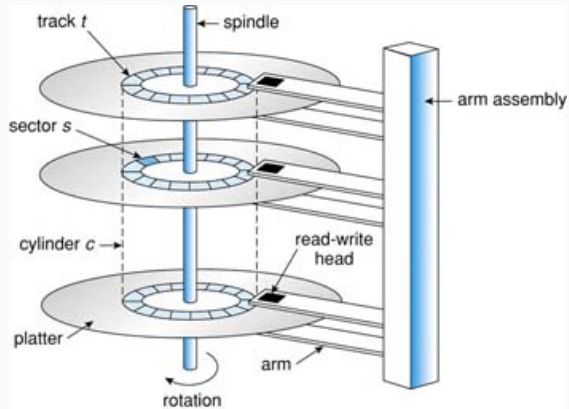
Hard Disk Layout



What do they do?

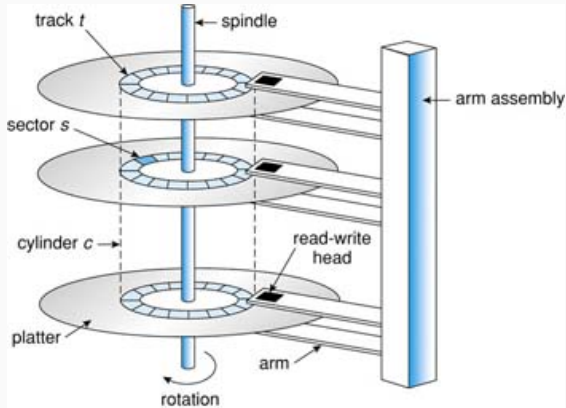
- Spindle: Spin connected platters!
- Head: Read/Write
- Arm: Move heads

Hard Disk Layout



How can you address data (512 byte blocks typically) on the disk?

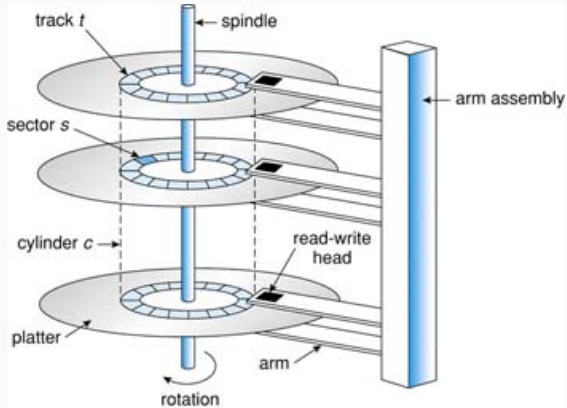
Hard Disk Layout



How can you address data (512 byte blocks typically) on the disk?

- Cylinder - Head - Sector (CHS). Limited to „small“ disks (< 8GB), rarely used these days

Hard Disk Layout



How can you address data (512 byte blocks typically) on the disk?

- Cylinder - Head - Sector (CHS). Limited to „small“ disks (< 8GB), rarely used these days
- Logical Block Addressing (LBA). Each data block has its own unique number.

How could you optimize the OS \Leftrightarrow Disk interface?

Native-Command-Queuing. OS sends reads and writes in batches and (the disk | the OS) reorders them based on internal geometry.

How could you optimize the OS \Leftrightarrow Disk interface?

Native-Command-Queuing. OS sends reads and writes in batches and *the disk* reorders them based on internal geometry.

What do you do when a sector is damaged?

What do you do when a sector is damaged?

Disk marks it as such and never uses it again \Rightarrow *Sector sparing*.

What adverse effect might this have?

What do you do when a sector is damaged?

Disk marks it as such and never uses it again \Rightarrow *Sector sparing*.

What adverse effect might this have? OS disk scheduler is unaware and optimizes for wrong geometry.

Hard Disk Layout - More Data!

Optimizing storage

Write area

Read area

A diagram illustrating storage optimization. It features a light gray background. In the center, there is a blue rectangle with the text "Read area" in white. This blue rectangle is enclosed within a larger red rectangle. To the left of the red rectangle, the text "Write area" is written in red.

What can you do with this?

Hard Disk Layout - More Data!

Optimizing storage

Write area

Read area



What can you do with this?

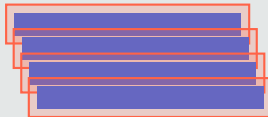
Conventional layout



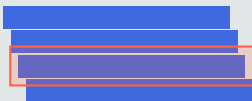
Shingled Magnet Recording

Hard Disk Layout - More Data!

Shingled Magnet Recording

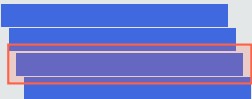


Shingled Magnet Recording



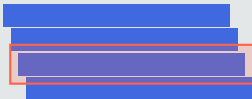
What happens when you write to this track?

Shingled Magnet Recording



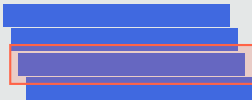
What happens when you write to this track? You overwrite the adjacent ones!

Shingled Magnet Recording



What happens when you write to this track? You overwrite the adjacent ones!
⇒ Append only and group shingled tracks

Shingled Magnet Recording



What happens when you write to this track? You overwrite the adjacent ones!

⇒ Append only and group shingled tracks

⇒ Can rewrite the whole group at once

How can such a device interface with the OS?

How can such a device interface with the OS?

- Pretend you are a normal disk. Buffer writes in a normal zone and flush them once they fill up a group.
⇒ *Device Managed*

How can such a device interface with the OS?

- Pretend you are a normal disk. Buffer writes in a normal zone and flush them once they fill up a group.
⇒ *Device Managed*
- Tell the OS where your shingled zones are. The OS needs to write carefully to not destroy data
⇒ *Host Managed*

How can such a device interface with the OS?

- Pretend you are a normal disk. Buffer writes in a normal zone and flush them once they fill up a group.
⇒ *Device Managed*
- Tell the OS where your shingled zones are. The OS needs to write carefully to not destroy data
⇒ *Host Managed*
- Compromise. Tell the OS where your shingled zones are and expose their tail. If the OS writes to the tail, directly commit it - else buffer.
⇒ *Host Aware*

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: 25 μ s

Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page: 25 μ s
- Writing a page: 250 μ s

Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page: 25 μ s
- Writing a page: 250 μ s
- Erasing a block:

NAND based flash memory

Rejoice, TI might be useful once

How long do writes/reads take normally?

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

What happens when you just write to a random page?

NAND based flash memory

Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page: 25 μ s
- Writing a page: 250 μ s
- 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?

NAND based flash memory

Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page: 25 μ s
- Writing a page: 250 μ s
- 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: 25 μ s
- Writing a page: 250 μ s
- 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?

NAND based flash memory

Rejoice, TI might be useful once

How long do writes/reads take normally?

- Reading a page: 25 μ s
- Writing a page: 250 μ s
- 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 Redundant Array of Independent/Inexpensive Disks

What is that?

A Redundant Array of Independent/Inexpensive Disks

Why would you use that?

What is that?

A Redundant Array of Independent/Inexpensive Disks

Why would you use that?

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

What is that?

A Redundant Array of Independent/Inexpensive Disks

Why would you use that?

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

What is that?

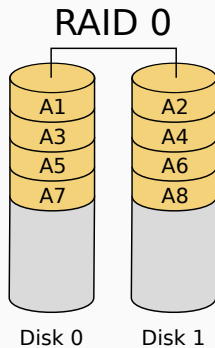
A Redundant Array of Independent/Inexpensive Disks

Why would you use that?

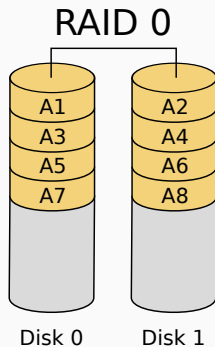
- Probably cheaper than a SLED (Single Large Expensive Disk)
- Might be more resilient
- Might be faster

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

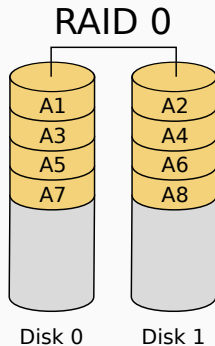


Benefits / Drawbacks?



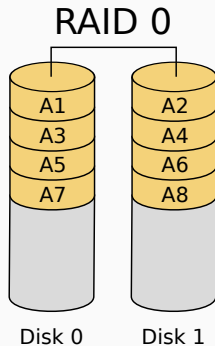
Benefits / Drawbacks?

- + Extremely fast (parallel reads and writes)



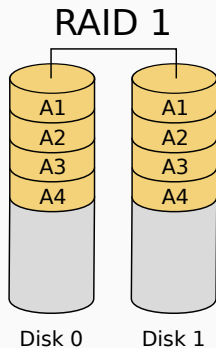
Benefits / Drawbacks?

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

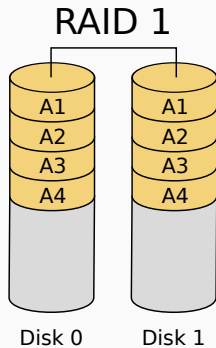


Benefits / Drawbacks?

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

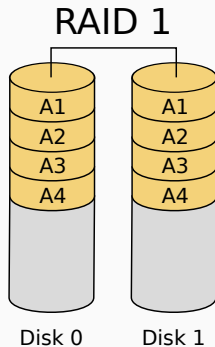


Benefits / Drawbacks?



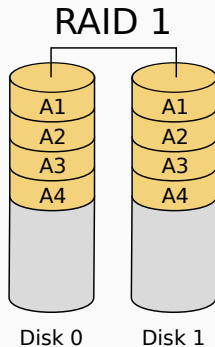
Benefits / Drawbacks?

- + You can lose all but one disk without losing data



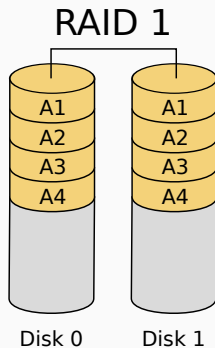
Benefits / Drawbacks?

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



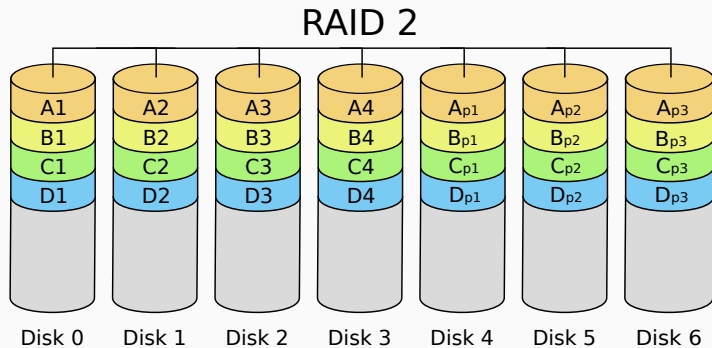
Benefits / Drawbacks?

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



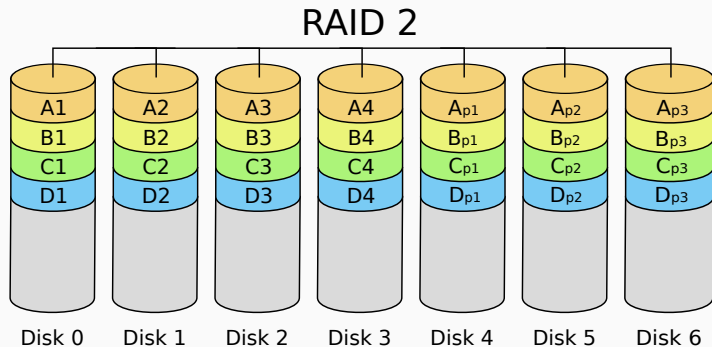
Benefits / Drawbacks?

- + 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



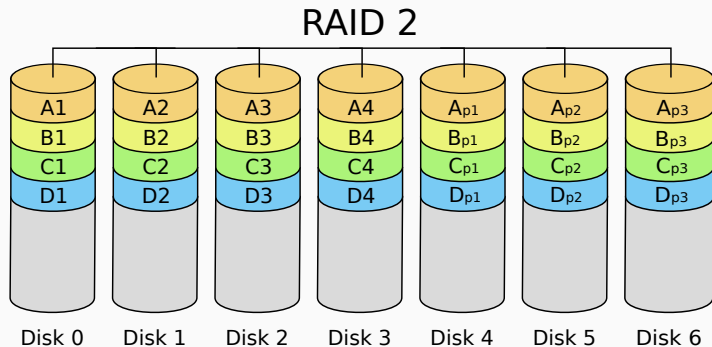
What is that?

- Have $\log_2(N)$ parity disk



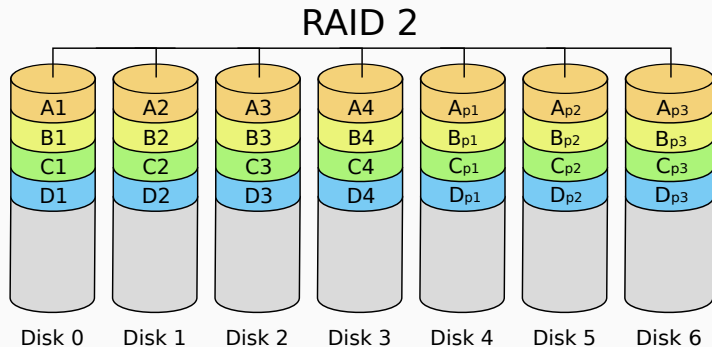
What is that?

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



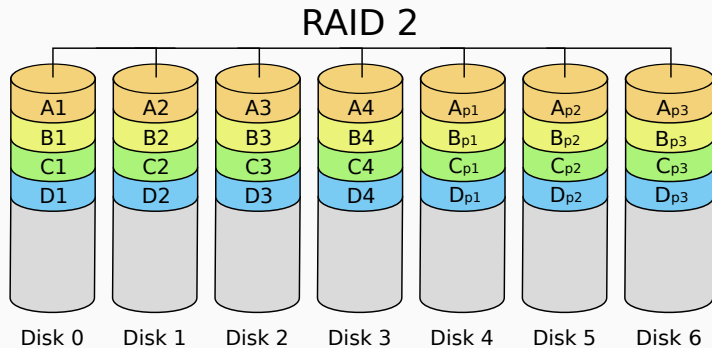
What is that?

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

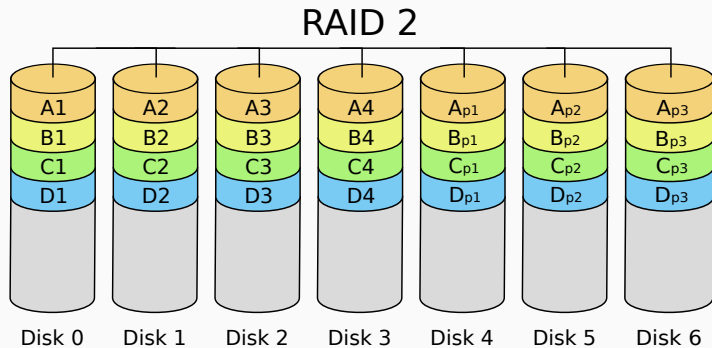


What is that?

- 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)

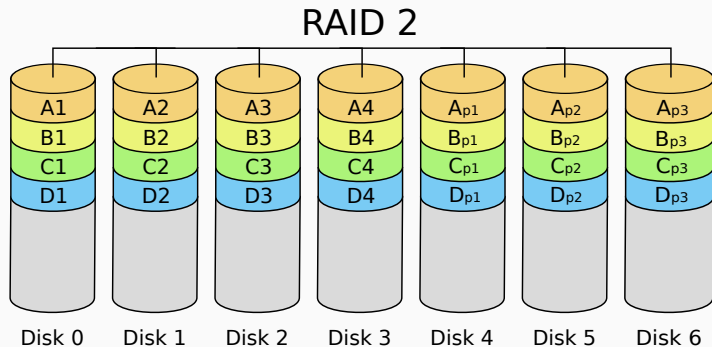


Benefits / Drawbacks



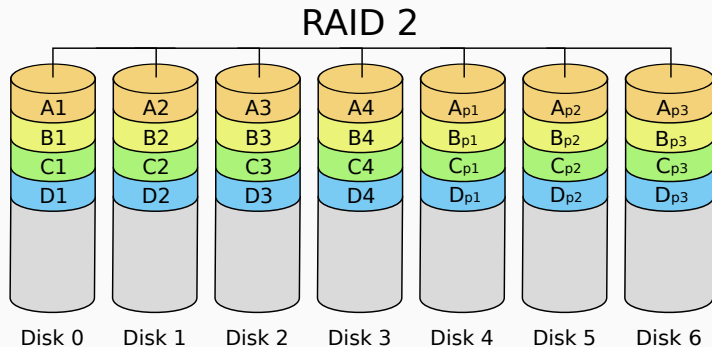
Benefits / Drawbacks

- + External error checking



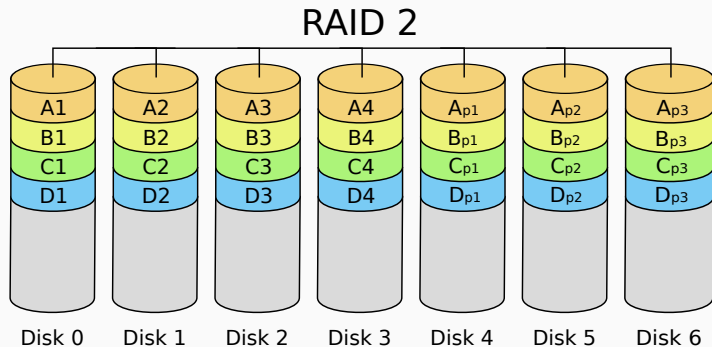
Benefits / Drawbacks

- + External error checking
- Really slow



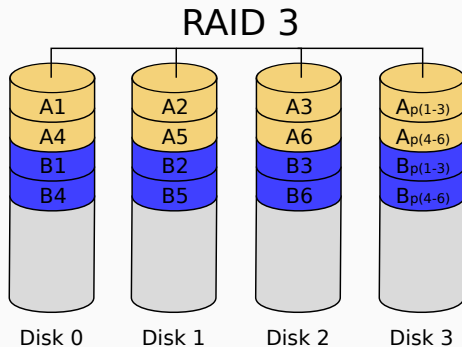
Benefits / Drawbacks

- + External error checking
- Really slow
- Not that useful as disks have internal error checking by now



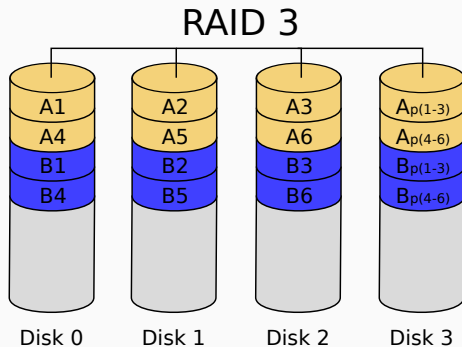
Benefits / Drawbacks

- + 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



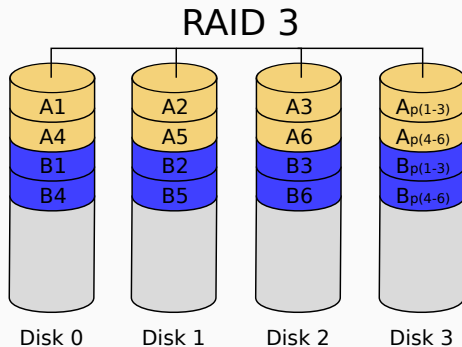
What is that?

- Have a dedicated parity disk



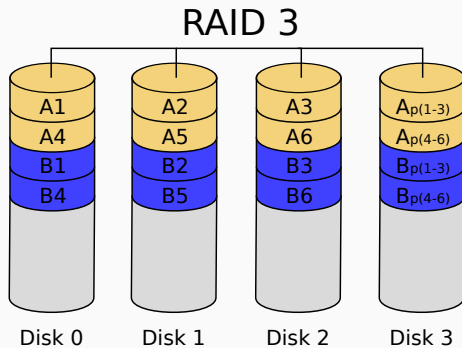
What is that?

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

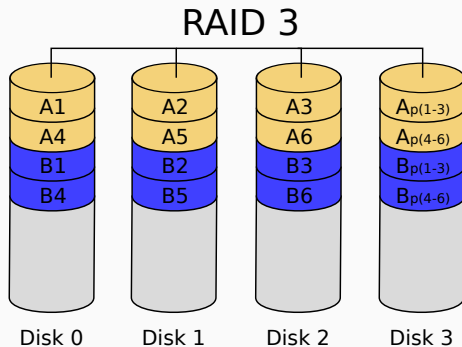


What is that?

- 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)

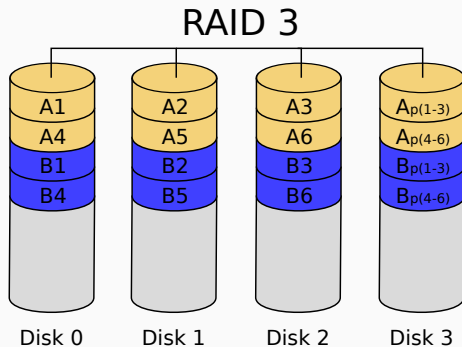


Benefits / Drawbacks



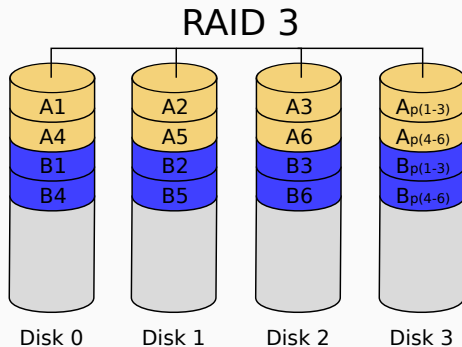
Benefits / Drawbacks

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



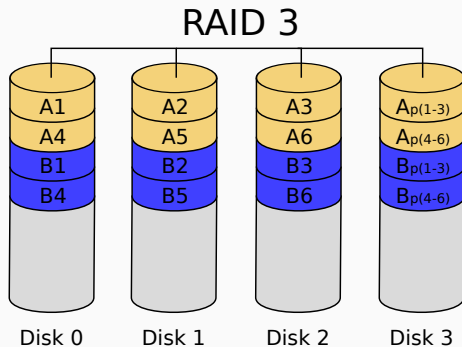
Benefits / Drawbacks

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



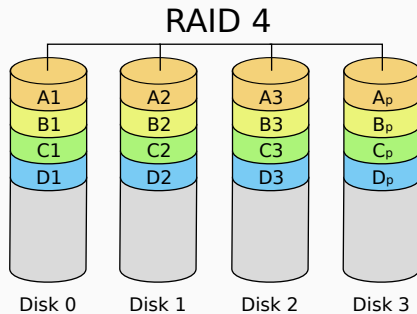
Benefits / Drawbacks

- + 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



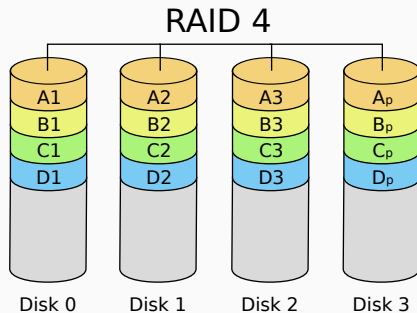
Benefits / Drawbacks

- + 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



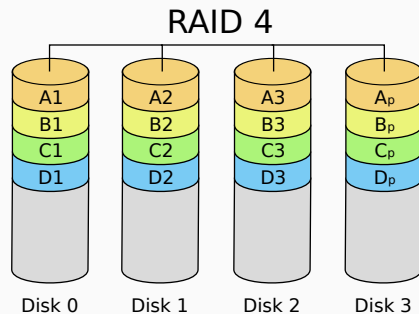
What is that?

- Have a dedicated parity disk

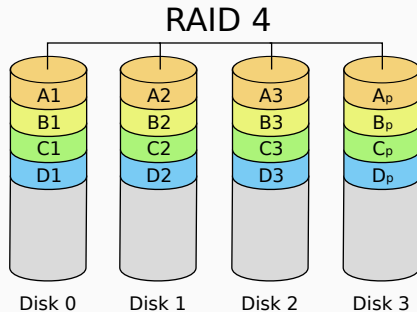


What is that?

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

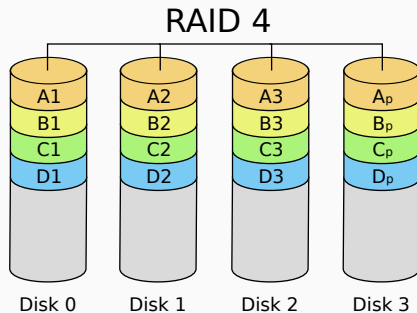


Benefits / Drawbacks



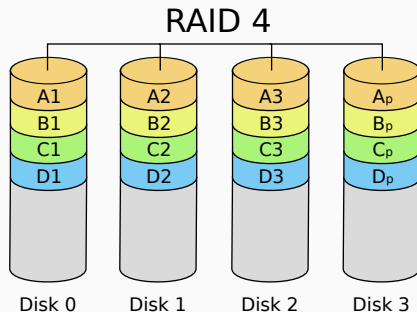
Benefits / Drawbacks

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



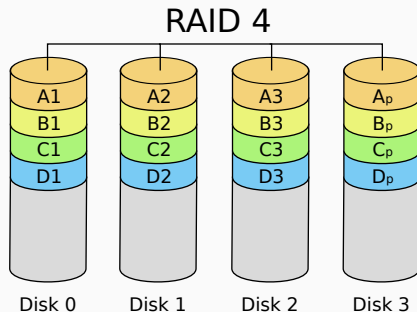
Benefits / Drawbacks

- + You can lose a disk and restore it using the parity
- + Good read performance



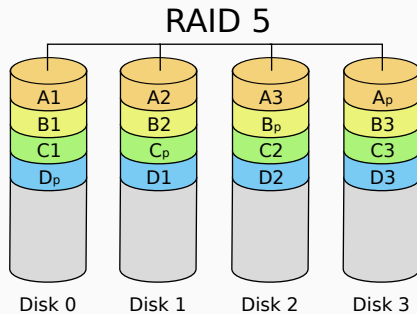
Benefits / Drawbacks

- + 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



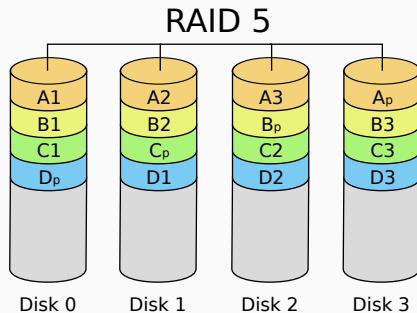
Benefits / Drawbacks

- + 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)



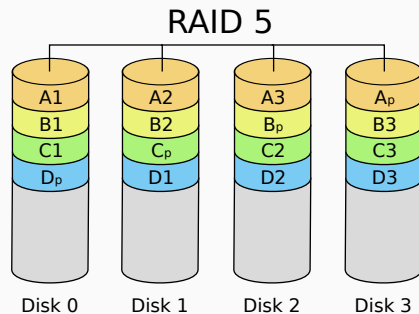
What is that?

- Stripe data at the *block* level

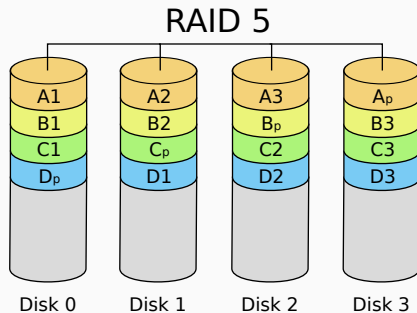


What is that?

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

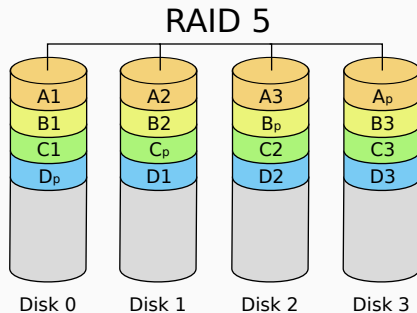


Benefits / Drawbacks



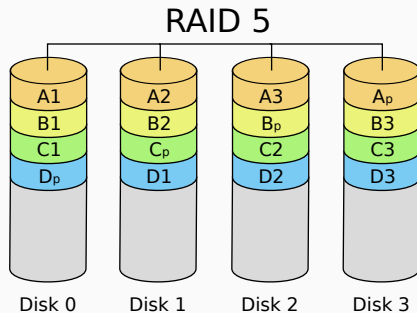
Benefits / Drawbacks

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



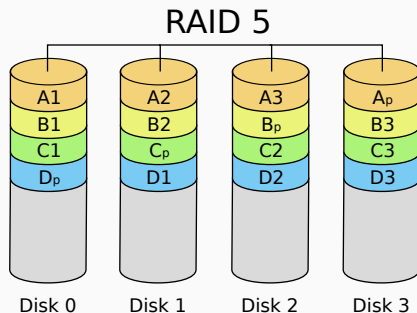
Benefits / Drawbacks

- + You can lose a disk and restore it using the parity
- + Good read performance



Benefits / Drawbacks

- + You can lose a disk and restore it using the parity
- + Good read performance
- + Okay write performance



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



XKCD 1360 - Old Files

FRAGEN?



<https://forms.gle/9CwJSKidKibubran9>

Bis nächste Woche