

# Betriebssysteme

## 11. Tutorium - Storage

---

Peter Bohner

24. Januar 2025

ITEC - Operating Systems Group

# 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?

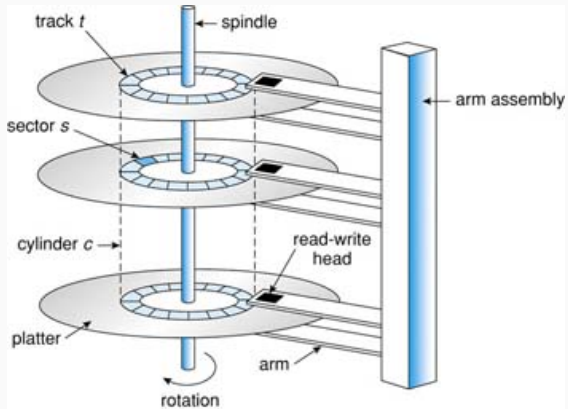
- Heads
- Arms
- Platters

## What parts can you find in a hard disk?

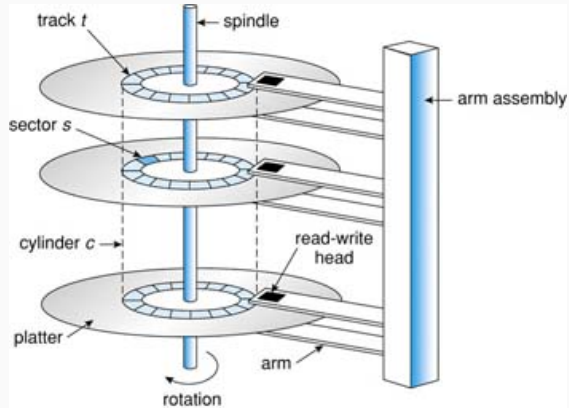
- Heads
- Arms
- Platters
- Spindle



# Hard Disk Layout

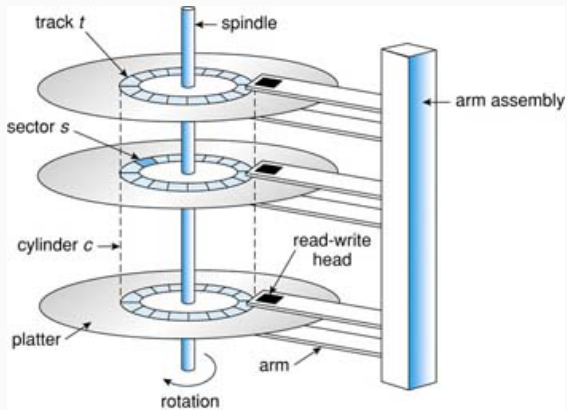


# Hard Disk Layout



What do they do?

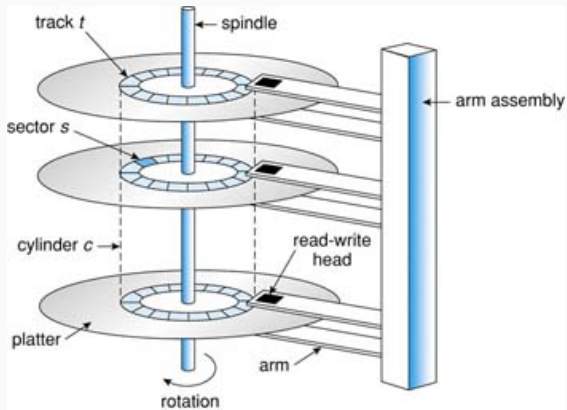
# Hard Disk Layout



## What do they do?

- Spindle: Spin connected platters!

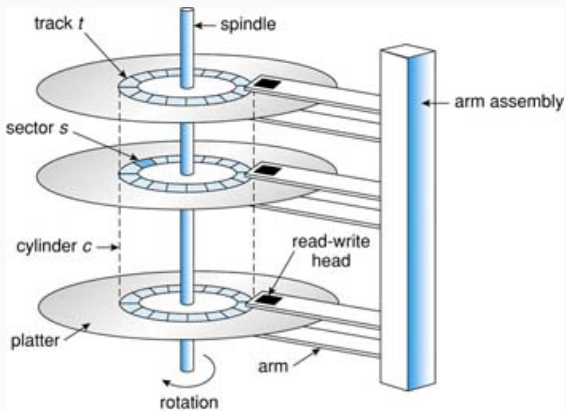
# Hard Disk Layout



## What do they do?

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

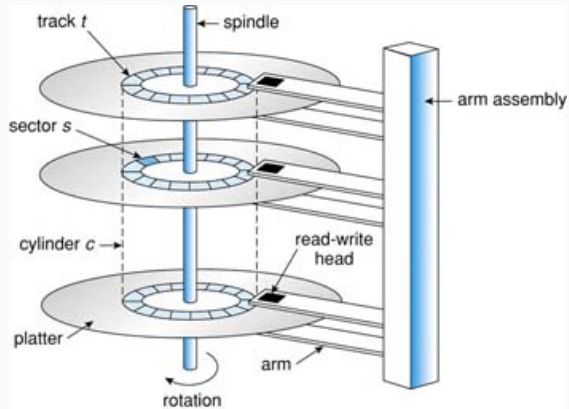
# Hard Disk Layout



## What do they do?

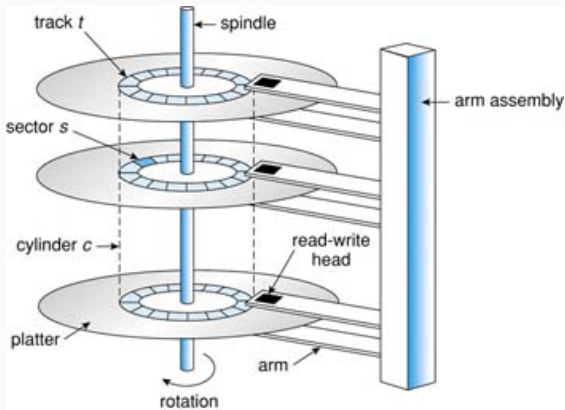
- Spindle: Spin connected platters!
- Head: Read/Write
- Arm: Move heads

# Hard Disk Layout



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

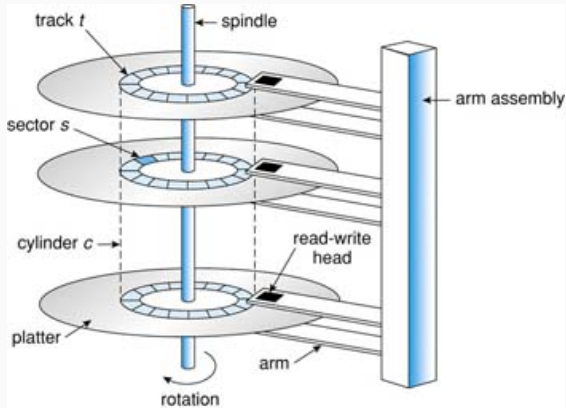
# Hard Disk Layout



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, thin red rectangular border. To the left of the red border, 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

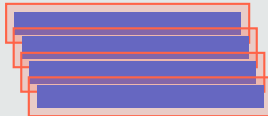


### Shingled Magnet Recording

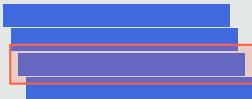


# Hard Disk Layout - More Data!

## Shingled Magnet Recording

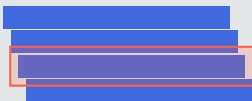


## Shingled Magnet Recording



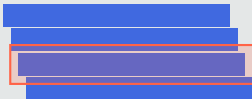
What happens when you write to this track?

## Shingled Magnet Recording



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

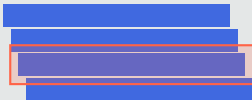
## Shingled Magnet Recording



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*

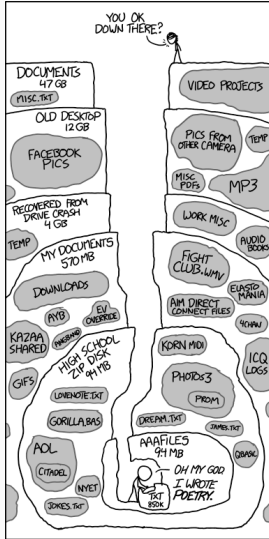
## How can such a device interface with the OS?

- 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*



XKCD 1360 - Old Files

## FRAGEN?



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

Bis nächste Woche