# Betriebssysteme

## 13. Tutorium - Files Systems

Peter Bohner

6. Februar 2025
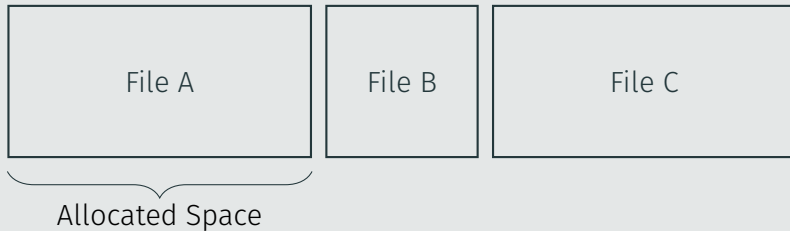
ITEC - Operating Systems Group

# Disk Space Allocation

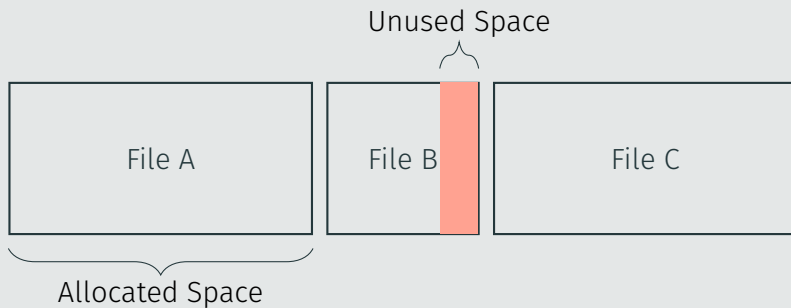# Contiguous Allocation

## What is that?

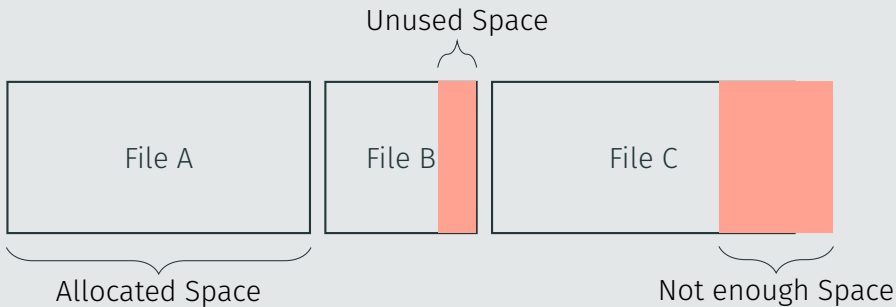## What is that?

File A

File B

File C

Allocated Space

## What is that?

Unused Space

File A

File B

File C

Allocated Space

## What is that?



Unused Space

| File A | File B | File C |

Allocated Space

Not enough Space

# Contiguous Allocation

## What is that?

Unused Space

File A

File B

File C

Allocated Space

Not enough Space

## Challenges

- Sizing your block (internal fragmentation, growth)
- External fragmentation

# Chained Allocation

## What's that?

## What's that?
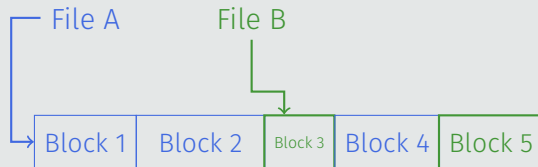
File A        File B

| Block 1 | Block 2 | Block 3 | Block 4 | Block 5 |

## What's that?

## What's that?



File A    File B

| Block 1 | Block 2 | Block 3 | Block 4 | Block 5 |

## What's that?

## What's that?

# Chained Allocation

## What's that?

File A   File B

| Block 1 | Block 2 | Block 3 | Block 4 | Block 5 |

## Benefits? Drawbacks?

# Chained Allocation

## What's that?



## Benefits? Drawbacks?

+ No longer need a contiguous chunk

# Chained Allocation

## What's that?



## Benefits? Drawbacks?

+ No longer need a contiguous chunk
- Only sequential access

# Chained Allocation

## What's that?



## Benefits? Drawbacks?

+ No longer need a contiguous chunk

- Only sequential access

- A single corrupted pointer is very bad news

## What is that?

## What is that?

## What is that?

# Linked List Allocation With FAT

## What is that?

## What is that?



⇐ File A

## What is that?

## What is that?

## What is that?

## What is that?

## What is that?

## What is that?



| 0 | |
|---|---|
| 1 | 7 | ⟸ File A
| 2 | -1 |
| 3 | 5 |
| 4 | -1 |
| 5 | 2 |
| 6 | 8 | ⟸ File B
| 7 | 3 |
| 8 | 4 |

3

Benefits? Drawbacks?

#### Benefits? Drawbacks?

+ Hopefully fits in RAM $\Rightarrow$ Iteration fast

### Benefits? Drawbacks?

+ Hopefully fits in RAM $\Rightarrow$ Iteration fast
+ Even if it doesn't fit: Less disk seeks

## Benefits? Drawbacks?

+ Hopefully fits in RAM ⇒ Iteration fast
+ Even if it doesn't fit: Less disk seeks
- Size depends on size of hard disk (one entry per Block)

### Benefits? Drawbacks?

- + Hopefully fits in RAM $\Rightarrow$ Iteration fast
- + Even if it doesn't fit: Less disk seeks
- - Size depends on size of hard disk (one entry per Block)
- - Might be too large to be cached in RAM (on large disk)

## What is that?

| | | |
|---|---|---|
| 0 | 1 | 2 |
| 3 | 4 | 5 |
| 6 | 7 | 8 |

## What is that?

## What is that?

## What is that?



Index Block **1**

## What is that?

Benefits? Drawbacks?

### Benefits? Drawbacks?

+ Hopefully fits in RAM $\Rightarrow$ Lookups fast

### Benefits? Drawbacks?

+ Hopefully fits in RAM $\Rightarrow$ Lookups fast
+ Random access possible (just fetch the $i$th block)

## Benefits? Drawbacks?

- + Hopefully fits in RAM $\Rightarrow$ Lookups fast
- + Random access possible (just fetch the $i$th block)
- - File size limited by how many references fit into one block

### Benefits? Drawbacks?

+ Hopefully fits in RAM $\Rightarrow$ Lookups fast
+ Random access possible (just fetch the $i$th block)
- File size limited by how many references fit into one block
- Wasteful for small files

### Benefits? Drawbacks?

+ Hopefully fits in RAM $\Rightarrow$ Lookups fast
+ Random access possible (just fetch the $i$th block)
- File size limited by how many references fit into one block
- Wasteful for small files

### How could you store *Huge* files?

### Benefits? Drawbacks?

+ Hopefully fits in RAM $\Rightarrow$ Lookups fast
+ Random access possible (just fetch the $i$th block)
- File size limited by how many references fit into one block
- Wasteful for small files

### How could you store *Huge* files?

- Increase the Block size

### Benefits? Drawbacks?

- + Hopefully fits in RAM $\Rightarrow$ Lookups fast
- + Random access possible (just fetch the *i*th block)
- - File size limited by how many references fit into one block
- - Wasteful for small files

### How could you store *Huge* files?

- - Increase the Block size
- + Indirection! Make the index block link to other index blocks (like pagetables)
  You can also mix that: First *N* pointers point to data blocks, next to indirect
  blocks, next to double-indirect blocks, etc.

## What does an inode look like?

| | data** | data* | data |
|---|---|---|---|
| Size | | | |
| Access flags | | | |
| ⋮ | | | |

Direct block ⟶ Data

Direct block 2

Indirect block ⟶ Pointer ⟶ Data

Indirect block 2

Double Ind. block ⟶ Pointer ⟶ Pointer ⟶ Data

Double Ind. block 2

### What is the maximum file size?

Assume that disk blocks are 8 KiB in size and a pointer to a disk block is 4 bytes long. An inode contains 12 pointers to direct blocks, and one pointer to a single, double, and triple indirect block, respectively.

## What is the maximum file size?

Assume that disk blocks are 8 KiB in size and a pointer to a disk block is 4 bytes long. An inode contains 12 pointers to direct blocks, and one pointer to a single, double, and triple indirect block, respectively.

## Calculating the amount of pointers to blocks

Pointers per Block: 8 *KiB* / 4 *bytes* $= 2^{13}$ / $2^2 = 2^{11} = 2048$
  Direct pointers

## Indexed Allocations With Inodes

### What is the maximum file size?

Assume that disk blocks are 8 KiB in size and a pointer to a disk block is 4 bytes long. An inode contains 12 pointers to direct blocks, and one pointer to a single, double, and triple indirect block, respectively.

### Calculating the amount of pointers to blocks

Pointers per Block: 8 *KiB* / 4 *bytes* $= 2^{13}$ / $2^2 = 2^{11} = 2048$
  Direct pointers              12 pointers to blocks
  One single-indirect block

### What is the maximum file size?

Assume that disk blocks are 8 KiB in size and a pointer to a disk block is 4 bytes long. An inode contains 12 pointers to direct blocks, and one pointer to a single, double, and triple indirect block, respectively.

### Calculating the amount of pointers to blocks

Pointers per Block: 8 *KiB* / 4 *bytes* $= 2^{13}$ / $2^2 = 2^{11} = 2048$

| | |
|---|---|
| Direct pointers | 12 pointers to blocks |
| One single-indirect block | 2048 pointers to blocks |
| One double-indirect block | |

### What is the maximum file size?

Assume that disk blocks are 8 KiB in size and a pointer to a disk block is 4 bytes long. An inode contains 12 pointers to direct blocks, and one pointer to a single, double, and triple indirect block, respectively.

### Calculating the amount of pointers to blocks

Pointers per Block: 8 $KiB$ / 4 $bytes = 2^{13}$ / $2^2 = 2^{11} = 2048$

| | |
|---|---|
| Direct pointers | 12 pointers to blocks |
| One single-indirect block | 2048 pointers to blocks |
| One double-indirect block | $2048 \cdot 2048$ pointers to blocks |
| One triple-indirect block | |

### What is the maximum file size?

Assume that disk blocks are 8 KiB in size and a pointer to a disk block is 4 bytes long. An inode contains 12 pointers to direct blocks, and one pointer to a single, double, and triple indirect block, respectively.

### Calculating the amount of pointers to blocks

Pointers per Block: 8 *KiB* / 4 *bytes* $= 2^{13}$ / $2^2 = 2^{11} = 2048$

| | |
|---|---|
| Direct pointers | 12 pointers to blocks |
| One single-indirect block | 2048 pointers to blocks |
| One double-indirect block | $2048 \cdot 2048$ pointers to blocks |
| One triple-indirect block | $2048 \cdot 2048 \cdot 2048$ pointers to blocks |

## Indexed Allocations With Inodes

### What is the maximum file size?

Assume that disk blocks are 8 KiB in size and a pointer to a disk block is 4 bytes long. An inode contains 12 pointers to direct blocks, and one pointer to a single, double, and triple indirect block, respectively.

### Calculating the amount of pointers to blocks

Pointers per Block: 8 *KiB* / 4 *bytes* $= 2^{13}$ / $2^2 = 2^{11} = 2048$

| | |
|---|---|
| Direct pointers | 12 pointers to blocks |
| One single-indirect block | 2048 pointers to blocks |
| One double-indirect block | $2048 \cdot 2048$ pointers to blocks |
| One triple-indirect block | $2048 \cdot 2048 \cdot 2048$ pointers to blocks |

### Final result

$(12 + 2048 + 2048^2 + 2048^3) \cdot 8 KiB \approx 64 TiB$

# File System Implementation

# File System Implementation

## Hard links

### Hard links

Pointer to the same inode.

⇒ *Everything* is identical: Size, content, access time, mode, ...

Can therefore *not* cross file system boundaries.

## File System Implementation

### Hard links

Pointer to the same inode.

⇒ *Everything* is identical: Size, content, access time, mode, ...

Can therefore *not* cross file system boundaries.

### Symlinks

## Hard links

Pointer to the same inode.
⇒ *Everything* is identical: Size, content, access time, mode, ...
Can therefore *not* cross file system boundaries.

## Symlinks

Are (mostly) normal files containing a filepath with their own access time, mode, .... When programs access the file the OS transparently:

1. Reads the contents of the symlink file
2. Resolves the file path it read
3. Performs the operation on that path instead

Can cross file system boundaries or point to non-existent files or change what they point to if the file is moved, ...

### Hard and symlink renames

```
1  echo "Hello" > test.txt
2  ln test.txt hardlink.txt
3  ln -s test.txt symlink.txt
4
5  mv test.txt renamed_test.txt
6  // Is the hardlink / symlink broken?
```

10

# Having Fun With Paths I

## Hard and symlink renames

```
1  echo "Hello" > test.txt
2  ln test.txt hardlink.txt
3  ln -s test.txt symlink.txt
4
5  mv test.txt renamed_test.txt
6  // Is the hardlink / symlink broken?
```

## Reaching for the stars

Does this work (on my machine ;)?

```
1  cd /tmp/
2  echo "Hello" > test.txt
3  ln test.txt ~/test_link.txt
```

# Having Fun With Paths I

## Hard and symlink renames

```
1  echo "Hello" > test.txt
2  ln test.txt hardlink.txt
3  ln -s test.txt symlink.txt
4
5  mv test.txt renamed_test.txt
6  // Is the hardlink / symlink broken?
```

## Reaching for the stars

Does this work (on my machine ;)?

```
1  cd /tmp/
2  echo "Hello" > test.txt
3  ln test.txt ~/test_link.txt
```

No, as **/tmp** is mounted as tempfs and **/** as ext4.

## Hard and symlink renames

```
1  echo "Hello" > test.txt
2  ln test.txt hardlink.txt
3  ln -s test.txt symlink.txt
4
5  cp test.txt renamed_test.txt
6  rm test.txt
7  // Is the hardlink / symlink broken?
8
```

## Hard and symlink renames

```
1  echo "Hello" > test.txt
2  ln test.txt hardlink.txt
3  ln -s test.txt symlink.txt
4
5  cp test.txt renamed_test.txt
6  rm test.txt
7  // Is the hardlink / symlink broken?
8
9  echo "Hello" > test.txt
10 // Is the hardlink / symlink still broken?
```

### Hard and symlink renames

```
1   echo "Hello" > test.txt
2   ln test.txt hardlink.txt
3   ln -s test.txt symlink.txt
4
5   cp test.txt renamed_test.txt
6   rm test.txt
7   // Is the hardlink / symlink broken?
8
9   echo "Hello" > test.txt
10  // Is the hardlink / symlink still broken?
```

If you delete a hardlink, how can the FS know when it can delete the file?

If you delete a hardlink, how can the FS know when it can delete the file?
If the refcount stored in the inode is zero

If you delete a hardlink, how can the FS know when it can delete the file?

If the refcount stored in the inode is zero

How can directories be implemented? What information is stored in them?

If you delete a hardlink, how can the FS know when it can delete the file?

If the refcount stored in the inode is zero

How can directories be implemented? What information is stored in them?

As a normal file with variable-length entries consisting of

- The filename
- The inode that represents that file

### Which of the following data are typically stored in an inode?

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

### Which of the following data are typically stored in an inode?

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

## Which of the following data are typically stored in an inode?

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

## Which of the following data are typically stored in an inode?

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

Which of the following data are typically stored in an inode?

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

13

**Which of the following data are typically stored in an inode?**

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

13

## Which of the following data are typically stored in an inode?

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

**Which of the following data are typically stored in an inode?**

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

Which of the following data are typically stored in an inode?

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

## Which of the following data are typically stored in an inode?

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

**Which of the following data are typically stored in an inode?**

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

**Which of the following data are typically stored in an inode?**

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

**Which of the following data are typically stored in an inode?**

1. filename
2. name of containing directory
3. file size
4. file type
5. number of symlinks to file
6. name/location of symlinks
7. number of hardlinks
8. name/location of hardlinks
9. access rights
10. timestamps (last access/modify)
11. file contents
12. ordered list of blocks occupied by file

# Virtual File System

What is the purpose of the VFS layer in an OS?

### What is the purpose of the VFS layer in an OS?

Abstraction. Provide a high-level API that works no matter what file system you use (XFS, ext4, NTFS, SMB, IliasFUSE, …).

What is the purpose of the VFS layer in an OS?

Abstraction. Provide a high-level API that works no matter what file system you use (XFS, ext4, NTFS, SMB, IliasFUSE, …).

What are some drawbacks / problems?

### What is the purpose of the VFS layer in an OS?

Abstraction. Provide a high-level API that works no matter what file system you use (XFS, ext4, NTFS, SMB, IliasFUSE, …).

### What are some drawbacks / problems?

- Lowest common denominator: Might hide special FS features

#### What is the purpose of the VFS layer in an OS?

Abstraction. Provide a high-level API that works no matter what file system you use (XFS, ext4, NTFS, SMB, IliasFUSE, …).

#### What are some drawbacks / problems?

- Lowest common denominator: Might hide special FS features
- $\Rightarrow$ Allow accessing underlying FS (e.g. via `ioctl`)
- $\Rightarrow$ Lose compatibility with other file systems

## Virtual File System

### What is the purpose of the VFS layer in an OS?

Abstraction. Provide a high-level API that works no matter what file system you use (XFS, ext4, NTFS, SMB, IliasFUSE, …).

### What are some drawbacks / problems?

- Lowest common denominator: Might hide special FS features
- ⇒ Allow accessing underlying FS (e.g. via `ioctl`)
- ⇒ Lose compatibility with other file systems

### What happens when you mount a filesystem?

E.g. `mount /dev/sda1 /mnt`.

## Virtual File System

### What is the purpose of the VFS layer in an OS?

Abstraction. Provide a high-level API that works no matter what file system you use (XFS, ext4, NTFS, SMB, IliasFUSE, …).

### What are some drawbacks / problems?

- Lowest common denominator: Might hide special FS features
- $\Rightarrow$ Allow accessing underlying FS (e.g. via `ioctl`)
- $\Rightarrow$ Lose compatibility with other file systems

### What happens when you mount a filesystem?

E.g. `mount /dev/sda1 /mnt`.

Makes files from that file system accessible. It will make the given path (e.g. `/mnt`) the *root* of the new file system. Any access to `/tmp/*` is stripped of the `/tmp/` prefix and then searched in the mounted file system.

## FUSE

File System in **USE**rspace.

- Allows users to write their own file system without writing kernel code

## FUSE

File System in **USE**rspace.

- Allows users to write their own file system without writing kernel code
- Runs as an unprivileged user process

## FUSE

File System in **USE**rspace.

- Allows users to write their own file system without writing kernel code
- Runs as an unprivileged user process
- Is really awsome! You suddenly can use all your normal tools on whatever the FUSE filesystem exposes!

# File System Cache

## File System Cache

Is data persisted after a call to `write`? If not, why and when?

Is data persisted after a call to `write`? If not, why and when?

- Data is probably buffered by libc first

# File System Cache

## Is data persisted after a call to `write`? If not, why and when?

- Data is probably buffered by libc first
- Even after that, it doesn't directly get written to disk. Enter...

# File System Cache

## Is data persisted after a call to `write`? If not, why and when?

- Data is probably buffered by libc first
- Even after that, it doesn't directly get written to disk. Enter...The *File System Cache*!

Is data persisted after a call to `write`? If not, why and when?

- Data is probably buffered by libc first
- Even after that, it doesn't directly get written to disk. Enter...The *File System Cache*!
- Why? Writing to disk every time is *reeeaaaally slow*, so it is probably cached in memory and flushed in chunks

Is data persisted after a call to `write`? If not, why and when?

- Data is probably buffered by libc first
- Even after that, it doesn't directly get written to disk. Enter...The *File System Cache*!
- Why? Writing to disk every time is *reeeaaaally slow*, so it is probably cached in memory and flushed in chunks

The standard question: Why is a cache even helpful?

### Is data persisted after a call to `write`? If not, why and when?

- Data is probably buffered by libc first
- Even after that, it doesn't directly get written to disk. Enter...The *File System Cache*!
- Why? Writing to disk every time is *reeeaaaally slow*, so it is probably cached in memory and flushed in chunks

### The standard question: Why is a cache even helpful?

- Temporal and spatial locality

## File System Cache

### Is data persisted after a call to `write`? If not, why and when?

- Data is probably buffered by libc first
- Even after that, it doesn't directly get written to disk. Enter...The *File System Cache*!
- Why? Writing to disk every time is *reeeaaaally slow*, so it is probably cached in memory and flushed in chunks

### The standard question: Why is a cache even helpful?

- Temporal and spatial locality
- When exactly is it flushed to disk then? Is that even important to know/control?

## File System Cache

### Is data persisted after a call to `write`? If not, why and when?

- Data is probably buffered by libc first
- Even after that, it doesn't directly get written to disk. Enter...The *File System Cache*!
- Why? Writing to disk every time is *reeeaaaally slow*, so it is probably cached in memory and flushed in chunks

### The standard question: Why is a cache even helpful?

- Temporal and spatial locality
- When exactly is it flushed to disk then? Is that even important to know/control?
- Using `flush` (for buffers) and `fsync` (for the page cache) and after some time by a daemon

## Let's talk about the page cache

| Userspace buffer | Page cache | Physical disk |
|---|---|---|
|  |  |  |
|  |  |  |
|  |  |  |

## Let's talk about the page cache

Userspace buffer

Page cache

Physical disk

# Page cache interactions

## Let's talk about the page cache

| Userspace buffer | | Page cache | Physical disk |



Userspace buffer →write→ Page cache → Physical disk

## Let's talk about the page cache

| Userspace buffer | | Page cache | | Physical disk |
|---|---|---|---|---|

write

## Let's talk about the page cache

| Userspace buffer | | Page cache | | Physical disk |
|---|---|---|---|---|

write → fsync / Time →

## Let's talk about the page cache

| Userspace buffer | | Page cache | | Physical disk |
|---|---|---|---|---|

write → fsync / Time →

How large would you make it? Do you give it a fixed size or let it grow?

How large would you make it? Do you give it a fixed size or let it grow?

Fix Easy to implement

How large would you make it? Do you give it a fixed size or let it grow?

Fix Easy to implement

Fix Can give hard guarantees (real-time systems anyone?)

How large would you make it? Do you give it a fixed size or let it grow?

Fix  Easy to implement

Fix  Can give hard guarantees (real-time systems anyone?)

Fix  Can't adapt to different workloads

# Sizing the page cache

How large would you make it? Do you give it a fixed size or let it grow?

Fix  Easy to implement

Fix  Can give hard guarantees (real-time systems anyone?)

Fix  Can't adapt to different workloads

When do you load data in the page cache?

18

How large would you make it? Do you give it a fixed size or let it grow?

Fix  Easy to implement

Fix  Can give hard guarantees (real-time systems anyone?)

Fix  Can't adapt to different workloads

When do you load data in the page cache?

Read-ahead: Read more data in the cache than requested

# Sizing the page cache

How large would you make it? Do you give it a fixed size or let it grow?

Fix  Easy to implement

Fix  Can give hard guarantees (real-time systems anyone?)

Fix  Can't adapt to different workloads

When do you load data in the page cache?

Read-ahead: Read more data in the cache than requested

+ Good sequential performance

# Sizing the page cache

## How large would you make it? Do you give it a fixed size or let it grow?

- Fix Easy to implement
- Fix Can give hard guarantees (real-time systems anyone?)
- Fix Can't adapt to different workloads

## When do you load data in the page cache?

Read-ahead: Read more data in the cache than requested

- + Good sequential performance
- + Improved throughput if files are sequential

How large would you make it? Do you give it a fixed size or let it grow?

Fix  Easy to implement

Fix  Can give hard guarantees (real-time systems anyone?)

Fix  Can't adapt to different workloads

When do you load data in the page cache?

Read-ahead: Read more data in the cache than requested

+ Good sequential performance

+ Improved throughput if files are sequential

- Wasted time and memory if not needed

How could you implement `mmap` with the file system cache?

How could you implement `mmap` with the file system cache?

Just map the file system cache in the process's memory (and handle faults!)
⇒ Acts as a shared memory segment

You synchronize accesses yourself :(

How could you implement `mmap` with the file system cache?

Just map the file system cache in the process's memory (and handle faults!)
⇒ Acts as a shared memory segment

You synchronize accesses yourself :(

`read()` and `write()` and the file system cache

The File System Cache and mmap

---

How could you implement `mmap` with the file system cache?

Just map the file system cache in the process's memory (and handle faults!)
⇒ Acts as a shared memory segment

You synchronize accesses yourself :(

---

`read()` and `write()` and the file system cache

- `read()`: Copy data from cache to your application buffers

How could you implement `mmap` with the file system cache?

Just map the file system cache in the process's memory (and handle faults!)
⇒ Acts as a shared memory segment

You synchronize accesses yourself :(

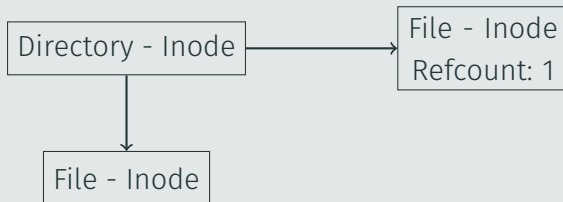`read()` and `write()` and the file system cache

- `read()`: Copy data from cache to your application buffers
- `write()`: Copy data from your application buffers to the file system cache

*File system synchronizes access!*

# Modern File Systems

Your computer crashes (e.g. your power fails). What horrible death does your file system die?

A fun game for the whole family and your „why would I need backups" crowd

Your computer crashes (e.g. your power fails). What horrible death does your file system die?

A fun game for the whole family and your „why would I need backups" crowd

Your computer crashes (e.g. your power fails). What horrible death does your file system die?

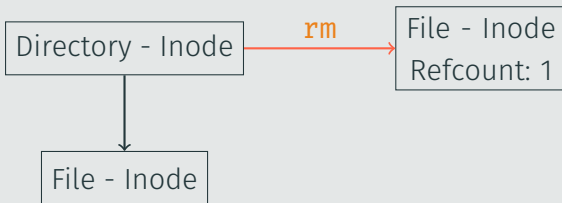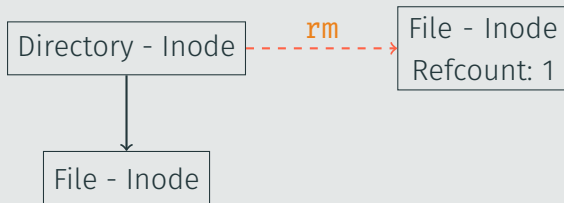A fun game for the whole family and your „why would I need backups" crowd

Your computer crashes (e.g. your power fails). What horrible death does your file system die?

A fun game for the whole family and your „why would I need backups" crowd



The inode has a refcount of 1 but is no longer referenced!

## Where is Data?

File 1

File 2

## Where is Data?

## Where is Data?



Truncate, delete → [blue block] ← Reassigned to this file

File 1

File 2

Both files might point to this block!

## Where is Data?

File 1

## Where is Data?



File 1

Truncate, delete

## Where is Data?



Block still referenced but marked free

## Where is Data?



Block still referenced but marked free

Or multiple directory entries with the same name, as another problem

What might be the outcome?

```
1  // create a file
2  echo "Hey" > a.txt
3  // And a second
4  echo "Hello" > b.txt
5  // CRASH
```

## What might be the outcome?

```
1  // create a file
2  echo "Hey" > a.txt
3  // And a second
4  echo "Hello" > b.txt
5  // CRASH
```

- Both files exist

## What might be the outcome?

```
1  // create a file
2  echo "Hey" > a.txt
3  // And a second
4  echo "Hello" > b.txt
5  // CRASH
```

- Both files exist
- No files exist

## What might be the outcome?

```
1  // create a file
2  echo "Hey" > a.txt
3  // And a second
4  echo "Hello" > b.txt
5  // CRASH
```

- Both files exist
- No files exist
- Only **a** exists

## What might be the outcome?

```
1  // create a file
2  echo "Hey" > a.txt
3  // And a second
4  echo "Hello" > b.txt
5  // CRASH
```

- Both files exist
- No files exist
- Only a exists
- Only b exists ‼

Reordering is *allowed* (unless you take precautions)

How could you detect those inconsistencies?

Use the **fsck** (**f**ile **s**ystem **c**onsistency chec**k**) program. And what could that do?

How could you detect those inconsistencies?

Use the **fsck** (**f**ile **s**ystem **c**onsistency chec**k**) program. And what could that do?

- Check whether all data blocks are referenced by *exactly one* inode

# Error Detection

## How could you detect those inconsistencies?

Use the **fsck** (**f**ile **s**ystem **c**onsistency chec**k**) program. And what could that do?

- Check whether all data blocks are referenced by *exactly one* inode
- Do all directories contain valid entries (and just one with a given name)

# Error Detection

## How could you detect those inconsistencies?

Use the **fsck** (**f**ile **s**ystem **c**onsistency chec**k**) program. And what could that do?

- Check whether all data blocks are referenced by *exactly one* inode
- Do all directories contain valid entries (and just one with a given name)
- Is every directory / file connected to the directory tree

### How could you detect those inconsistencies?

Use the **fsck** (**f**ile **s**ystem **c**onsistency chec**k**) program. And what could that do?

- Check whether all data blocks are referenced by *exactly one* inode
- Do all directories contain valid entries (and just one with a given name)
- Is every directory / file connected to the directory tree
- Is the refcount consistent with the number of hardlinks? What do you do if not?
  *refcount* > 0 but no links?

How could you detect those inconsistencies?

Use the **fsck** (**f**ile **s**ystem **c**onsistency chec**k**) program. And what could that do?

- Check whether all data blocks are referenced by *exactly one* inode
- Do all directories contain valid entries (and just one with a given name)
- Is every directory / file connected to the directory tree
- Is the refcount consistent with the number of hardlinks? What do you do if not?
  *refcount* $> 0$ but no links? $\Rightarrow$ Attach to `lost+found`
  Update the refcount to be consistent

How could you detect those inconsistencies?

Use the **fsck** (**f**ile **s**ystem **c**onsistency chec**k**) program. And what could that do?

- Check whether all data blocks are referenced by *exactly one* inode
- Do all directories contain valid entries (and just one with a given name)
- Is every directory / file connected to the directory tree
- Is the refcount consistent with the number of hardlinks? What do you do if not?
  *refcount* > 0 but no links? ⇒ Attach to `lost+found`
  Update the refcount to be consistent

Can this fix your application data?

How could you detect those inconsistencies?

Use the **fsck** (**f**ile **s**ystem **c**onsistency chec**k**) program. And what could that do?

- Check whether all data blocks are referenced by *exactly one* inode
- Do all directories contain valid entries (and just one with a given name)
- Is every directory / file connected to the directory tree
- Is the refcount consistent with the number of hardlinks? What do you do if not?
  *refcount* > 0 but no links? ⇒ Attach to `lost+found`
  Update the refcount to be consistent

Can this fix your application data?

No! It just tries to keep the filesystem internally consistent, it won't find corrupted data blocks

### General principle

Walk over the journal and execute any outstanding entries.

### Let's crash

| Journal | Disk | Journal |
|---|---|---|
| Journal Entry | Disk Entry | Done marker |

### What to do

The happy path, everything's nice

## General principle

Walk over the journal and execute any outstanding entries.

## Let's crash

| Journal | Disk | Journal |
|---------|------|---------|
| Journal Entry | Disk Entry | Done marker |

## What to do

$\Rightarrow$ We didn't write anything!
$\Rightarrow$ Operation failed

## General principle

Walk over the journal and execute any outstanding entries.

## Let's crash

| Journal | Disk | Journal |
|---------|------|---------|

⟶ Journal %§/! ----→ Disk Entry ----→ Done marker ----→

## What to do

⇒ Invalid checksum
⇒ Skip entry
⇒ Operation failed

## General principle

Walk over the journal and execute any outstanding entries.

## Let's crash

Journal                     Disk                     Journal

⟶ | Journal Entry | - - - ⇢ ⌐ Disk %§/! ¬ - - - ⇢ ⌐ Done marker ¬ - - ⇢

## What to do

⇒ Non-terminated journal entry
⇒ Retry and complete operation
⇒ Operation successful

### General principle

Walk over the journal and execute any outstanding entries.

### Let's crash

| Journal | Disk | Journal |
|---------|------|---------|
| Journal Entry | Disk Entry | Done marker |

### What to do

⇒ Consistent state
⇒ Execute operation again
⇒ Operation successful

### General principle

Walk over the journal and execute any outstanding entries.

### Let's crash

| Journal | Disk | Journal |
|---|---|---|
| Journal Entry | Disk Entry | Done marker |

### What to do

$\Rightarrow$ Consistent state
$\Rightarrow$ Don't do anything
$\Rightarrow$ Operation successful

## Physical vs logical logging

What data could you log in the journal?

### Physical vs logical logging

What data could you log in the journal?

- **Logical logging**: Store a high level entry (like: „rename a to b")

#### Physical vs logical logging

What data could you log in the journal?

- **Logical logging**: Store a high level entry (like: „rename a to b")
- **Physical logging**: Store the file system blocks that will be modified

What might be problems with the journal presented before?

What might be problems with the journal presented before?

You write *every block twice*!

What might be problems with the journal presented before?

You write *every block twice*! How could you make that less painful?

# Log Structured File System

What might be problems with the journal presented before?

You write *every block twice*! How could you make that less painful? Only journal *metadata*.

What might be problems with the journal presented before?

You write *every block twice*! How could you make that less painful? Only journal *metadata.*

How could you design a safe system that only needs to write once?

What might be problems with the journal presented before?

You write *every block twice*! How could you make that less painful? Only journal *metadata*.

How could you design a safe system that only needs to write once?

A Log Structured File System (Copy on Write). What is the core idea there?

# Log Structured File System

What might be problems with the journal presented before?

You write *every block twice*! How could you make that less painful? Only journal *metadata*.

How could you design a safe system that only needs to write once?

A Log Structured File System (Copy on Write). What is the core idea there?

*Only write to unused blocks.*

What might be problems with the journal presented before?

You write *every block twice*! How could you make that less painful? Only journal *metadata*.

How could you design a safe system that only needs to write once?

A Log Structured File System (Copy on Write). What is the core idea there?

*Only write to unused blocks.*

Then you can *atomically* update indexing datastructures to point to the new blocks. If it crashes, you either have *all* of the new state or *exactly* the old state.

**Das letzte Tut ist *nächste Woche***

Das wird ein Wiederholungstut von Sachen, die euch noch irgendwo unklar sind oder über die ihr nochmal reden wollt.

⇒ Daher bitte bis nächste Woche Fragen an mich per Mail
peter.bohner@student.kit.edu