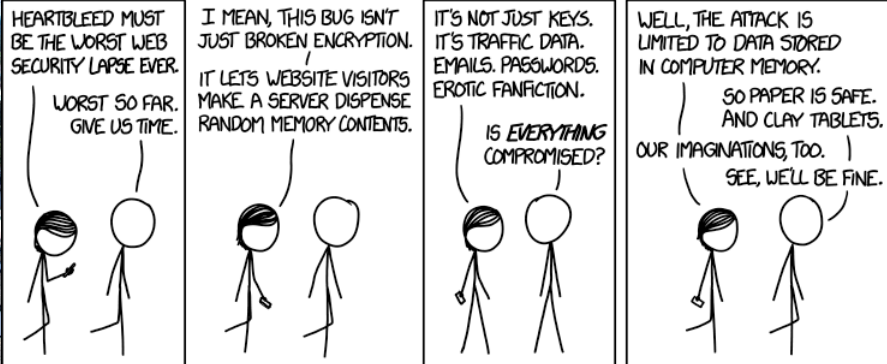**ETH**zürich
*spcl.inf.ethz.ch*
🐦 *@spcl_eth*

**ADRIAN PERRIG & TORSTEN HOEFLER**

# Networks and Operating Systems (252-0062-00)
# Chapter 8: Filesystem Implementation



---

**ETH**zürich
*spcl.inf.ethz.ch*
🐦 *@spcl_eth*

# Thanks for the feedback! ☺

- **Some answers:**
  - I'll provide references to books (I'm not **only** teaching from books)
    *Yesterday: "Operating Systems - An Advanced Course", Chapter 3a*
    *Today: Anderson/Dahlin: 11/12 (partial), 13 (full)*
  - Everything that's mentioned on the slides is essential
    *You should make sure you understand it*
    *This may require listening ☺*
    *Everything else and the stories I tell are optional*
  - Why are your slides not self-contained?
    *See Rebecca Schumann "Digital Slideshows are the scourge of education"*
  - I talked to the assistants to deepen parts of the lecture in the exercise
    *I \*hope\* that works --- they're open for additional feedback!*
  - Next Friday is Karfreitag (probably free?)
    *Exercises will not be replaced but skipped! You may go to Thursday's exercise!*

# File system operations

**We've already seen the file system as a naming scheme.**

**Directory (name space) operations:**
- **Link (bind a name)**
- **Unlink (unbind a name)**
- **Rename**
- **List entries**

---

# Acyclic-Graph Directories

- **Two different names (aliasing)**

- **If *dict* deletes *list* ⇒ dangling pointer**
  **Solutions:**
  - Backpointers, so we can delete all pointers
    Variable size records a problem
  - Backpointers using a daisy chain organization
  - Entry-hold-count solution

- **New directory entry type**
  - **Link** – another name (pointer) to an existing file
  - **Resolve the link** – follow pointer to locate the file

## Slide 1

### General Graph Directory

- **How do we guarantee no cycles? Options:**
  - Allow only links to files and not directories
  - Garbage collection (with cycle collector)
  - Check for cycles when every new link is added
  - Restrict directory links to parents

    *E.g., "." and "..".*

    *All cycles are therefore trivial*



## Slide 2

### Access Control

# Protection

- **File owner/creator should be able to control:**
  - what can be done
  - by whom

- **Types of access**
  - **Read**
  - **Write**
  - **Execute**
  - **Append**
  - **Delete**
  - **List**

# Access control matrix

For a single file or directory:

Principals

| Rights | A | B | C | D | E | F | G | H | J | … |
|---|---|---|---|---|---|---|---|---|---|---|
| Read | ☑ | ☑ | ☑ | | | ☑ | ☑ | | | |
| Write | ☑ | ☑ | | ☑ | | | ☑ | | | |
| Append | ☑ | | | | ☑ | | | | | |
| Execute | ☑ | ☑ | ☑ | ☑ | | | | | | |
| Delete | ☑ | | | | | | | | | |
| List | ☑ | | | | ☑ | | | | | |
| … | | | | | | | | | | |

Problem: how to scalably represent this matrix?

# Row-wise: ACLs

- **Access Control Lists**
  - For each right, list the principals
  - Store with the file
- **Good:**
  - Easy to change rights quickly
  - Scales to large numbers of files
- **Bad:**
  - Doesn't scale to large numbers of principals

# Column-wise: Capabilities

- **Each principal with a right on a file holds a *capability* for that right**
  - Stored with principal, not object (file)
  - Cannot be forged or (sometimes) copied
- **Good:**
  - Very flexible, highly scalable in principals
  - Access control resources charged to principal
- **Bad:**
  - Revocation: hard to change access rights
    (need to keep track of who has what capabilities)

**ETH**zürich

# POSIX (Unix) Access Control

- **Simplifies ACLs: each file identifies 3 principals:**
  - Owner (a single user)
  - Group (a collection of users, defined elsewhere)
  - The World (everyone)
- **For each principal, file defines 3 rights:**
  - Read (or traverse, if a directory)
  - Write (or create a file, if a directory)
  - Execute (or list, if a directory)

---

**ETH**zürich

# Example

```
drwx--x--x  9 htor htor     4096 May  9 13:14 pagai
htor@lenny ~ > ls -l projekte/llvm/llvm-svn              < 09.05.13 19:08:49 >
total 860
drwx--x--x  3 htor htor     4096 Jan 29 15:58 autoconf
drwx--x--x  4 htor htor     4096 Dec 25 13:20 bindings
drwx--x--x  4 htor htor     4096 Jan 29 15:57 cmake
-rw-------  1 htor htor    16401 Dec 25 13:20 CMakeLists.txt
-rw-------  1 htor htor     2782 Jan 29 15:57 CODE_OWNERS.TXT
-rwx------  1 htor htor   658352 Jan 29 15:57 configure
-rw-------  1 htor htor    10048 Dec 25 13:20 CREDITS.TXT
drwxr-xr-x 11 htor htor     4096 Apr  4 11:13 Debug
drwx--x--x 10 htor htor     4096 Jan 29 15:57 docs
drwx--x--x 10 htor htor     4096 Dec 25 13:20 examples
drwx--x--x  4 htor htor     4096 Dec 25 13:20 include
drwx--x--x 18 htor htor     4096 Jan 29 15:58 lib
-rw-------  1 htor htor     3254 Jan 29 15:57 LICENSE.TXT
-rw-------  1 htor htor      752 Dec 25 13:20 LLVMBuild.txt
-rw-------  1 htor htor     1865 Dec 25 13:20 llvm.spec.in
-rw-------  1 htor htor     8618 Jan 29 15:58 Makefile
-rw-------  1 htor htor     2599 Dec 25 13:20 Makefile.common
-rw-------  1 htor htor    12068 Jan 29 15:57 Makefile.config.in
-rw-------  1 htor htor    79586 Jan 29 15:57 Makefile.rules
drwx--x--x  4 htor htor     4096 Dec 25 13:21 projects
-rw-------  1 htor htor      687 Jan 29 15:58 README.txt
drwx--x--x  3 htor htor     4096 Dec 25 13:20 runtime
drwx--x--x 27 htor htor     4096 Jan 29 15:57 test
drwx--x--x 35 htor htor     4096 Dec 25 13:21 tools
drwx--x--x 11 htor htor     4096 Jan 29 15:57 unittests
drwx--x--x 32 htor htor     4096 Jan 29 15:57 utils
```

# Full ACLs

- **POSIX now supports full ACLs**
  - Rarely used, interestingly
  - setfacl, getfacl, …
- **Windows has very powerful ACL support**
  - Arbitrary groups as principals
  - Modification rights
  - Delegation rights

# Our Small Quiz

- **True or false (raise hand)**
  - A file name identifies a string of data on a storage device
  - The file size is part of the file's metadata
  - Names provide a means of abstraction through indirection
  - Names are always assigned at object creation time
  - A context is implicit to a name
  - A context is implicit to an object
  - Name resolve may be specific to a context
  - Each file has exactly one name
  - The call "unlink file" always removes the contents of "file"
  - A fully qualified domain name is resolved recursively starting from the left
  - A full (absolute) path identifies a unique file (piece of data)
  - A full (absolute) path identifies a unique name
  - Stable bindings can be changed with bind()
  - Each name identifies exactly one object in a single context

14

# File types

---

## Is a directory a file?

- **Yes…**
  - Allocated just like a file on disk
  - Has entries in other directories like a file
- **…and No…**
  - Users can't be allowed to read/write to it
    *Corrupt file system data structures*
    *Bypass security mechanisms*
  - File system provides special interface
    *opendir, closedir, readdir, seekdir, telldir, etc.*

# Directory Implementation

- **Linear list of (file name, block pointer) pairs**
  - Simple to program
  - Lookup is slow for lots of files (linear scan)
- **Hash Table – linear list with closed hashing.**
  - Fast name lookup
  - **collisions**
  - fixed size
- **B-Tree – name index, leaves are block pointers**
  - Increasingly common
  - Complex to maintain, but scales well

# File types

- **Other file types treated "specially" by the OS**
- **Simple, common cases:**
  - Executable files
  - Directories, symbolic links, other file system data
- **Some distinguish between text and binary**
- **Some have many types**
  - "Document" or "media" types
  - Used to select default applications, editors, etc.

---

# Unix devices and other file types

- **Unix also uses the file namespace for**
  - Naming I/O devices (/dev)
  - Named pipes (FIFOs)
  - Unix domain sockets
- **More recently:**
  - Process control (/proc)
  - OS configuration and status (/proc, /sys)
- **Plan 9 from Bell Labs**
  - Evolution of Unix: almost *everything* is a file

**ETH**zürich

# Executable files

- **Most OSes recognize binary executables**
  - Sometimes with a "magic number"
  - Will load, dynamically link, and execute in a process
- **Other files are sometimes recognized**
  - E.g. "#!" script files in Unix
    "#!/usr/bin/python"

**ETH**zürich

# File system operations

**File operations:**
- **Create and variants**
  - Unix: `mknod`, `mkfifo`, `ln -s`, …
- **Change access control**
  - Unix: `chmod`, `chgrp`, `chown`, `setfacl`, …
- **Read metadata**
  - Unix: `stat`, `fstat`, …
- **Open**
  - Operation: *file → open file handle*

**ETH**_zürich_

## "Files" vs. "Open Files"

- **Typical operations on files:**
  - Rename, stat, create, delete, etc.
  - **Open**

- **Open creates an "open file handle"**
  - Different class of object
  - Allows reading and writing of file data

---

**ETH**_zürich_

## Open File Interface

**ETH**zürich

# Kinds of files

1. **Byte sequence**
   - The one you're probably familiar with
2. **Record sequence**
   - Fixed (at creation time) records
   - Mainframes or minicomputer OSes of the 70s/80s
3. **Key-based, tree structured**
   - E.g. IBM Indexed Sequential Access Method (ISAM)
   - Mainframe feature, now superseded by databases
   - In other words, moved into libraries

**ETH**zürich

# Byte-sequence files

- **File is a vector of bytes**
  - Can be appended to
  - Can be truncated
  - Can be updated in place
  - Typically no "insert".
- **Accessed as:**
  - Sequential files (rare these days)
  - Random access

2014-04-10

zürich

*spcl.inf.ethz.ch*
*@spcl_eth*

# Random access

- **Support read, write, seek, and tell**
  - State: current position in file
  - **Seek** absolute or relative to current position.
  - **Tell** returns current index
- **Index units:**
  - For byte sequence files, offset in bytes

---

zürich

*spcl.inf.ethz.ch*
*@spcl_eth*

# Record-sequence files

- **File is now a vector of fixed-size records**
  - Can be appended to
  - Can be truncated
  - Can be updated in place
  - Typically no "insert".
- **Record size (and perhaps format) fixed at creation time**
  - Read/write/seek operations take records and record offsets instead of byte addresses

> Compare with databases!

**ETH**zürich

## Memory-mapped files

- **Basic idea: use VM system to cache files**
  - Map file content into virtual address space
  - Set the backing store of region to file
  - Can now access the file using load/store
- **When memory is paged out**
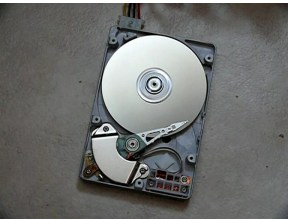  - Updates go back to file instead of swap space

**ETH**zürich

## On-disk data structures

# Disk addressing

- **Disks have tracks, sectors, spindles, etc.**
    - And bad sector maps!
- **More convenient to use *logical block addresses***
    - Treat disk as compact linear array of usable blocks
    - Block size typically 512 bytes
    - Ignore geometry except for performance (later!)
- **Also abstracts other block storage devices**
    - Flash drives (load-levelling, etc.)
    - Storage-area Networks (SANs)
    - Virtual disks (RAM, RAID, etc.)

# Implementation aspects

- **Directories and indexes**
    - Where on the disk is the data for each file?
- **Index granularity**
    - What is the unit of allocation for files?
- **Free space maps**
    - How to allocate more sectors on the disk?
- **Locality optimizations**
    - How to make it go fast in the common case

**ETH**zürich

# File system implementations

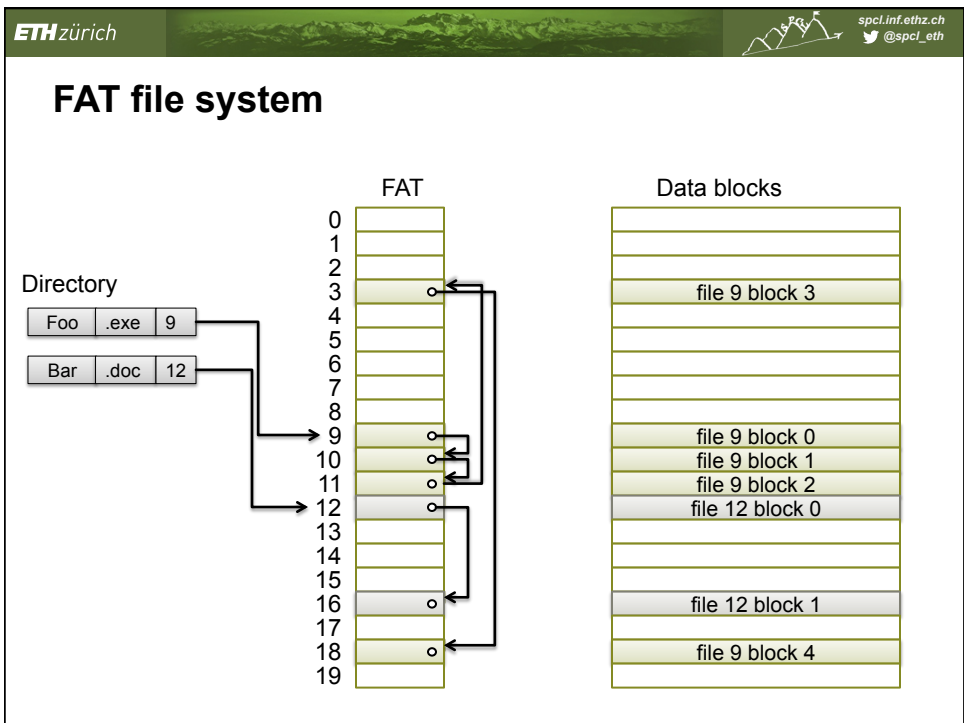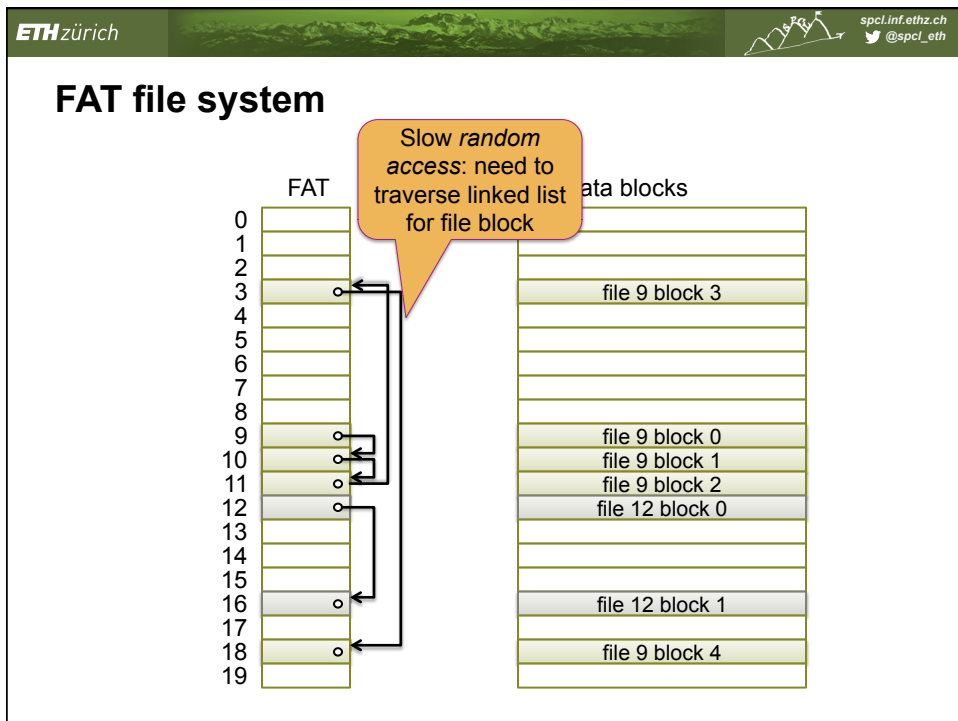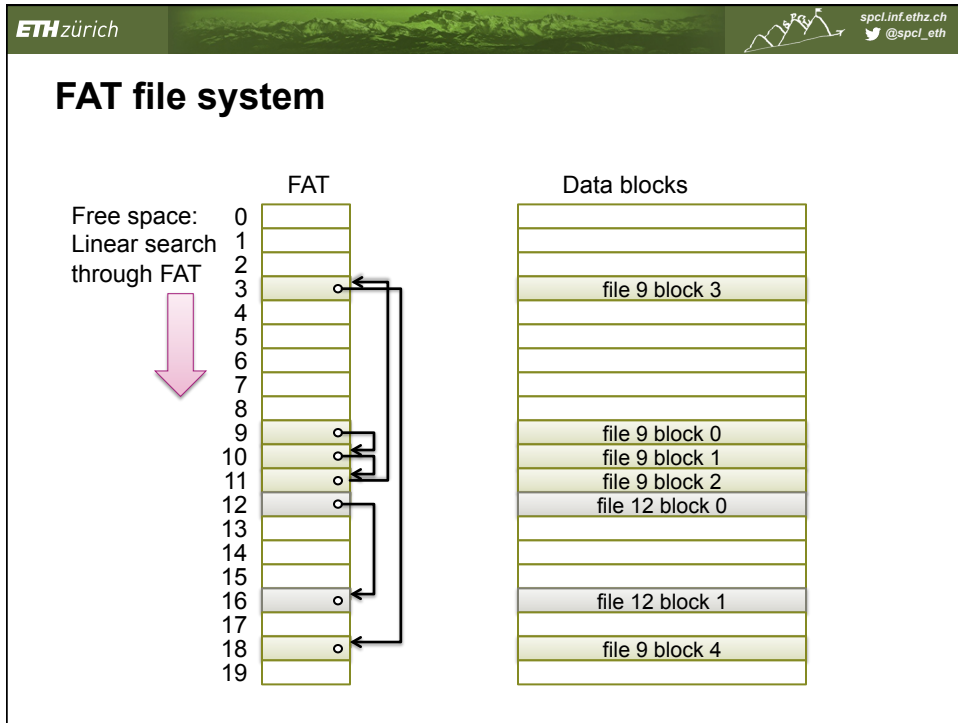|  | FAT | FFS | NTFS | ZFS |
|---|---|---|---|---|
| **Index structure** | Linked list | Fixed, asymmetric tree | Dynamic tree | Dynamic COW tree |
| **Index granularity:** | Block | Block | Extent | Block |
| **Free space management** | FAT Array | Fixed bitmap | Bitmap in file | Log-structured space map |
| **Locality heuristics** | Defragmentation | Block groups, Reserve space | Best fit, Defragmentation | Write anywhere, Block groups |

See book
for details

---

**ETH**zürich

**FAT-32**

# FAT background

- **Very old – dates back to 1970s!**
- **No access control**
- **Very little metadata**
- **Limited volume size**
- **No support for hard links**
- **BUT still extensively used** ☹
  - Flash devices, cameras, phones



# FAT file system

| | FAT | | Data blocks |
|---|---|---|---|
| 0 | | | |
| 1 | | | |
| 2 | | | |
| 3 | ○ | | file 9 block 3 |
| 4 | | | |
| 5 | | | |
| 6 | | | |
| 7 | | | |
| 8 | | | |
| 9 | ○ | | file 9 block 0 |
| 10 | ○ | | file 9 block 1 |
| 11 | ○ | | file 9 block 2 |
| 12 | ○ | | file 12 block 0 |
| 13 | | | |
| 14 | | | |
| 15 | | | |
| 16 | ○ | | file 12 block 1 |
| 17 | | | |
| 18 | ○ | | file 9 block 4 |
| 19 | | | |

Directory

| Foo | .exe | 9 |
|---|---|---|

| Bar | .doc | 12 |
|---|---|---|

# FAT file system



# FAT file system

**FFS**

## Unix Fast File System (FFS)

- **First appeared in BSD in the mid 1980's**
- **Based on original Unix FS, with performance optimizations**
- **Basis for Linux ext{2,3} file systems**

**FFS uses indexed allocation**

Inode array

File inode number from directory entry

Inode

Metadata

Block pointers

Data block
Data block
Data block
Data block

File is represented by an *index block* or ***inode***
- File metadata
- List of blocks for each part of file
- Directory contains pointers to inodes

---

**Inode and file size in FFS**

- **Example:**
  - Inode is 1 block = 4096 bytes
  - Block addresses = 8 bytes
  - Inode metadata = 512 bytes
- **Hence:**
  - (4096-512) / 8 = 448 block pointers
  - 448 * 4096 = 1792kB max. file size

ETH*zürich*

## Unix file system inode format (simplified)

Inode:

(all blocks 4kB)

| File mode |
| --- |
| Owner/group |
| Timestamps |
| Size |
| <other metadata> |
| 12 direct block pointers  … |

Data block
Data block
Data block
…
Data block

---

ETH*zürich*

## Unix file system inode format (simplified)

Inode:

(all blocks 4kB)

| File mode |
| --- |
| Owner/group |
| Timestamps |
| Size |
| <other metadata> |
| 12 direct block pointers  … |

Data block
Data block
Data block
…
Data block

Question:
How to extend file size if there are no more block pointers in the Inode?

Unix file system inode format (**simplified**)



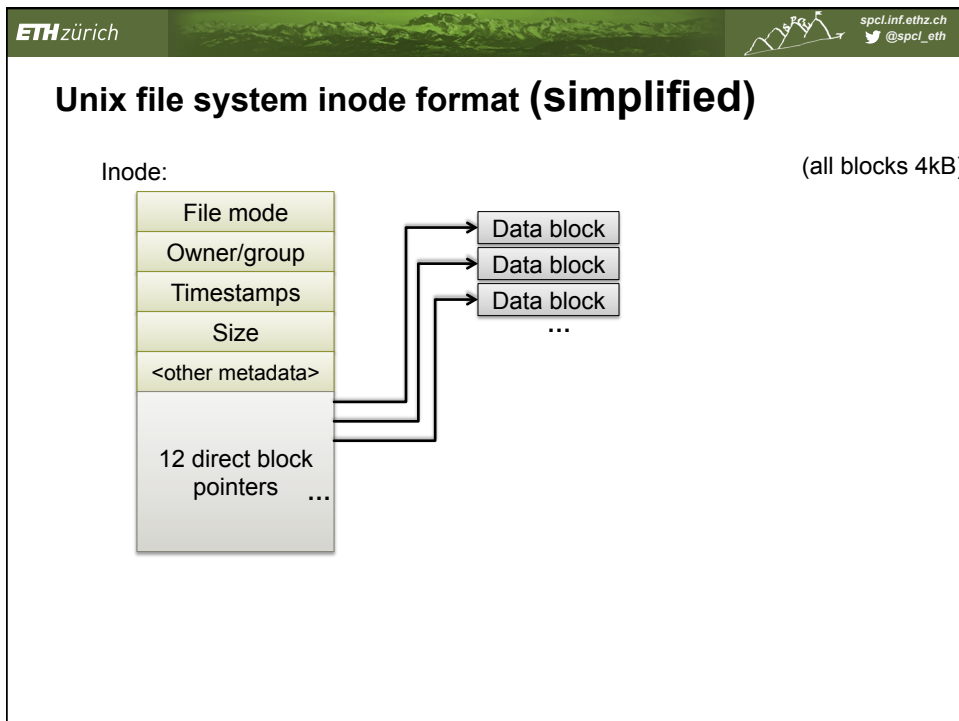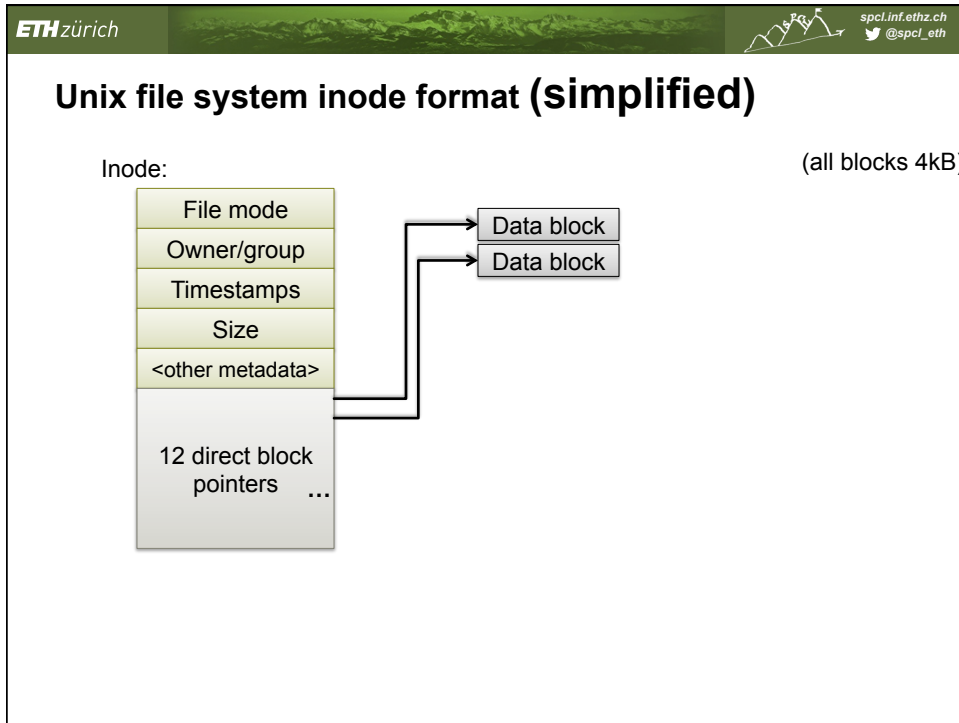Unix file system inode format (**simplified**)
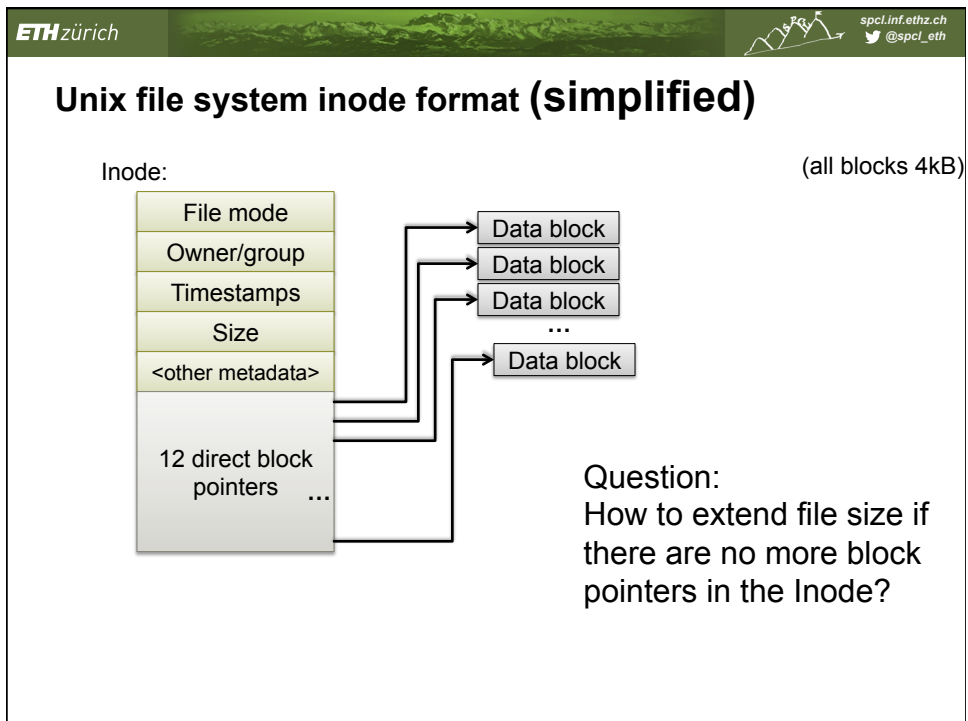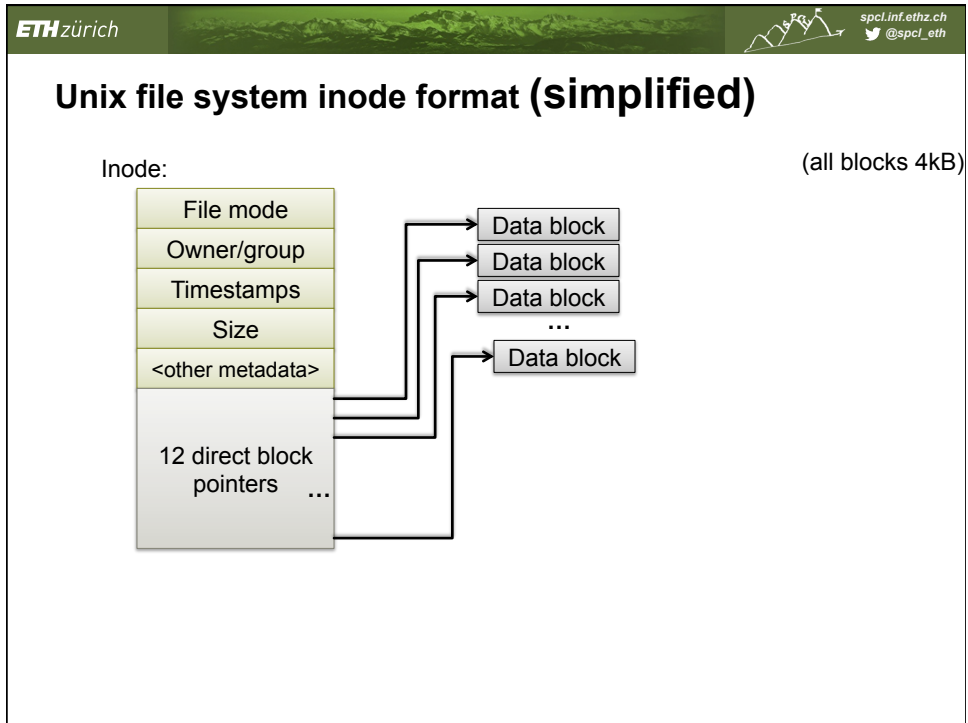
ETH *zürich*

*spcl.inf.ethz.ch*
*@spcl_eth*
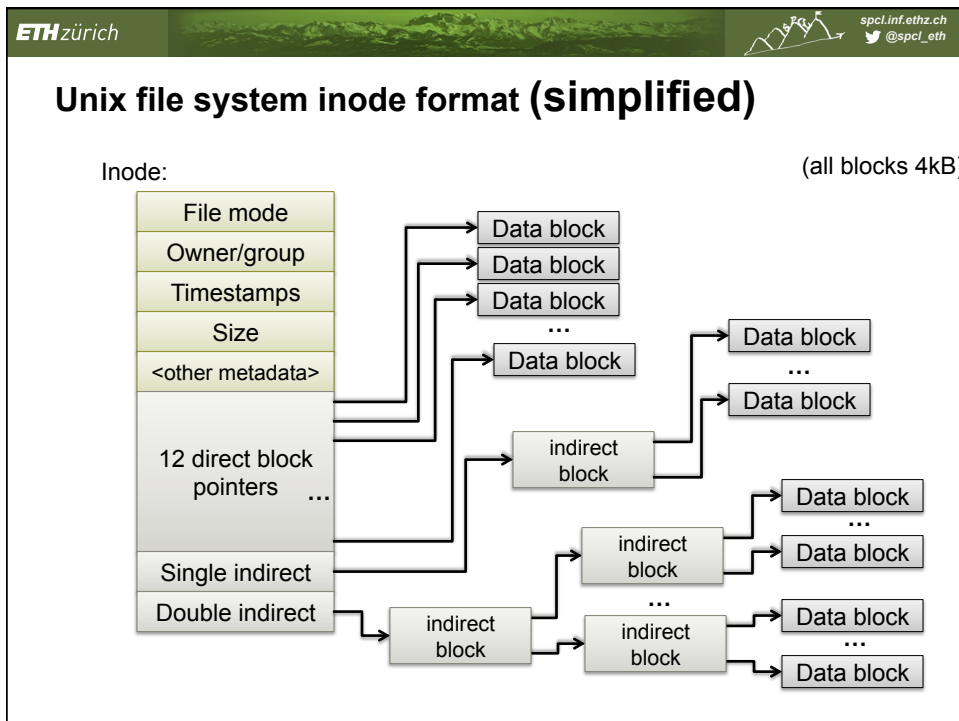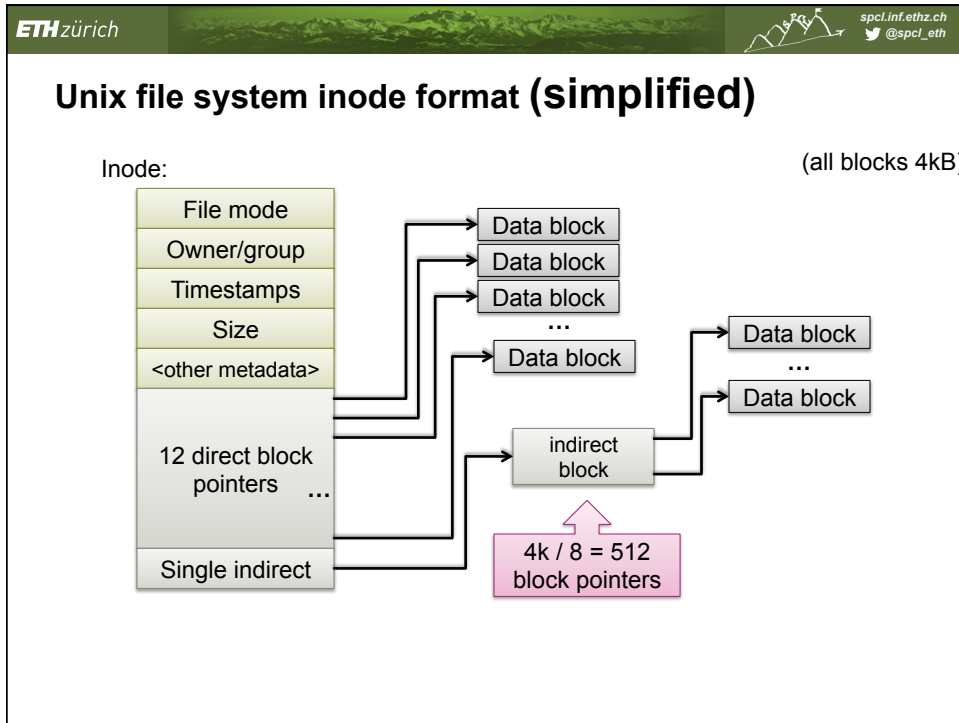
# Unix file system inode format (**simplified**)

Inode:

(all blocks 4kB)

| File mode |
| Owner/group |
| Timestamps |
| Size |
| <other metadata> |
| 12 direct block pointers ... |
| Single indirect |
| Double indirect |
| Triple indirect |

Data block
Data block
Data block
...
Data block
Data block

indirect block

Data block
...
Data block

indirect block

Data block
...
Data block

indirect block
indirect block

Data block
...
Data block

---

ETH *zürich*

*spcl.inf.ethz.ch*
*@spcl_eth*

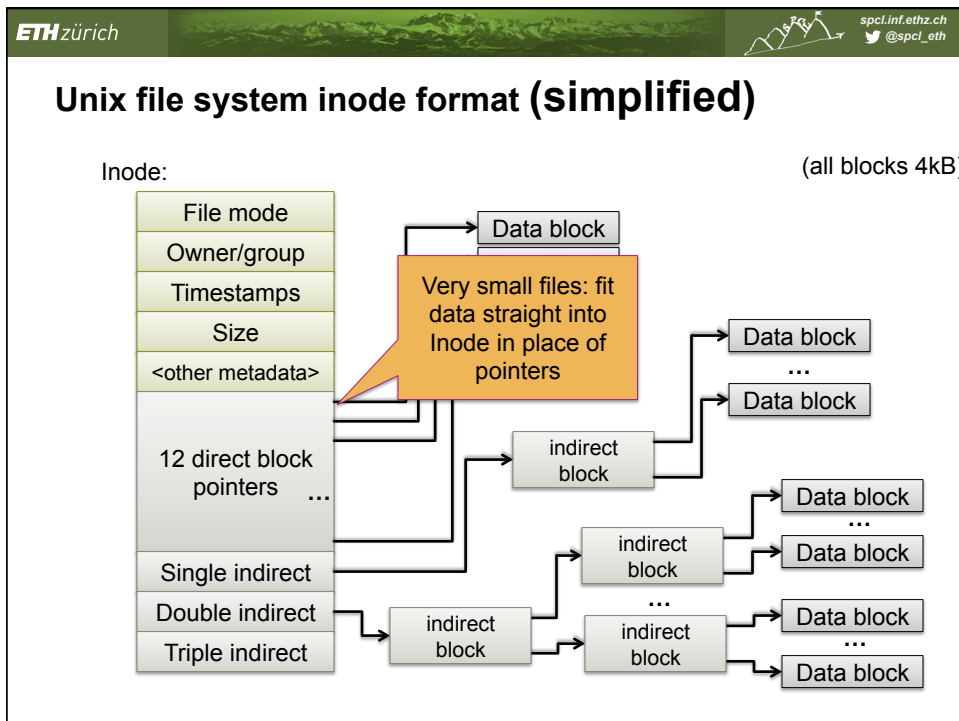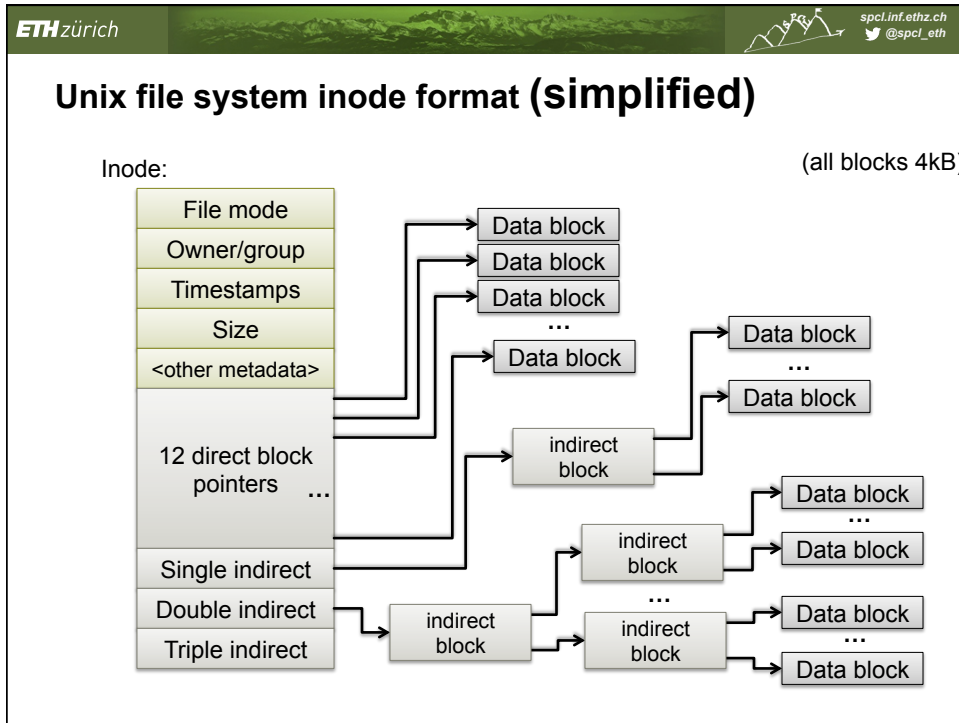# Unix file system inode format (**simplified**)

Inode:

(all blocks 4kB)

| File mode |
| Owner/group |
| Timestamps |
| Size |
| <other metadata> |
| 12 direct block pointers ... |
| Single indirect |
| Double indirect |
| Triple indirect |

Data block

Very small files: fit data straight into Inode in place of pointers

Data block
...
Data block

indirect block

Data block
...
Data block

indirect block

Data block
...
Data block

indirect block
indirect block

Data block
...
Data block

Unix file system inode format (simplified)



Unix file system inode format (simplified)

# Free space map

- **FFS uses a simple bitmap**
  - Initialized when the file system is created
  - One bit per disk (file system) block
- **Allocation is reasonably fast**
  - Scan though lots of bits at a time
  - Bitmap can be cached in memory

# Block groups

1. Optimize disk performance by keeping together related:
- Files
- Metadata (inodes)
- Free space map
- Directories

spcl.inf.ethz.ch
@spcl_eth

*ETH*zürich

# Block groups

2. Use *first-fit* allocation within a block group to improve disk locality

Block Group 0

Block Group 1

Block Group 2

Inodes

Free Space Bitmap

Data Blocks for files in directories /d/q, /c and /b/c

Data Blocks for files in directories /a, /d, and /b/c

Inodes

Data Blocks for files in directories /b, /a/g, /z

Free Space Bitmap

Inodes

Free Space Bitmap

---

spcl.inf.ethz.ch
@spcl_eth

*ETH*zürich

# Block groups

Block Group 0

Block Group 1

Block Group 2

Inodes

Free Space Bitmap

Data Blocks for files in directories /d/q, /c and /b/c

Data Blocks for files in directories /a, /d, and /b/c

Inodes

Data Blocks for files in directories /b, /a/g, /z

Free Space Bitmap

Inodes

Free Space Bitmap

3. Layout and block groups defined in the *superblock* (not shown); Replicated several times.

**NTFS**

---



**NFTS Master file table**

MFT

MFT record:

| Std. info | Attributes, data, metadata | *free* |

Lots of options for what goes in here

1 kB fixed size

**NTFS normal files**

- **MFT holds list of *extents*:**



**Too many attributes?**

- **Attribute list holds list of attribute locations**



In addition, attributes can also be stored in extents ⇒ very large scaling (see book)

# Metadata files

- **File system metadata in NTFS is held *in files!***

| File num. | Name | Description |
|---|---|---|
| 0 | $MFT | Master file table |
| 1 | $MFTirr | Copy of first 4 MFT entries |
| 2 | $Logfile | Transaction log of FS changes |
| 3 | $Volume | Volume information & metadata |
| 4 | $AttrDef | Table mapping numeric IDs to attributes |
| 5 | . | Root directory |
| 6 | $Bitmap | Free space bitmap |
| 7 | $Boot | Volume boot record |
| 8 | $BadClus | Bad cluster map |
| 9 | $Secure | Access control list database |
| 10 | $UpCase | Filename mappings to DOS |
| 11 | $Extend | Extra file system attributes (e.g. quota) |

# Metadata files

▪ **File system metadata in NTFS is held *in files!***

| File num. | Name | Description |
|---|---|---|
| 0 | $MFT | Master file table |
| 1 | $MFTirr | Copy of first 4 MFT entries |
| 2 | $Logfile | Transaction log of FS changes |
| 3 | $Volume | Volume information & metadata |
| 4 | $AttrDef | Table mapping numeric IDs to attributes |
| 5 | . | Root directory |
| 6 | $Bitmap | Free space bitmap |
| 7 | $Boot | Volume boot record |
| 8 | $BadClus | Bad cluster map |
| 9 | $Secure | Access control list database |
| 10 | $UpCase | Filename mappings to DOS |
| 11 | $Extend | Extra file system attributes (e.g. quota) |

# Metadata files

▪ **File system metadata in NTFS is held *in files!***

| File num. | Name | Description |
|---|---|---|
| 0 | $MFT | Master file table |
| 1 | $MFTirr | Copy of first 4 MFT entries |
| 2 | $Logfile | Transaction log of FS changes |
| 3 | $Volume | Volume information & metadata |
| 4 | $AttrDef | Table mapping numeric IDs to attributes |
| 5 | . | Root directory |
| 6 | $Bitmap | Free space bitmap |
| 7 | $Boot | Volume boot record |
| 8 | $BadClus | Bad cluster map |
| 9 | $Secure | Access control list database |
| 10 | $UpCase | Filename mappings to DOS |
| 11 | $Extend | Extra file system attributes (e.g. quota) |

# In-memory data structures

---

# Opening a file

- **Directories translated into kernel data structures on demand:**

open("foo");

directory

file inode

directory structure

User space          Kernel          Disk

# Reading and writing

- **Per-process open file table → index into…**
- **System open file table → cache of inodes**



read(5,…)

5

Per-process open file table

System open file table

file inode

File blocks

User space | Kernel | Disk

---

# Efficiency and Performance

- **Efficiency dependent on:**
  - disk allocation and directory algorithms
  - types of data kept in file's directory entry

- **Performance**
  - disk cache – separate section of main memory for frequently used blocks
  - free-behind and read-ahead – techniques to optimize sequential access
  - improve PC performance by dedicating section of memory as virtual disk, or RAM disk

**ETH**zürich

## Addendum

1. **FAT file size limit (4 GB) – due to the file length entry in the FAT table (guess what, 32 bits)**
   - Volume size is up to 16 TB ($2^{32}$ * 4k blocks)

3. **Insert data into the middle of a file**
   - May allocate a block (true)
   - But who wants to insert only block sizes?
   - Solution: Streams (D'oh, not in FFS/ext…?)

**ETH**zürich

## Addendum

# 3. Block Groups
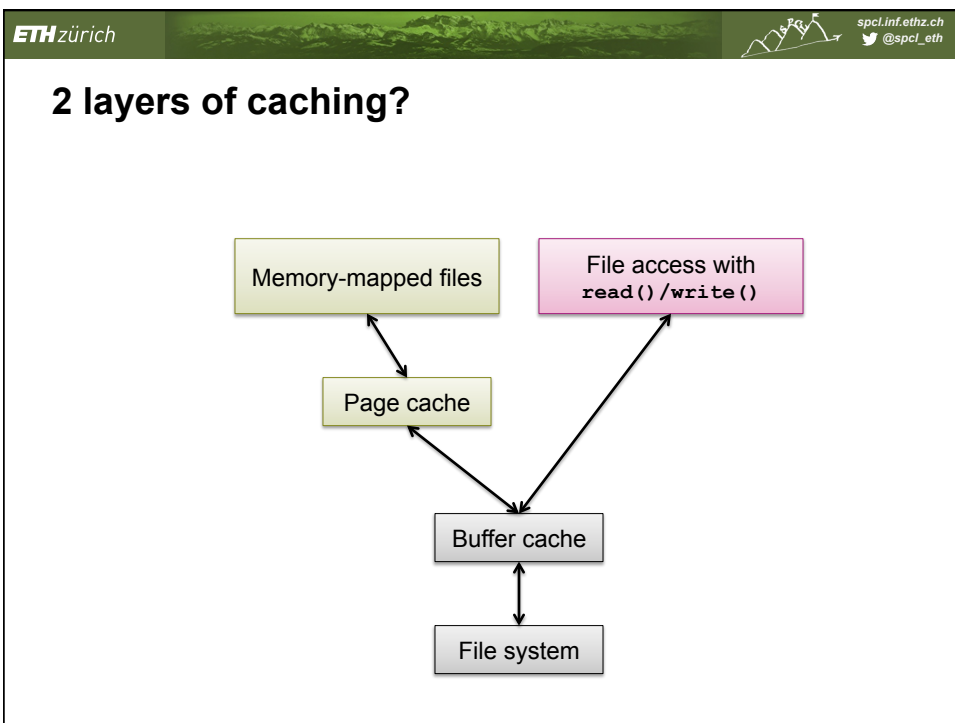
**ETH**zürich
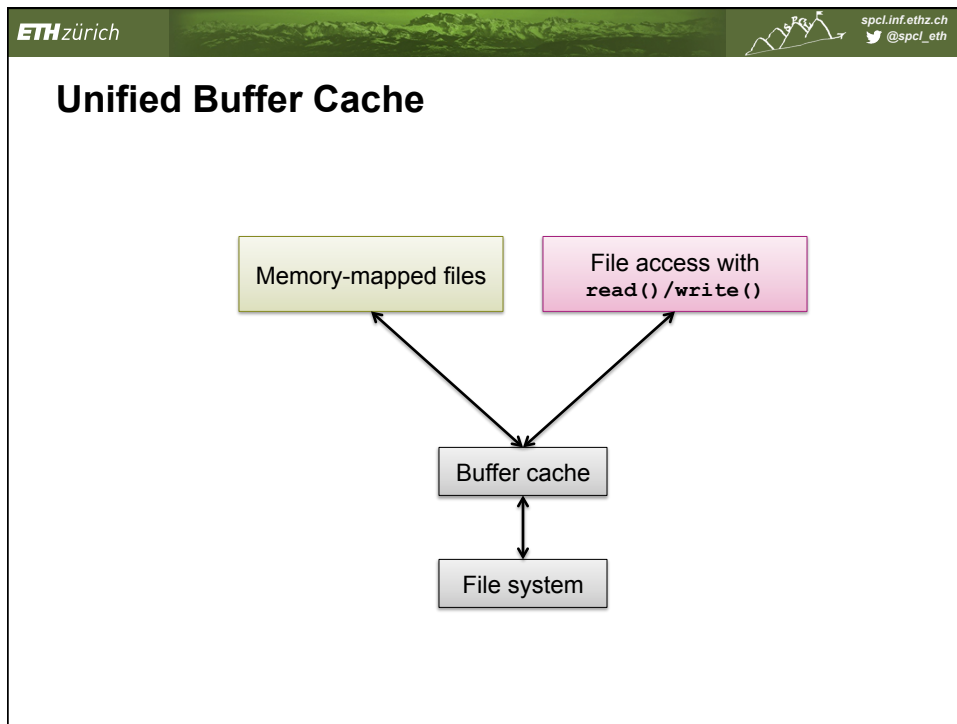
# Page Cache

- **A page cache caches pages rather than disk blocks using virtual memory techniques**

- **Memory-mapped I/O uses a page cache**

- **Routine I/O through the file system uses the buffer (disk) cache**

- **This leads to the following figure**

---

**ETH**zürich

# 2 layers of caching?

**ETH**zürich

# Unified Buffer Cache

| Memory-mapped files | File access with `read()/write()` |
|---|---|

Buffer cache

File system

---

**ETH**zürich

# Recovery

- **Consistency checking – compares data in directory structure with data blocks on disk, and tries to fix inconsistencies**

- **Use system programs to back up data from disk to another storage device (floppy disk, magnetic tape, other magnetic disk, optical)**

- **Recover lost file or disk by restoring data from backup**