

# Dateisysteme

## Hochleistungs-Ein-/Ausgabe

Michael Kuhn

Wissenschaftliches Rechnen  
Fachbereich Informatik  
Universität Hamburg

2017-04-21



Universität Hamburg  
DER FORSCHUNG | DER LEHRE | DER BILDUNG

- 1 Dateisysteme
  - Orientierung
  - Dateisysteme
  - ext4
  - Object Stores
  - Datenstrukturen
  - Leistungsbewertung
  - Ausblick und Zusammenfassung

- 2 Quellen



# Aufgabe

- **Strukturierung**
  - Üblicherweise Dateien und Verzeichnisse
  - Hierarchische Organisation
  - Andere Ansätze: Tagging
- **Verwaltung von Daten und Metadaten**
  - Blockallokation
  - Zugriffsrechte, Zeitstempel etc.
- **Dateisysteme nutzen ein darunter liegendes Speichergerät**
  - Oder einen Speicherverbund
  - Logical Volume Manager (LVM) und/oder mdadm

# Beispiele

- Linux: ext4, XFS, btrfs, ZFS, ...
- Windows: FAT, exFAT, NTFS
- OS X: HFS+, APFS
- Universal: ISO9660, UDF, ...

# Beispiele...

- Netzwerk: NFS, AFS, Samba
- Kryptographisch: EncFS, eCryptfs
- Parallel verteilt: GPFS, Lustre, ...
- Pseudo: procfs, ...
- Setzen häufig auf darunterliegenden Dateisystemen auf

# E/A-Schnittstellen

- Anfragen werden über E/A-Schnittstellen realisiert
  - Weiterleitung an das Dateisystem
  - Unterschiedliche Abstraktionsebenen
- Low-Level-Funktionalität
  - POSIX, MPI-IO, ...
- High-Level-Funktionalität
  - HDF, NetCDF, ...









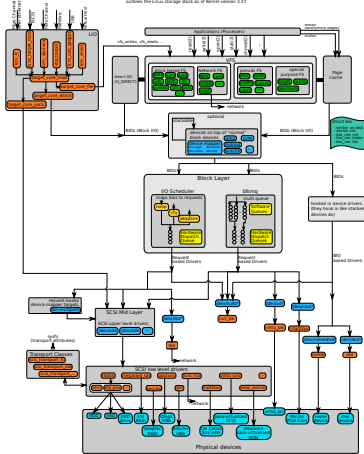




# VFS... [3]

### The Linux Storage Stack Diagram

version 3.17, 2015-07-17  
outlines the Linux storage stack as of kernel version 3.17



# Dateisystemobjekte

- Unterscheidung in Benutzer- und Systemsicht
  - Benutzer sehen Dateien und Verzeichnisse
  - Das System kennt zusätzlich Inodes
    - Relevant für stat etc.
- Inodes
  - Enthalten Metadaten
  - Eigentliche Basisobjekte des Dateisystems
    - Jeder Datei und jedem Verzeichnis ist ein Inode zugeordnet
  - Üblicherweise eindeutige IDs

# Dateisystemobjekte...

- Dateien
  - Enthalten Daten in Form eines Byte-Arrays
  - Können gelesen/geschrieben werden (explizit)
  - Können in den Speicher gemappt werden (implizit)
- Verzeichnisse
  - Enthalten Dateien und Verzeichnisse
  - Zur Organisation des Namensraumes

# Dateien

```
1 nb = pwrite(fd, data, sizeof(data), 42);
2 nb = pread(fd, data, sizeof(data), 42);
```

## Listing 2: Expliziter Zugriff

- `fwrite` und `pread` verhalten sich wie `write` bzw. `read`
  - Explizite Angabe des Offsets und damit threadsicher
- Zugriff über File Descriptor
  - Kann von mehreren Threads parallel genutzt werden





# Dateien...

- Beide Zugriffsarten haben jeweils Vor- und Nachteile
  - Beide Modi profitieren vom Caching durch das Betriebssystem
- Expliziter Zugriff
  - Vorteile: genaue Kontrolle über E/A
  - Nachteile: separate Puffer notwendig, Kopiervorgänge zwischen Kernel- und Userspace
- Impliziter Zugriff
  - Vorteile: keine separaten Puffer notwendig, effiziente E/A durch das Betriebssystem, keine unnötigen Kopiervorgänge
  - Nachteile: keine genaue Kontrolle, kompliziertere Fehlerbehandlung (Signale)

# Verzeichnisse

Inode	Größe	Namenslänge	Dateityp	Name
23	10	2	2	.\0
24	11	3	2	..\0
⋮	⋮	⋮	⋮	⋮
42	14	6	1	hello\0
43	14	6	2	world\0

**Abbildung:** ext4-Verzeichniseintrag [1]

- Traditionell lineares Array
  - Langsam, da über das komplette Array iteriert werden muss
- Heutzutage eher Baumstrukturen
  - Deutlich komplexer, dafür geringere Zugriffszeiten
- Name wird nicht im Inode gespeichert
  - Mehrere Namen können auf denselben Inode zeigen

# Inodes

Feldgröße	Inhalt
2 Bytes	Berechtigungen
2 Bytes	Benutzer-ID
4 Bytes	Dateigröße
4 Bytes	Zugriffszeit
4 Bytes	Inode-Änderungszeit
4 Bytes	Datenänderungszeit
4 Bytes	Löschzeit
2 Bytes	Gruppen-ID
2 Bytes	Linkzahl
⋮	⋮
60 Bytes	Blockzeiger, Extent-Baum oder Inline-Daten
⋮	⋮
4 Bytes	Versionsnummer
100 Bytes	Freier Speicher

**Abbildung:** ext4-Inode (256 Bytes) [1]

# Inodes...

- Kompliziert durch Rückwärtskompatibilität
  - On-Disk-Format kann nur schwer geändert werden
- Viele Felder sind aus Kompatibilitätsgründen aufgeteilt
  - Zeitstempel: 4 Bytes für Sekunden seit 1970, 4 Bytes für Nanosekundenauflösung
  - Größe: Obere und untere 4 Bytes
- Felder mehrfach überladen
  - Blockzeiger, Extent-Baum oder Inline-Daten (falls Datei kleiner als 60 Bytes)
  - 100 Bytes am Inode-Ende für erweiterte Attribute



# POSIX-Schnittstelle

- Syntax
  - open, close, creat
  - read, write, lseek
  - chmod, chown, stat
  - link, unlink
- Semantik
  - Spezifiziert auch wie sich E/A-Operationen verhalten sollen
  - *w r i t e: "POSIX requires that a read(2) which can be proved to occur after a write() has returned returns the new data. Note that not all filesystems are POSIX conforming."*

# ext4

- Standard-Dateisystem in vielen Linux-Distributionen
  - Eingeführt 2006, stabil 2008
  - Vorgänger: ext, ext2, ext3
- Statische Festlegung bei Dateisystemerzeugung
  - Inode-Zahl
  - Blockgröße
- Traditionelles Dateisystem
  - Daten werden direkt geändert (kein Copy on Write)
  - Keine Prüfsummen für Daten



# ext

- Erstes Dateisystem speziell für Linux
  - Nutzte als erstes Dateisystem die VFS-Schicht
- Inspiriert vom Unix File System (UFS)
- Beseitigte Beschränkungen des MINIX-Dateisystems
  - Dateigrößen bis 2 GB
  - Dateinamen bis 255 Zeichen



# ext2

- Separate Zeitstempel für Zugriff und Inode-/Datenänderung
- Datenstrukturen für zukünftige Erweiterungen ausgelegt
- Testumgebung für neue VFS-Funktionen
  - Access Control Lists (ACLs)
  - Erweiterte Attribute

# ext3

- Journaling
  - Erklärung folgt später
- Dateisystemvergrößerung zur Laufzeit
  - Nützlich für LVM-Umgebungen
- H-Baum für größere Verzeichnisse
  - Verkürzt die Suchzeiten im Verzeichnis

# ext4

- Größere Dateisysteme, Dateien und Verzeichnisse
- Extents
- Preallokation, verzögerte Allokation und verbesserte Multiblockallokation
- Journal-Prüfsummen
- Schnellere Dateisystemüberprüfung
- Nanosekunden-Zeitstempel
- Unterstützung für TRIM

## ext4...

Inhalt	Größe
Padding (Blockgruppe 0)	1.024 Bytes
Superblock	1 Block
Gruppenbeschreibung	n Blöcke
Reservierte GDT-Blöcke	m Blöcke
Daten-Bitmap	1 Block
Inode-Bitmap	1 Block
Inode-Tabelle	k Blöcke
Daten-Blöcke	l Blöcke

Abbildung: ext4-Blockgruppe [1]

- Das Speichergerät ist in mehrere Blockgruppen unterteilt
  - Flexible Blockgruppen fassen mehrere Blockgruppen zusammen

## ext4...

Blockgröße	1 KiB	2 KiB	4 KiB	64 KiB
Blöcke	$2^{64}$	$2^{64}$	$2^{64}$	$2^{64}$
Inodes	$2^{32}$	$2^{32}$	$2^{32}$	$2^{32}$
Dateisystemgröße	16 ZiB	32 ZiB	64 ZiB	1 YiB
Dateigröße (Extents)	4 TiB	8 TiB	16 TiB	256 TiB
Dateigröße (Blöcke)	16 GiB	256 GiB	4 TiB	256 PiB

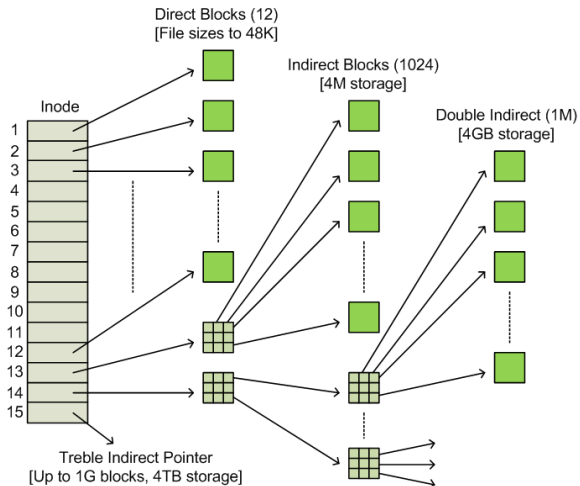
Abbildung: ext4-Limits im 64-Bit-Modus [1]

- Standardgröße ist 4 KiB (und offizielles Maximum)
  - Sollte nicht größer als Seitengröße gewählt werden

# Allokation

- Blockbasiert
  - Viele Blöcke gleicher Größe (üblicherweise 4 KiB)
  - Zeiger auf Blöcke im Inode
    - Direkt, indirekt, doppelt indirekt, dreifach indirekt
  - Hoher Overhead bei großen Dateien
    - Beispiel: 1 TiB große Datei benötigt 268.435.456 Zeiger
  - Beschränkt maximale Dateigröße

## Allokation... [2]



# Allokation...

- Extentbasiert
  - Wenige möglichst große Extents
    - Vier Extents können im Inode gespeichert werden
    - Mehr in einer Baumstruktur und zusätzlichen Blöcken
  - Zeiger auf Startblock und Länge
    - Maximale Länge: 32.768 Blöcke
    - Entspricht 128 MiB bei einer Blockgröße von 4 KiB
  - Ermöglicht größere Dateien



# Allokation...

- Blockallokation
  - Versuche zusammenhängende Blöcke zu allokiieren
  - Versuche Blöcke in derselben Blockgruppe zu allokiieren
- Multiblockallokation und verzögerte Allokation
  - Spekulativ 8 KiB bei Dateierzeugung allokiieren
  - Allokation wird erst durchgeführt, wenn Blöcke auf das Speichergerät geschrieben werden müssen

# Allokation...

- Dateien und Verzeichnisse
  - Blöcke möglichst in der Blockgruppe des Inodes allokkieren
  - Dateien möglichst in der Blockgruppe des Verzeichnisses allokkieren
- Ziele der Allokationsstrategien
  - Möglichst große Zugriffe
    - Festplatten erreichen nur geringe IOPS-Werte
  - Zugriffe nahe beieinander
    - Reduziert Kopfbewegungen bei Festplatten
    - Metadaten der Blockgruppe eventuell schon im Cache
- Optimierungen bei SSDs weniger von Bedeutung

# Sparse-Dateien und Preallokation

- Sparse-Dateien: Dateien mit „Löchern“
  - Z.B. mit `lseek` oder `truncate`
  - Effiziente Speicherung von Dateien mit vielen 0-Bytes

```
1 $ truncate --size=1G dummy
2 $ ls -lh dummy
3 -rw-r--r--. 1 u g 1,0G 18. Apr 23:49 dummy
4 $ du -h dummy
5 0    dummy
```

Listing 5: Erzeugung einer Sparse-Datei

# Sparse-Dateien und Preallokation...

- Preallokation: Speicher vorallokieren
  - Mit `fallocate` bzw. `posix_fallocate`
  - Verhindert Fragmentierung bei vielen Dateivergrößerungen

```
1 $ fallocate --length $((1024 * 1024 * 1024)) dummy
2 $ ls -lh dummy
3 -rw-r--r--. 1 u g 1,0G 19. Apr 19:14 dummy
4 $ du -h dummy
5 1,1G      dummy
```

Listing 6: Preallokation einer Datei

- Unterschiedliche Basen je nach Werkzeug

# Journaling

- Journaling zur Sicherung der Konsistenz des Dateisystems
- Dateisystemoperationen benötigen mehrere Schritte
- Z.B. das Löschen einer Datei
  - 1 Entfernen des Verzeichniseintrags
  - 2 Freigeben des Inodes
  - 3 Freigeben der Datenblöcke
- Problematisch im Fall eines Absturzes

# Journaling...

- Geplante Änderungen werden ins Journal eingetragen
  - Entfernen wenn Operation vollständig durchgeführt
- Bei der anschließenden Dateisystemüberprüfung
  - Änderungen wiederholen oder
  - Änderungen verwerfen
- Unterschiedliche Modi
  - Metadaten-Journaling und volles Journaling

# Journaling...

- **Journal:** Alle Änderungen werden ins Journal geschrieben
  - Deaktiviert verzögerte Allokation und O\_DIRECT
- **Ordered:** Metadaten werden ins Journal geschrieben
  - Zugehörige Daten werden vor Metadaten geschrieben
  - Problematisch mit verzögerter Allokation
  - Ist die Standardeinstellung
- **Writeback:** Metadaten werden ins Journal geschrieben
  - Bietet höchste Leistung aber geringste Sicherheit

# Funktionen

- „Dateisystem light“
  - Dünne Abstraktionsschicht über Speichergeräten
  - Objektbasierter Zugriff auf Daten
- Nur Grundoperationen
  - Erstellen, Öffnen, Schließen, Lesen, Schreiben
- Manchmal Object Sets
  - Können benutzt werden um verwandte Objekte zu gruppieren



# Funktionen...

- Üblicherweise keine Pfade
  - Zugriff über eindeutige IDs
  - Kein Overhead durch Pfadauflösung
  - Dadurch auch flacher Namensraum
- Block-/Extent-Allokation
  - Einer der komplexesten und leistungsrelevantesten Aspekte
- Auf unterschiedlichen Abstraktionsebenen verfügbar
  - Cloudspeicher, Festplatte

# Schichtung

- Können als Unterbau für Dateisysteme genutzt werden
  - Erlaubt Konzentration auf Dateisystemfunktionalität
  - Speicherverwaltung durch separate Schicht
- Bei lokalen Dateisystemen nicht sinnvoll
  - Funktionalität größtenteils durch POSIX vorgegeben
  - Hauptunterschied ist Blockallokation
- Sehr sinnvoll für parallele verteilte Dateisysteme
  - Kein redundanter Dateisystem-Overhead

## B-Baum vs. B+-Baum

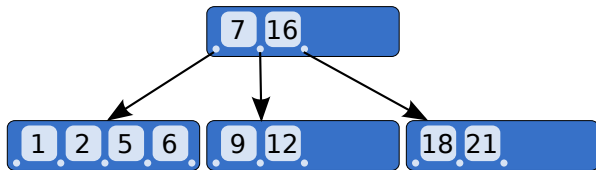


Abbildung: B-Baum [4]

- Verallgemeinerter Binärbaum
- Optimiert für Systeme, die große Blöcke lesen/schreiben
- Zeiger und Daten gemischt

# B-Baum vs. B+-Baum...

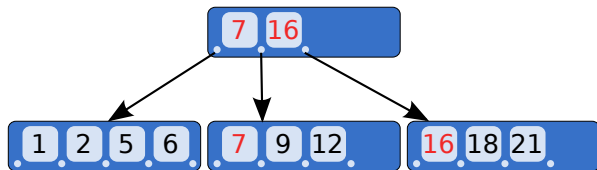


Abbildung: B+-Baum [4]

- Daten nur in Blättern
- Vorteilhaft für Caching, da einfacher alle Knoten zu cachen
- Benutzt in NTFS, XFS, ...

# Alternativen

- H-Baum
  - Basiert auf B-Baum
  - Andere Behandlung von Hash-Kollisionen
  - Benutzt in ext3 und ext4
  
- B<sup>ε</sup>-Baum
  - Optimiert für Schreibvorgänge
  - Verbesserte Leistung für Einfügeoperationen, Bereichsabfragen und Aktualisierungen

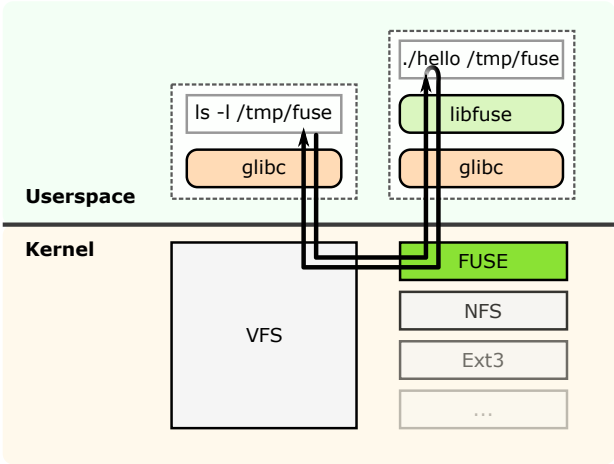
# Leistungsbewertung

- Dateisystemleistung ist schwierig zu bewerten
  - Viele unterschiedliche Faktoren
  - Daten- vs. Metadatenleistung
  - Leistung unterschiedlicher Funktionen
  - Leistung für spezifische Anforderungen messen
- Datensicherheit kostet üblicherweise Leistung
  - Volles Journaling, Prüfsummen etc.

# Kernel- vs. Userspace

- Dateisysteme üblicherweise direkt im Kernel implementiert
  - Hoher Wartungsaufwand
  - Komplexere Implementierung
- Alternative: Filesystem in Userspace (FUSE)
  - Besteht aus Kernelmodul und Bibliothek
  - Entwicklung von Dateisystemen als normale Prozesse
  - Umleitung in Userspace durch VFS und FUSE-Modul
  - Geringere Leistung durch Kontextwechsel

# Kernel- vs. Userspace... [5]





# Ausblick

- Moderne Dateisysteme integrieren zusätzliche Funktionen
  - Volumenverwaltung, Prüfsummen, Schnappschüsse, ...
  - Komfort vs. Datensicherheit
- Basis für parallele verteilte Dateisysteme
  - Existierende und optimierte Blockallokation etc.
  - Object Stores häufig besser geeignet

# Zusammenfassung

- Dateisysteme organisieren Daten und Metadaten
  - Üblicherweise standardisierte Schnittstelle
- Hauptobjekte sind Dateien und Verzeichnisse
  - Inodes speichern Metadaten
- Neue Techniken zur Effizienzsteigerung
  - Journaling um Konsistenz sicherzustellen
  - Speicherallokation mit Hilfe von Extents
  - Baumstrukturen für skalierbaren Zugriff

- 1 Dateisysteme
  - Orientierung
  - Dateisysteme
  - ext4
  - Object Stores
  - Datenstrukturen
  - Leistungsbewertung
  - Ausblick und Zusammenfassung

## 2 Quellen

# Quellen I

- [1] djwong. Ext4 Disk Layout. [https://ext4.wiki.kernel.org/index.php/Ext4\\_Disk\\_Layout](https://ext4.wiki.kernel.org/index.php/Ext4_Disk_Layout).
- [2] Hal Pomeranz. Understanding Indirect Blocks in Unix File Systems. <http://digital-forensics.sans.org/blog/2008/12/24/understanding-indirect-blocks-in-unix-file-system>
- [3] Werner Fischer and Georg Schönberger. Linux Storage Stack Diagramm. [https://www.thomas-krenn.com/de/wiki/Linux\\_Storage\\_Stack\\_Diagramm](https://www.thomas-krenn.com/de/wiki/Linux_Storage_Stack_Diagramm).
- [4] Wikipedia. B-tree. <http://en.wikipedia.org/wiki/B-tree>.

## Quellen II

- [5] Wikipedia. Filesystem in Userspace. [http://en.wikipedia.org/wiki/Filesystem\\_in\\_Userspace](http://en.wikipedia.org/wiki/Filesystem_in_Userspace).