

Dateisysteme

Hochleistungs-Ein-/Ausgabe

Michael Kuhn

Wissenschaftliches Rechnen
Fachbereich Informatik
Universität Hamburg

2015-04-20

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

- 2 Quellen

ext4

ext3

- Journaling
- Dateisystemvergrößerung zur Laufzeit
- H-Baum für größere Verzeichnisse

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]

Allokation

- Blockbasiert
 - Viele Blöcke gleicher Größe (üblicherweise 4 KiB)
 - Zeiger auf Blöcke
 - Direkt, indirekt, doppelt indirekt, dreifach indirekt

ext4

Allokation...

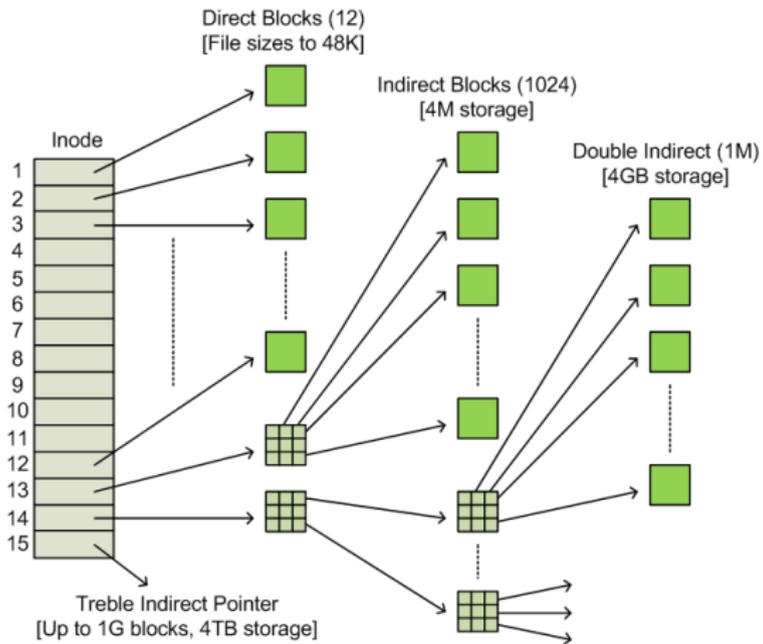


Abbildung: Block-Zeiger [2]

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

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

ext4

Allokation...

- Dateien und Verzeichnisse
 - Blöcke möglichst in der Blockgruppe des Inodes allokiieren
 - Dateien möglichst in der Blockgruppe des Verzeichnisses allokiieren

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

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
- Ordered: Metadaten werden ins Journal geschrieben
 - Daten vor Metadaten
 - Problematisch mit verzögerter Allokation
- Writeback: Metadaten werden ins Journal geschrieben
 - Reihenfolge beliebig

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
- Block-/Extent-Allokation
- Auf unterschiedlichen Abstraktionsebenen verfügbar
 - Cloudspeicher, Festplatte

Schichtung

- Object Stores 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

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
- Alternative: Filesystem in Userspace (FUSE)
 - Besteht aus Kernelmodul und Bibliothek
 - Entwicklung von Dateisystemen als normale Prozesse

Ausblick

- Moderne Dateisysteme integrieren zusätzliche Funktionen
 - Volumenverwaltung, Prüfsummen, Schnappschüsse, ...
- Basis für parallele verteilte Dateisysteme
 - Object Stores besser geeignet

- 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.
- [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-systems>.
- [3] Werner Fischer and Georg Schönberger. 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.