C IV Dateisystem

C | IV Dateisystem

Überblick

- Medien
- Speicherung von Dateien
- Freispeicherverwaltung
- Beispiele: Dateisysteme unter UNIX und Windows

1 Überblick

- Dateisysteme mit Fehlererholung
- Datensicherung

Systemprogrammierung

Dateisystem

Jürgen Kleinöder

Lehrstuhl Informatik 4

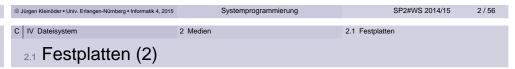
13. Januar 2015

20. Januar 2015

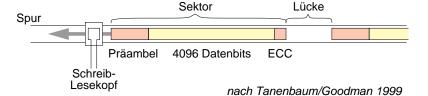


-Zylinder Spur Plattenarme

◆ Kopf schwebt auf Luftpolster



Sektoraufbau



♦ Breite der Spur: 5–10 μm

◆ Spuren pro Zentimeter: 800-2000 ♦ Breite einzelner Bits: 0,1–0,2 μm

Zonen

◆ Mehrere Zylinder (10–30) bilden eine Zone mit gleicher Sektorenanzahl (bessere Plattenausnutzung)

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 3/56 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 4/56

2.1 Festplatten (3)

■ Datenblätter von drei (alten) Beispielplatten

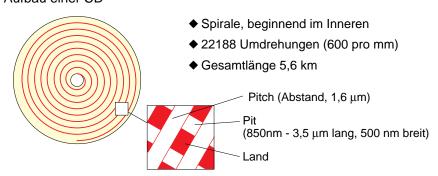
Plattentyp		Fujitsu M2344 (1987)	Seagate Cheetah	Seagate Barracuda	
Kapazität		690 MB	300 GB	400 GB	
Platten/Köpfe		8 / 28	4/8	781.422.768 Sektoren	
Zylinderzahl		624	90.774		
Cache		-	4 MB	8 MB	
Posititionier- zeiten	Spur zu Spur	4 ms	0,5 ms	-	
	mittlere	16 ms	5,3 ms	8 ms	
	maximale	33 ms	10,3 ms	-	
Transferrate		2,4 MB/s	320 MB/s	-150 MB/s	
Rotationsgeschw.		3.600 U/min	10.000 U/min	7.200 U/min	
eine Plattenumdrehung		16 ms	6 ms	8 ms	
Stromaufnahme		?	16-18 W	12,8 W	

09.2014: Kapazität bis 8 TB bei 7.200 U/min oder 600 GB bei 15.000 U/min, Zugriffszeit ab 2,7 ms, Transferrate bis 200 MB/s

SSD: Kapazität bis 4 TB, Zugriffszeit ca. 0,1 ms, Transferrate bis 3 GB/s

temprogrammierung SP2#WS 2014/15 5 / 56
2.2 CD-ROM / DVD

Aufbau einer CD



- ◆ Pit: Vertiefung, wird von Laser (780 nm Wellenlänge) abgetastet
- DVD
 - ◆ gleiches Grundkonzept, Wellenlänge des Lasers 650 nm
 - ◆ Pits und Spurabstand weniger als halb so groß

2.1 Festplatten (4)

- Zugriffsmerkmale
 - ◆ blockorientierter und wahlfreier Zugriff
 - ◆ Blockgröße zwischen 32 und 4096 Bytes (typisch 512 Bytes)
 - ◆ Zugriff erfordert Positionierung des Schwenkarms auf den richtigen Zylinder und Warten auf den entsprechenden Sektor
 - ◆ heutige Platten haben internen Cache und verbergen die Hardware-Details
- Blöcke sind üblicherweise nummeriert
 - ◆ früher getrennte Nummerierung: Zylindernummer, Sektornummer
 - ◆ heute durchgehende Nummerierung der Blöcke
 - ➤ Kompatibilität zu alten Betriebssystemen wird durch logical CHS (Cylinder/Head/Sector)-Umrechnung hergestellt

© Jürgen Kleinöder • Univ. Erlangen-Nümberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 6 / 56

C | IV Dateisystem 2 Medien 2.2 CD-ROM / DVD

2.2 CD-ROM / DVD (2)

- Kodierung einer CD
 - ◆ **Symbol**: ein Byte wird mit 14 Bits kodiert (kann bereits bis zu zwei Bitfehler korrigieren)
 - ◆ Frame: 42 Symbole (192 Datenbits, 396 Fehlerkorrekturbits)
 - ◆ Sektor : 98 Frames werden zusammengefasst (16 Bytes Präambel, 2048 Datenbytes, 288 Bytes Fehlerkorrektur)
 - ◆ Effizienz: 7203 Bytes transportieren 2048 Nutzbytes (28,4 %)
- Kodierung einer DVD
 - ◆ Codierung mit Reed-Solomon-Product-Code, 8/16-Bit-Modulation, 43,2 % Nutzdaten
- Transferrate
 - ◆ CD-Single-Speed-Laufwerk: 75 Sektoren/Sek. (153.600 Bytes/Sek.)
 - ◆ CD-72-fach-Laufwerk: 11,06 MB/Sek.
 - ◆ DVD 1-fach: 1.3 MB/sec, 24-fach: 33.2 MB/sec

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 7 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 8 / 56

C IV Dateisystem 2 Medien 2.2 CD-ROM / DVD

2.2 CD-ROM / DVD (3)

Kapazität

◆ CD: ca. 650 MB

◆ DVD single layer: 4.7 GB

◆ DVD dual layer: 8.5 GB, beidseitig: 17 GB

Varianten

◆ DVD/CD-R (Recordable): einmal beschreibbar

◆ DVD/CD-RW (Rewritable): mehrfach beschreibbar

Speicherung von Dateien

Dateien benötigen oft mehr als einen Block auf der Festplatte

3 Speicherung von Dateier

◆ Welche Blöcke werden für die Speicherung einer Datei verwendet?

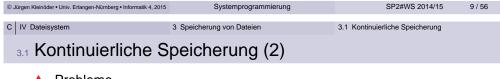
3.1 Kontinuierliche Speicherung

3.1 Kontinuierliche Speicherung

- Datei wird in Blöcken mit aufsteigenden Blocknummern gespeichert
 - ◆ Nummer des ersten Blocks und Anzahl der Folgeblöcke muss gespeichert werden
- Vorteile

C IV Dateisystem

- ◆ Zugriff auf alle Blöcke mit minimaler Positionierzeit des Schwenkarms
- Schneller direkter Zugriff auf bestimmter Dateiposition
- ◆ Einsatz z. B. bei Systemen mit Echtzeitanforderungen



Probleme

- ◆ Finden des freien Platzes auf der Festplatte (Menge aufeinanderfolgender und freier Plattenblöcke)
- ◆ Fragmentierungsproblem (Verschnitt: nicht nutzbare Plattenblöcke; siehe auch Speicherverwaltung)
- ◆ Größe bei neuen Dateien oft nicht im Voraus bekannt
- Erweitern ist problematisch
 - ➤ Umkopieren, falls kein freier angrenzender Block mehr verfügbar



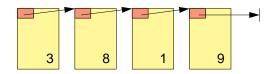
- Variation
 - ◆ Unterteilen einer Datei in Folgen von Blöcken (Chunks, Extents)
 - ◆ Blockfolgen werden kontinuierlich gespeichert
 - ◆ Pro Datei muss erster Block und Länge jedes einzelnen Chunks gespeichert werden
- Problem
 - ◆ Verschnitt innerhalb einer Folge (siehe auch Speicherverwaltung: interner Verschnitt bei Seitenadressierung)

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 SP2#WS 2014/15 11 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 SP2#WS 2014/15 12/56 Systemprogrammierung Systemprogrammierung

C IV Dateisystem 3 Speicherung von Dateien 3.2 Verkettete Speicherung

3.2 Verkettete Speicherung

■ Blöcke einer Datei sind verkettet



- ◆ z. B. Commodore Systeme (CBM 64 etc.)
 - Blockgröße 256 Bytes
 - die ersten zwei Bytes bezeichnen Spur- und Sektornummer des nächsten Blocks
 - wenn Spurnummer gleich Null: letzter Block
 - 254 Bytes Nutzdaten
- ★ Datei kann wachsen und verlängert werden

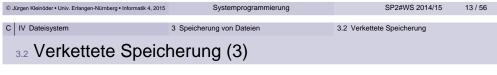
c | IV Dateisystem 3 Speicherung von Dateien 3.2 Verkettete Speicherung (2)

Probleme

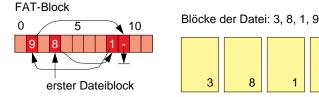
 Speicher für Verzeigerung geht von den Nutzdaten im Block ab (ungünstig im Zusammenhang mit Paging: Seite würde immer aus Teilen von zwei Plattenblöcken bestehen)

3.2 Verkettete Speicherung

- Fehleranfälligkeit: Datei ist nicht restaurierbar, falls einmal Verzeigerung fehlerhaft
- ◆ schlechter direkter Zugriff auf bestimmte Dateiposition
- häufiges Positionieren des Schreib-, Lesekopfs bei verstreuten Datenblöcken



- Verkettung wird in speziellen Plattenblocks gespeichert
 - ◆ FAT-Ansatz (FAT: File Allocation Table), z. B. MS-DOS, Windows 95



- ★ Vorteile
 - ♦ kompletter Inhalt des Datenblocks ist nutzbar (günstig bei Paging)
 - mehrfache Speicherung der FAT möglich: Einschränkung der Fehleranfälligkeit



Probleme

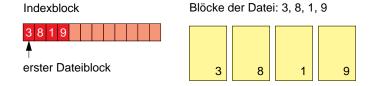
- ◆ mindestens ein zusätzlicher Block muss geladen werden (Caching der FAT zur Effizienzsteigerung nötig)
- ◆ FAT enthält Verkettungen für alle Dateien: das Laden der FAT-Blöcke lädt auch nicht benötigte Informationen
- aufwändige Suche nach dem zugehörigen Datenblock bei bekannter Position in der Datei
- ♦ häufiges Positionieren des Schreib-, Lesekopfs bei verstreuten Datenblöcken

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 15 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 16 / 56

C V Dateisystem 3 Speicherung von Dateien 3.3 Indiziertes Speicherung

3.3 Indiziertes Speichern

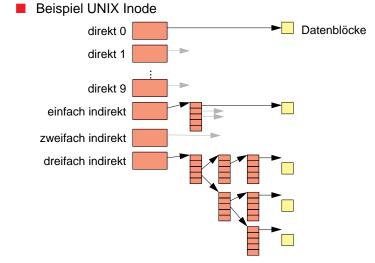
 Spezieller Plattenblock enthält Blocknummern der Datenblocks einer Datei



Problem

- ◆ feste Anzahl von Blöcken im Indexblock
 - · Verschnitt bei kleinen Dateien
 - · Erweiterung nötig für große Dateien



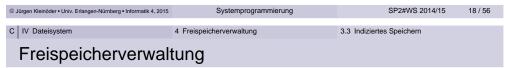




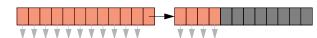
- ★ Einsatz von mehreren Stufen der Indizierung
 - ◆ Inode benötigt sowieso einen Block auf der Platte (Verschnitt unproblematisch bei kleinen Dateien)
 - ♦ durch mehrere Stufen der Indizierung auch große Dateien adressierbar

Nachteil

◆ mehrere Blöcke müssen geladen werden (nur bei langen Dateien)



- Prinzipiell ähnlich wie Verwaltung von freiem Hauptspeicher
 - ◆ Bitvektoren zeigen für jeden Block Belegung an
 - ◆ verkettete Listen repräsentieren freie Blöcke
 - Verkettung kann in den freien Blöcken vorgenommen werden
 - Optimierung: aufeinanderfolgende Blöcke werden nicht einzeln aufgenommen, sondern als Stück verwaltet
 - Optimierung: ein freier Block enthält viele Blocknummern weiterer freier Blöcke und evtl. die Blocknummer eines weiteren Blocks mit den Nummern freier Blöcke



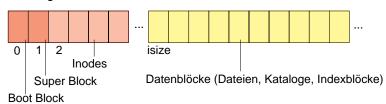
© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 19 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 20 / 56

C IV Dateisystem 5 Beispiel: UNIX File Systems 5.1 System V File System

Beispiel: UNIX File Systems

5.1 System V File System

Blockorganisation



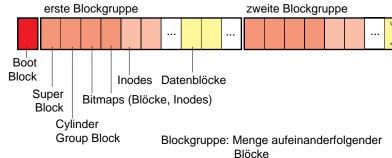
- ◆ Boot Block enthält Informationen zum Laden eines initialen Programms
- ◆ Super Block enthält Verwaltungsinformation für ein Dateisystem
 - · Anzahl der Blöcke, Anzahl der Inodes
 - Anzahl und Liste freier Blöcke und freier Inodes
 - Attribute (z.B. Modified flag)

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 21 / 56

C | IV Dateisystem 5 Beispiel: UNIX File Systems 5.3 Linux EXT2 File System

5.3 Linux EXT2 File System

Blockorganisation



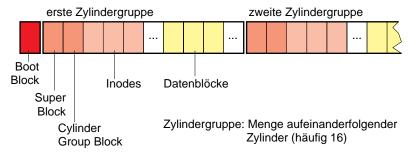
- ♦ Ähnliches Layout wie BSD FFS
- ◆ Blockgruppen unabhängig von Zylindern

5.2 BSD 4.2 (Berkeley Fast File System)

5 Beispiel: UNIX File Systems

Blockorganisation

C IV Dateisystem



5.2 BSD 4.2 (Berkeley Fast File System)

- ◆ Kopie des Super Blocks in jeder Zylindergruppe
- ♦ freie Inodes u. freie Datenblöcke werden im *Cylinder Group Block* gehalten
- ♦ eine Datei wird möglichst innerhalb einer Zylindergruppe gespeichert
- ★ Vorteil: kürzere Positionierungszeiten

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 22 / 56

C | IV Dateisystem 6 Beispiel: Windows NT (NTFS) 5.3 Linux EXT2 File System

Beispiel: Windows NT (NTFS)

Deispiel. Williauws IVI (IVIF3

- Dateisystem für Windows NT
- Datei
 - ◆ beliebiger Inhalt; für das Betriebssystem ist der Inhalt transparent
 - ◆ Rechte verknüpft mit NT-Benutzern und -Gruppen
 - ◆ Datei kann automatisch komprimiert oder verschlüsselt gespeichert werden
 - ◆ große Dateien bis zu 2⁶⁴ Bytes lang
 - ◆ Hard links: mehrere Einträge derselben Datei in verschiedenen Katalogen möglich
- Dateiinhalt: Sammlung von Streams
 - ◆ Stream: einfache, unstrukturierte Folge von Bytes
 - ◆ "normaler Inhalt" = unbenannter Stream (default stream)
 - dynamisch erweiterbar
 - ◆ Syntax: dateiname:streamname

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 23 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 24 / 56

C | IV Dateisystem 6 Beispiel: Windows NT (NTFS) 6.1 Dateiverwaltung

6.1 Dateiverwaltung

- Basiseinheit "Cluster"
 - ◆ 512 Bytes bis 4 Kilobytes (beim Formatieren festgelegt)
 - ◆ wird auf eine Menge von hintereinanderfolgenden Blöcken abgebildet
 - ◆ logische Cluster-Nummer als Adresse (LCN)
- Basiseinheit "Strom"
 - ◆ jede Datei kann mehrere (Daten-)Ströme speichern
 - ♦ einer der Ströme wird für die eigentlichen Daten verwendet
 - ◆ Dateiname, MS-DOS Dateiname, Zugriffsrechte, Attribute und Zeitstempel werden jeweils in eigenen Datenströmen gespeichert (leichte Erweiterbarkeit des Systems)

6.1 Dateiverwaltung (2)

■ File-Reference

nummer

C IV Dateisystem

◆ Bezeichnet eindeutig eine Datei oder einen Katalog



6 Beispiel: Windows NT (NTFS)

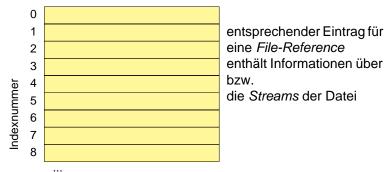
• Dateinummer ist Index in eine globale Tabelle (MFT: Master File Table)

6.1 Dateiverwaltung

 Sequenznummer wird hochgezählt, für jede neue Datei mit gleicher Dateinummer



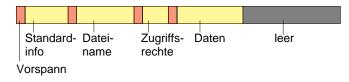
- Rückgrat des gesamten Systems
 - ◆ große Tabelle mit gleich langen Elementen (1KB, 2KB oder 4KB groß, je nach Clustergröße)
 - ◆ kann dynamisch erweitert werden



◆ Index in die Tabelle ist Teil der File-Reference



Eintrag für eine kurze Datei



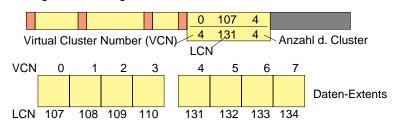
- Streams
 - ◆ Standard-Information (immer in der MFT)
 - enthält Länge, Standard-Attribute, Zeitstempel, Anzahl der Hard links, Sequenznummer der gültigen File-Reference
 - ◆ Dateiname (immer in der MFT)
 - kann mehrfach vorkommen (Hard links)
 - ◆ Zugriffsrechte (Security Descriptor)
 - ◆ Eigentliche Daten

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 27 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 28 / 56

C | IV Dateisystem 6 Beispiel: Windows NT (NTFS) 6.2 Master-File-Table

6.2 Master-File-Table (3)

■ Eintrag für eine längere Datei



- ◆ Extents werden außerhalb der MFT in aufeinanderfolgenden Clustern gespeichert
- ◆ Lokalisierungsinformationen werden in einem eigenen Stream gespeichert

6.2 Master-File-Table (4)

- Mögliche weitere Streams (Attributes)
 - ◆ Index

C IV Dateisystem

 Index über einen Attributschlüssel (z.B. Dateinamen) implementiert Katalog

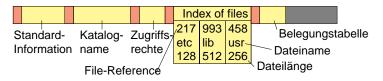
6 Beispiel: Windows NT (NTFS)

6.2 Master-File-Table

- ◆ Indexbelegungstabelle
 - Belegung der Struktur eines Index
- ◆ Attributliste (immer in der MFT)
 - wird benötigt, falls nicht alle Streams in einen MFT Eintrag passen
 - referenzieren weitere MFT Einträge und deren Inhalt
- ◆ Streams mit beliebigen Daten
 - wird gerne zum Verstecken von Viren genutzt, da viele Standard-Werkzeuge von Windows nicht auf die Bearbeitung mehrerer Streams eingestellt sind (arbeiten nur mit dem unbenannten Stream)



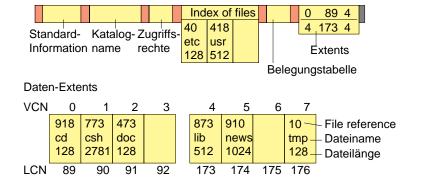
■ Eintrag für einen kurzen Katalog



- ◆ Dateien des Katalogs werden mit File-References benannt
- Name und Standard-Attribute (z.B. Länge) der im Katalog enthaltenen Dateien und Kataloge werden auch im Index gespeichert (doppelter Aufwand beim Update; schnellerer Zugriff beim Kataloglisten)



Eintrag für einen längeren Katalog

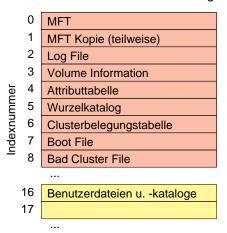


- ◆ Speicherung als B+-Baum (sortiert, schneller Zugriff)
- ♦ in einen Cluster passen zwischen 3 und 15 Dateien (im Bild nur eine)

© Jürgen Kleinöder • Univ. Erlangen-Nürmberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 31 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürmberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 32 / 56

c | IV Dateisystem 6 Beispiel: Windows NT (NTFS) 6.3 Metadaten c | IV Dateisystem 6.3 Metadaten 6.3 Metadaten

Alle Metadaten werden in Dateien gehalten



Feste Dateien in der MFT

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015	Systemprogrammierung	SP2#WS 2014/15	33 / 56
C IV Dateisystem	6 Beispiel: Windows NT (NTFS)	6.4 Fehlererholung	
6.4 Fehlererholung			

- NTFS ist ein Journal-File-System
 - ◆ Änderungen an der MFT und an Dateien werden protokolliert.
 - ◆ Konsistenz der Daten und Metadaten kann nach einem Systemausfall durch Abgleich des Protokolls mit den Daten wieder hergestellt werden.
- Nachteile
 - etwas ineffizienter
 - ◆ nur für Volumes >400 MB geeignet

6.3 Metadaten (2)

- Bedeutung der Metadateien
 - ◆ MFT und MFT Kopie: MFT wird selbst als Datei gehalten (d.h. Cluster der MFT stehen im Eintrag 0)
 MFT Kopie enthält die ersten 16 Einträge der MFT (Fehlertoleranz)
 - ◆ Log File: enthält protokollierte Änderungen am Dateisystem

6 Beispiel: Windows NT (NTFS)

◆ Volume Information: Name, Größe und ähnliche Attribute des Volumes

6.3 Metadater

- ◆ Attributtabelle: definiert mögliche Ströme in den Einträgen
- ♦ Wurzelkatalog
- ◆ Clusterbelegungstabelle: Bitmap für jeden Cluster des Volumes
- Boot File: enthält initiales Programm zum Laden, sowie ersten Cluster der MFT
- ◆ Bad Cluster File: enthält alle nicht lesbaren Cluster der Platte NTFS markiert automatisch alle schlechten Cluster und versucht die Daten in einen anderen Cluster zu retten

- Metadaten und aktuell genutzte Datenblöcke geöffneter Dateien werden im Hauptspeicher gehalten (Dateisystem-Cache)
 - ◆ effizienter Zugriff
 - ♦ Konsistenz zwischen Cache und Platte muss regelmäßig hergestellt werden
 - ➤ synchrone Änderungen: Operation kehrt erst zurück, wenn Änderungen auf der Platte gespeichert wurden
 - ➤ asynchrone Änderungen: Änderungen erfolgen nur im Cache, Operation kehrt danach sofort zurück, Synchronisation mit der Platte erfolgt später
- Mögliche Fehlerursachen
 - ◆ Stromausfall (dummer Benutzer schaltet einfach Rechner aus)
 - ◆ Systemabsturz

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 35 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 36 / 56

c N Dateisystem 7 Dateisysteme mit Fehlererholung 7.1 Konsistenzprobleme 7 Dateisysteme mit Fehlererholung 7.2 Journaling-File-Systems 7.1 Konsistenzprobleme 7.2 Journaling-File-Systems

- Fehlerursachen & Auswirkungen auf das Dateisystem
 - ◆ Cache-Inhalte und aktuelle E/A-Operationen gehen verloren
 - ◆ inkonsistente Metadaten
 - z. B. Katalogeintrag fehlt zur Datei oder umgekehrt
 - z. B. Block ist benutzt aber nicht als belegt markiert
- * Reparaturprogramme
 - Programme wie chkdsk, scandisk oder fsck k\u00f6nnen inkonsistente Metadaten reparieren
- Datenverluste bei Reparatur möglich
- Große Platten bedeuten lange Laufzeiten der Reparaturprogramme

- Zusätzlich zum Schreiben der Daten und Meta-Daten (z. B. Inodes) wird ein Protokoll der Änderungen geführt
 - ◆ Grundidee: Log-based Recovery bei Datenbanken
 - ◆ alle Änderungen treten als Teil von Transaktionen auf.
 - ◆ Beispiele für Transaktionen:
 - Erzeugen, Löschen, Erweitern, Verkürzen von Dateien
 - Dateiattribute verändern
 - Datei umbenennen
 - ◆ Protokollieren aller Änderungen am Dateisystem zusätzlich in einer Protokolldatei (*Log File*)
 - beim Bootvorgang wird Protokolldatei mit den aktuellen Änderungen abgeglichen und damit werden Inkonsistenzen vermieden.

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 37 / 56

C | IV Dateisystem 7 Dateisysteme mit Fehlererholung 7.2 Journaling-File-Systems

7.2 Journaling-File-Systems (2)

- Protokollierung
 - ◆ für jeden Einzelvorgang einer Transaktion wird zunächst ein Logeintrag erzeugt und
 - ◆ danach die Änderung am Dateisystem vorgenommen
 - ◆ dabei gilt:
 - ➤ der Logeintrag wird immer vor der eigentlichen Änderung auf Platte geschrieben
 - wurde etwas auf Platte geändert, steht auch der Protokolleintrag dazu auf der Platte
- Fehlererholung
 - Beim Bootvorgang wird überprüft, ob die protokollierten Änderungen vorhanden sind:
 - ➤ Transaktion kann wiederholt bzw. abgeschlossen werden (Redo) falls alle Logeinträge vorhanden
 - ➤ angefangene, aber nicht beendete Transaktionen werden rückgängig gemacht (Undo).

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 38 / 56

C | IV Dateisystem 7 Dateisysteme mit Fehlererholung 7.2 Journaling-File-Systems

7.2 Journaling-File-Systems (3)

- Beispiel: Löschen einer Datei im NTFS
 - ◆ Vorgänge der Transaktion
 - ➤ Beginn der Transaktion
 - ➤ Freigeben der Extents durch Löschen der entsprechenden Bits in der Belegungstabelle (gesetzte Bits kennzeichnen belegten Cluster)
 - ➤ Freigeben des MFT-Eintrags der Datei
 - Löschen des Katalogeintrags der Datei (evtl. Freigeben eines Extents aus dem Index)
 - ➤ Ende der Transaktion
 - ◆ Alle Vorgänge werden unter der File-Reference im Log-File protokolliert, danach jeweils durchgeführt.
 - ➤ Protokolleinträge enthalten Informationen zum *Redo* und zum *Undo*

© Jürgen Kleinöder • Univ. Erlangen-Nürmberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 39 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürmberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 40 / 56

C IV Dateisystem 7 Dateisysteme mit Fehlererholung 7.2 Journaling-File-Systems

7.2 Journaling-File-Systems (4)

- ◆ Log vollständig (Ende der Transaktion wurde protokolliert und steht auf Platte):
 - ➤ Redo der Transaktion: alle Operationen werden wiederholt, falls nötig
- ◆ Log unvollständig (Ende der Transaktion steht nicht auf Platte):
 - ➤ *Undo* der Transaktion: in umgekehrter Reihenfolge werden alle Operation rückgängig gemacht
- Checkpoints
 - Log-File kann nicht beliebig groß werden
 - ◆ gelegentlich wird f
 ür einen konsistenten Zustand auf Platte gesorgt (Checkpoint) und dieser Zustand protokolliert (alle Protokolleinträge von vorher können gelöscht werden)
 - ◆ ähnlich verfährt NTFS, wenn Ende des Log-Files erreicht wird.

7.2 Journaling-File-Systems (5)

★ Ergebnis

C IV Dateisystem

◆ eine Transaktion ist entweder vollständig durchgeführt oder gar nicht

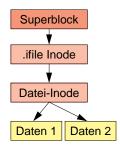
7.2 Journaling-File-Systems

7 Dateisysteme mit Fehlererholung

- ◆ Benutzer kann ebenfalls Transaktionen über mehrere Dateizugriffe definieren, wenn diese ebenfalls im Log erfasst werden
- ◆ keine inkonsistenten Metadaten möglich
- ◆ Hochfahren eines abgestürzten Systems benötigt nur den relativ kurzen Durchgang durch das Log-File.
 - Alternative chkdsk benötigt viel Zeit bei großen Platten
- Nachteile
 - ◆ ineffizienter, da zusätzliches Log-File geschrieben wird
- Beispiele: NTFS, EXT3, EXT4, ReiserFS



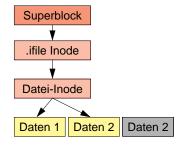
- Alternatives Konzept zur Realisierung von atomaren Änderungen
- Alle Änderungen im Dateisystem erfolgen auf Kopien
 - ◆ Der Inhalt veränderter Blöcke wird in einen neuen Block geschrieben



◆ Beispiel LinLogFS: Superblock einziger nicht ersetzter Block



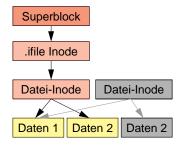
- Alternatives Konzept zur Realisierung von atomaren Änderungen
- Alle Änderungen im Dateisystem erfolgen auf Kopien
 - ◆ Der Inhalt veränderter Blöcke wird in einen neuen Block geschrieben



◆ Beispiel LinLogFS: Superblock einziger nicht ersetzter Block

SP2#WS 2014/15 43 / 56 SP2#WS 2014/15 44 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung

- Alternatives Konzept zur Realisierung von atomaren Änderungen
- Alle Änderungen im Dateisystem erfolgen auf Kopien
 - ◆ Der Inhalt veränderter Blöcke wird in einen neuen Block geschrieben



◆ Beispiel LinLogFS: Superblock einziger nicht ersetzter Block

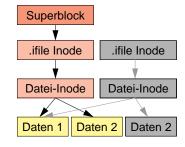
7.3 Copy-on-Write- / Log-Structured-File-Systems

Alternatives Konzept zur Realisierung von atomaren Änderungen

7 Dateisysteme mit Fehlererholung

- Alle Änderungen im Dateisystem erfolgen auf Kopien
 - ◆ Der Inhalt veränderter Blöcke wird in einen neuen Block geschrieben

7.3 Copy-on-Write- / Log-Structured-File-Systems

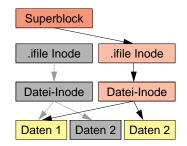


C IV Dateisystem

◆ Beispiel LinLogFS: Superblock einziger nicht ersetzter Block



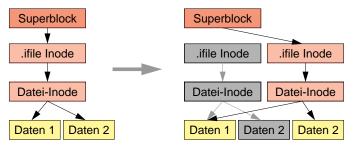
- Alternatives Konzept zur Realisierung von atomaren Änderungen
- Alle Änderungen im Dateisystem erfolgen auf Kopien
 - ◆ Der Inhalt veränderter Blöcke wird in einen neuen Block geschrieben



◆ Beispiel LinLogFS: Superblock einziger nicht ersetzter Block



- Alternatives Konzept zur Realisierung von atomaren Änderungen
- Alle Änderungen im Dateisystem erfolgen auf Kopien
 - ◆ Der Inhalt veränderter Blöcke wird in einen neuen Block geschrieben



Beispiel LinLogFS: Superblock einziger statischer Block (Anker im System)

SP2#WS 2014/15 47 / 56 SP2#WS 2014/15 48 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung

C IV Dateisystem C IV Dateisystem 7 Dateisysteme mit Fehlererholung 7.4 Copy-on-Write- / Log-Structured-File-Systems 8 Fehlerhafte Plattenblöcke 7.4 Copy-on-Write- / Log-Structured-File-Systems

7.4 Copy-on-Write- / Log-Structured-File-Systems (2)

- ★ Vorteile
 - ◆ Datenkonsistenz bei Systemausfällen
 - ein atomare Änderung macht alle zusammengehörigen Änderungen
 - ◆ Schnappschüsse / Checkpoints einfach realisierbar
- Nachteile
 - ◆ Gesamtperformanz geringer
- Unterschied zwischen Copy-in-Write- und Log-Structured-File-Systems
 - ◆ Log-Structured-File-Systems schreiben kontinuierlich an's Ende des belegten Plattenbereichs und geben vorne die Blöcke wieder frei
 - ➤ Gute Schreibeffizienz
 - ➤ Annahme: Lesen kann primär aus dem Cache erfolgen
 - ◆ Beispiele: Log-Structured: LinLogFS, BSD LFS Copy-on-Write: Btrfs (Oracle)

SP2#WS 2014/15 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015

C IV Dateisystem

9 Datensicherung

9.1 Sichern der Daten auf Tertiärspeiche

Datensicherung

- Schutz vor dem Totalausfall von Platten
 - z. B. durch Head-Crash oder andere Fehler

9.1 Sichern der Daten auf Tertiärspeicher

- ➤ Bänder
- ➤ WORM-Speicherplatten (Write Once Read Many)
- Sichern großer Datenbestände
 - ◆ Total-Backups benötigen lange Zeit
 - ◆ Inkrementelle Backups sichern nur Änderungen ab einem bestimmten Zeitpunkt
 - ◆ Mischen von Total-Backups mit inkrementellen Backups

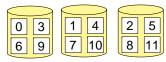
Fehlerhafte Plattenblöcke

Blöcke, die beim Lesen Fehlermeldungen erzeugen

- ◆ z.B. Prüfsummenfehler
- Hardwarelösung
 - ◆ Platte und Plattencontroller bemerken selbst fehlerhafte Blöcke und maskieren diese aus
 - ◆ Zugriff auf den Block wird vom Controller automatisch auf einen "gesunden" Block umgeleitet
- Softwarelösung
 - ◆ File-System bemerkt fehlerhafte Blöcke und markiert diese auch als belegt



- Gestreifte Platten (Striping; RAID 0)
 - ◆ Daten werden über mehrere Platten gespeichert

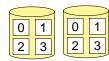


- ◆ Datentransfers sind nun schneller, da mehrere Platten gleichzeitig angesprochen werden können
- Nachteil
 - ◆ keinerlei Datensicherung: Ausfall einer Platte lässt Gesamtsystem ausfallen

C IV Dateisystem 9 Datensicherung 9.2 Einsatz mehrerer (redundanter) Platten

9.2 Einsatz mehrerer redundanter Platten (2)

- Gespiegelte Platten (*Mirroring*; RAID 1)
 - ◆ Daten werden auf zwei Platten gleichzeitig gespeichert



- Implementierung durch Software (File-System, Plattentreiber) oder Hardware (spez. Controller)
- ◆ eine Platte kann ausfallen
- schnelleres Lesen (da zwei Platten unabhängig voneinander beauftragt werden können)
- Nachteil
 - doppelter Speicherbedarf
- wenig langsameres Schreiben durch Warten auf zwei Plattentransfers
- Verknüpfung von RAID 0 und 1 möglich (RAID 0+1)

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015

Systemprogrammierung

SP2#WS 2014/15

53 / 56

C | IV Dateisystem

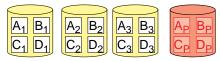
9 Datensicherung

9.2 Einsatz mehrerer (redundanter) Platten

9.2 Einsatz mehrerer redundanter Platten (4)

- Nachteil von RAID 4
 - ◆ jeder Schreibvorgang erfordert auch das Schreiben des Paritätsblocks
 - ◆ Erzeugung des Paritätsblocks durch Speichern des vorherigen Blockinhalts möglich: P_{neu} = P_{alt} ⊕ B_{alt} ⊕ B_{neu} (P=Parity, B=Block)
 - ◆ Schreiben eines kompletten Streifens benötigt nur einmaliges Schreiben des Paritätsblocks
 - ◆ Paritätsplatte ist hoch belastet

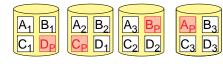
- Paritätsplatte (RAID 4)
 - ◆ Daten werden über mehrere Platten gespeichert, eine Platte enthält Parität



- ◆ Paritätsblock enthält byteweise XOR-Verknüpfungen von den zugehörigen Blöcken aus den anderen Streifen
- ◆ eine Platte kann ausfallen
- ◆ schnelles Lesen
- ◆ prinzipiell beliebige Plattenanzahl (ab drei)



- Verstreuter Paritätsblock (RAID 5)
 - ◆ Paritätsblock wird über alle Platten verstreut



- zusätzliche Belastung durch Schreiben des Paritätsblocks wird auf alle Platten verteilt
- ♦ heute g\u00e4ngigstes Verfahren redundanter Platten
- ♦ Vor- und Nachteile wie RAID 4

© Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 55 / 56 © Jürgen Kleinöder • Univ. Erlangen-Nürnberg • Informatik 4, 2015 Systemprogrammierung SP2#WS 2014/15 56 / 56