Kapitel 6

Dateiverwaltung

- <u>FAT-Dateisystem</u> (MS-DOS, Win95, Win98)
- NTFS-Dateisystem (NT, W2K, WinXP)
- Journaling Dateisysteme
- RAID

- Unix-Dateisystem
- CD-ROM Dateisystem
- Plattenzugriffsverwaltung

Betriebssystem e

235 Prof. Dr. U. Wienkop

Dateiverwaltung

Systemaufrufschnittstelle Betriebssystemkern Dateiverwaltungskomponente Plattentreiber Geräte-Hardware

O Begriffe

- > Datei oder File benannte logischer Abschnitt der Platte
- Dateisystem Gesamtheit der Dateien
- Directories, Subdirectories, root, working directory
- > flache vs. hierarchische Dateisysteme

Plattenbuchhaltung

- > Merken der physikalischen Lage einer benannten Datei
- Umsetzen von Schreib- und Leseaufträge in physikalisch adressierte Transferaufträge

O Von der Plattenbuchhaltung benötigte Daten

- > freie Allokationseinheiten (z. B. Allokationsvektor)
- Verzeichnis(se) der auf der Platte vorhandenen Dateien
- > für jede Datei die belegten Allokationseinheiten
- > ggf. Informationen über die Plattenpartitionierung
- > ggf. Dateiattribute und Eigentums- und Zugriffsrechte

Betriebssystem e

Der interne Aufbau des FAT-Datei-Systems

Plattenaufteilung

- Partitionssektor
- Urlade-Bereich (Boot sector)
- > FAT 1 (File Allocation Table), FAT 2 (Sicherheitskopie von FAT 1)
- Stamm-Katalog (\)Nutzdatenbereich (inklusive aller anderen Kataloge)

Aufbau der FAT

- > 1½ Byte=12 Bit (~FAT12) bzw. 2 Byte (~FAT16) > FAT(i) = 0 --> frei, FAT(i) = FFF7H --> defekt, FAT(i) > FFF7H --> Dateiende, sonst durch Datei belegt

O Directory-Einträge

jeweils 32 Byte pro Eintrag

Betriebssystem e

237 Prof. Dr. U. Wienkop

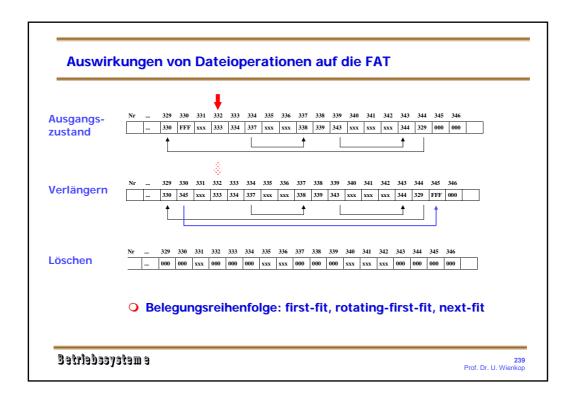
Aufbau der FAT-Directory Struktur

0								8			11				15
		Name			Extension Attr			reserviert							
16						22		24		26		28			31
	reserviert					Uhr	zeit	Datum		FAT-	Eintr.		Datei	größe	!

1	R	read only file	schreibgeschützte Datei
2	Н	hidden file	versteckte Datei
4	S	system file	System-Datei
8	٧	volume label	Partitions- oder Datenträgernamen
0X10	D	directory	Katalogdatei
0X20	Α	archive bit	Archiv-Bit (Datei wurde geändert)

Attribut-Flags

Betriebssystem e



Größe der Partition	Sektoren pro Clu	ısterClustergröße	
0-31 MB	1	0.5K	
-63 MB	2	1K	
-127 MB	4	2K	
-255 MB	8	4K	
-511 MB	16	8K	
-1023 MB	32	16K	
-2047 MB	64	32K	
 Durch den Overh größer als 511ME 		ı das Dateisystem FAT16 nic	ht für Datenträge

Lange Dateinamen bei Win95 und NT

0x42	,	v	I	1			:	f	()	0x0F	0x00	Prüf.	x
0x0	000	0xF	FFF	0xF	FFF	0xF	FFF	0xF	FFF	0x0	000	0xF	FFF	0xFFFF
0x01	7	Γ	ł	ı	(е			(1	0x0F	0x00	Prüf.	u
	i c		2	1	1			ı)	0x0	0000 r		o	
T	Н	Е	Q	U	I	~	1	F	О	X	0x20	NT		Uhrzeit
Dat	tum		er Zu- iff	0x0	000		derung: zeit	Veränd Dat	derung:	erster (Cluster		Datei	größe

- **O** Verwendung von sekundären Ordnereinträgen
- Gekennzeichnet durch das Attribut 0x0F = "Datenträgerbezeichnung", "Schreibgeschützt", "System" und "versteckt"
- O Nummerierung: 01, 02, ... Ende+0x40

Betriebssystem e

241 Prof. Dr. U. Wienkop

NTFS-Dateisystem

Betriebssystem e

Dateisysteme FAT und NTFS im Vergleich

Vorteile des Dateisystems FAT

➤ Bessere Eignung bei kleinen Datenträgern (~500MB), da nur geringer Overhead

O Vorteile des Dateisystems NTFS

- Zugriffsberechtigungen
- > Effiziente Verwaltung großer Datenträger bis 2 TBytes (derzeit), FAT16 nur bis 2 GB
- > Kleinere Clustergrößen durch 64 Bit-Adressen möglich --> geringerer Verschnitt
- > Wiederherstellbarkeit des Dateisystems
- > Einbindung eines Fehlertoleranztreibers möglich
- > Schnellerer Zugriff auf Dateien in großen Ordnern (B-Tree)
- Kleine Dateien und Ordner können direkt im MFT-Eintrag abgelegt sein --> weniger Plattenzugriffe

Betriebssystem e

243 Prof. Dr. U. Wienkop

Windows NT Dateisystem (NTFS)

O NTFS-Entwurfsziele:

- > Wiederherstellbarkeit des Dateisystems
- Sicherheit (Zugriffsschutz)
- Datenredundanz und Fehlertoleranz
- ➤ Große Festplatten (2⁶⁴ Cluster mit max. Größe zu 64KB)
- > Multiple Datenströme: Dateiaufbau aus Attributen und Datenströmen

➤ max. Dateilänge: 2⁶⁴ Byte



Betriebssystem e

Windows NT Dateisystem (NTFS)

- O Datenträger enthält nur Dateien Selbst die Metadaten des Dateisystems sind Teil einer Datei
- O Dateien bestehen nur aus Attributen, z.B.
 - Datenattribut
 - Sicherheitsattribut
 - Dateinamenattribut
- O Aufbau eines NTFS-Datenträgers

Boot-Sektor der	Master-Dateitabelle	Systemdateien	Dateibereich
Partition		-	

Betriebssystem e

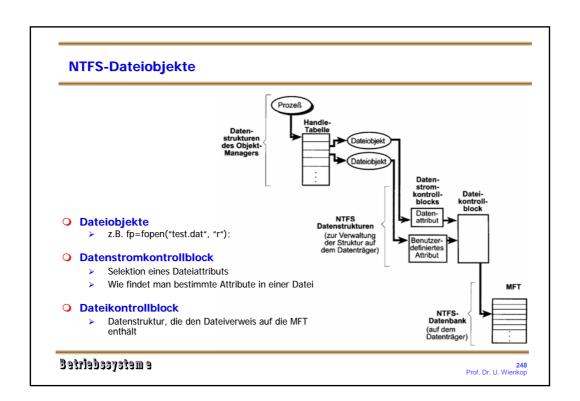
245 Prof. Dr. U. Wienkop

NTFS-Systemdateien

Systemdatei	D-Name	Nr	Zweck der Datei
Masterdateitabelle	\$Mft	0	Auflistung des Inhalts des NTFS-Datenträgers
Master-Dateitabelle2	\$MftMirr	1	Spiegelung der ersten drei Datensätze der MFT
Logdatei	\$LogFile	2	Schritte für die Wiederherstellung von Dateien
Datenträger	\$Volume	3	Name, NTFS-Version und andere Informationen
Attribut-Definitionen	\$AttrDef	4	Tabelle mit Attributnamen, -ziffern und – beschreibungen
Stammindex	\$.	5	Stammverzeichnis
Cluster-Bitmuster	\$Bitmap	6	Zeigt, welche Zuordnungseinheiten belegt sind
Boot-Datei	\$Boot	7	Enthält den Bootstrap des Boot-Datenträgers
Beschädigtr Cluster	\$BadClus	8	Enthält die Orte beschädigter Cluster
Kontingent-Tabelle	\$Quota	9	Auslastung der Festplatte für jeden Benutzer
Großbuchstaben-Tab.	\$Upcase	10	Konvertierung der Kleinbuchstaben in die entsprechenden Großbuchstaben im Unicode
		11-	15 Reserviert für zukünftige Verwendung

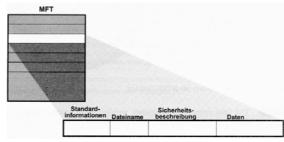
6

```
NTFS-Systemdateien
 NTFS File Sector Information Utility (nfi)
  File 0
  Master File Table ($Mft)
      $STANDARD_INFORMATION (resident)
      $FILE_NAME (resident)
      $DATA (nonresident)
          logical sectors 32-102207 (0x20-0x18f3f)
      $BITMAP (nonresident)
          logical sectors 16-23 (0x10-0x17)
logical sectors 2048272-2048279 (0x1f4110-0x1f4117)
 Master File Table Mirror ($MftMirr)
      $STANDARD_INFORMATION (resident)
      $FILE_NAME (resident)
      $DATA (nonresident)
          logical sectors 12289688-12289695 (0xbb8698-0xbb869f)
 File 2
 Log File ($LogFile)
      $STANDARD_INFORMATION (resident)
      $FILE_NAME (resident)
      $DATA (nonresident)
          logical sectors 12289696-12416703 (0xbb86a0-0xbd76bf)
                                                                             247
Prof. Dr. U. Wienkop
Betriebssystem e
```



Zuordnung MFT zur Datei





- Eine Datei wird mit einem 64-Bit-Wert, dem sogenannten Dateiverweis identifiziert
 - Dateinummer = Position der Datei in der MFT-1
 - Sequenznummer = fortlaufende Nummer bei Wiederverwendung von Dateinummern

Betriebssystem e

249 Prof. Dr. U. Wienkop

Dateiattribute

Attributtyp Beschreibung

StandardinformationenZeitpunkt der letzten Speicherung, Bindungszähler, usw.AttributlisteOrte aller Attributeinträge, wenn diese nicht alle in den MFT-

Datensatz passen.

Dateiname Wiederholbares Attribut sowohl für lange / kurze Dateinamen

Sicherheitsdeskriptor Eigentums und Zugriffsrechte

DatenDaten der Datei. NTFS erlaubt mehrere DatenattributeIndex-WurzelWird für die Implementierung von Ordnern benötigt (B-Tree)

Index-Zuordnung Wird für die Implementierung von Ordnern benötigt

Datenträgerinformation Wird nur bei der Systemdatei des Datenträgers verwendet und

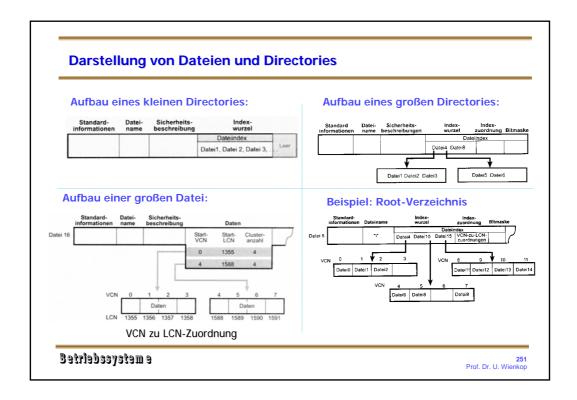
enthält u.a. die Version und den Namen des Datenträgers Fine Karte der in der MFT oder im Ordner belegten Finträge

Bitmuster Eine Karte der in der MFT oder im Ordner belegten Einträge **Erweiterte Attributinformation** Wird von Datei-Servern verwendet, die mit OS/2-

Systemen verbunden sind

Erweiterte Attribute dto.

Betriebssystem e



Dateisystemwiederherstellung

Wiederherstellung des Dateisystems nach

- Stromausfall (vorzeitigem Abschalten)
- Systemfehlern (Absturz)

Dateisystemarten:

Careful-Write-Dateisysteme

- Schreiboperationen auf das Dateisystem werden so angeordnet, daß ein Systemabsturz im schlimmsten Fall eine vorhersagbare, unkritische Inkonsistenz erzeugt
- Konsequenz: strenge Sequenzialisierung der Schreiboperationen
 -> eher langsam

O Lazy-Write-Dateisysteme

- Veränderungen von Dateien finden zunächst nur im Cache statt und werden dann in optimierter Weise (im Hintergrund) zurückgeschrieben
- Verringerung der physikalischen Plattenzugriffe
 - -> schnell, aber eher unsicher

Betriebssystem e

Dateisystemwiederherstellung (2)

- Wiederherstellbares Dateisystem
 - > Erreichen der Sicherheit von Careful-Write-Systemen
 - > Erzielen der Geschwindigkeit von Lazy-Write-Systemen
 - > Konsistenzgarantie durch Protokollierungstechniken
- Transaktionsbasiertes Protokollierungsschema
 - Protokolldatei (Log File)
 - > Vermerk: Aktion abgeschlossen oder ...
 - > aufgezeichnete Teilaktionen wieder rückgängig machen
 - → Journaling File Systems

Betriebssystem e

of Dr. II. Wionke

Prof. Dr. U. Wienko

Journaling-Dateisysteme am Beispiel von NTFS

- Ursprünglich bei DEC für das Betriebssystem VMS entwickelt, später Übernahme nach Unix und auch NTFS
- Maßnahmen von NTFS sichern nur die Konsistenz des Dateisystems, verhindern jedoch keinen Nutzdatenverlust
 - Transaktionskonzept: Alle Veränderungen an Dateien werden durch einen Transaktionsbeginn und ein Transaktionsende (Zusicherung, committment) geklammert, die beide zusammen mit den geplanten Datei-Operationen in der Protokolldatei \$LogFil (sog. Journal) festgehalten
 - Die angegebenen Datei-Operationen k\u00f6nnen in einem beliebigen Zeitraum bis zum Transaktionsende abgewickelt und dabei insbesondere mit Hilfe eines Platten-Caches im Hauptspeicher zwecks Durchsatzerh\u00f6hung aufgeschoben werden.
 - Der Cache-Manager muss nur dafür sorgen, dass der Transaktionsbeginn als erste Teiloperation physikalisch auf die Platte geschrieben wird und das Transaktionsende als letzte
 - Durch einen Ausfall des Rechners unvollständig gebliebene Transaktionen werden später mit Hilfe der Protokolldatei entweder vervollständigt oder ungeschehen gemacht

Betriebssystem e

Journaling-Dateisysteme Protokolldatei \$LogFil

O Protokolldatei \$LogFil ist von ihrer Lage und Größe her konstant.

- > Veränderungen an dieser Datei betreffen daher nur ihren Nutzdatenbereich
- Schreibaufträge, die die Buchhaltungsdaten von \$LogFil betreffen, sind somit überflüssig und können, sofern sie überhaupt ausgeführt werden, keine Konsistenzprobleme verursachen
- Die Protokolldatei wird als (scheinbar endloser) Ringpuffer geführt Sinnvoll, weil die älteren Einträge in der Protokolldatei nach dem Abschluss der zugehörigen Transaktion ohnehin veraltet sind und überflüssig werden.

Verwaltung der Protokolldatei: Protokolldateidienst (LFS, log file service)

- Nummerierung der Protokolleinträge mit Folgenummern (LSNs, logical sequence numbers) und
- Führt vor dem eigentlichen Protokollbereich einen Zeiger auf den Beginn des aktuellen Bereichs, ab dem die Protokolldatei zur Wiederherstellung des Dateisystems nach einem Ausfall des Rechners gelesen werden muss
- Auch alle Änderungen in der Protokolldatei werden vom LFS veranlasst.

Betriebssystem e

255 Prof. Dr. U. Wienko

Journaling-Dateisysteme Einträge in die Protokolldatei

Einträge in die Protokolldatei

Aktualisierungseintrag (update entry)

Beschreibt eine Teiloperation und enthält zweierlei Informationen (s. Beispiel unten):

Informationen für die Wiederholung geben an, wie eine Teiloperation erneut durchgeführt werden soll

- □ Informationen für die Rückgängigmachung geben an, wie eine Teiloperation rückgängig gemacht wird
- > Aufsetzpunkteintrag (checkpoint entry)

Definiert einen Wiederaufsetzpunkt (checkpoint) und wird alle 5 Sekunden in die Protokolldatei geschrieben

Mit den Informationen, die in diesem Eintrag gespeichert sind, kann NTFS entscheiden, wie weit es in der Protokolldatei bei einer Wiederherstellungsaktion zurückgehen muss. Der obengenannte Zeiger auf den Beginn des aktuellen Bereichs zielt auf den letzten Aufsetzpunkteintrag.

Zusicherungseintrag (committment entry)

Ein solcher Eintrag schließt eine Transaktion ab. Die zu einer Transaktion gehörigen Einträge sind in der Protokolldatei in einer verzeigerten Liste miteinander verbunden, zu deren letztem Element der Zusicherungseintrag wird

Transaktionsanfang

Markierung durch den ersten Aktualisierungseintrag der Transaktion

Betriebssystem e

Journaling-Dateisysteme Übertragungsaspekte

O Schreiben auf die Festplatte

- Nicht sofort, wenn ein Zusicherungseintrag in die Protokolldatei geschrieben wird, wird diese auch auf die Platte hinausgeschrieben, sondern mehrere zugesicherte Transaktionen werden erst einmal gesammelt und dann in einem Rutsch physikalisch übertragen
- Jedoch spätestens dann, wenn ein Aufsetzpunkteintrag in die Protokolldatei geschrieben wird; dies löst immer auch ein Hinausschreiben der Protokolldatei aus.

O Gestaltung der Teiloperationen

- Mehrfache Ausführung bewirkt dasselbe wie nur eine einfache Ausführung (Idempotenz)
 - Beispiel: bei der Allokation neuer Sektorgruppen müssen die zugehörigen Bits im Allokationsvektor (Sektorgruppen-Bitmuster) in der Datei \$Bitmap gesetzt werden. Eine Wiederholung dieser Teiloperation setzt die bereits gesetzten Bits eben ein zweites Mal

Betriebssystem e

Prof. Dr. U. Wienko

Journaling-Dateisysteme Ablauf einer Transaktion

Ablauf einer Transaktion

- Mit dem Einfügen eines Aktualisierungseintrags in die Protokolldatei, der keiner anderen Transaktion angehört (also mit deren Einträgen verkettet ist), beginnt eine neue Transaktion
- > Es werden danach verschiedene Teiloperationen ausgeführt
 - jeweils vorher mit einem Eintrag in der Protokolldatei vermerkt und
 - untereinander verkettet werden
- Wegen der Parallelität der Dateizugriffe verschiedener Prozesse sind übrigens die verketteten Listen mehrerer Transaktionen ineinander verschachtelt
- Das alles findet zunächst im Hauptspeicher, also im Platten-Cache, statt. Wenn zu diesem Zeitpunkt der Rechner ausfällt, geht diese Transaktion zwar komplett verloren, aber das Dateisystem verbleibt in einem konsistenten (dem alten) Zustand.

O Hinausschreiben auf die Platte

- z. B., wenn der Platten-Cache voll ist und Platz geschaffen werden muss
- > spätestens aber an einem Wiederaufsetzpunkt

Betriebssystem e

Journaling-Dateisysteme Ablauf einer Transaktion, Log File Service

O LFS (log file service)

- LFS & Cache-Manager: Immer erst die Änderungen in der Protokolldatei vor den Auswirkungen der zugehörigen Teiloperationen hinausschreiben
- Wenn nach der Änderung der Protokolldatei auf der Platte der Rechner ausfällt, können anhand der Einträge in der Protokolldatei die betroffenen Teiloperationen entweder wiederholt oder, falls nötig, rückgängig gemacht werden

O Abschluß einer Transaktion:

Schreiben eines Zusicherungseintrags in die Protokolldatei geschrieben (nicht sofort auch auf die Platte)

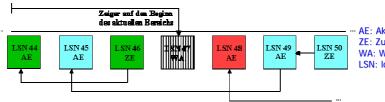
Wiederaufsetzpunkt

- Alle 5 Sekunden wird ein Wiederaufsetzpunkt definiert
 - Zunächst Schreiben aller Änderungen in der Protokolldatei
 - □ Dann alle übrigen Änderungen im Platten-Cache auf die Platte schreiben
 - □ Schließlich einen Aufsetzpunkteintrag in die Protokolldatei (und auch auf die Platte) schreiben
- Außerdem wird in der Protokolldatei der Zeiger auf den Beginn des aktuellen Bereichs auf den neuen Aufsetzpunkteintrag gesetzt.
- Abgeschlossene Transaktionen brauchen später bei einer Wiederherstellung nicht mehr berücksichtigt werden

Betriebssystem e

Prof Dr II Wienk

Journaling-Dateisysteme Beispiel für den Ablauf einer Transaktion – 1



--- AE: Aktualisierungseintrag ZE: Zusicherungseintrag WA: Wiederaufsetzpunkt

LSN: logical sequence number

Neustart bzw. Wiederanlauf von Windows NT

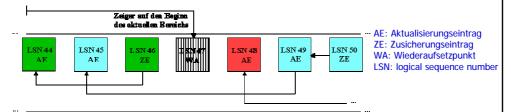
Veranlassung einer Wiederherstellung (die im Normalfall natürlich überflüssig ist). Dazu wird in der Protokolldatei der letzte Aufsetzpunkteintrag bestimmt, und es werden alle Transaktionen, die am Wiederaufsetzpunkt noch nicht zugesichert waren oder die später begonnen (und möglicherweise schon zugesichert) wurden, untersucht.

Nicht abgeschlossene Transaktionen (kein Zusicherungseintrag vorhanden)

werden mit Hilfe der Informationen für die Rückgängigmachung ungeschehen gemacht (man weiß ja nicht, welche Teiloperationen noch bis zur Zusicherung gefehlt haben). In der obigen Abbildung betrifft dies die mit rot gekennzeichnete Transaktion.

Betriebssystem e

Journaling-Dateisysteme Beispiel für den Ablauf einer Transaktion – 2



Zugesicherte Transaktionen

- werden hingegen mit Hilfe der Informationen für die Wiederholung in den Aktualisierungseinträgen wiederholt (denn es können ja zum Zeitpunkt des Ausfalls noch Teilergebnisse im Cache gewesen sein) und damit komplett wiederhergestellt. In der obigen Abbildung betrifft dies die mit blau gekennzeichnete Transaktion, an der man auch erkennt, dass die Wiederherstellungsaktionen teilweise auch noch Einträge vor dem Wiederaufsetzpunkt betreffen.
- Anschließend ist das Dateisystem wieder konsistent, auch wenn die nicht zugesicherten Transaktionen dabei verloren gegangen sind!

Betriebssystem e

261 Prof. Dr. U. Wienkop

Journaling-Dateisysteme vs. Standard-Dateisysteme

Geschwindigkeit

- Im Gegensatz zu Unix läuft eine Wiederherstellung des Dateisystems in NT wesentlich schneller ab, da im wesentlichen ja nur die Transaktionen seit dem letzten Wiederaufsetzpunkt berücksichtigt werden müssen.
- In Unix wird, sofern das Dateisystem nicht sauber heruntergefahren wurde (erkenntlich an einem Bit im Superblock), das Programm fsck gestartet, das den ganzen Dateisystembaum nach Inkonsistenzen absucht. Das kann bei großen Festplatten einige Minuten dauern.

Sicherheit

 Es werden tatsächlich alle Operationen auf dem Dateisystem erfasst. Es bleiben keine etwaigen Inkonsistenzen – wie sie trotz fsck oder scandisk auftreten können – zurück (s. Übungsaufgabe)

Betriebssystem e

Weitere NTFS-Eigenschaften

Fehlertoleranztreiber

Funktionen zur redundanten Datenspeicherung und zur dynamischen Wiederherstellung von Daten

O Datenträgersätze

> Mehrere Partitionen können zu einem Dateisystem konsolidiert werden

Stripe Sets

- > Folge von Partitionen, von denen jede auf einer anderen Platte liegt
- Zusammenfassung zu einem logischen Dateisystem
- > Dateien werden über mehrere Platten hinweg verteilt -> Geschwindigkeitssteigerung



Betriebssystem e

Prof. Dr. U. Wier

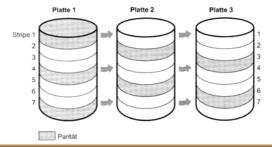
NTFS-Fehlertoleranz-Eigenschaften

Spiegelsätze

- Inhalt einer Partition wird dupliziert und in eine Partition der gleichen Größe auf einer anderen Platte geschrieben
- Im Fehlerfall: Zugreifen auf die andere Platte
- Möglichkeit: Lastverteilung durch parallelen Lesezugriff auf beide Platten

O Stripe-Sets mit Parität

Parität: Byte-für-Byte gebildetes log. XOR



Betriebssystem e

RAID

Redundant Array of Independent Disks

Betriebssystem e

265 Prof. Dr. U. Wienkop

RAID – Redundant Array of Independent Disks

(früher auch: Redundant Array of Inexpensive Disks)

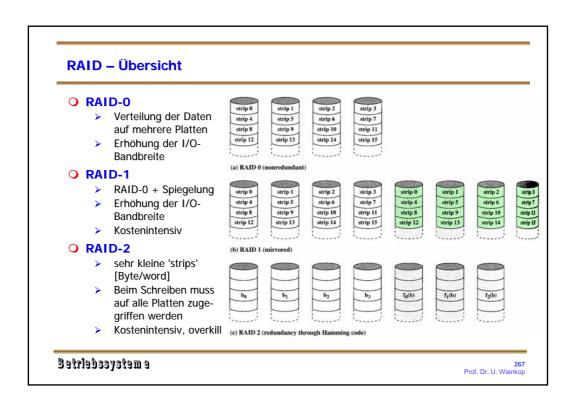
Charakteristika

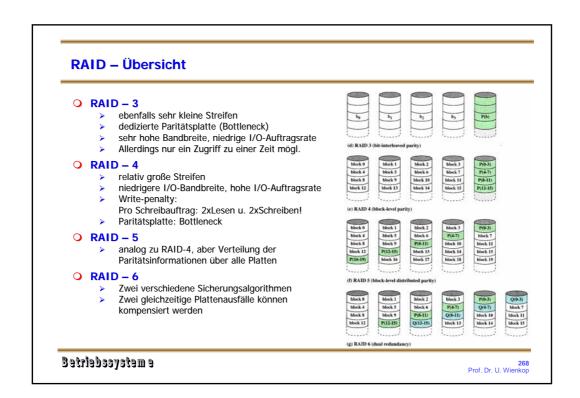
- RAID: Menge von physikalischen Laufwerken, die vom BS als ein einzelnes (logisches) Laufwerk gesehen werden
- > Die Daten sind über die physikalischen Laufwerke verteilt
- Redundante Datenkapazität wird verwendet, um Paritätsinformationen abzuspeichern, die beim Ausfall einer Platte ermöglicht, die Daten wieder herzustellen
- > RAID-Ebenen: 0 ... 6

O Ausfallsicherheit

- Verwendung mehrerer unabhängiger Platten erhöht die Fehlerwahrscheinlichkeit n Geräte haben i.a. 1/n der Ausfallsicherheit eines Laufwerks!
- > Redundante Informationen können die Ausfallsicherheit erhöhen
- MTTF (mean time to fail) von Platten [10 Jahre], MTTR (mean time to repair) [Stunden], d.h. kritischer Zeitraum lediglich: Während eine Platte repariert wird, fällt noch eine zweite aus.
- O Hot-spares: Platten, die nur als Ersatz für den Ausfall anderer Platten bereitstehen
- Hot-swapping: Platten können ausgetauscht werden, ohne dass das System heruntergefahren werden muss

Betriebssystem e





Redundanz durch Parität, Zusatzaufwendungen

O RAID - 3

 \rightarrow X4(i) = X3(i) \oplus X2(i) \oplus X1(i) \oplus X0(i)

X0 .. X3 - Datenplatten

X4 – Paritätsplatte

- ➤ Annahme X1 fällt aus. Wir addieren X4(i) ⊕ X1(i) auf beide Seiten der Gleichung
- X1(i) = X4(i) ⊕ X3(i) ⊕ X2(i) ⊕ X0(i)
- > Fehlende Daten können somit aus den verbleibenden Platten rekonstruiert werden

O RAID - 4/5

- Schreiben von Daten in den Streifen auf Platte X1; d.h. X1 → X1¹
- Initial: X4(i) = X3(i) ⊕ X2(i) ⊕ X1(i) ⊕ X0(i)

> X4'(i) = X3(i) ⊕ X2(i) ⊕ X1'(i) ⊕ X0(i)

= $X3(i) \oplus X2(i) \oplus X1'(i) \oplus X0(i) \oplus X1(i) \oplus X1(i)$ Einsetzen von X1(i) s.o.

 $= X3(1) \oplus X2(1) \oplus X1'(1) \oplus X0(1) \oplus X1(1) \oplus X4(1) \oplus X3(1) \oplus X2(1) \oplus X0(1)$

= $X4(i) \oplus X1(i) \oplus X1'(i)$

- Write-penalty: Zwei Lesezugriffe (X1-alt, X4-Parität) und zwei Schreibzugriffe (X1'-neu, X4-neue Parität) erforderlich!
- Im Falle größerer Schreiboperationen, die alle Platten betreffen, ist ein einfaches Neuerrechnen der Parität ausreichend

Betriebssystem e

269 Prof. Dr. U. Wienkop

RAID - Ebenen im Vergleich

Category	Level	Description	I/O Request Rate (Read/Write)	Data Transfer Rate (Read/Write)	Typical Application
Striping 0		Non-redundant	Large strips: Excellent	Small strips: Excellent	Applications requiring high performance for noncritical data
Mirroring	1	Mirrored	Good/fair	Fair/fair	System drives; critical files
	2	Redundant via Hamming code	Poor	Excellent	
Parallel access	3	Bit-interleaved parity	Poor	Excellent	Large I/O request size applications such as imaging, CAD
	4	Block-interleaved parity	Excellent/fair	Fair/poor	
Independent access	5	Block-interleaved distributed parity	Excellent/fair	Fair/poor	High request rate, read-intensive, data lookup
	6	Block-interleaved dual distributed parity	Excellent/poor	Fair/poor	Applications requiring extremely high availability

Betriebssystem e

Unix-Dateisystem

Betriebssystem e

271 Prof. Dr. U. Wienkop

Unix - Dateisystem



O Urladeblock

enthält Code für den Urladevorgang

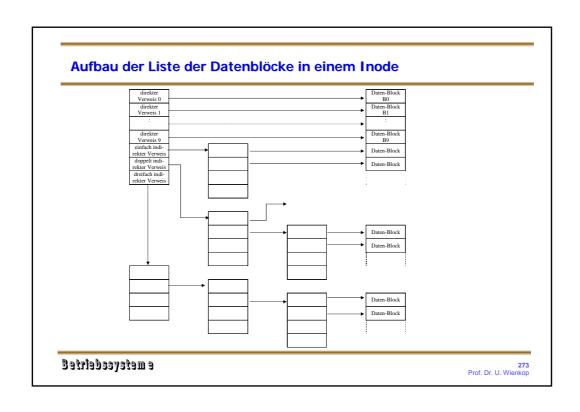
Superblock

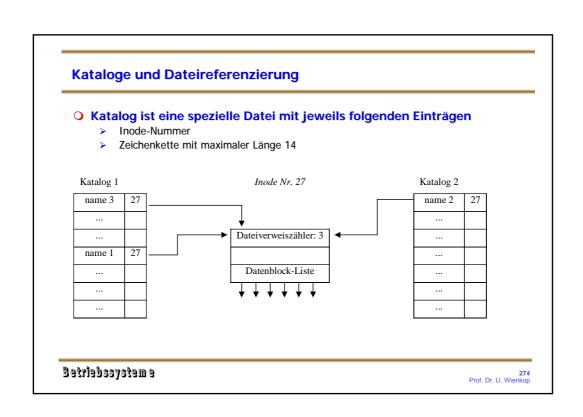
- Größe des Dateisystems (# Blöcke)
- > # freie Blöcke des Dateisystems
- > Liste der freien Blöcke
- > Zeiger auf nächsten Block der Freiliste
- Größe der Inode-Tabelle
- # freie Inodes
- > Liste der freien Inodes
- Zeiger auf nächsten Inode der Freiliste
- > Semaphorfelder für die zwei Freilisten
- Update-Flag

O Inode (index node)

- Dateityp (regulär, Katalog, Zeichenoder Block-Spezial, FIFO)
- Zugriffsrechte (rwxrwxrwx)
- Dateiverweiszähler
- Datei-Eigentümer und dessen Benutzer-Gruppe
- Dateigröße in Bytes
- Zeitstempel für die letzte Datei-Änderung
- > den letzten Datei-Zugriff
- die letzte Inode-Änderung
- Liste der Datenblöcke der Datei

Betriebssystem e





CD-ROM Dateisystem

Betriebssystem e

275 Prof. Dr. U. Wienkop

CD ROM - Dateisystem

O Unterteilung in Sektoren à 3234 Byte

- 2048 Byte Nutzdaten
- > Byte wird durch 14 Bits + 3 Koppelbits dargestellt
- > Synchronisations-, Typ-, Steuer- und Fehlerkorrektur-Informationen
- > Sektor ist in 1,2 oder 4 Blöcke unterteilbar

Adressierung

- > in Form einer Zeitangabe MM:SS:DD (DD ~ 1/75)
- fortlaufende Sektornumerierung
- fortlaufende Blocknumerierung

O Audio-CDs

- > Verzicht auf eine der Fehlerkorrektur-Ebenen
- > 2352 Byte Nutzdaten pro Sektor
- ➤ Fehlerwahrscheinlichkeit 10-8 statt 10-12

Betriebssystem e

Grundstruktur einer CD-ROM

- O Die ersten 16 Sektoren (0 bis 15) bleiben bei CD-ROMs frei
- O Sektor 16: primärer CD-Deskriptor
 - > die Größe der logischen Blöcke (512, 1024 oder 2048 byte)
 - der Name der CD-ROM
 - belegter Platz in Blöcken
 - > Zeiger auf ersten Sektor und Größe der Pfadtabellen
 - Katalogeintrag für den Stammkatalog (mit Zeiger auf dessen ersten Sektor)
 - > Zeitpunkte für früheste und letztmalige Nutzung der CD-ROM (!)
 - verwendeter Standard (normal oder XA, ISO-9660-Version)
 - > Hersteller- und urheberrechtliche Informationen
 - Datum und Zeitpunkt der Herstellung
 - > Ende durch Terminatorsektor
- O Anschließend Pfadtabellen und der Datenbereich des Stammkatalogs
- O Dann Datenbereiche der Unterkataloge und Dateien des Stammkatalogs

Betriebssystem e

27 Prof. Dr. U. Wienko

Kataloge

- O Kataloge ~ spezielle Dateien, Einträge lexikographisch geordnet
- Zusätzlich: gesamte Katalog-Hierarchie in eigenem Datenbereich (Pfadtabelle) am Anfang des Systembereichs der CD ROM
- O Datei-/Katalogeintrag unterschiedlicher Länge

```
Länge des vorliegenden Eintrags in byte (maximal 256)
            Länge des optionalen Attributbereichs in byte (oder Blöcken?), ≤ 256
 2
            des Blocks, mit dem der Datenbereich der Datei beginnt
3 – 10
11 – 18
            Länge des Datenbereichs in byte
19 – 25
            Erstellungsdatum und -zeit
26
            Dateiflaggen
            Bei Verschachtelung: Größe einer Dateneinheit (Sektorgruppe)
27
            Lücke zw. 2 Dateneinheiten in Anzahl Sektoren bei Verschachtelung
28
29 - 32
            Folgenummer der CD-ROM in einem Satz von CD-ROMs
            Länge des Dateinamens
33
34 - x
            ISO-konformer Dateiname der Länge x-33 byte (max. 31 byte)
            falls x ungerade: Füllbyte zur Ausrichtung auf Wortadressen.
(x+1)
            Systemdaten: frei für Systemnutzung, z. B. bei XA und RRIP belegt.
ab x+1|2
```

Betriebssystem e

Dateidarstellung O Die Dateien selbst sind zusammenhängend abgelegt Adressierung durch Angabe des Anfangssektors und der Dateilänge in Byte O Bei XA Möglichkeit zur Verschachtelung (Interleaving) von Binärdaten mit "nichtbinären" Daten Datei durch Dateinummer von 1 bis 255 (= 0, wenn nicht verschachtelt) identifiziert Abstand der Dateisektoren als Anzahl von Sektoren im Katalogeintrag O Beispiel einer verschachtelten Datei zwei Sektor langer Attributbereich Sektorgruppengröße (Byte Nr. 27) von zwei Sektoren Lücke (Byte Nr. 28) von drei Sektoren zwischen zwei Sektorgruppen 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 Attribut-bereich 0 1 Sektoren mit "nichtbinären" Daten Sektoren mit "binären Daten" Byteposition eines Bytes im 4. logischen Nutzdatensektor Betriebssystem e Prof. Dr. U. Wie

Auflistung von Dateisystemen affs - Amiga fast file system ext / ext2 - extended Filesystem (Linux) hpfs - high performance file system (OS/2) o iso9660 - CDROMs minix - Lehrdateisystem (A. Tanenbaum) ncpfs - Einbinden von Novell-Volumes onfs - Network file system smbfs - Server Message Block (Windows Protokoll für das Anbinden von Netzlaufwerken) - SCO-Unix, Xenix und Coherent (Unix-Derivate) sysv ufs - BSD, SunOS und NeXTstep umsdos - Unix on MSDOS vfat - virtual fat (FAT32 und lange Dateinamen von FAT16) Betriebssystem e 280 Prof. Dr. U. Wienkop

Plattenzugriffsverwaltung

Betriebssystem e

281 Prof. Dr. U. Wienkop

Ein-/Ausgabe-Behandlung Wiederholung: Zweck der Geräteverwaltung durch das BS

- O BS bietet über Systemaufrufe beim HW-Zugriff mehr Komfort und vor allem Geräteunabhängigkeit
 - Programmierer braucht sich nicht mit der Programmierung einzelner Geräte zu befassen, Geräte-Adressen und die Länge von Transferabschnitten (z. B. Sektoren) kennen, Zeitabhängigkeiten beachten, die Daten in eine gerätespezifische Form bringen
- O reale physikalische Geräte werden zu virtuellen Geräten mit idealisierten **Eigenschaften**
 - virtueller Drucker
 - > Plattendateien (Abstraktion von Festplatten, Disketten, CD-ROMs...)
- O Betriebssystem beansprucht das Monopol auf alle Gerätezugriffe
 - > notwendig für die Koordination solcher Zugriffe zwischen verschiedenen Prozessen
 - Verhinderung des Verstellens von Geräteparametern durch einzelne Prozesse
 - Alle Ein-/Ausgabebefehle des Prozessors müssen privilegiert sein
 - Abbruch als Privilegverletzung
 - □ Betriebssystem behandelt den Befehl als E/A-Auftrag und führt die Ein-/Ausgabe nach eigenen Gesichtspunkten durch.
 - Betriebssystem kann besser als einzelne Prozesse Optimierungsmaßnahmen zur Effizienzsteigerung durchführen

Betriebssystem e

Block- und Zeichengeräte O Datenübertragung an Zeichengerät in streng sequentieller Reihenfolge Bsp.: Drucker, Terminals oder serielle Schnittstellen > Datenpufferung nur insoweit, wie die vorgegebene Reihenfolge nicht verletzt wird O Blockgeräte > typischerweise Massenspeichergeräte mit wahlfreier Zugriffspositionierung Buchhaltung (Dateisystem) zum Merken der räumlichen Verteilung. Zur Vereinfachung der Buchhaltung Datenablage mit einheitlicher Länge Optimierungsmöglichkeiten: □ Weiterleiten in anderer Reihenfolge als in der Transfer-Warteliste ☐ Wichtig: Blockpufferung mit entsprechendem Vorspann Vorgegebene Reihenfolge Buchhaltung: AAAAAA BBBBBB CCCCC Block auf Position y Block auf Position x Ablage im 3. Block auf Position z AAAAAA CCCCCC BBBBBB Blockpuffer Schreibreihen-BBBBBB AAAAAA CCCCCC folge Betriebssystem e Prof. Dr. U. Wie

Optimierungsstrategien Kriterien

O Kriterien zur Beurteilung verschiedener Strategien:

- Mittlerer Durchsatz (abgearbeitete Aufträge pro Zeiteinheit)
- > mittlere Zugriffszeit T (wird durch lange Kopfbewegungen vergrößert)
- (Theoretisch) maximale Zugriffszeit T_{max} einer Auftragsfolge
- Varianz v der Wartezeit (vom Eintreffen bis zur Abarbeitung eines Auftrags)
- Fairness (werden bestimmte Aufträge bzw. Datenbereiche im Gerät benachteiligt oder gar ausgehungert?)

O Aspekte:

- hoher Durchsatz --> evtl. Benachteiligung einzelner Aufträge (hohes T_{max}).
- > Niedriges v --> Aufträge werden mit einer Zugriffszeit in der Nähe von T abgearbeitet
- Hohes v --> viele Aufträge weichen stark vom Mittelwert T ab

O Wahl der günstigsten Strategie abhängig von:

- > Typ des bedienten Geräts
- Zugriffsverhalten der Anwendungsprozesse (sequentiell/wahlfrei positionierend)
- Verteilung der Zugriffspositionen (Fragmentierungsgrad)
- Systemlast

Betriebssystem e

Disk-Scheduling Strategien

Beispielauftragsfolge: 23 111 44 66 105 11 66

- Ankunftsreihenfolge (FCFS)
 - ightharpoonup wenig effizient, sehr gute Fairness, T besonders schlecht, insbes. bei hoher Systemlast Ausführung: 23 111 44 66 105 11 66 Start: 100: $\Sigma = 442$
- O Einfache Fahrstuhlstrategie (Pickup)
 - > ungefähr genauso fair wie FCFS, deutliche Durchsatzverbesserung Ausführung: 66 66 44 23 105 111 11 Start: 100: $\Sigma = 265$
- O Volle Fahrstuhlstrategie (Scan, Look)
 - noch fair, Tendenz zu hohen Wartezeiten in den äußeren Spuren (hohes v)

 Ausführung: 66 66 44 23 11 105 111 Start: 100ψ : $\Sigma = 189$ Ausführung: 105 111 66 66 44 23 11 Start: 100ψ : $\Sigma = 111$
- O Zirkuläre Fahrstuhlstrategie (Circular Scan / Look)
 - Niedrigeres v (als bei der vollen Fahrstuhlstrategie), etwas h\u00f6heres T. Praxis: niedrige Systemlast: volle Fahrstuhlstrategie verwenden, sonst: zirkuläre Ausführung: 66 44 23 111 Start: 100ψ : $\Sigma = 195$ Ausführung: 105 23 Start: $100 \uparrow$: $\Sigma = 166$ 111 11 44 66 66

Betriebssystem e

285 Prof. Dr. U. Wienkop

Disk-Scheduling Strategien (2)

- O Beispielauftragsfolge: 23 111 44 66 105 11 66
- Kürzeste Positionierungszeit (SSTF)
 - Minimierung der Positionierungsvorgänge, hoher Durchsatz, Gefahr des Aushungerns (besonders hohes v, für Dialogbetrieb unzumutbar)

 Ausführung: 105 111 66 66 44 23 11 Start: 100: Σ = 111
- Prioritätsgesteuerte Strategie
 - > Blöcken erhalten Prioritäten (z.B. Prozeßpriorität)
 - > Gesamtsystemdurchsatz???, Gefahr des Aushungerns für Blöcke schlechter Priorität
 - Meist: Kombination mit anderer Strategie zur Durchsatzoptimierung

Betriebssystem e

Vorauslesen und aufgeschobenes Schreiben

O Vorauslesen

- > nicht nur den gewünschten Block, sondern alle Blöcke auf derselben Spur lesen
- Lohnend, wenn Spurwechselzeit deutl. höher als die Zeit für eine Umdrehung ist
- Variante: Mit dem Lesen der Spur sofort nach Erreichen beginnen, OHNE erst auf den gewünschten Sektor zu warten.
 - Die zu lesende Datei sollte nicht zu klein sein (Umfang von mehr als einer Spur)
 - weitgehend zusammenhängende Plattenablage (Fragmentierungsgrad ?) überwiegend sequentielles Lesen (Programmladen vs. Datenbankzugriff)

Aufgeschobenes (verzögertes) Schreiben (deferred / delayed write, write behind)

- > Schreibauftrag nicht sofort, sondern zu einem geeigneteren Zeitpunkt ausführen
 - System (Prozessor bzw. Gerät) wird nicht anderweitig beansprucht
 - Blockpuffer voll
 - anach Ablauf einer Zeitschranke, so mit Sicherheit irgendwann tatsächlich geschrieben wird
- Nebeneffekt: schreibender Prozeß kann sofort nach Auftragsabgabe fortgesetzt werden
- Problem, wenn der Schreibauftrag aus irgendwelchen Gründen später nicht ausgeführt werden kann (Benachrichtigung des Auftraggebers? Evtl. bereits beendet!)

Betriebssystem e

Prof Dr II Wie

Datenmodifikation

O Für den Anwendungsprozeß transparentes Modifizieren von Daten:

- - Zur Reduzierung der Datenmenge werden die Daten vor dem Transfer komprimiert.
- Verschlüsselung:
 - Zur Wahrung des Datenschutzes werden die Daten vor dem Transfer verschlüsselt (chiffriert).
- Fehlererkennung und -korrektur:
 - Die Daten werden mit Redundanz (Paritätsbits, Prüfsummen usw.) versehen.
- speziell bei Tastaturen:
 - einfache Editiermöglichkeiten, bei denen keine Rücktransformation erforderlich ist (Einige Beispiele: das eingegebene Zeichen ^H wird zusammen mit dem vorangehenden nicht transferiert, Tabulatorzeichen werden expandiert, ein Zeilenende dient als Abschlußzeichen (Begrenzer) für die aktuelle Eingabe, wird selbst aber nicht übertragen).

Betriebssystem e

Systemaufrufschnittstelle zur Dateiverwaltung O Datei einrichten create O Datei (Gerät) eröffnen open O Datei (Gerät) schließen close O Datei löschen unlink O Zugriffsposition bestimmen / positionieren seek Lesezugriff read Schreibzugriff write Parametrieraufruf control + Erweiterungen Betriebssystem e 289 Prof. Dr. U. Wienkop

