

## J UNIX-Dateisystem

### J.1 Funktionalität

- Abstraktionen für Benutzersicht:
  1. Pfade, Dateinamen
  2. Dateibaum verdeckt mehrere Platten (bzw. Partitionen)
  3. Dateien = unstrukturierte Byteströme
  4. *special files* fast wie reguläre Dateien zugreifbar
- diese Abstraktionen müssen vom Dateisystem auf konkrete Objekte abgebildet werden
- den obigen Abstraktionen entsprechen etwa folgende Objekte:
  1. / 2. Plattenpositionen, Adressen von Gerätecontrollern, etc.
  3. Platten verwalten Sektoren, Spuren und Zylinder
  4. unterschiedliche Gerätetreiber werden für E/A-Operationen auf verschiedenen Peripheriegeräten benötigt

AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

J-Filesystem.doc 1999-02-03 08.50

J.1

Reproduktion jeder Art oder Vervielfältigung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## 1 Umwandlung: Pfad : Inode

J.3 Inodes

- Root = Inode 2
- Inode-Nummer ist nur in dem Directory zu finden, in dem die Datei eingetragen ist
  - ↳ die Pfadkomponenten müssen einzeln bearbeitet werden
  - ↳ um den Inode der jeweils nächsten Pfadkomponente zu bestimmen, muß der Inhalt des voranstehenden Directories eingelesen werden
- Systemkern-Funktion: **lookup()**
- Directory Lookup Cache zur Beschleunigung
  - ◆ LRU Cache mit Einträgen  
Directory-Inode-Nummer / Datei(name) in Directory / Inode der Datei

AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

J-Filesystem.doc 1999-02-03 08.50

J.3

Reproduktion jeder Art oder Vervielfältigung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

### J.2 Directories (Kataloge)

- wichtigster Bestandteil zum Aufbau eines hierarchischen Dateisystems
- Funktion: Strukturierung des Namensraums
- *Directories* enthalten Verweise auf Inodes (von Dateien oder weiteren Directories)

### J.3 Inodes

- ein Inode enthält alle Informationen, die zur Verwaltung einer Datei und zum Auffinden ihres Inhalts benötigt werden
- beim Öffnen einer Datei werden die Inode-Daten von der Platte (**disk inode**) in die Inode-Tabelle des Systemkerns (**in-core inode**) gelesen

AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

J-Filesystem.doc 1999-02-03 08.50

J.2

Reproduktion jeder Art oder Vervielfältigung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## 1 Umwandlung: Pfad : Inode (2)

J.3 Inodes

- Berücksichtigt werden hierbei:
  - ◆ Symbolic Links
    - werden beim Absteigen im Baum verfolgt
  - ◆ Mount-Punkte
    - beim Absteigen im Baum am Mount-Punkt Wechsel zur *root* des eingehängten Dateisystems
    - beim Aufsteigen im Baum (über "..") an der *Root* eines Dateisystems Wechsel auf den Mount-Punkt
  - ◆ Current Root
    - Überschreiten beim Aufsteigen im Baum (über "..") wird verhindert

AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

J-Filesystem.doc 1999-02-03 08.50

J.4

Reproduktion jeder Art oder Vervielfältigung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## 2 Datei öffnen

J.3 Inodes

### Aufruf:

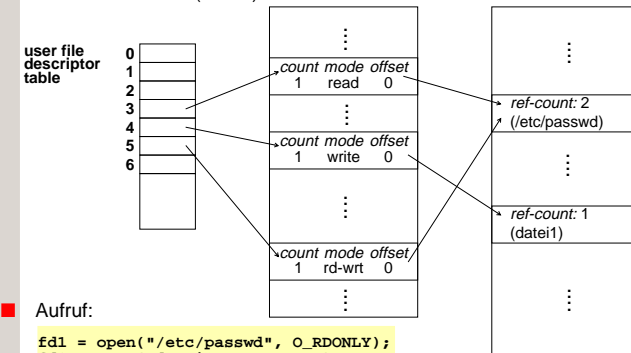
```
fd = open(path, flags, modes);
```

<b>fd</b>	Filedescriptor (kleiner Integer-Wert)
<b>path</b>	Pfadname der Datei
<b>flags</b>	Angabe über Zugriffsform (lesen, schreiben, lesen+schreiben), ob Datei neu erzeugt / vorher gelöscht wird
<b>modes</b>	Zugriffsrechte, falls die Datei neu erzeugt wird

## 2 Datei öffnen (3)

J.3 Inodes

### Kontrollstrukturen (früher)



### Aufruf:

```
fd1 = open("/etc/passwd", O_RDONLY);
fd2 = open("datei1", O_WRONLY);
fd3 = open("/etc/passwd", O_RDWR);
-> fd1 == 3, fd2 == 4, fd3 == 5
```

## 2 Datei öffnen (2)

J.3 Inodes

### open-Algorithmus

in-core inode zu Pfadnamen erzeugen ( <i>namei, lookup</i> ) (Inode suchen und einlesen)	
Datei existiert nicht oder keine Zugriffsrechte	
ja	nein
in-core Inode freigeben	<b>file table</b> entry für Inode belegen
return(Fehler)	<b>count</b> und <b>offset</b> initialisieren
	<b>user file descriptor</b> belegen, Zeiger auf <b>file table entry</b> setzen
	falls <b>truncate</b> , alte Blöcke freig.
	return( <b>user file descriptor</b> )

## 2 Datei öffnen (4)

J.3 Inodes

### user file descriptor table

- ◆ Feld von Zeigern auf Einträge in der **file table**
- ◆ Teil der **user-area** des Prozesses  

```
struct file *u_ofile[NOFILE]
```

  
siehe `<sys/user.h>`
- ◆ Filedescriptor (aus **open(2)**) ist ein Index in diese Tabelle  
↳ Entkopplung von Systemdatenstrukturen und Benutzerdaten  
Filedeskriptor = Capability (Sicherheit!)

### file table

- ◆ Feld von **file**-Strukturen (siehe `<sys/file.h>`), enthält
  - Zeiger auf Eintrag in der **inode table**
  - Eintrag für Zugriffsmodus
  - Schreib/Lesezeiger (offset)
  - Referenzzähler

## 2 Datei öffnen (5)

J.3 Inodes

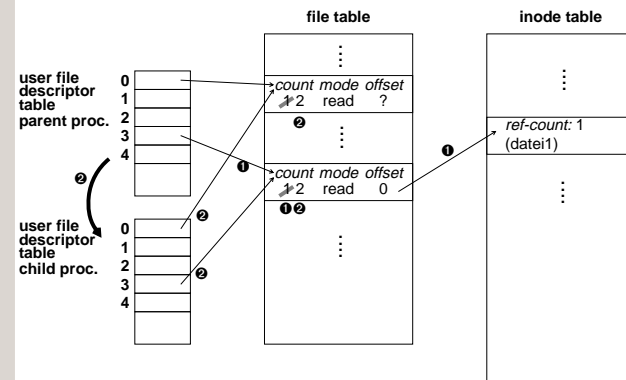
### inode table

- ◆ Feld von *in-core inodes* (siehe `<ufs/inode.h>`), enthält
  - Referenzzähler
  - Info über physikalisches Gerät (Platte, Partition)
  - *disk inode* Struktur
  - Informationen für *lock*-Funktionen (*flock(2)*, *lockf(2)*)

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## 4 Kontrollstrukturen und Systemaufruf *fork()*

J.3 Inodes

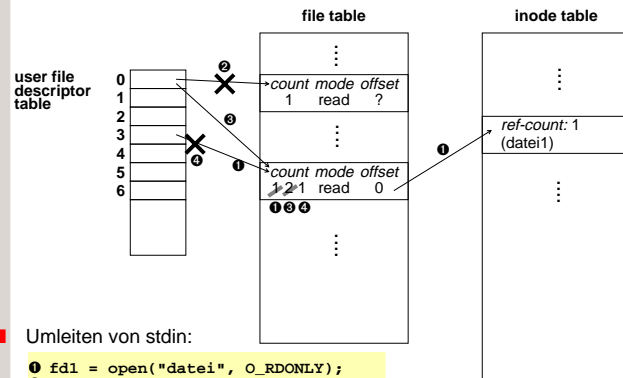


```
1 fd1 = open("datei", O_RDONLY);
2 fork( );
```

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## 3 Kontrollstrukturen und Systemaufruf *dup()*

J.3 Inodes



### Umleiten von stdin:

```
1 fd1 = open("datei", O_RDONLY);
2 close(0);
3 dup(fd1);
4 close(fd1);
```

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## J.4 Vnodes

### 1 Motivation

- Inode = konkrete Implementierung eines Dateiverwaltungsblocks eines speziellen Dateisystems
- Problem: unterschiedliche Dateisysteme in einem UNIX-System
  - ◆ S5 – SystemV-File-System
  - ◆ UFS - UNIX-File-System
  - ◆ NFS - Network File-System
  - ◆ specFS - special File-System (Geräte-dateien, FIFO-Dateien)
- ➔ unterschiedliche Implementierung (Inode, *lookup*-Funktion, ...)

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

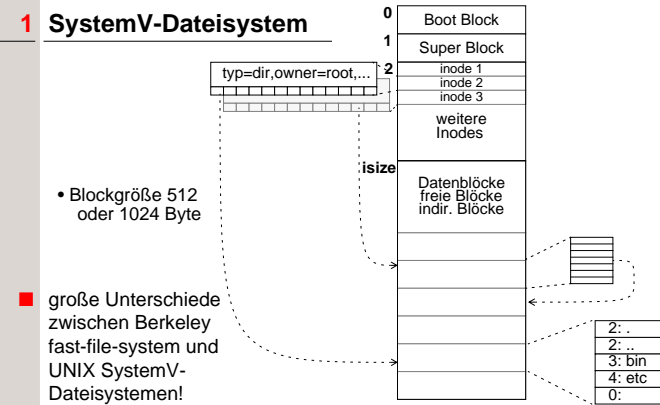
# 1 Motivation (2)

- Lösung: **Vnode** = abstrakte Darstellung einer Datei
  - ↳ definiert Operationen auf Dateien
    - *open, close, read/write, lookup, ...*
  - ↳ + allgemeine Zustandsinformationen
    - Dateityp (*regular file/directory*)
    - Referenzzähler
    - Verwaltungsinformationen für *locking* etc.
    - Zeiger auf konkrete Operationen und konkreten Zustand
- konkrete Darstellung einer Datei
  - ↳ Implementierung der Operationen
    - *ufs\_open, ufs\_close, ufs\_read, ufs\_lookup, ...*
  - ↳ + konkreter Zustand
    - *in-core Inode*

# J.5 Dateisysteme

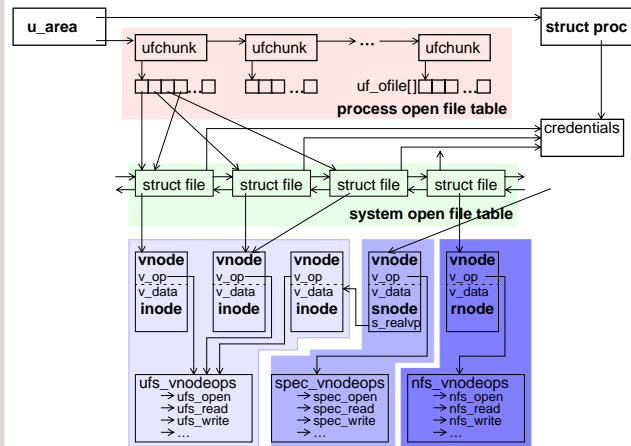
## Aufbau eines Dateisystems auf einer Disk

### 1 SystemV-Dateisystem



■ große Unterschiede zwischen Berkeley fast-file-system und UNIX SystemV-Dateisystemen!

# 2 Kontrollstrukturen (SystemV.4)



# 1 SystemV-Dateisystem (2)

- Super-Block
  - ◆ Der *Super-Block* enthält die Verwaltungsinformationen für ein Dateisystem
    - Größe
    - Anzahl der freien Blöcke
    - Liste freier Blöcke (*free block list*)
    - Zeiger auf nächsten freien Block in der *free block list*
    - Größe der Inode-Liste
    - Anzahl freier Inodes
    - Liste freier Inodes (*free inode list*)
    - Zeiger auf ersten freien Inode
    - Lock-Felder für *free block* und *free inode list*
    - *modified flag* für Super-Block

## 2 Berkeley Fast File System (UNIX File System — UFS)

J.5 Dateisysteme

- für UNIX 4.2bsd entwickelt - effizienter als SystemV-Dateisysteme
- Blockgröße minimal 4k, meist 8k
  - ⇒ über zweifach indirekte Blöcke können bis zu 64 GB adressiert werden
- Directories anders aufgebaut - Dateinamen bis zu 255 Zeichen lang
- Dateisystem in **cylinder groups** aufgeteilt (*cylinder group* = Gruppe von aufeinanderfolgenden Zylindern auf Platte, häufig 16)
- pro *cylinder group* werden folgende Daten gehalten:
  - ◆ Backup-Kopie des Superblocks
  - ◆ Bitmap für freie Blöcke (statt *free list*)
  - ◆ Inodes
  - ◆ Datenblöcke

AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

J-Filesystem.doc 1999-02-03 08.50

J.17

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## J.6 Virtuelles Dateisystem

### 1 Überblick

- Vnodes bilden Schnittstelle zu dateisystemunabhängigen Operationen auf Dateien
- das *Virtual File System (VFS)* stellt eine vergleichbare Schnittstelle für Operationen auf Dateisystemen zur Verfügung — z. B.
  - ◆ *vfs\_mount, vfs\_unmount*
  - ◆ *vfs\_statvfs*                      Statistikdaten ermitteln
  - ◆ *vfs\_root*                            root-vnode ermitteln
  - ◆ *vfs\_sync*                            gepufferte Daten schreiben
- Implementierung durch VFS-switch-Tabelle
  - ◆ Dateisystemname, Initialisierungsfunktion, *vfsops* — Feld mit Zeigern auf die FS-Funktionen

AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

J-Filesystem.doc 1999-02-03 08.50

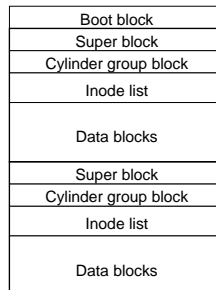
J.19

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## 2 Berkeley Fast File System (2)

J.5 Dateisysteme

- Um Verschnitt gering zu halten, können Blöcke in 2, 4 oder 8 **Fragmente** aufgeteilt werden, die einzeln adressierbar sind
- Eine Datei wird - wenn möglich - in einer *cylinder group* gehalten



AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

J-Filesystem.doc 1999-02-03 08.50

J.18

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## 2 Mount-Mechanismus

J.6 Virtuelles Dateisystem

- über die mount-Funktion wird der root-Vnode eines Dateisystems auf einen Directory-Vnode des bestehenden Dateibaums montiert
  - ◆ Systemkern verwaltet die montierten Dateisysteme in verketteten *vfs*-Strukturen
  - ◆ über die *vfs*-Struktur kann ermittelt werden:
    - auf welchen Vnode das Dateisystem montiert ist
    - der root-Vnode des Dateisystems
  - ◆ Vnode des Mount-Punkts (Directory) enthält Eintrag mit Zeiger auf *vfs*-Datenstruktur des montierten Dateisystems

AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

J-Filesystem.doc 1999-02-03 08.50

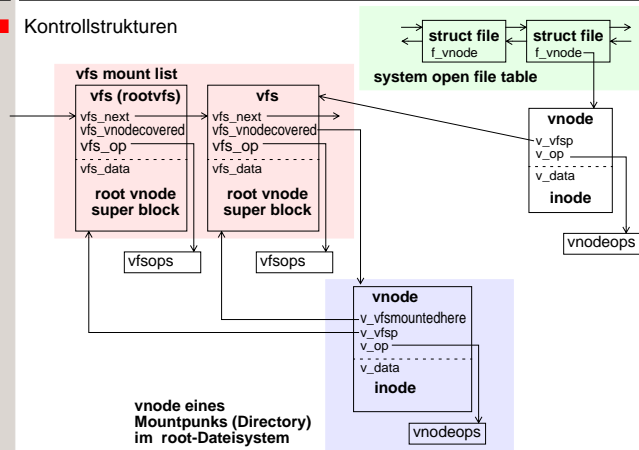
J.20

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## 2 Mount-Mechanismus (2)

J.6 Virtuelles Dateisystem

### Kontrollstrukturen



AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

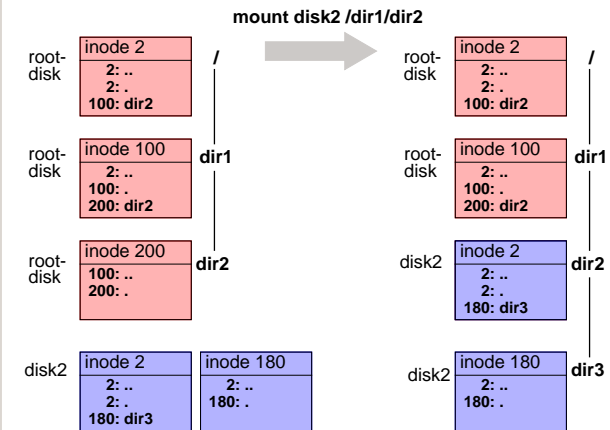
J-Filesystem.doc 1999-02-03 08.50

J.21

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## 2 Mount-Mechanismus (4)

J.6 Virtuelles Dateisystem



AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

J-Filesystem.doc 1999-02-03 08.50

J.23

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## 2 Mount-Mechanismus (3)

J.6 Virtuelles Dateisystem

### lookup-Funktion des UNIX-Kerns berücksichtigt bei Pfad-Vnode-Umwandlung Mount-Punkte

- ◆ ist in einem vnode `vfs_mountedhere` gesetzt, wird der Vnode beim Abstieg im Pfad durch den `root-vnode` des `VFS` ersetzt
- ◆ ein `root-vnode` eines `VFS` wird bei Verfolgung von `".."` durch `vfs_vnodecovered` des `VFS` ersetzt
- ◆ an Directory-Einträgen ändert sich nichts
  - ➔ an Mount-Punkten Diskrepanz zwischen Directory-Einträgen und Ergebnis von `stat(2)`

AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

J-Filesystem.doc 1999-02-03 08.50

J.22

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.

## J.7 Network File System (NFS)

### 1 Überblick

- 1984 von Sun Microsystems Inc. entwickelt
- NFS ermöglicht, Teilbäume von Dateisystemen anderer Rechner in das lokale Dateisystem einzuhängen (**mounten**)
- ein Teilbaum eines *Remote*-Dateisystems verhält sich dadurch genau so, wie eine Partition einer lokalen Platte, die mit einem **mount**-Befehl in den Dateibaum eingehängt wurde
- die Dateien des *Remote*-Dateisystems können mit den gleichen Kommandos und Systemaufrufen wie lokale Dateien bearbeitet werden

AKBP I

Ausgewählte Kapitel der praktischen Betriebsprogrammierung I  
© Jürgen Kleinöder, Universität Erlangen-Nürnberg, IMMD IV, 1999

J-Filesystem.doc 1999-02-03 08.50

J.24

Reproduktion jeder Art oder Verwendung dieser Unterlagen, außer zu Lehrzwecken an der Universität Erlangen-Nürnberg, bedarf der Zustimmung des Autors.