

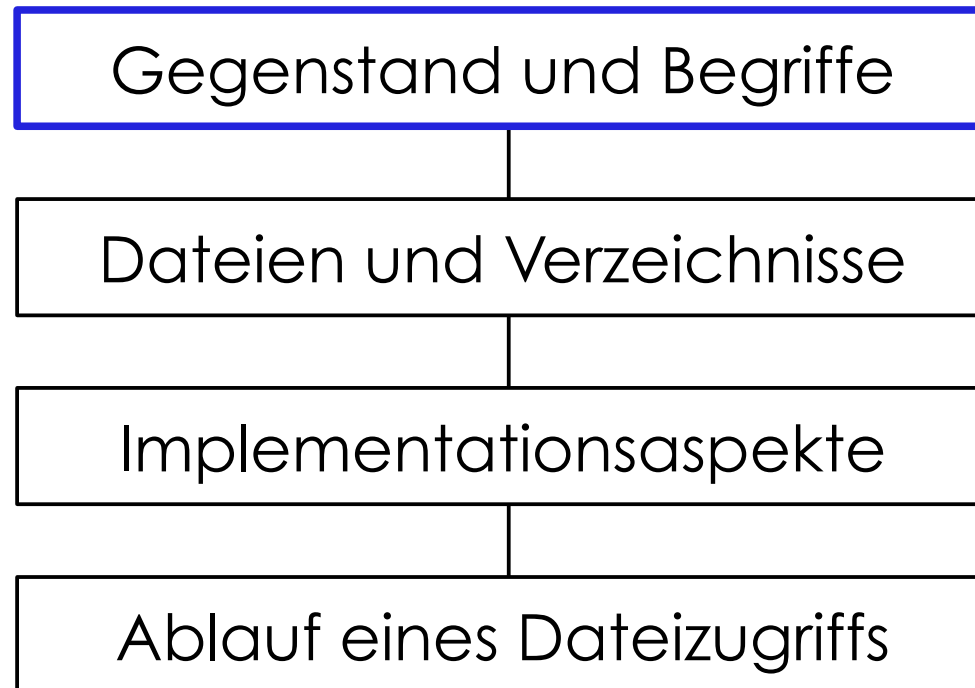
Ein einfaches Dateisystem

Betriebssysteme

Hermann Härtig
TU Dresden



Wegweiser



Gegenstand

Aufgaben eines Dateisystems

- **Persistenz** von Daten, d. h. Daten sollen Prozess- und Systembeendigungen überstehen (auch Abstürze)
- Umgang mit großen Datenmengen (z.B. $> 4\text{GB} = 2^{32}$ auf IA32)
- Schutz (separate Kapitel dieser Vorlesung)
- Basismechanismen für Datenbanken
- Organisation des Zugriffs paralleler Prozesse auf Dateien

Zentrale Herausforderungen

- Hardwareeigenschaften des persistenten Speichermediums
- Robustheit (→ separate Kapitel in dieser Vorlesung)

Begriffe (nach NEHMER)

Dateisystem

- BS-Komponente, die Anwendungsprogrammen einen effizienten Zugriff auf persistent gespeicherte Daten ermöglicht

Datei

- „Behälter“ für die dauerhafte Speicherung von Informationen (gleicher oder ähnlicher Struktur) unter einem inhaltlichen Gesichtspunkt
- codiert in Bytes, durch einen „Bezeichner“ repräsentiert

Verzeichnis

- besondere, vom Dateisystem verwaltete Datei zur Strukturierung der auf Externspeicher abgelegten Dateien

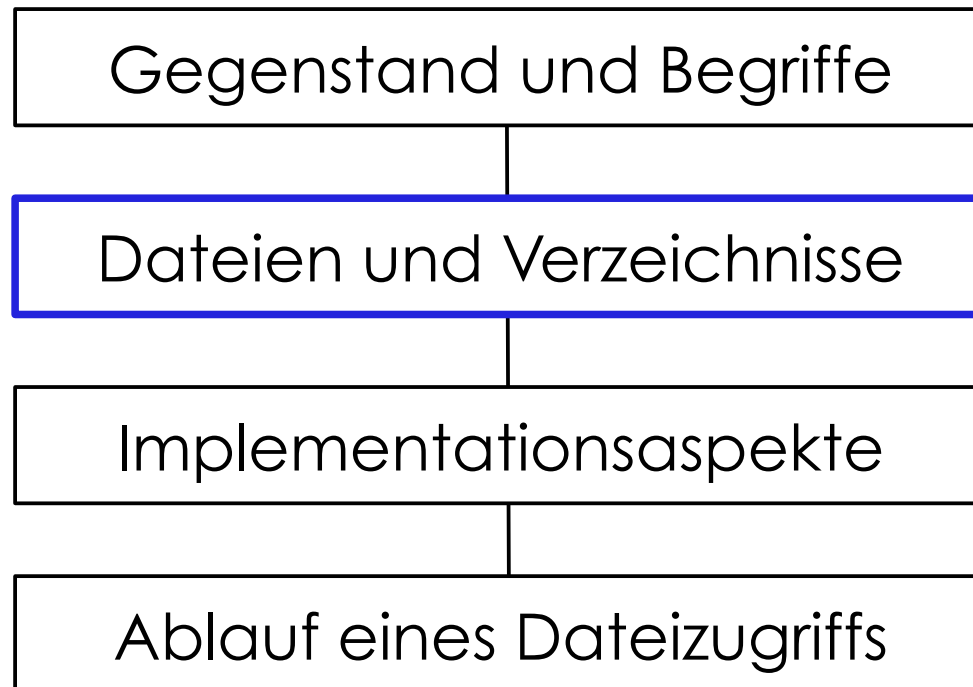
Operation

- Erzeugen (Create)
- Entfernen (Delete, Unlink)
- Öffnen (Open)
- Schließen (Close)
- Lesen (Read)
- Schreiben (Write)
- Append
- Seek
- Lesen/Setzen von Attributen
- Umbenennen
- „Einblenden“

typische Parameter

Name, Attribute, Zugriffsart → Handle
Handle
Name, Zugriffsart → Handle
Handle
Handle, Anzahl Zeichen, Puffer
Handle, Anzahl Zeichen, Puffer
Handle, Anzahl Zeichen, Puffer
Handle, Position

Wegweiser



Benennung (von Dateien)

Aufgaben:

- Benennen: symbolisch, (auch) Mensch als Benutzer
- Identifizieren: Programme, Datenbanken als Benutzer
- (schnelles) Lokalisieren
- Zugreifen

Beispiele

- /home/hh7/tmp/Dateien02S.ppt
- 4, 19, 329 (I-Knoten-Nummer); 17 (Filedescriptor)
- canaletto.inf.tu-dresden.de; 141.76.20.10

Benutzer legen Dateinamen fest ...

- innerhalb bestimmter Grenzen, je nach System
- nach Konventionen:
 - ♦ Folgen lesbarer Zeichen
 - ♦ Groß-/Kleinschreibung signifikant (oder nicht)
 - ♦ kontextabhängig (oder nicht)
 - ♦ Extensionen, z. B. komprimiertes Quell-Programm: `.c.gz`
 - ♦ Längenbeschränkungen
 - ♦ mit Rechnernamen, z. B. `erwin:/home/hh7/tmp/x`
 - ♦ Hierarchiebildung
 - ♦ mehrere Namen für ein Objekt
 - ♦ inhaltsbezogene Benennung
 - ♦ Möglichkeiten zur Abkürzung

Datei-Attribute

Zugriffschutz

- Eigentümer, Besitzer, Rechte des Besitzers und anderer

Zeiten

- Erzeugung, letzte Modifikation, letzter Zugriff, ...

Organisatorisches

- aktuelle/maximale Größe
- Verwaltungsinformationen
- Strukturinformationen

Datei-Typen

- in Unix-Welt:
- alle (cum grano salis) Objekte werden auch als Dateien behandelt

Beispiele

- normale Dateien (regular file)
- Verzeichnisse (directories)
- E/A-Geräte (special files)
- Datenströme (pipes)
- Symbolic Links

Zugriffsstrukturen

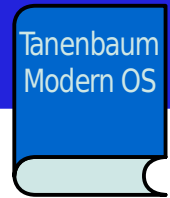
flach (Unix)

- ARRAY of BYTE
- Lesen/Schreiben in beliebigen Einheiten an beliebigen Stellen
- Aufprägung von Strukturen durch Anwendungen

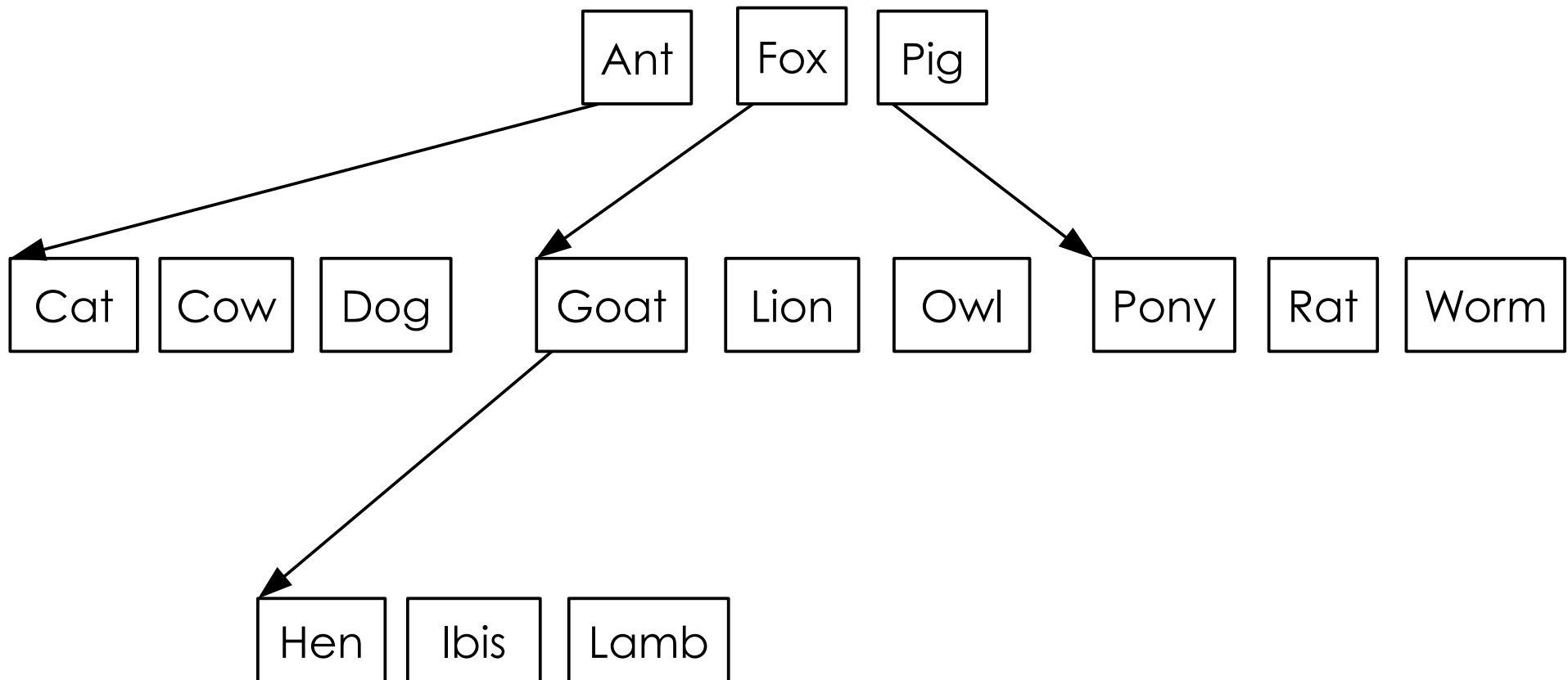
„Record“-Strukturen (satzstrukturierte Dateien)

- Länge konstant oder variabel
- ein Datensatz (record) pro Zugriff (Lesen/Schreiben)
 - ♦ sequentiell
 - ♦ direkt (B- bzw. B*-Bäume), Zugriff über Schlüssel
 - ♦ indexsequentiell

Indexsequentielle Dateien



- Datensätze beliebiger Länge mit Schlüssel
- Zugriffsreihenfolge: sequentiell oder über „Schlüssel“



Verzeichnisse und Pfadnamen

Problem:

- systemweite Eindeutigkeit von Dateinamen

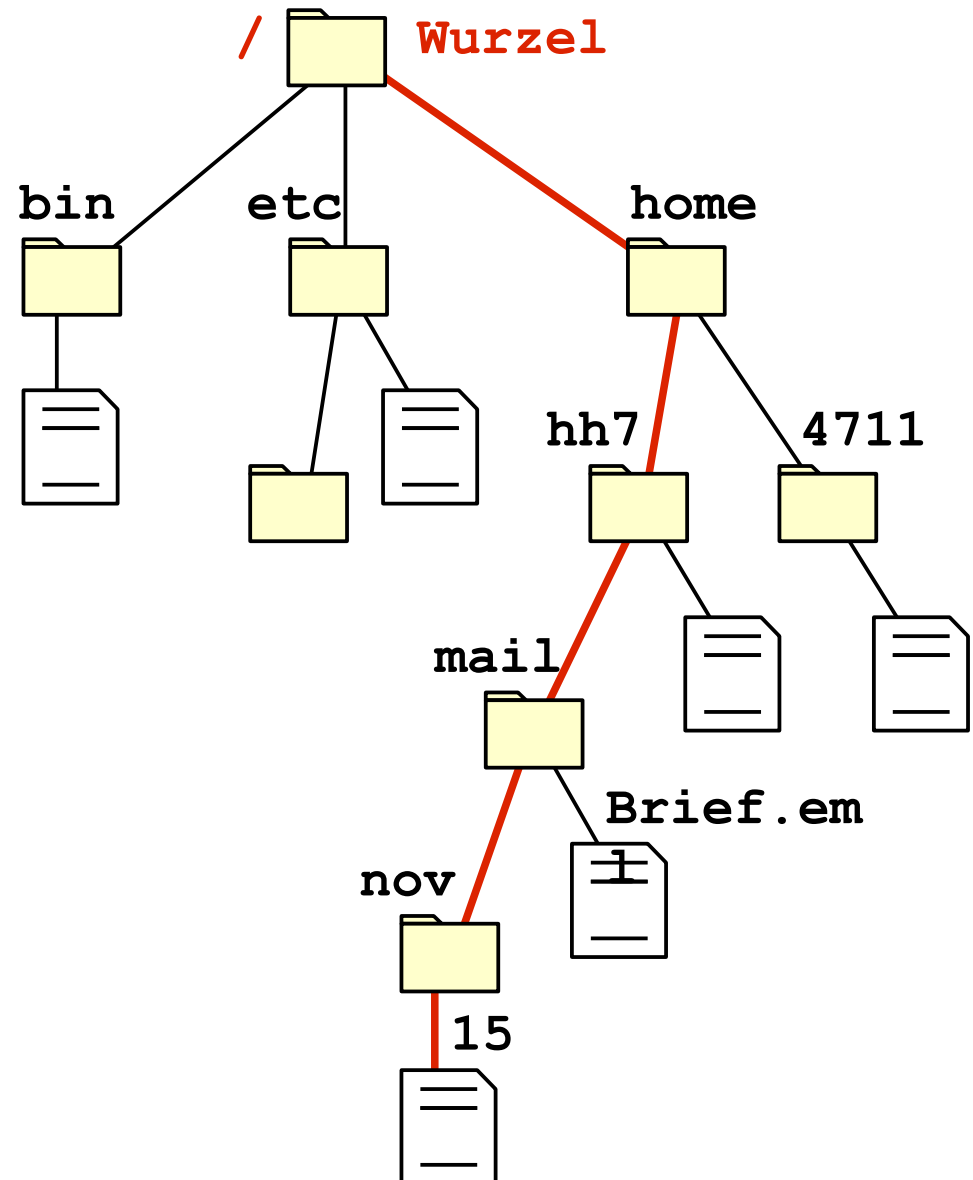
Lösung:

- Name eindeutig in einem „Kontext“
- auch Kontexte haben Namen

Pfadname:

- Folge von Teilnamen
- Namens-Auflösung:
von einem Ausgangs-Kontext ausgehend wird jeder Teilname in seinem Kontext interpretiert

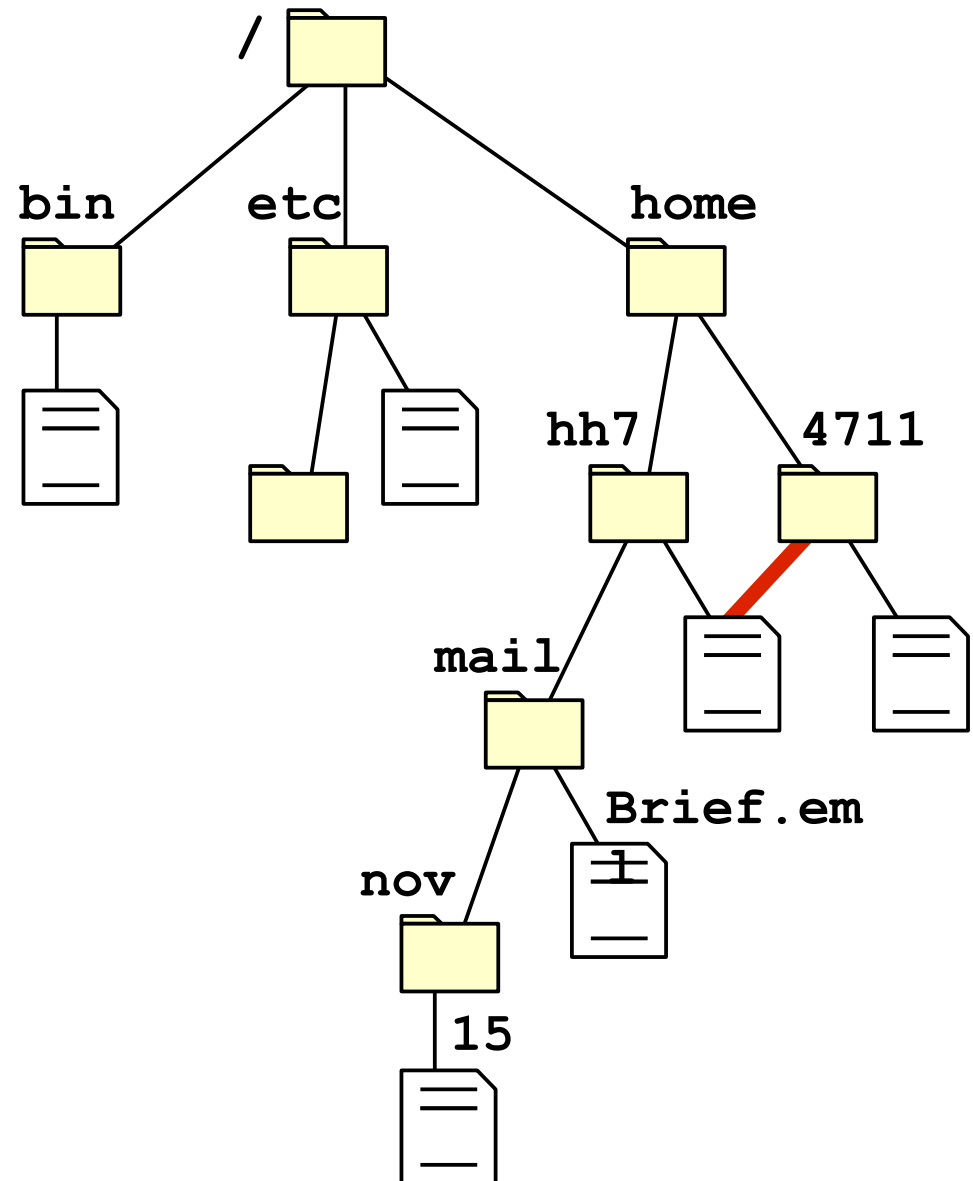
Beispiel: `/home/hh7/mail/nov/15`



Mehrere Namen für eine Datei: Harte Links

„Harte Links“

- Identifikatoren in mehreren Verzeichnissen
- Limitation: Gültigkeit Identifikatoren über Rechnergrenzen hinweg ?



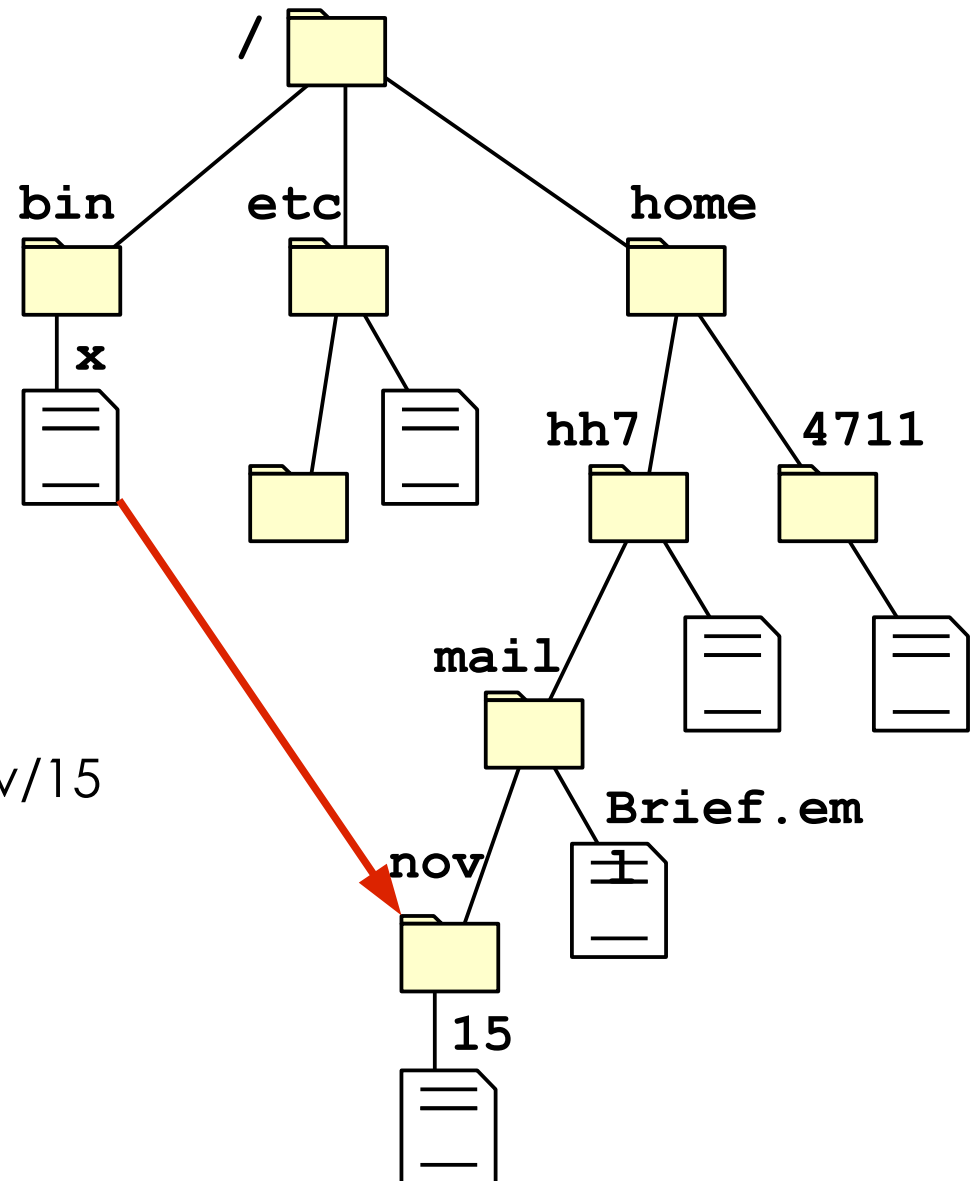
Mehrere Namen für ein Objekt: Weiche Links

„Symbolische Links“

- Namen als Objekte
- Zyklen !

Beispiel

- $x \rightarrow /home/hh7/mail/nov$
- $/bin/x/15 = /home/hh7/mail/nov/15$



Beispiel für Zyklen

/etc/bin/XYUtils → /bin/XYUtils

/bin/XYUtils → /usr/local/bin/XYUtils

/usr/local/bin/XYUtils → /etc/bin/XYUtils

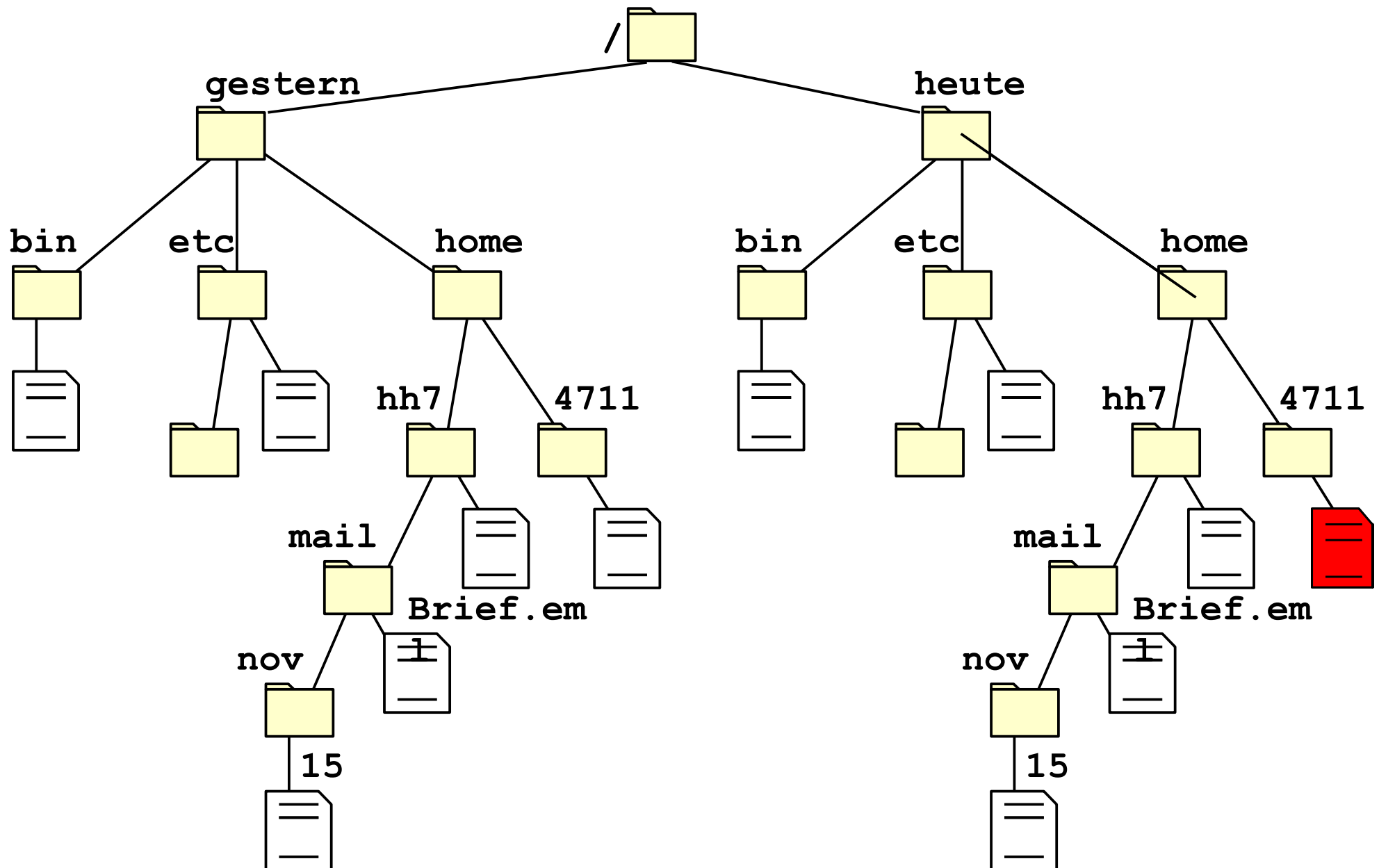
Auflösung von /bin/XYUtils/example

- /usr/local/bin/XYUtils/example
- /etc/bin/XYUtils/example
- /bin/XYUtils/example
- /usr/local/bin/XYUtils/example
- /etc/bin/XYUtils/example

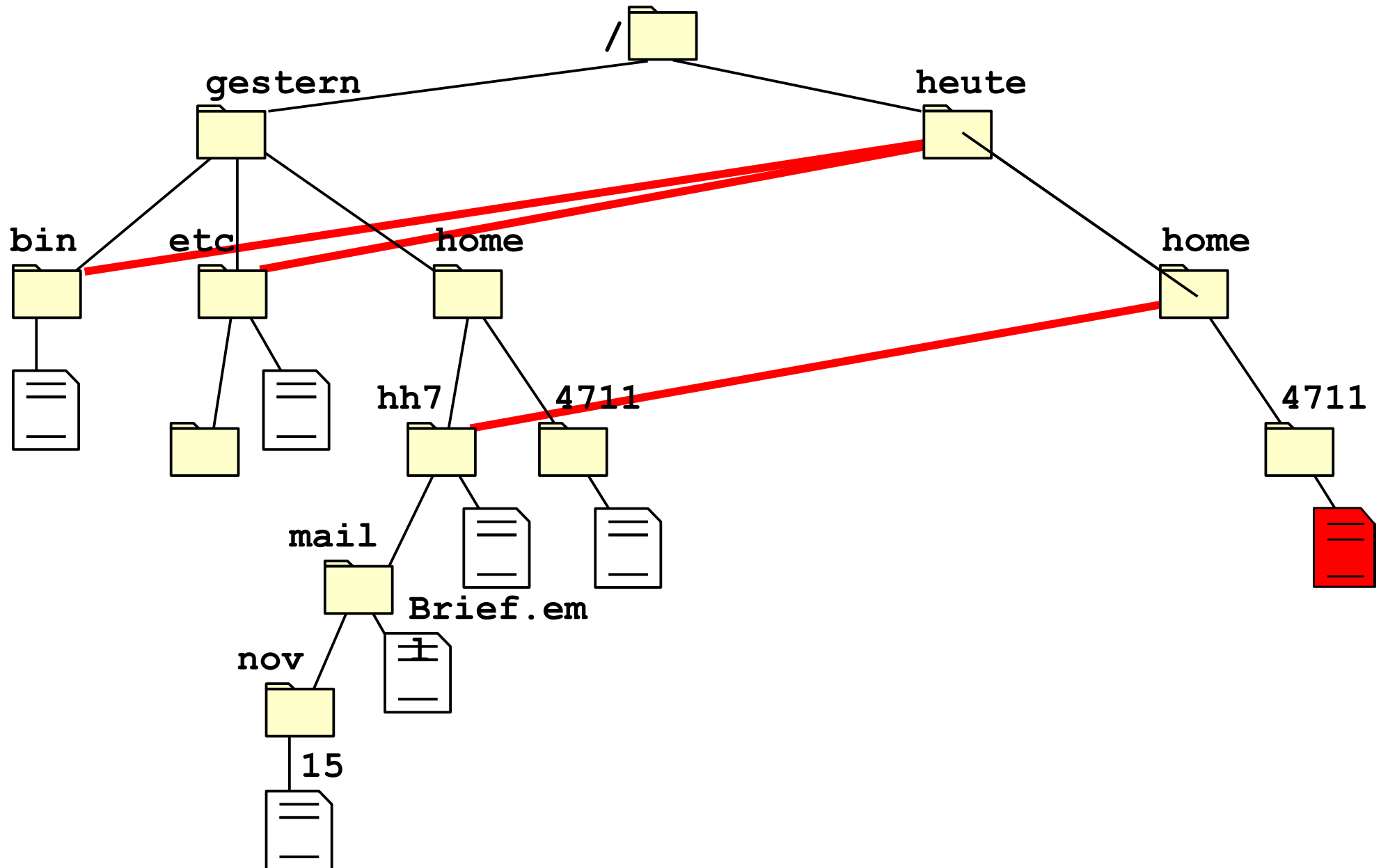
Beispiel: Backup

- Alle X Stunden vollständige Kopie des Dateisystems
- Nur das Nötige Kopieren
- Mehrfache harte Links auf alles, was sich nicht änderte

Beispiel: Backup



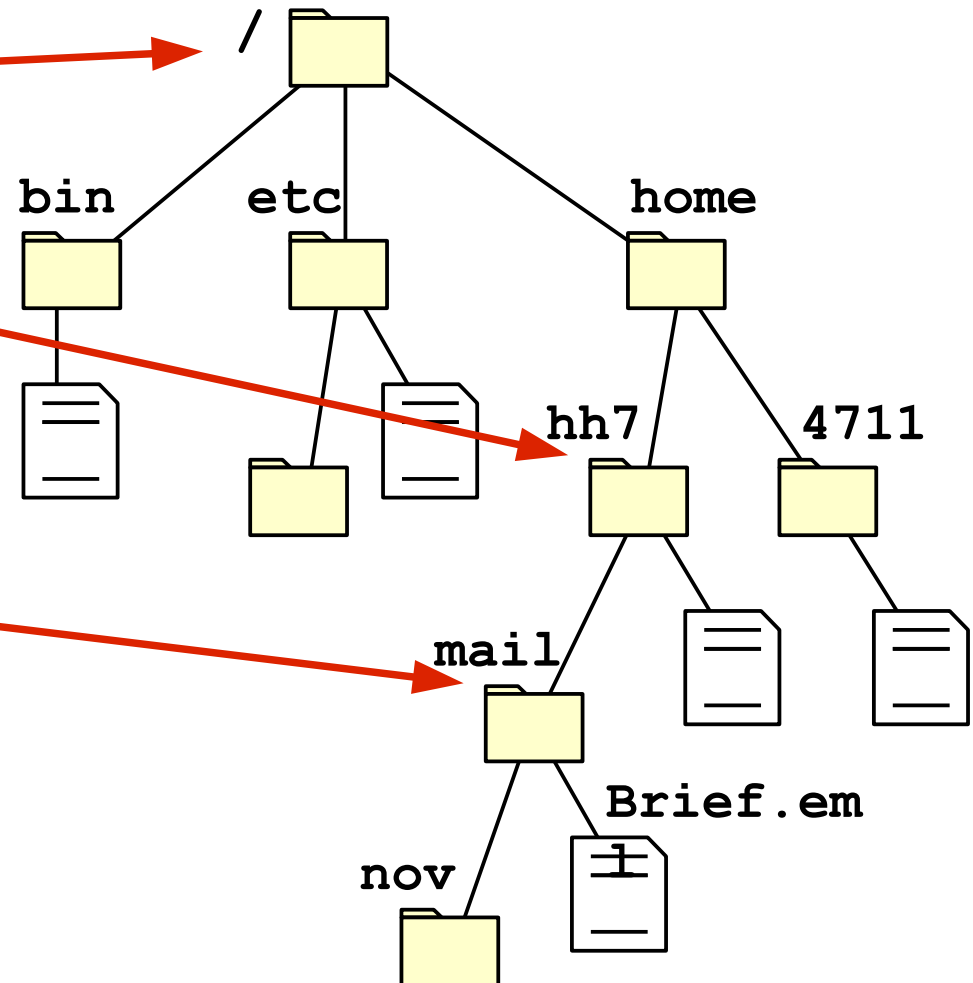
Beispiel: Backup



Kontexte für Namen

Ausgangs-Kontexte durch Prozesse festgelegt, z. B. in Unix :

- Root directory:
Kontext für globale Namen
- home directory:
current directory bei login
und durch parameterloses
„cd“
- current directory:
Ausgangskontext für relative
Pfadnamen



Prozesse und Pfadnamen

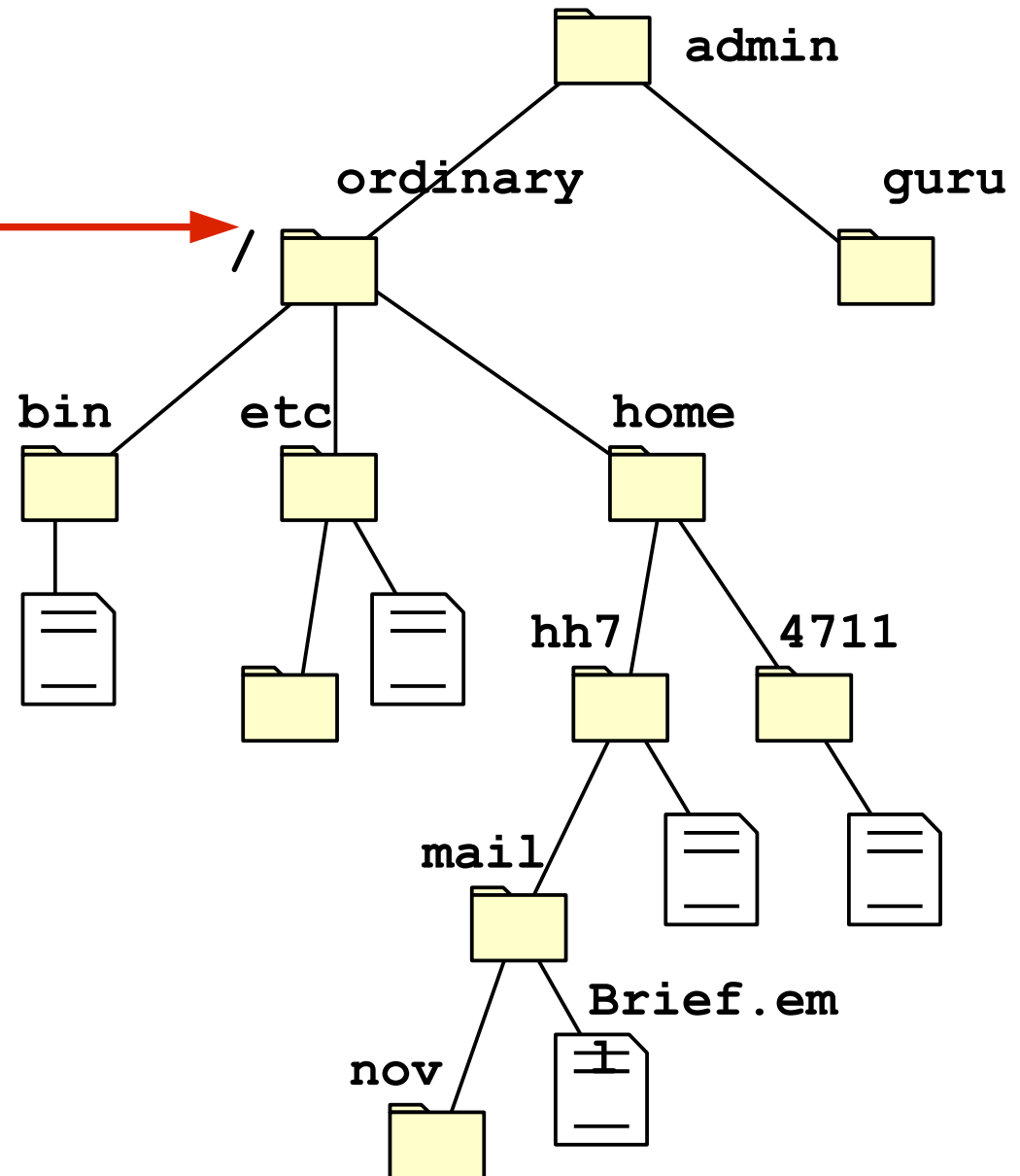
Ausgangs-Kontexte durch Prozesse festgelegt, z. B. in Unix :

- current root:
Ausgangskontext für absolute Pfadnamen

„chroot“ Kommando

mehrere „Globale“ Namensräume

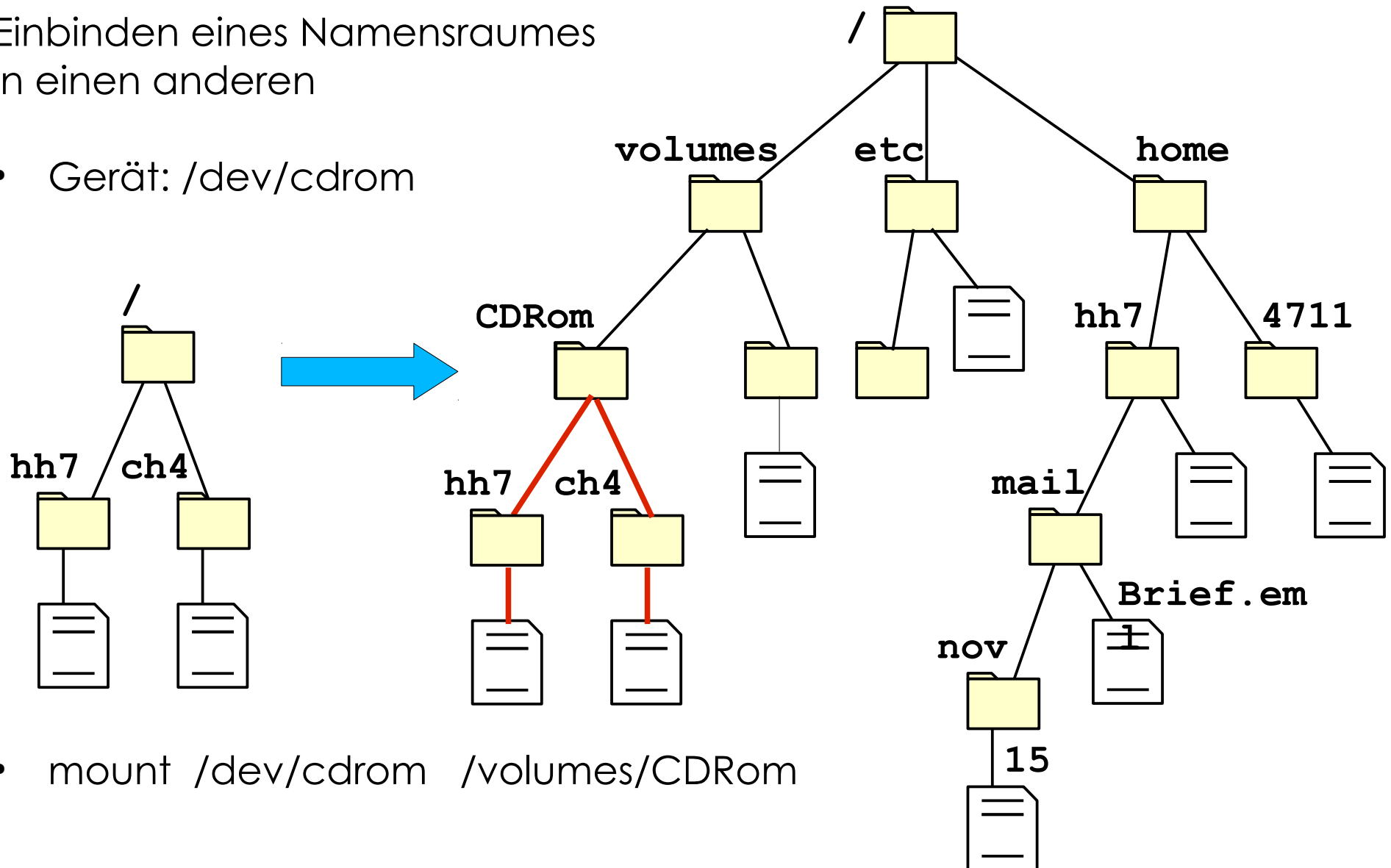
OS Virtualisierung



Verbinden von Namensräumen: Mounting

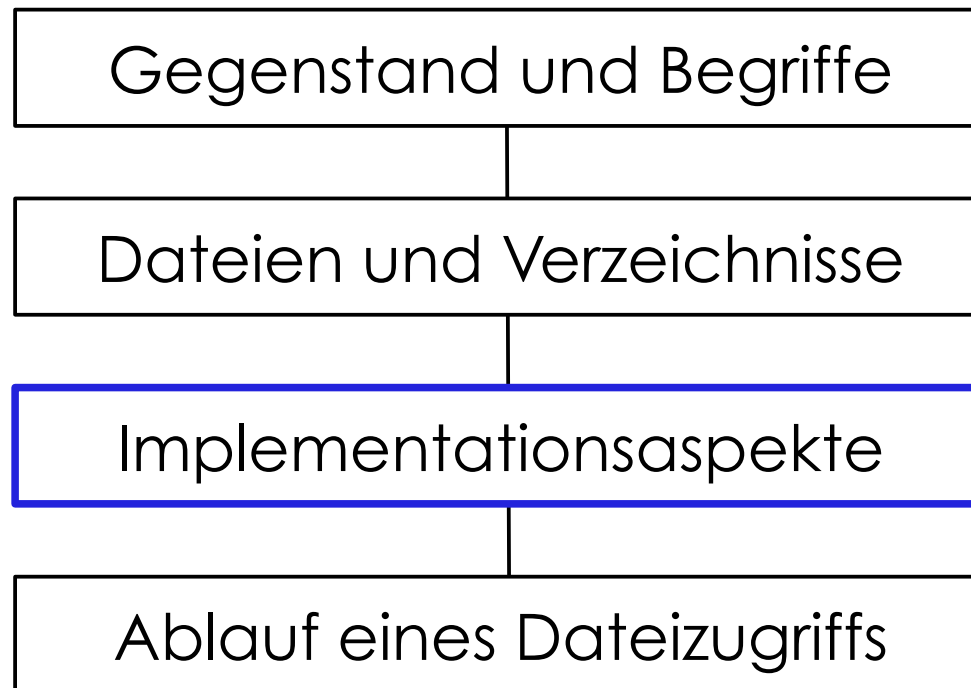
Einbinden eines Namensraumes
in einen anderen

- Gerät: `/dev/cdrom`



- `mount /dev/cdrom /volumes/CDRom`

Wegweiser



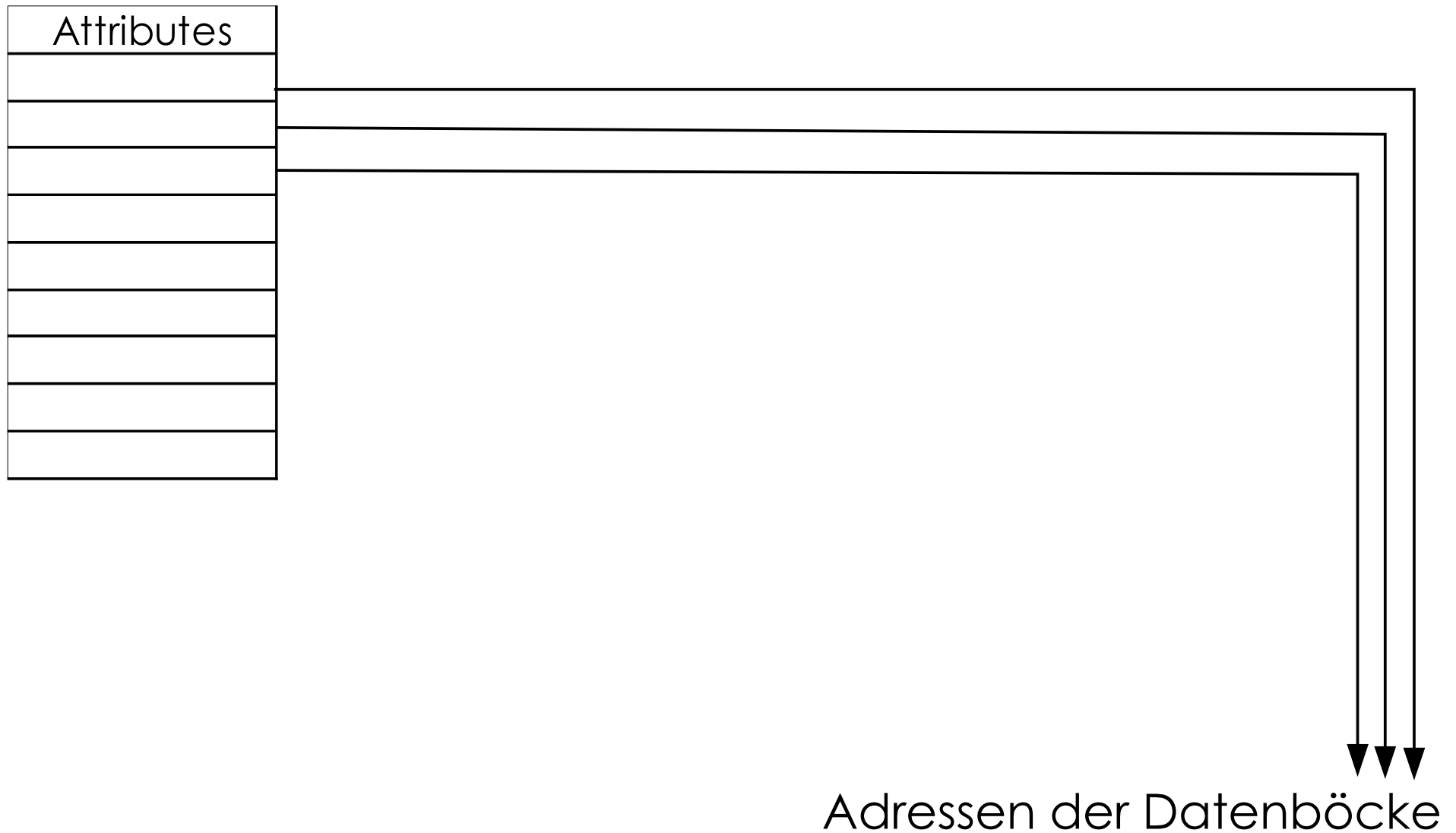
Verzeichnisse

Implementierung und Zugriff analog normale Datei mit:

- interner Struktur
- speziellen Operationen
 - ◆ opendir / closedir
 - ◆ readdir
 - ◆ lookup
 - ◆ rename

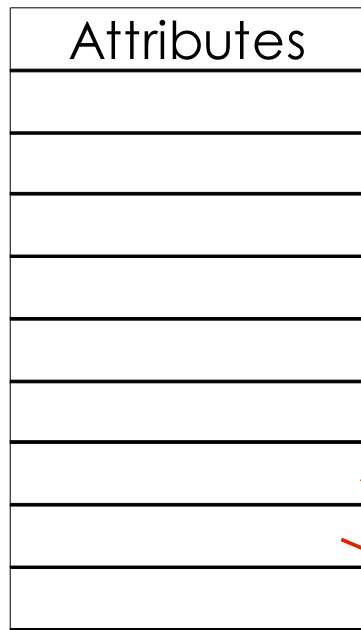
i-nodes (I-Knoten)

i-node



i-nodes (I-Knoten)

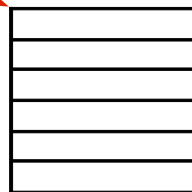
i-node



Single indirect block



Double indirect block



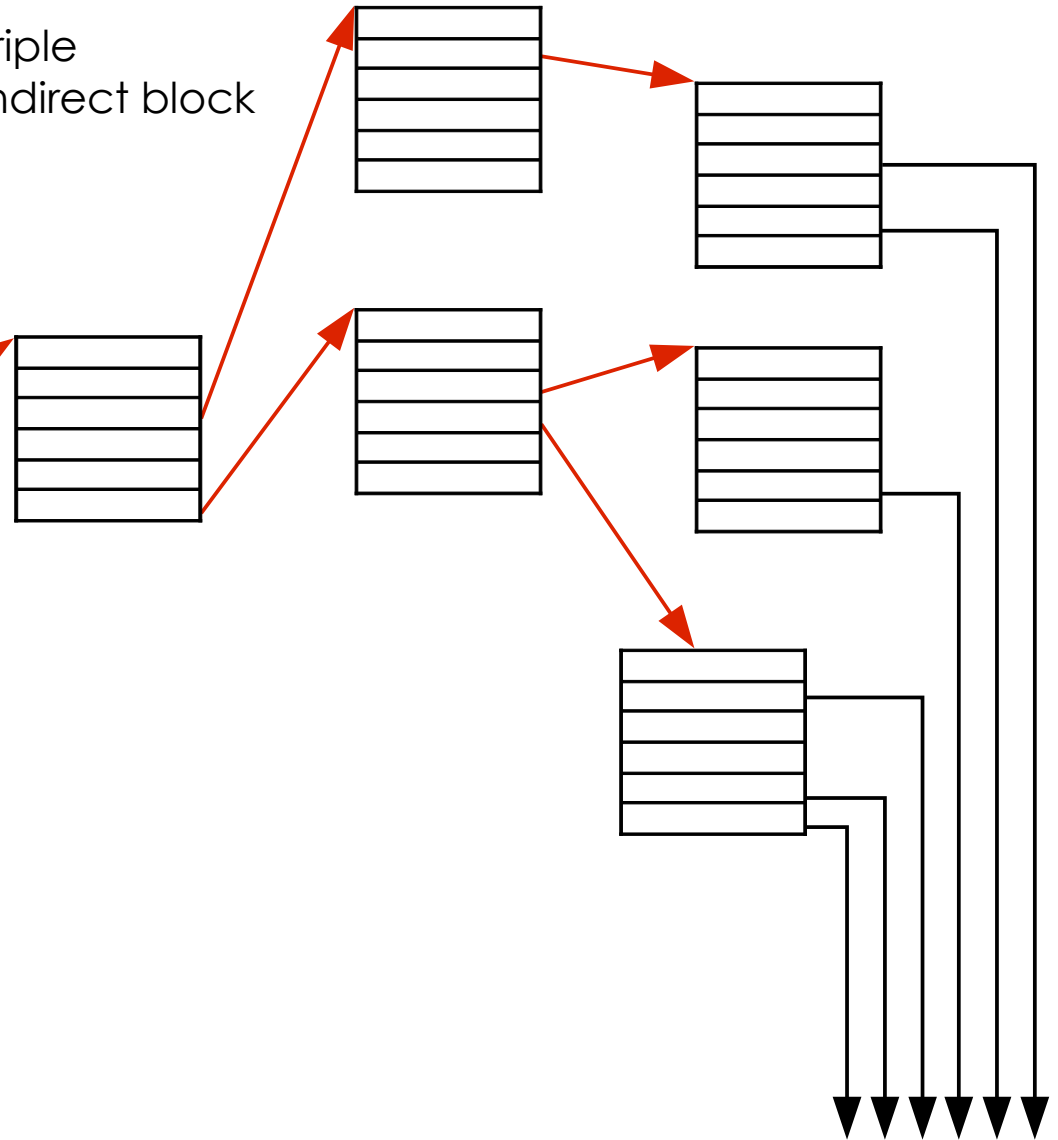
Adressen der Datenblöcke

i-nodes (I-Knoten)

i-node



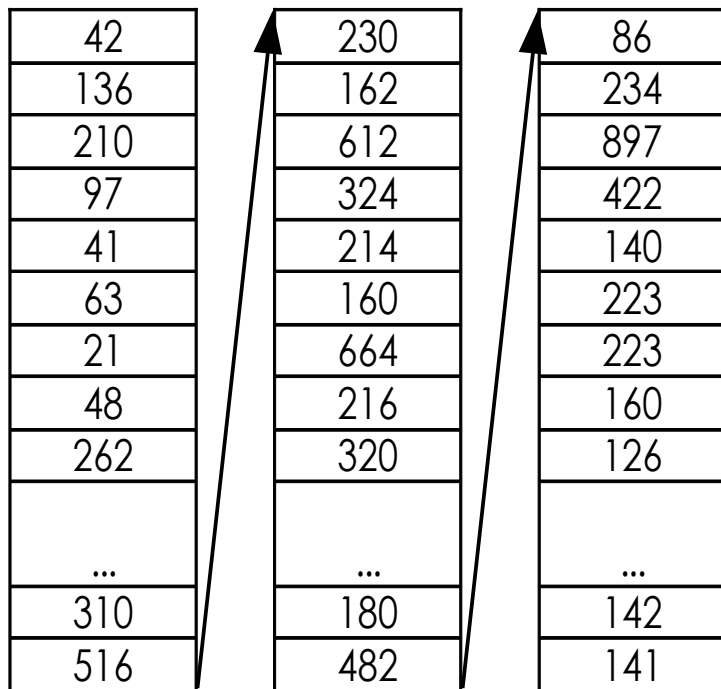
Triple indirect block



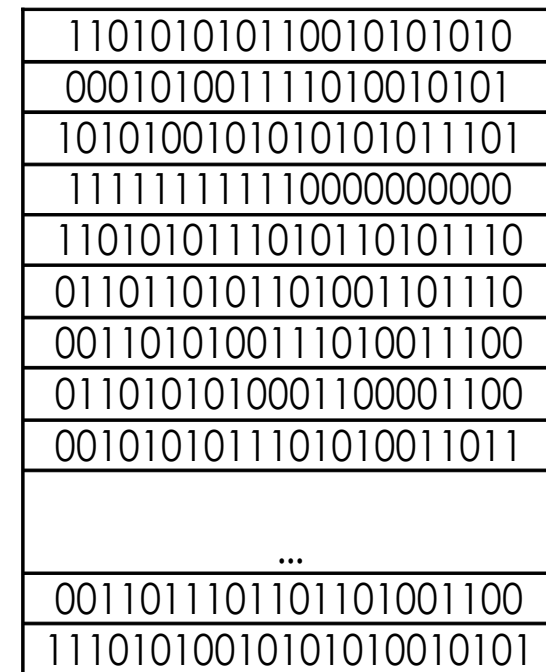
Adressen der Datenböcke

Verwaltung des freien Plattenspeichers

Verwendung von Listen



Verwendung von Bitmaps



Wegweiser



Schnittstelle

Namensdienste:

symbolische Namen → Identifikatoren

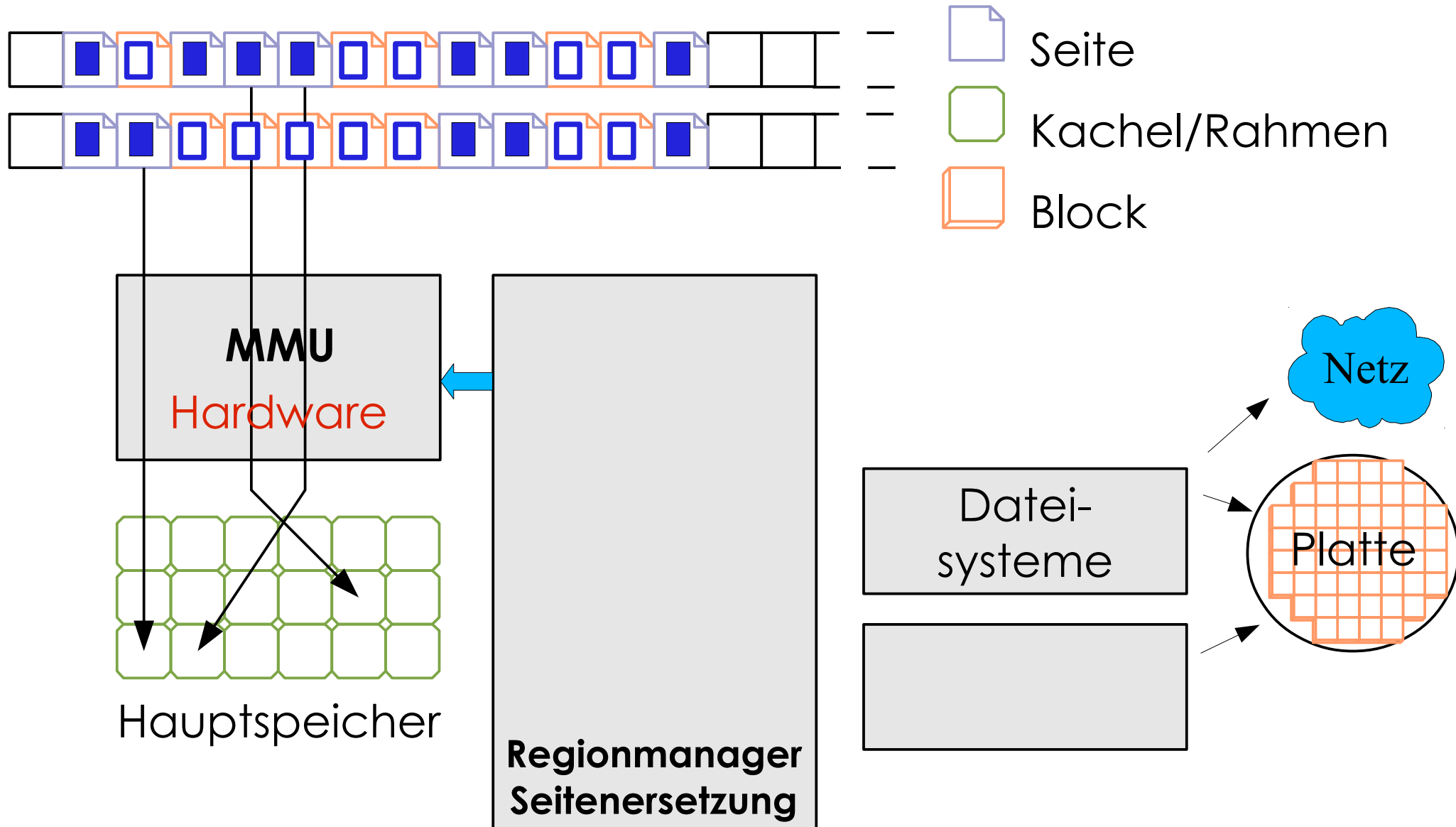
Zugriffe auf (Daten in) Dateien:

- häufig: **open** (Zugriff)* **close**
- Zugriff über
 - Kopieren aus Datei in Adressraum und zurück:
read/write (Datei, Position in Datei, Länge,
Adresse in Adressraum)
keine Ausrichtung der Adressen notwendig
seek-Operation ersetzt „Position in Datei“ als Parameter
 - Einblenden in Adressraum:
map (Datei, Adresse in Adressraum, Offset, Länge)
dann Zugriff über normale Maschinen-Instruktionen
Offset und Adresse müssen an Seitengrenzen ausgerichtet sein

Zugriffe auf Dateiattribute ...

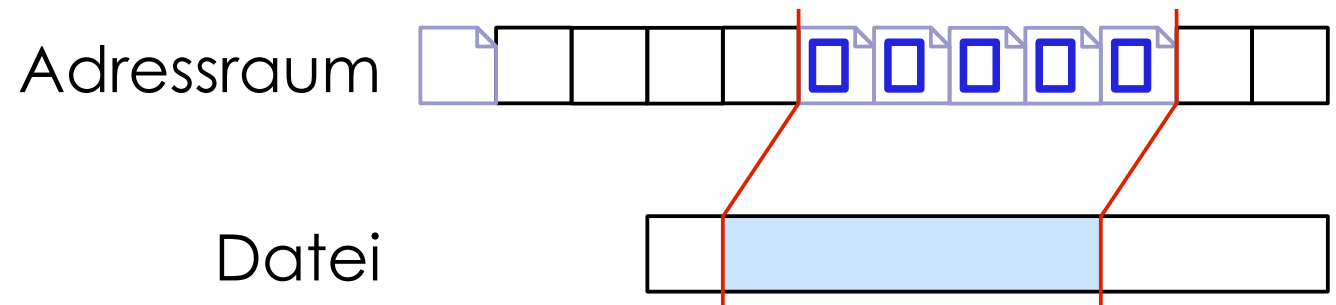
Integration in Systemarchitektur

Adressräume z. B. 4 GB

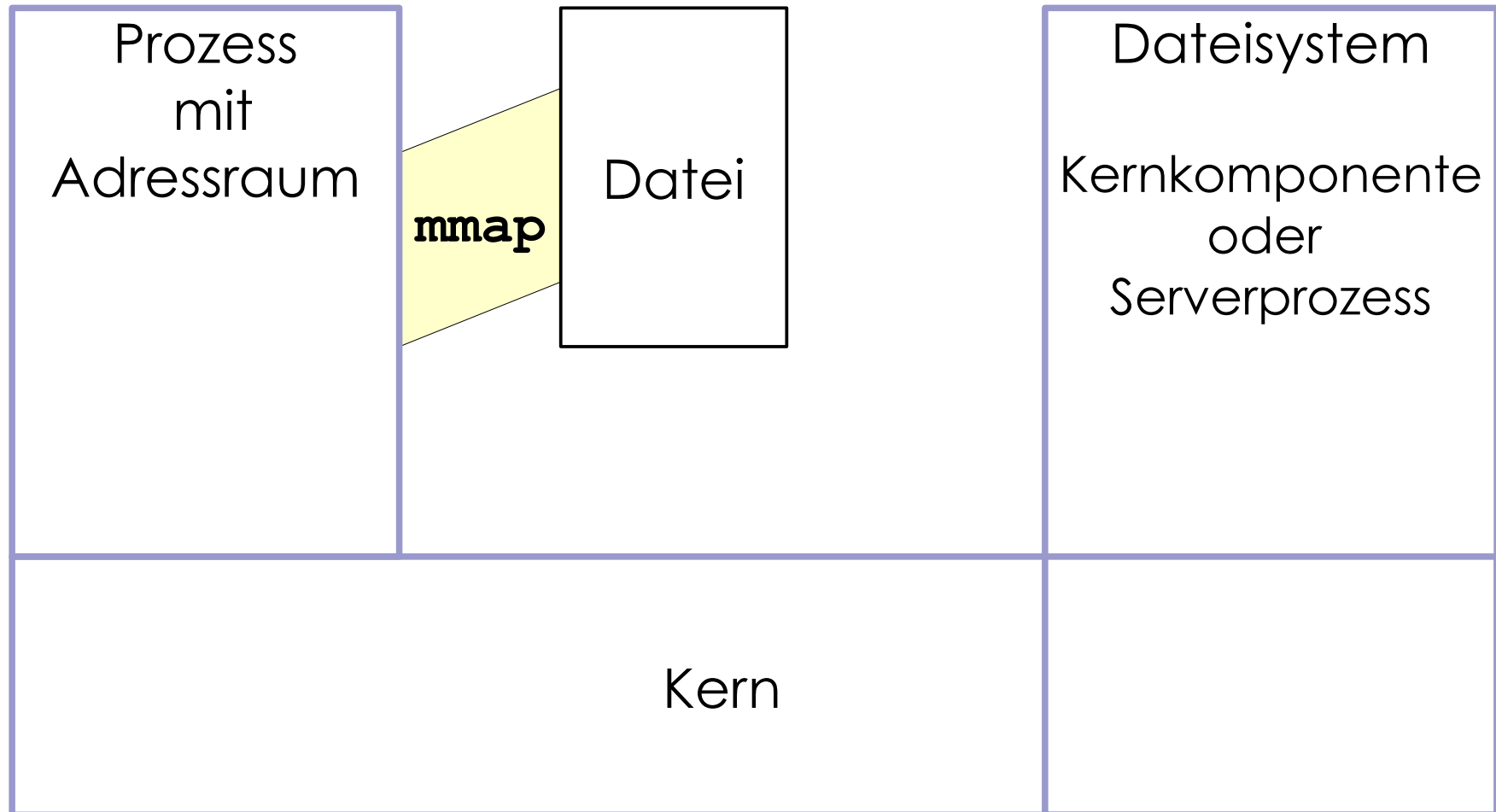


`mmap` (Datei, Adresse, Länge, Offset)

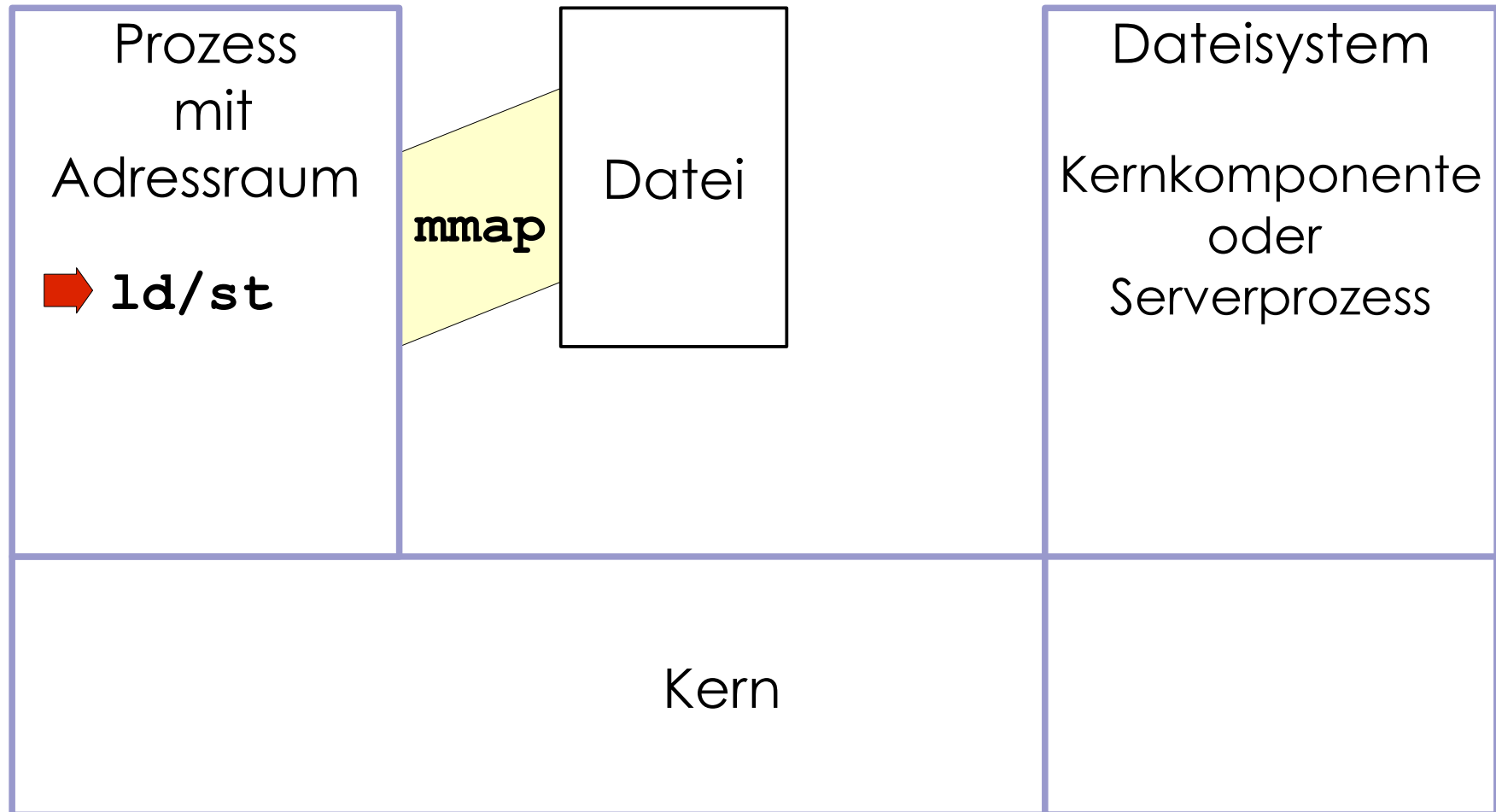
- bei Seitenfehler:
 - Zugriff auf Datei über Seitenfehlerbehandlung
 - dann Zugriff über normale HW-Instruktionen
- „Pager“ muss Struktur der Dateiimplementierung kennen!



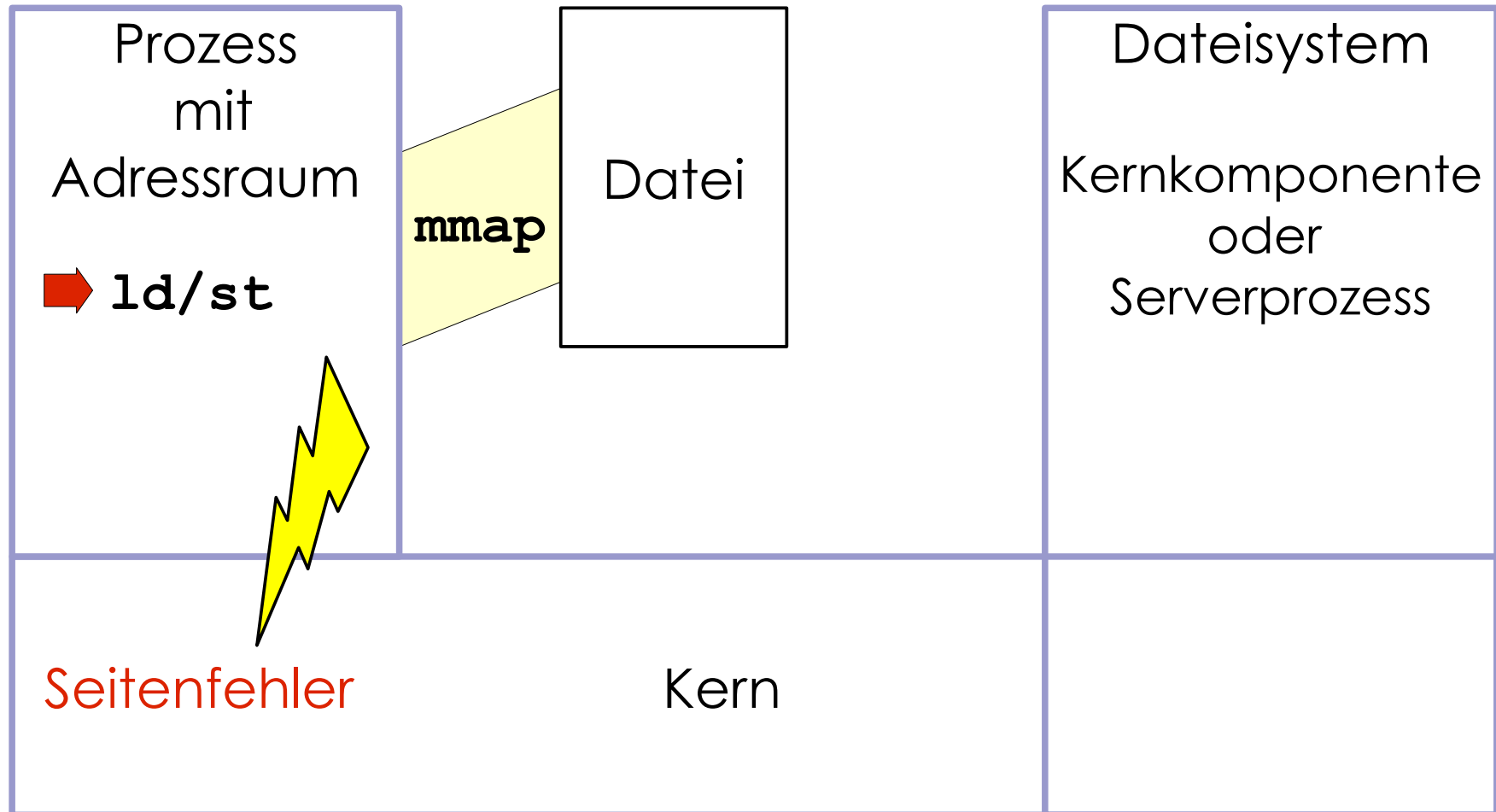
Zugriff mittels Einblendung



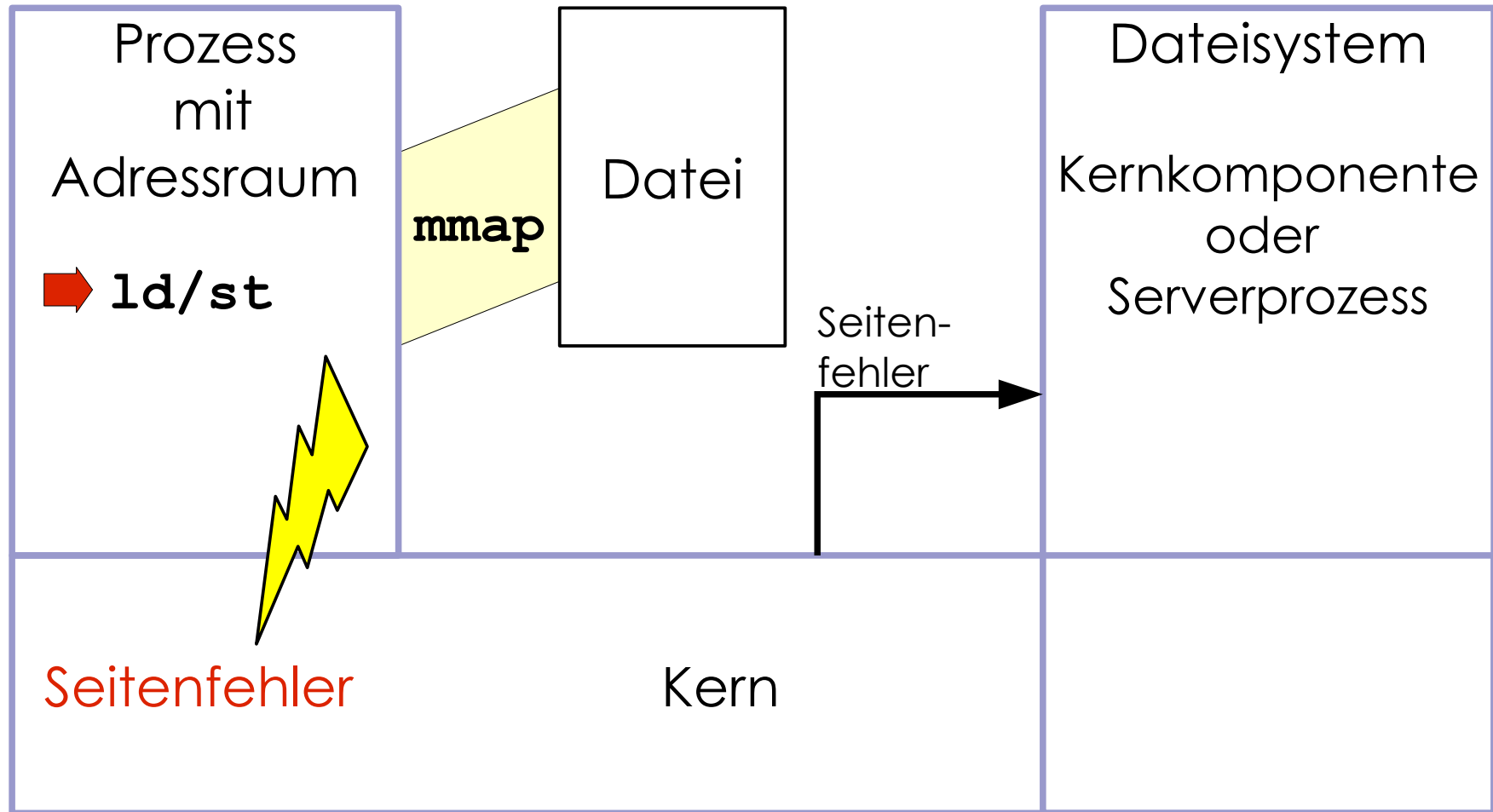
Zugriff mittels Einblendung



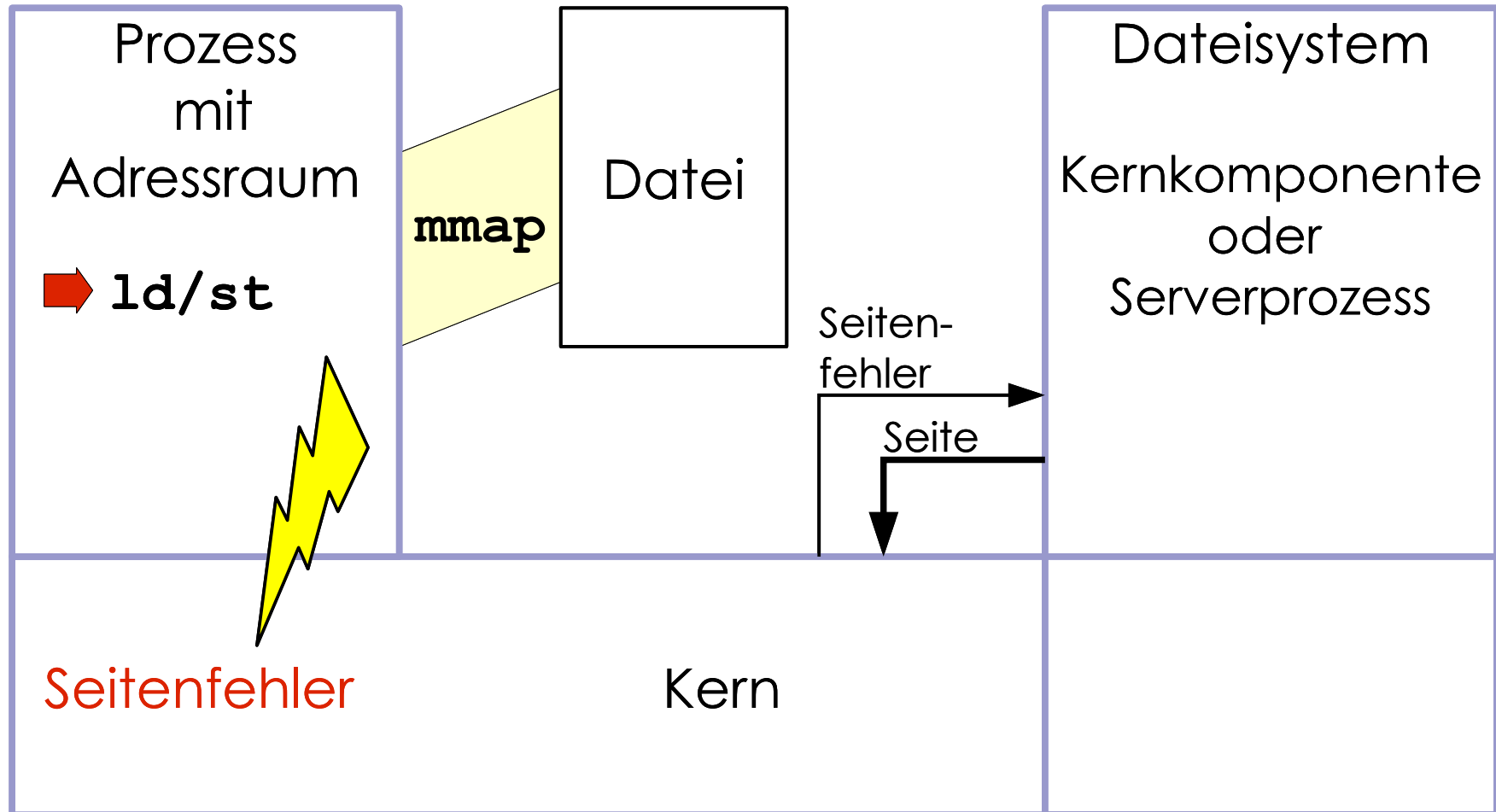
Zugriff mittels Einblendung



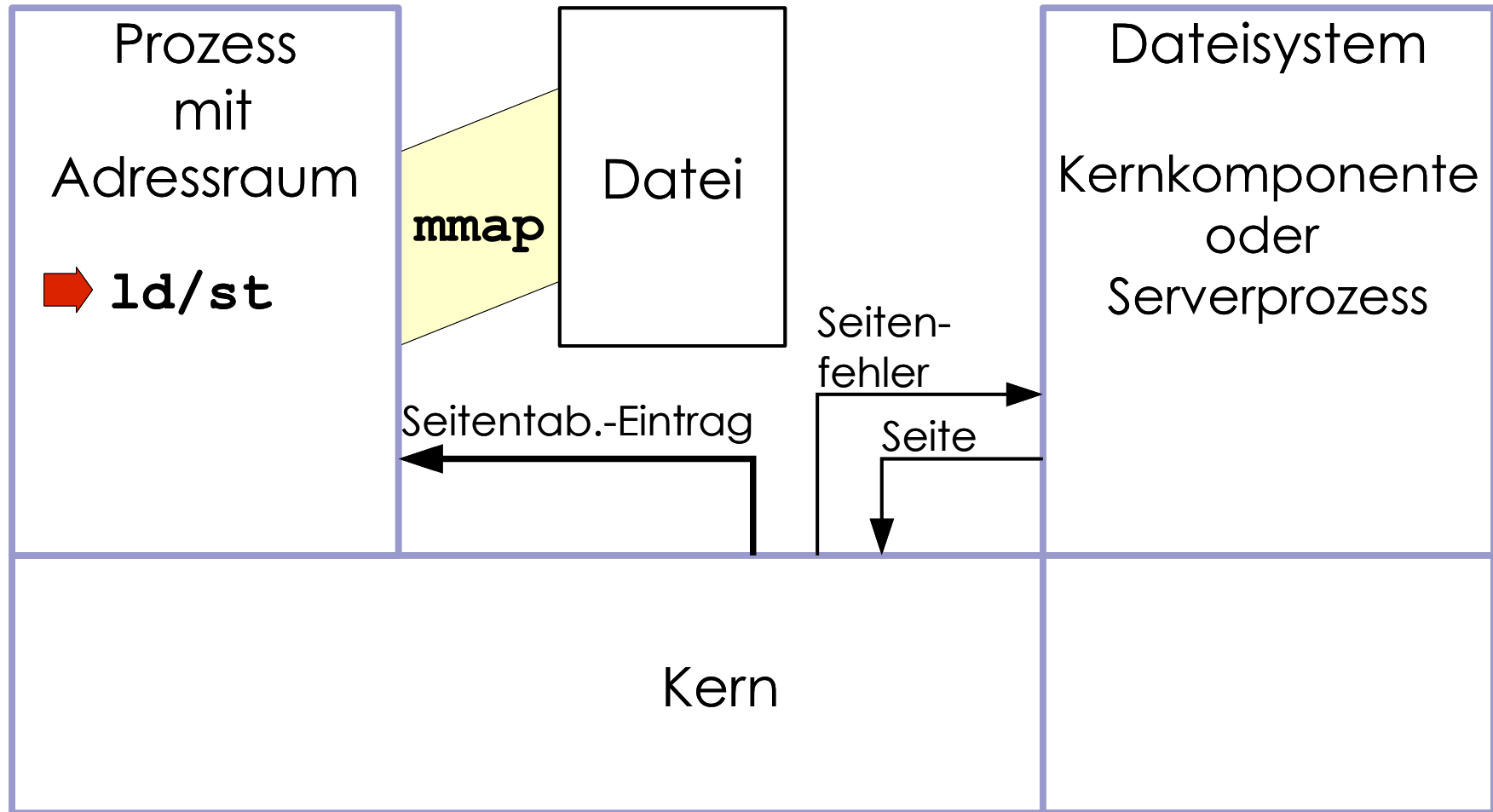
Zugriff mittels Einblendung



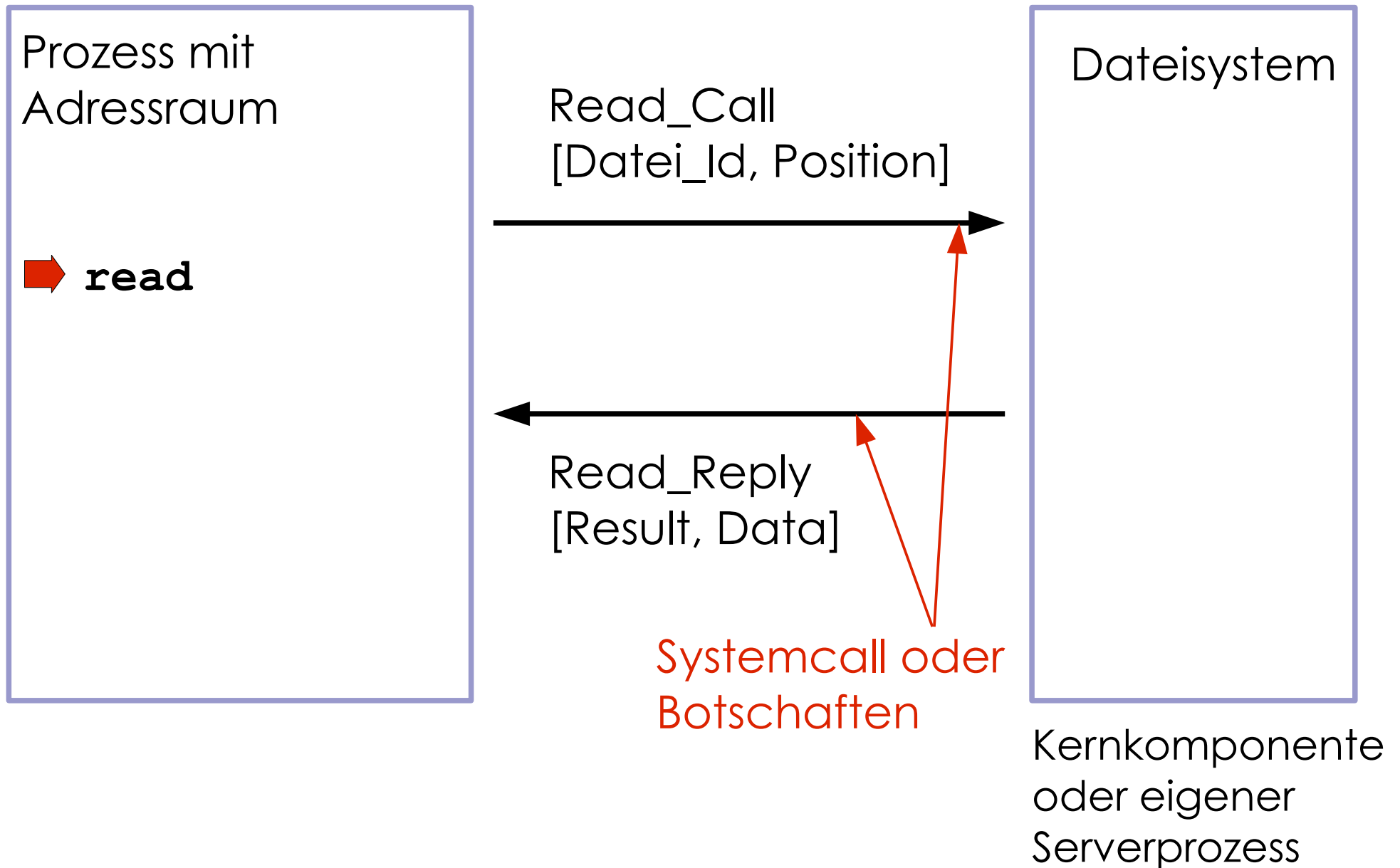
Zugriff mittels Einblendung



Zugriff mittels Einblendung



Zugriff mittels Kopieren (read/write)



Ablauf eines „write“ Dateizugriffs (Unix)

```
fd = open („drops/papers/xy.tex“, RW, ...)
```

Pfadname → i-node number (Namensauflösung)

i-node number → fd

```
seek (fd, position)
```

```
write (fd, address, length)
```

Im folgenden

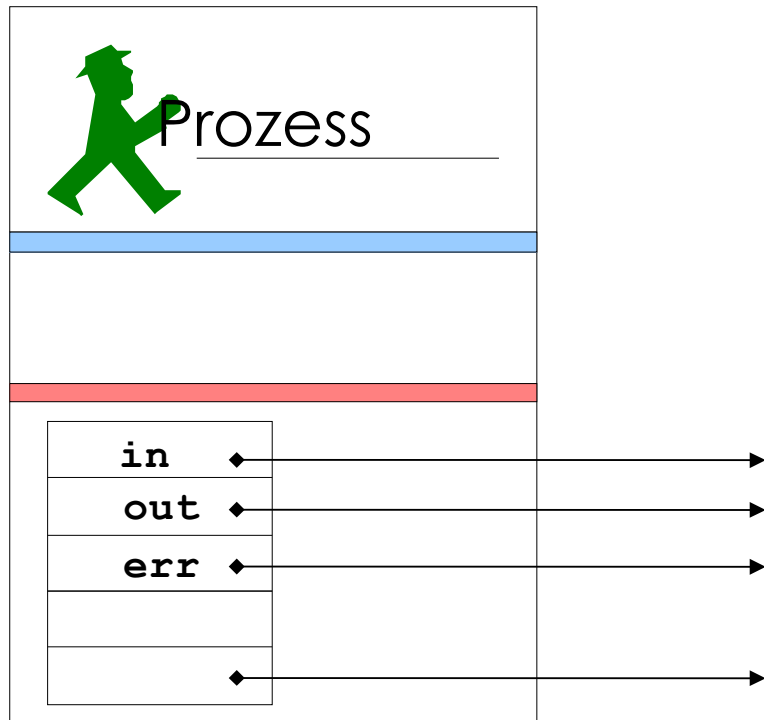
- beteiligte Datenstrukturen
- Ablauf

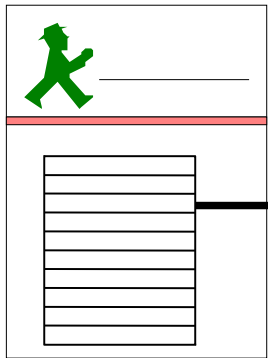
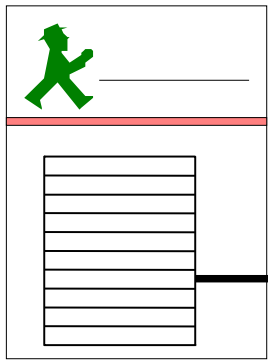
vereinfacht: ignoriert werden

- Fehlerbehandlung
- Attribute, Rechte
- Vorhandensein mehrerer Dateisysteme pro Rechner
- Plattentreiberdetails

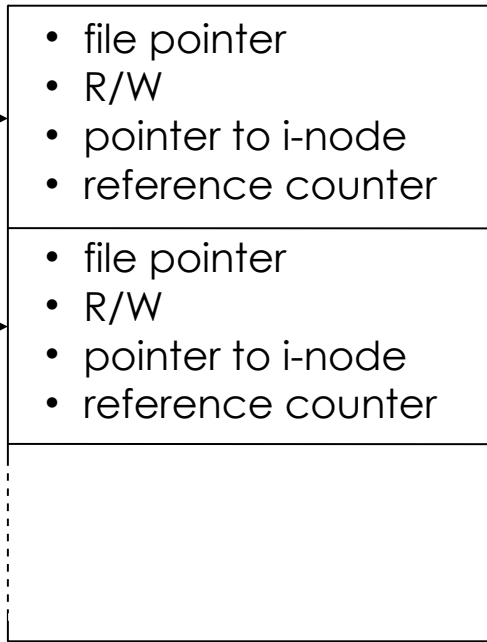
Beteiligte Datenstrukturen (des Dateisystems)

- Prozessteuerblock: current directory, Filedescriptor fd
- Verzeichnisse
- Tabelle geöffneter Dateien, pro Datei:
 - file pointer (gegenwärtige Schreib-/Leseposition), ...
 - i-node number
- Dateitabelle auf Platte (Tabelle der I-Knoten)
 - Index in dieser Tabelle (i-node number) identifiziert Datei
- für alle (auch nicht geöffnete) Dateien: i-node
 - Allokation (Zuordnung von Dateiblöcken zu Plattenblöcken)
 - Attribute
- Puffer-Verwaltung

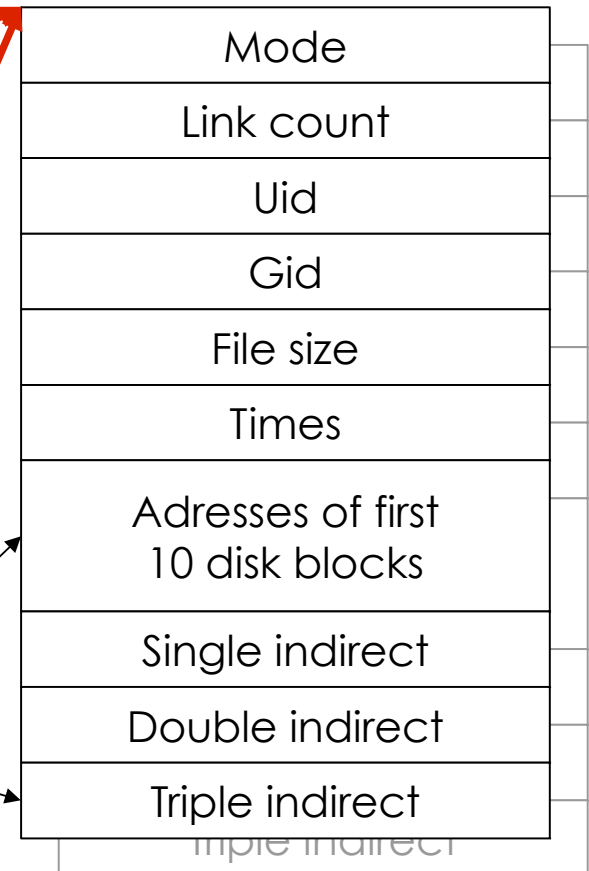




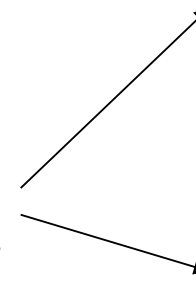
Open File Table



I-Node



Allokation
Plattenblöcke



Pufferverwaltung und Ablauf read/write

Block-Kachel-Tabelle

(„Buffer Cache“)

- Struktur im Kernadressraum
- Zuordnung:
Plattenadresse \leftrightarrow Puffer
- Puffer allokkieren/laden
(resident im Hauptspeicher)

Ablauf (read/write)

- fd \rightarrow file pointer, i-node
 \rightarrow Plattenadresse
- Puffer
 - evtl. Puffer besorgen
 - evtl. Laden von Platte
- Kopieren in Adressraum

Ablauf write

```
write (fd, address, length);
```

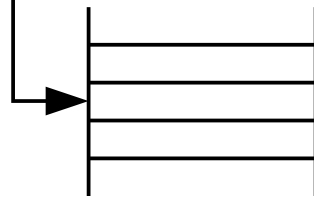
Nutzer
Adressraum



Datei



filepointer



Adressraum



Datei



Ablauf write

```
write (fd, address, length);
```

Nutzer
Adressraum

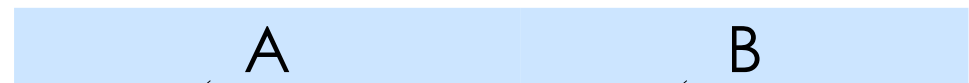


Datei



filepointer

Adressraum

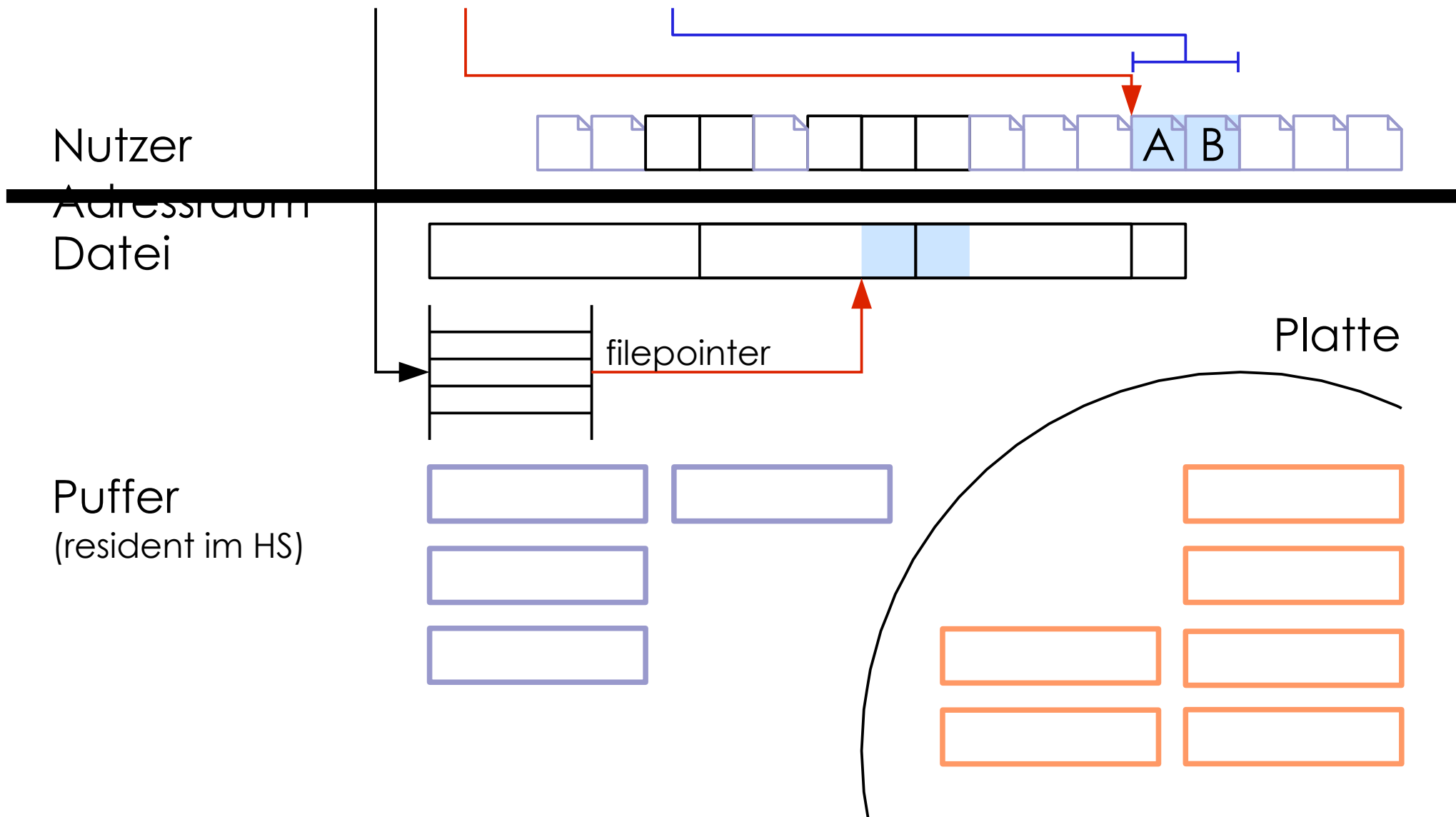


Datei



Ablauf write

```
write (fd, address, length);
```



Ablauf write

```
write (fd, address, length);
```

Adressraum



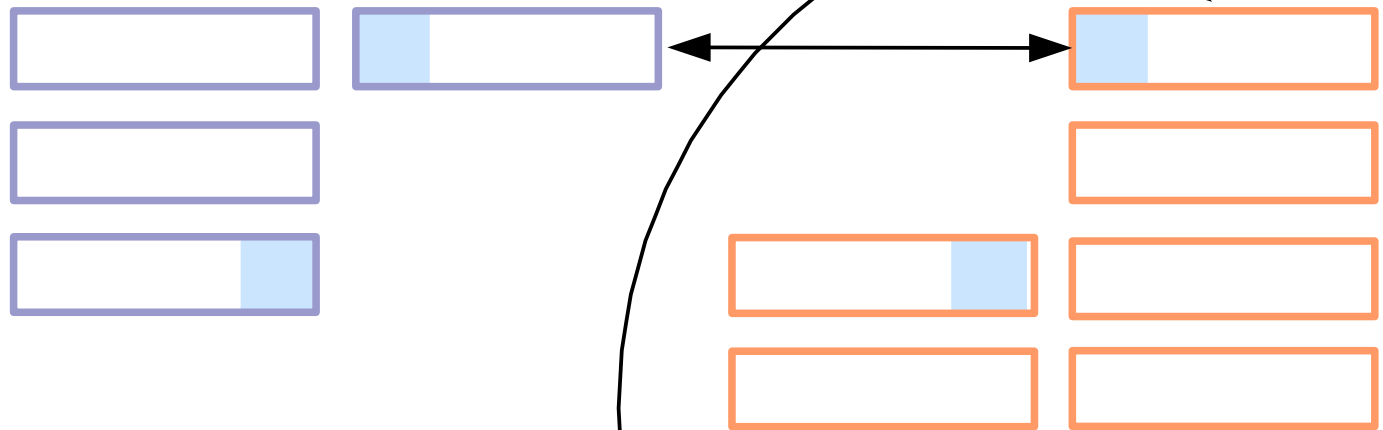
Datei



filepointer

Platte

Puffer
(resident im HS)



← → Datei ↔ Platte

↔ Puffer ↔ Platte

Ablauf write

```
write (fd, address, length);
```

Adressraum



Datei

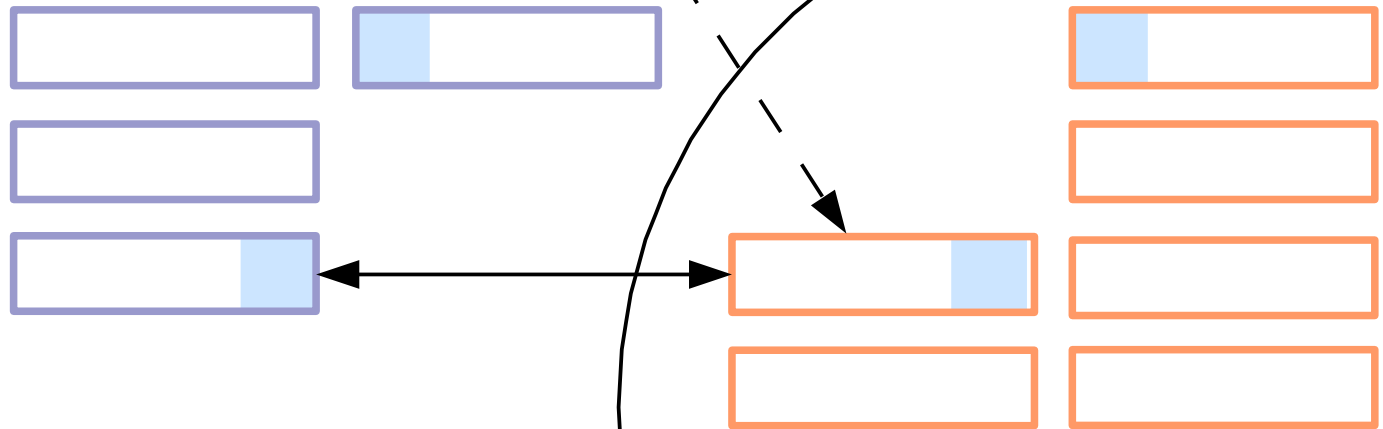


filepointer

Platte

Puffer

(resident im HS)

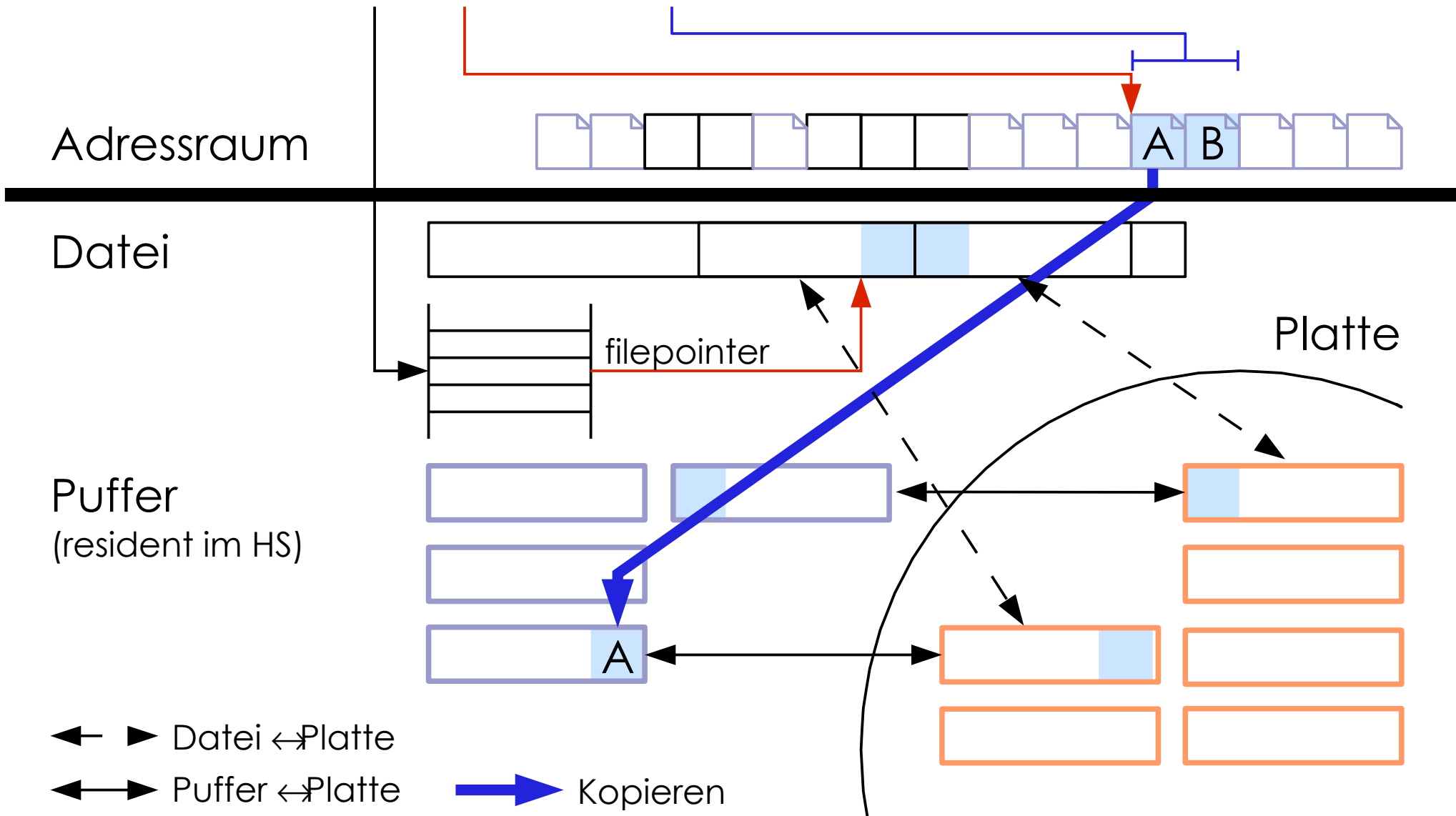


← → Datei ↔ Platte

↔ Puffer ↔ Platte

Ablauf write

```
write (fd, address, length);
```



Ablauf write

```
write (fd, address, length);
```

Adressraum



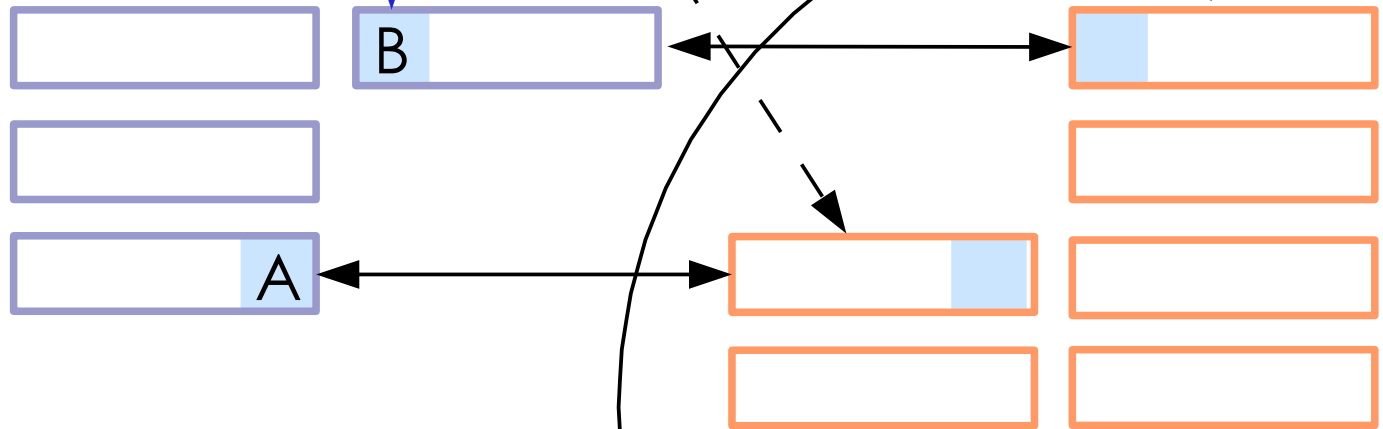
Datei



filepointer

Platte

Puffer
(resident im HS)



← → Datei ↔ Platte

← → Puffer ↔ Platte → Kopieren

Ablauf write

```
write (fd, address, length);
```

Adressraum



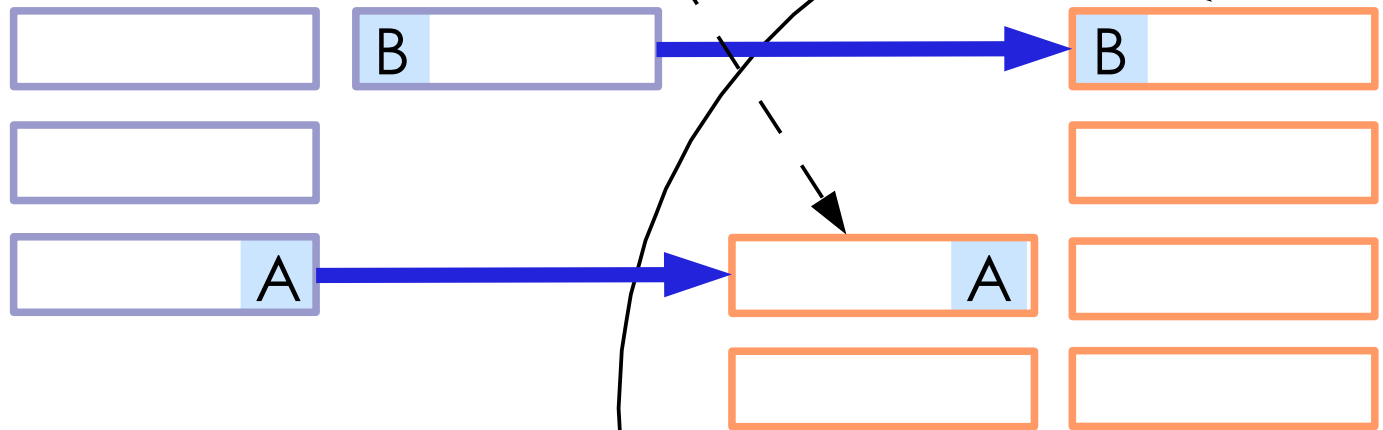
Datei



filepointer

Platte

Puffer
(resident im HS)



← → Datei ↔ Platte

← → Puffer ↔ Platte → Kopieren

Gegenüberstellung: read/write vs. mmap

read/write

- Kopie vor jeder Manipulation
- mehrere Prozesse können sich per read eine Kopie holen und dann den vom anderen Prozess vorher geschriebenen Wert überschreiben

mmap

- Einsparen von Kopieroperationen !!!
- Synchronisation wie bei gemeinsamem Speicher

Gegenüberstellung: read/write vs. map

read/write

- Kopie vor jeder Manipulation
- mehrere Prozesse können sich per read eine Kopie holen und dann den vom anderen Prozess vorher geschriebenen Wert überschreiben

mmap

- Einsparen von Kopieroperationen !!!
- Synchronisation wie bei gemeinsamem Speicher

Beispiel: Konto

```
open (Kontodatei, ... );
```

```
KontoAdrInDatei =  
    Find_Konto;  
read (Kontodatei,  
    KontoAdrInDatei,  
    &Konto, Länge);  
  
Konto.Stand += Betrag;  
  
write (Kontodatei,  
    KontoAdrInDatei,  
    &Konto, Länge);
```

```
mmap (Kontodatei, ...);  
  
KontoAdr = Find_Konto;  
  
KontoAdr->Stand += Betrag;
```

Beispiel: Konto

```
open (Kontodatei, ... );
```

```
KontoAdrInDatei =  
    Find_Konto;  
read (Kontodatei,  
    KontoAdrInDatei,  
    &Konto, Länge);  
  
Konto.Stand += Betrag;  
  
write (Kontodatei,  
    KontoAdrInDatei,  
    &Konto, Länge);
```

```
mmap (Kontodatei, ...);  
  
KontoAdr = Find_Konto;  
  
KontoAdr->Stand += Betrag;
```

Zusammenfassung

- Integration von Dateisystemen in Systemarchitektur
- Prinzipielle Teilaufgaben
- Ein einfaches altes Beispiel

ABER: Limitationen

Kleine Blöcke, riesige Dateien

Ausfall von Platten

Absturz eines BS → Verlust von Pufferinhalten,

Inkonsistenzen

Platte ./ Flash-Speicher

→ spätere Kapitel in dieser Vorlesung