

Betriebssystemabstraktionen

Adressraum ➡ physikalische, logische oder virtuelle Adressen

Speicher ➡ im Vorder- bzw. Hintergrund

Datei ➡ langfristige (permanente) Speicherung von Informationen

Prozess ➡ feder-, leicht- oder schwergewichtige Aktivitätsträger

Darüberhinaus sind oftmals noch **Koordinations-** und **Kommunikationsmittel** zur Unterstützung der Interaktion nebenläufiger Prozesse verfügbar.

Adressraum

physikalischer Adressraum ➡ Hardware (Ebene₂)

- die Größe entspricht der Adressbreite der CPU: N Bit $\Rightarrow 2^N$ Bytes
- nicht alle Adressen sind gültig und zur Programmspeicherung verwendbar

logischer Adressraum ➡ Kompilierer, Binder, Betriebssystem (Ebene_{5/4/3})

- definiert einen zusammenhängenden, linear adressierbaren Programmbereich
- alle Adressen sind gültig und zur Programmspeicherung verwendbar
- ist typischerweise (sehr viel) kleiner als die Adressbreite der CPU hergibt

virtueller Adressraum ➡ Betriebssystem (Ebene₃)

- ein logischer Adressraum, dessen Größe der Adressbreite der CPU entspricht

Physikalischer Adressraum

Toshiba Tecra 730CDT, 1996:

Adressbereich	Belegung
00000000–0009ffff	RAM
000a0000–000c7fff	System
000c8000–000dffff	keine
000e0000–000fffff	System
00100000–090fffff	RAM
09100000–fffdffff	keine
fffe0000–ffffffff	System



Logischer Adressraum (1)

- jedes Programm wird in einem eigenen logischen Adressraum ausgeführt
 - die $\left\{ \begin{array}{c} \text{Anfangs-} \\ \text{End-} \end{array} \right\}$ Adressen aller logischen Adressräume sind (meist) gleich
 - logische Adressen sind auf die gültigen physikalischen Adressen abzubilden
- die erforderliche *Adressabbildung* erfolgt (typischerweise) mehrstufig:

Programmadresse \rightsquigarrow logische Adresse
logische Adresse \rightsquigarrow physikalische Adresse

- im Gegensatz zu physikalischen Adressen sind logische Adressen mehrdeutig

Logischer Adressraum (2)

- in der Adress(raum)abbildung sind verschiedene Ebenen (X Kap. 5) involviert:

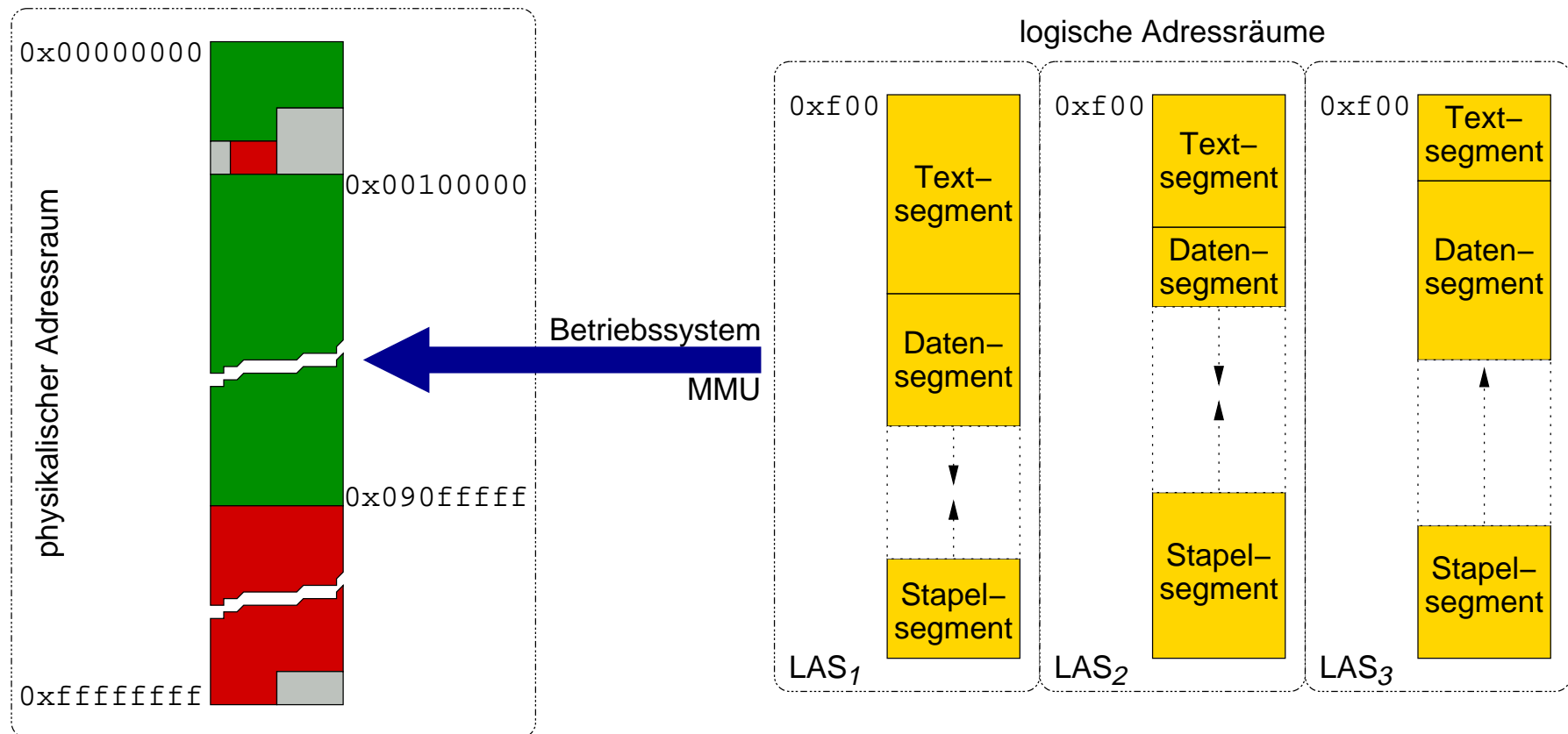
Entwicklungszeit	➡ Programmierer	➡ Ebene ₆
Übersetzungszeit	➡ Kompilierer, Assembler	➡ Ebene _{5/4}
Bindezeit	➡ Binder	➡ Ebene ₄
Ladezeit	➡ Lader	➡ Ebene ₃
Laufzeit	➡ bindender Lader, MMU	➡ Ebene _{3/2}

- je später diese Abbildung durchgeführt wird, desto. . .
 - ➡ höher das Abstraktionsniveau und geringer die Hardwareabhängigkeit
 - ➡ höher der Systemaufwand und geringer der Spezialisierungsgrad
- der Zeitpunkt unterliegt einem „*trade off*“ zwischen Flexibilität und Effizienz

Adressraumsegmentierung

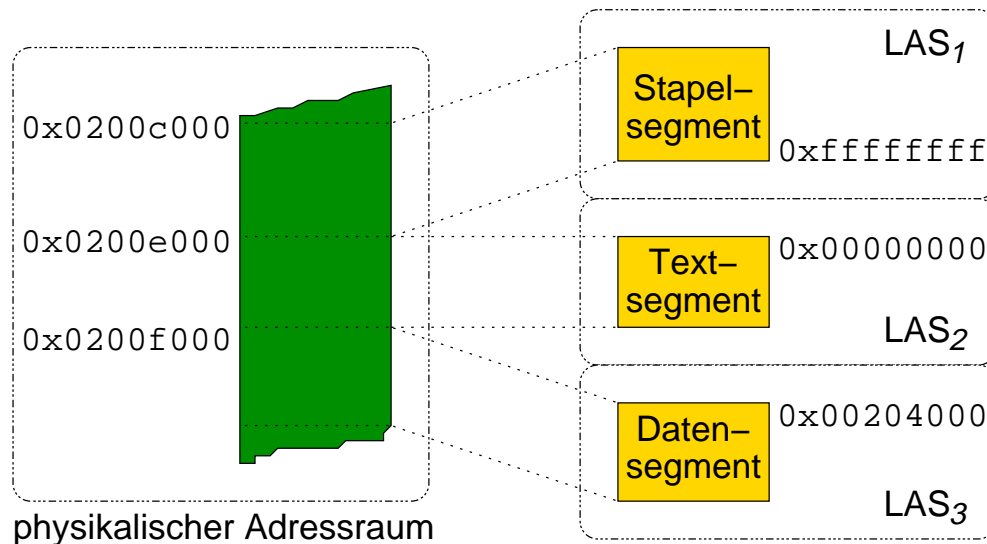
- gebräuchlich ist die logische Unterteilung des Adressraums in drei Segmente:
 - Textsegment** (*text segment*) enthält die Maschinenabweisungen der Ebene₂ (CPU) und andere Programmkonstanten ➡ **statisch/dynamisch**
 - Datensegment** (*data segment*) kapselt initialisierte Daten, globale Variablen und ggf. eine Halde (*heap*) ➡ **statisch/dynamisch**
 - Stapelsegment** (*stack segment*) beherbergt lokale Variablen, Hilfsvariablen und aktuelle Parameter ➡ **dynamisch**
- typischerweise werden mehrere logische Adressräume nebeneinander verwaltet
 - Ebene₃ (Betriebssystem) sorgt für die *Adressraumabbildung*
 - Ebene₂ (*memory management unit*, MMU) sorgt für die *Adressumsetzung*

Adressraumabbildung durch Betriebssystem/MMU (1)



Relokation zur Laufzeit

- die Abbildung gleichzeitig vorhandener logischer Adressräume ist disjunktiv
☞ alle Segmente werden überschneidungsfrei im phys. Adressraum angeordnet



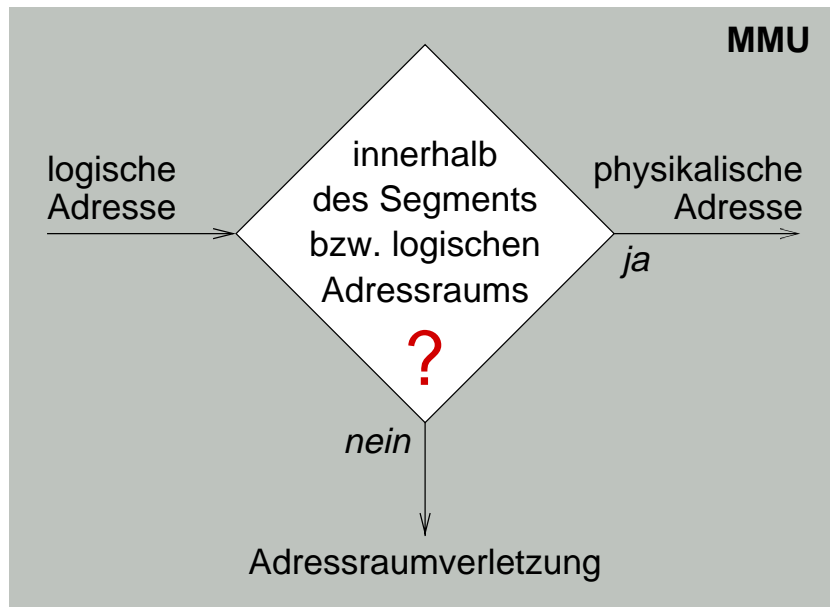
zur **Ladezeit** lokalisiert das BS den passenden Bereich im phys. Adressraum und programmiert die MMU


zur **Laufzeit** setzt die MMU jede logische Adresse um in die korrespondierende physikalische Adresse

- ein „Ausbruch“ aus einem Segment/logischen Adressraum ist zu verhindern

Isolation



- die MMU (Ebene₂) kapselt Adressräume und sichert Zugriffsschutz



- Adressraumverletzung führt zum Fehler
 - synchrone Programmunterbrechung
 -  *segmentation violation*.....
 - Programmabbruch ist die Folge
- Spiritus rector „Betriebssystem“
 - korrektes Funktionieren ist zwingend

- das Betriebssystem ist verantwortlich für die Integrität *aller* Adressräume

Schutzdomäne

- Hardware-gestützte Adressraumisolation erhöht die *Sicherheit* von Systemen
safety (Sicherheit, Gefahrlosigkeit, Ungefährlichkeit)  Schutz von Menschen und Sachwerten vor dem Versagen technischer Systeme
 - eine Fehlerausbreitung durch „Bitkipper“ z. B. im Speicher ist eingrenzbar
 - gleiches gilt für die Auswirkung von Berechnungsfehlern in Programmen
- security* (Sicherheit, Schutz, Sorglosigkeit, Zuversicht)  Schutz von Informationen und Informationsverarbeitung vor intelligenten Angreifern
 - „Eindringlinge“ können von den Programmen fern gehalten werden
 - Programmen wird ein Ausbruch aus ihren Adressräumen erschwert
- bei fehlerhaftem Betriebssystem nützt die beste Hardware/MMU nichts — letztlich ist ihr korrektes Funktionieren auch überhaupt nicht garantierbar !!!

Virtueller Adressraum (1)

- die Anzahl gültiger phys. Adressen dimensioniert einen logischen Adressraum
 - in Realität entspricht dies der Größe des gesamten Arbeitsspeichers (RAM)
 - der Speicherbedarf eines Programms kann diese Größe leicht übersteigen
 - umso problematischer: gleichzeitig nebeneinander bestehende Programme
- Überbelegung des Arbeitsspeichers im Mehrprogrammbetrieb ist Normalität
 - Einsatz von Hintergrundspeicher (Platte) entschärft die Engpasssituation:
 - ☞ zur Zeit nicht benötigte Programmbereiche liegen auf der Platte
 - ☞ nur die zur „Arbeitsmenge“ gehörenden Bereiche sind im RAM
- ein virtueller Adressraum ist dimensioniert durch die Adressbreite der CPU

Adressbreite (N Bits)	Adressraumgröße (2^N Bytes)	Dimension		
16	65 536	64	Kilo	2^{10}
20	1 048 576	1	Mega	2^{20}
32	4 294 967 296	4	Giga	2^{30}
48	281 474 976 710 656	256	Tera	2^{40}
64	18 446 744 073 709 551 616	16 384	Peta	2^{50}

☞ Ein Rechner ist nur mit einem Bruchteil des von einer (zukünftigen) CPU adressierbaren Arbeitsspeichers wirklich bestückt!

Die Aufgabe soll sein, den gesamten Adressraum byteweise `:-()` zu löschen. Dabei wird pro Byte eine Zugriffszeit von 1 Nanosekunde (ns) angenommen:

```
void clear () {
    char* p = 0;
    do *p++ = 0;
    while (p);
}
```

`gcc -O6 -S`

```
_clear:
    li    r2,0
    li    r0,0
L2:
    stb    r0,0(r2)
    addic. r2,r2,1
    bne+   cr0,L2
    blr
```

PowerPC G4. Jeder Befehl ist vier Bytes lang. Die Löschschleife (L2) umfasst drei Befehle, die von der CPU aus dem Speicher zu lesen sind. Der Löschbefehl (`stb r0,0(r2)`) schreibt ein Byte mit dem Wert 0 in die nächste Speicherzelle. Jeder Schleifendurchlauf greift somit auf $4 \times 3 + 1 = 13$ Bytes zu.

[Warum funktioniert das so nicht?]

☞ Der Löschvorgang eines Bytes kostet (schlimmstenfalls) 13 ns!

Adressraumgröße (Bytes)	Löschdauer	Einheit
2^{16}	0,000851968	Sekunden
2^{20}	0,013631488	Sekunden
2^{32}	55,834574848	Sekunden
2^{48}	42,351558995816	Tage
2^{64}	7604,251425615937	Jahre (ohne Schaltjahre)

Virtueller Speicher. Die zur Zeit nicht benötigten Programmbereiche bzw. Abschnitte des logischen Adressraums liegen im Hintergrundspeicher. Bei Bedarf werden diese Abschnitte „seitenweise“ in den Vordergrundspeicher geholt, d. h., eingelagert. Angenommen, jede Seite ist 4 KB groß und die mittlere Zugriffszeit auf den Hintergrundspeicher (Platte), um eine Seite einzulagern, liegt bei 5 ms. Wie hoch ist dann die mittlere Zugriffszeit auf jedes einzelne Byte? $1,220703125 \mu s$. Wie lange dauert dann jeweils der Löschvorgang für die Adressräume? $> 1,5$ Stunden (2^{32}).

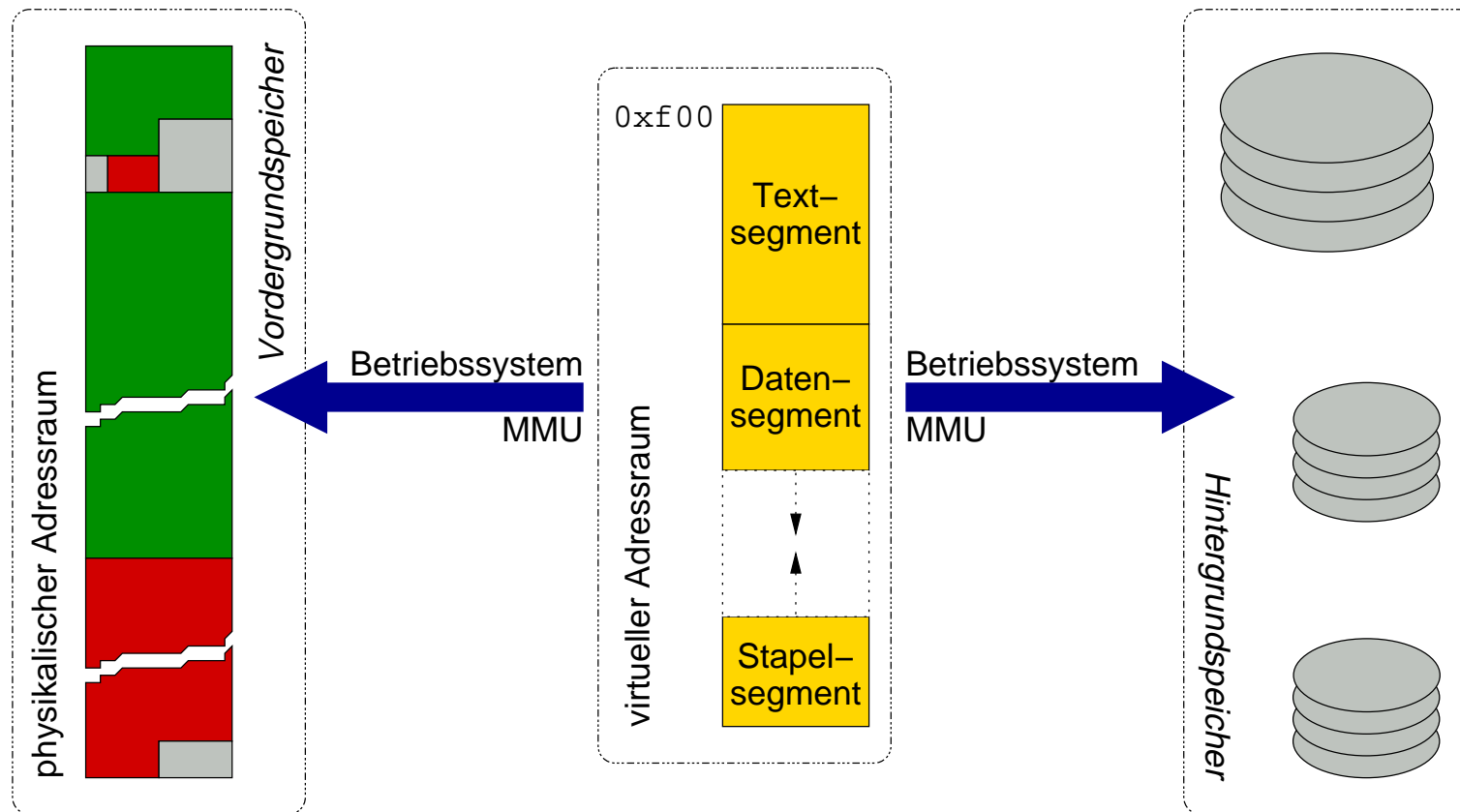
Virtueller Adressraum (2)

- die *Adressabbildung* (logischer Adressräume) ist um eine Stufe zu erweitern:

Programmadresse	\leadsto	logische Adresse
logische Adresse	\leadsto	virtuelle Adresse
virtuelle Adresse	\leadsto	physikalische Adresse

- Zugriffe über virtuelle Adressen können implizit Ein-/Ausgabe zur Folge haben
 - die MMU lässt nur gültige Zugriffe auf den *Vordergrundspeicher* zu
 - ggf. werden Zugriffe dann partiell interpretiert vom Betriebssystem (☞ *trap*)
 - Ergebnis ist die Einlagerung des „Operanden“ vom *Hintergrundspeicher*
 - dies kann zur Auslagerung anderer Vordergrundspeicherinhalte führen
- die Ein-/Auslagerung erfolgt typischerweise auf Seitenbasis (☞ *demand paging*)

Adressraumabbildung durch Betriebssystem/MMU (2)



`pa = mmap (addr, len, prot, flags, fd, offset)` legt ein Segment über einen (gemeinsamen) Speicherbereich oder eine (gemeinsame) Datei

- für den aufrufenden Prozess ändert sich seine *Adressraumabbildung*
- die zurück gelieferte Adresse, `pa`, identifiziert den Segmentanfang
- mit `addr = 0` erhält das System Freiheit über die Bestimmung von `pa`²⁴
- ein an `pa` ggf. vorher liegendes Segment wird ausgeblendet (`munmap(2)`)

`ok = munmap (addr, len)` entfernt die Abbildung für das Segment

²⁴Als „Nebeneffekt“ ist so die Erzeugung eines Speichersegments vorgegebener Länge (`len`) möglich.

Speicher

- die Abstraktion „Speicher“ kommt mit zwei fundamentalen Ausprägungen:

Vordergrundspeicher auch *Hauptspeicher* (RAM)

- der entsprechend bestückte (RAM-) Bereich im physikalischen Adressraum
- in ihm erfolgt die Ausführung der Programme („von Neumann Rechner“)
- kann größer sein als der physikalische Adressraum (*bank switching*)

Hintergrundspeicher auch *Massenspeicher* (Band, Platte, CD, DVD)

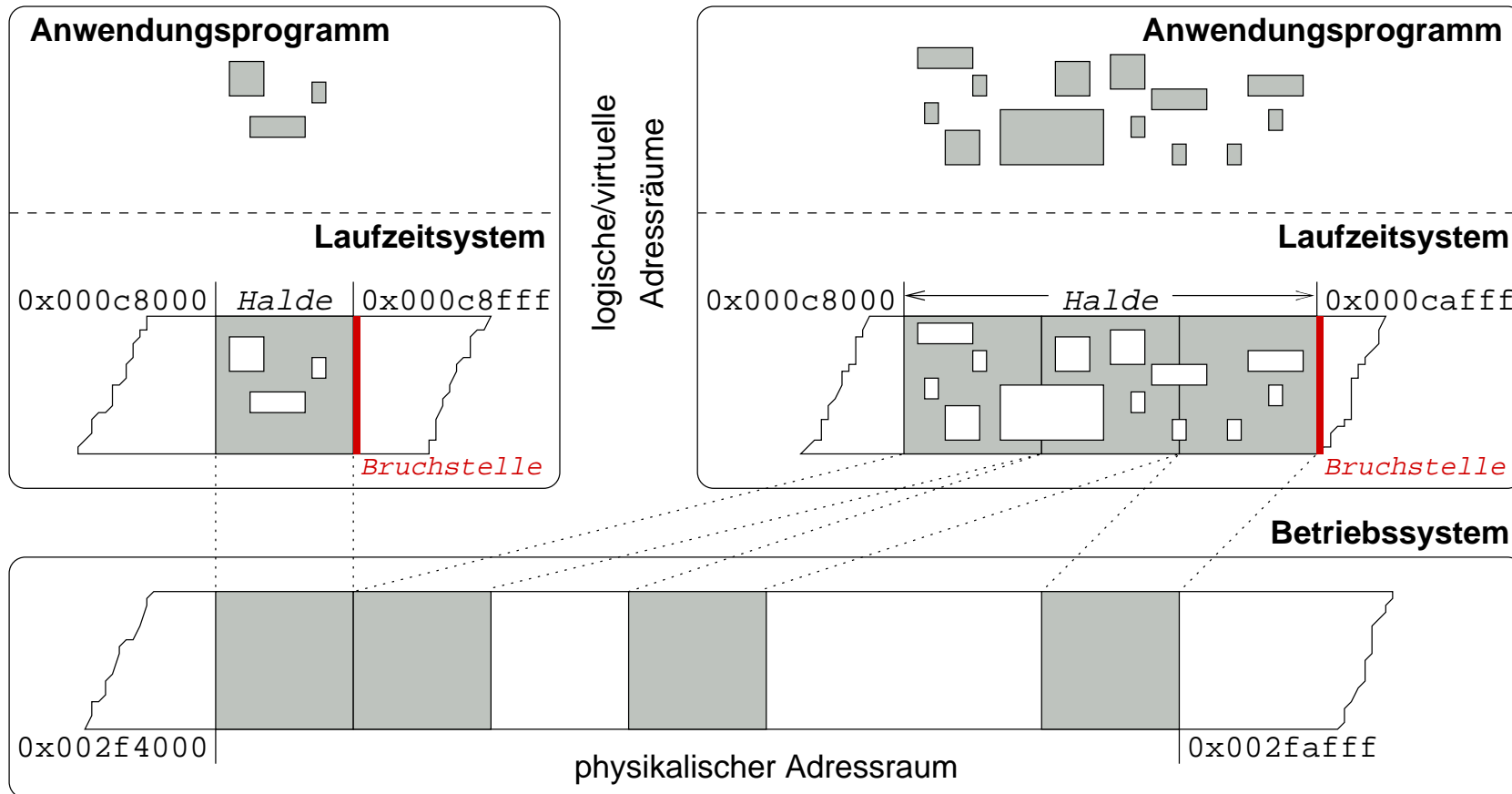
- die über die *Rechnerperipherie* (E/A-Geräte) angeschlossenen Bereiche
- dient der Datenablage und Implementierung virtueller Adressräume
- ist größer als der physikalische Adressraum: *Petabytes* (2^{50} bzw. 10^{15})

- die Differenzierung kann „nach aussen“ unsichtbar sein (☞ Multics [8])

Speicherverwaltung

- die Verwaltung von *Hauptspeicher* erfolgt typischerweise auf zwei Ebenen:
 - das **Laufzeitsystem** bzw. die Bibliotheksebene verwaltet den Speicher (lokal) innerhalb eines logischen/virtuellen Adressraums
 - Speicherblöcke können von sehr feinkörniger Struktur/Größe sein
 - Bedürfnisse von Programmiersprachen bestimmen die Verfahrensweisen
 - das **Betriebssystem** verwaltet den im physikalischen Adressraum (global) vorrätigen Speicher
 - Speicherblöcke sind üblicherweise von grobkörniger Struktur/Größe
 - Arten des Rechnerbetriebs üben starken Einfluss auf Verfahrensweisen aus
- „*separation of concerns*“ [13] — beide Ebenen/Systeme ergänzen sich einander

Synergie bei der Speicherverwaltung



`ptr = malloc (size)` legt einen (ausgerichteten) Speicherblock an

- Resultat ist die Anfangsadresse eines Blocks von mindestens size Bytes

`ptr = realloc (addr, size)` verkleinert/vergrößert den Speicherblock

- bei Verkleinerung steht der Rest ggf. der Wiedervergabe zur Verfügung

`free (ptr)` macht den Speicherblock zur Wiedervergabe verfügbar

- der Speicherbereich wird nicht ans Betriebssystem zurück gegeben !

Die Festlegung einer neuen „Bruchstelle“²⁵ für das Datensegment eines Prozesses bewirkt die Veränderung der diesem Segment zugeordneten Speichermenge. Die Bruchstelle kann dabei eine vom System vorgegebene Größe nicht überschreiten.

`addr = brk (brkval)` setzt die Bruchstelle auf den angegebenen Wert

`addr = sbrk (incr)` addiert den angegebenen Wert zur Bruchstelle

- beim negativen Summanden (`incr`) wird das Datensegment verkleinert
- der vor der Veränderung gültige Wert der Bruchstelle wird zurück geliefert

☞ `getrlimit(2), mmap(2)`

²⁵Der „*break value*“, d. h. die der gegenwärtigen Endadresse des Datensegments folgende Speicheradresse.

Datei

Da'tei allgemein eine Sammlung von Daten; im Kontext von Rechensystemen:

- ☞ eine *zusammenhängende, abgeschlossene Einheit von Daten*
- ☞ eine „beliebige“ Anzahl eindimensional adressierter Bytes
- die „Datei“ als Synonym für anhaltende (langfristige) Datenspeicherung
 - *persistente Speicherung* der Daten ist allerdings nicht zwingend:
 - ☞ „*RAM Disk*“, „*Rohr*“ (*pipe*), „*raw device*“
- nicht-flüchtige, blockorientierte Speichermedien (Platten, Magnetbänder) kommen bevorzugt zum Einsatz



Arten von Dateien

ausführbare Dateien Binär- und Skriptprogramme

- Dateien, die von einem Prozessor ausführbaren *Programmtext* enthalten

Binärprogramm	☞	CPU, FPU, MCU, JVM, . . . , Basic, Lisp, Prolog
Skriptprogramm	☞	perl(1), python(1), {a,ba,c,tc}sh(1), tc1(n)

- der Prozessor kann in Hard-, Firm- und/oder Software realisiert sein

nicht-ausführbare Dateien Text-, Bild- und Tondateien

- Dateien, die von einem Prozessor verarbeitbare *Programmdaten* enthalten

☞ `.{doc, fig, gif, jpg, mp3, pdf, tex, txt, wav, xls, . . . }`

☞ `.{a, c, cc, f, F, h, l, o, p, r, s, S, y, . . . }`

- der Prozessor liegt als Programmtext (zumeist als ausführbare Datei) vor



Bezeichnung von Dateien

- „von aussen“ ist eine Datei über eine **symbolische Adresse** bekannt/erreichbar
 - ein *benutzerdefinierter Name* von beliebiger aber maximaler Länge
 - dieser *Dateiname* ist für das Betriebssystem allgemein „Schall und Rauch“
 - je nach System (☞ Windows) erfolgt jedoch eine Interpretation des *Suffix*
- „nach innen“ besitzt die Datei eine logische bzw. **numerische Adresse**
 - diese Adresse identifiziert den die Datei beschreibenden „*Dateikopf*“
 - zwischen Dateiname und -kopf besteht eine *1:1*- oder *N:1*-Beziehung
 - je nach System (☞ UNIX) hat ein Dateikopf daher ggf. mehrere Namen
- symbolische und numerische Dateiadresse werden als Paar „verzeichnet“

Dateierweiterung

Dateinamensuffix oder *Extension*: eine üblicherweise durch einen Punkt vom Dateinamen abgegrenzte *symbolische Erweiterung* des Dateinamens

- liefert einen Hinweis auf das Dateiformat bzw. den Dateitypen

	$\left\{ \begin{array}{l} .doc \\ .fm \\ .tex \end{array} \right\}$	Textdokumente	$\left\{ \begin{array}{ll} \text{MS-Word} & \\ \text{Framemaker} & \text{maker(1)} \\ \text{\LaTeX} & \text{latex(1)} \end{array} \right\}$
	$\left\{ \begin{array}{l} .h \\ .c \\ .s \\ .o \end{array} \right\}$	Programme	$\left\{ \begin{array}{ll} \text{Präprozessor} & \text{cpp(1)} \\ \text{Kompilierer} & \text{cc(1)} \\ \text{Assemblerer} & \text{as(1)} \\ \text{Binder} & \text{ld(1)} \end{array} \right\}$

- ist anwendungsspezifisch ausgelegt und ggf. dem Betriebssystem bekannt

Dateikopf

Indexknoten (*index node*, *inode*) im UNIX-Jargon, enthält *Dateiattribute*:

- ☞ Eigentümer (*user ID*)
- ☞ Gruppenzugehörigkeit (*group ID*)
- ☞ Typ (reguläre/spezielle Datei)
- ☞ Zugriffsrechte (lesen, schreiben, ausführen; Eigentümer, Gruppe, „Welt“)
- ☞ Zeitstempel (letzter Zugriff, letzte [*mode*-]²⁶ Änderung)
- ☞ Anzahl der Verweise („*hard links*“)
- ☞ Größe (in Bytes)
- ☞ Adresse(n) der Daten auf dem Speichermedium

Indexknotennummer (*inode number*): die **numerische Adresse** der Datei

²⁶Typ und Zugriffsrechte.

Dateityp

reguläre Datei (*regular file, ordinary file*) eindimensionales Bytefeld

Vom System vordefinierte, spezielle Dateien (*special files*):

Verzeichnis (*directory*) Katalog von regulären und/oder speziellen Dateien

Gerätedatei (*device file*) Zugang zu zeichen-/blockorientierten Geräte(treiber)n

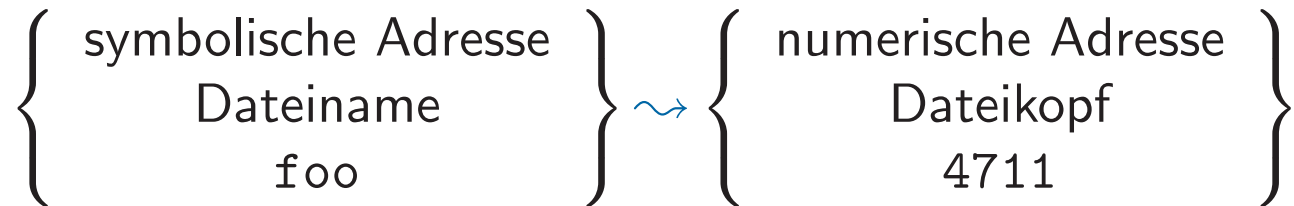
„benanntes Rohr“ (*named pipe*) zwischen unverwandten lokalen Prozessen

„Sockel“ (*socket*) Kanal zwischen lokalen/entfernten Prozessen

Dateiverzeichnis

das **Verzeichnis** (*directory*)²⁷ dient der Gruppierung von Dateinamen

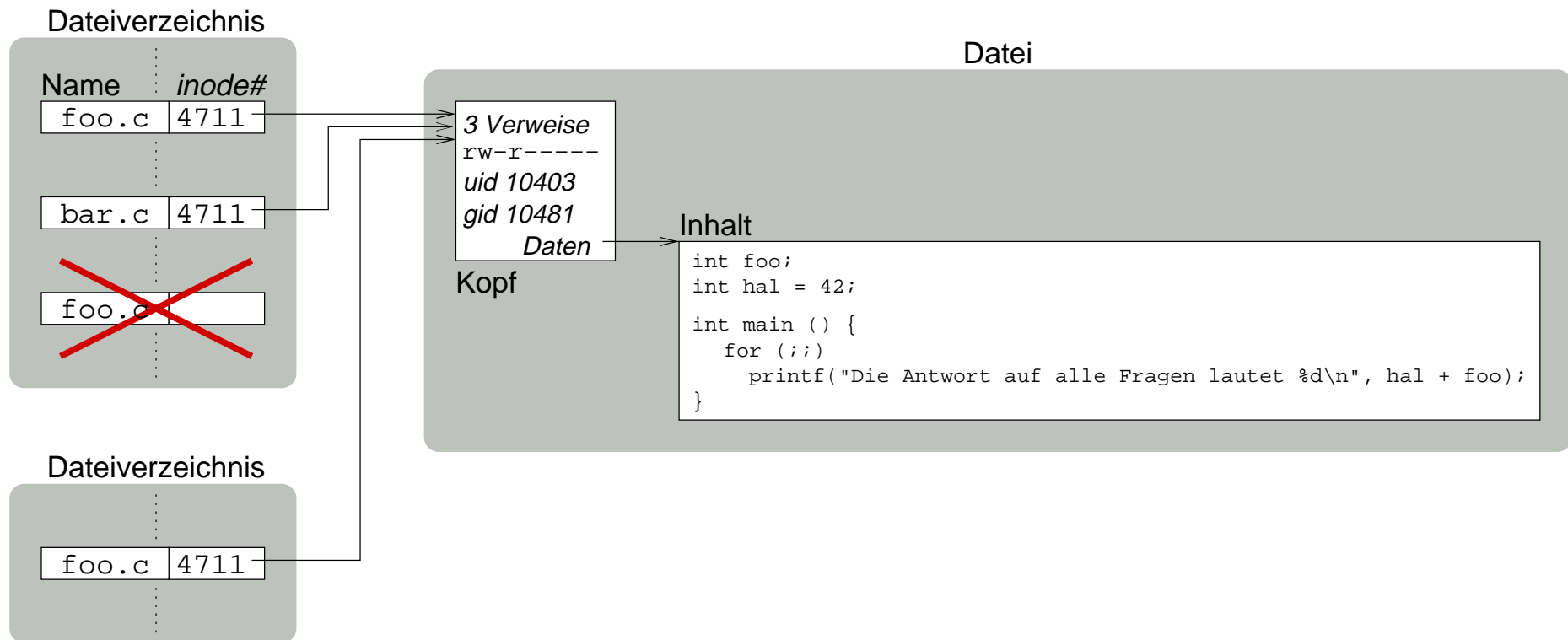
- definiert einen gemeinsamen *Kontext* für die symbolischen Adressen
☞ symbolische Adressen sind nur innerhalb ihrer Kontexte eindeutig
- implementiert eine „Umsetzungstabelle“ für symbolische Adressen:



- führt Buch über die Beziehung zwischen Dateinamen und Dateikopf

²⁷Auch als „Katalog“ (*catalogue*) bezeichnet.

Dateiname vs. Dateikopf vs. Datei



Blockorientierte Datenspeicherung

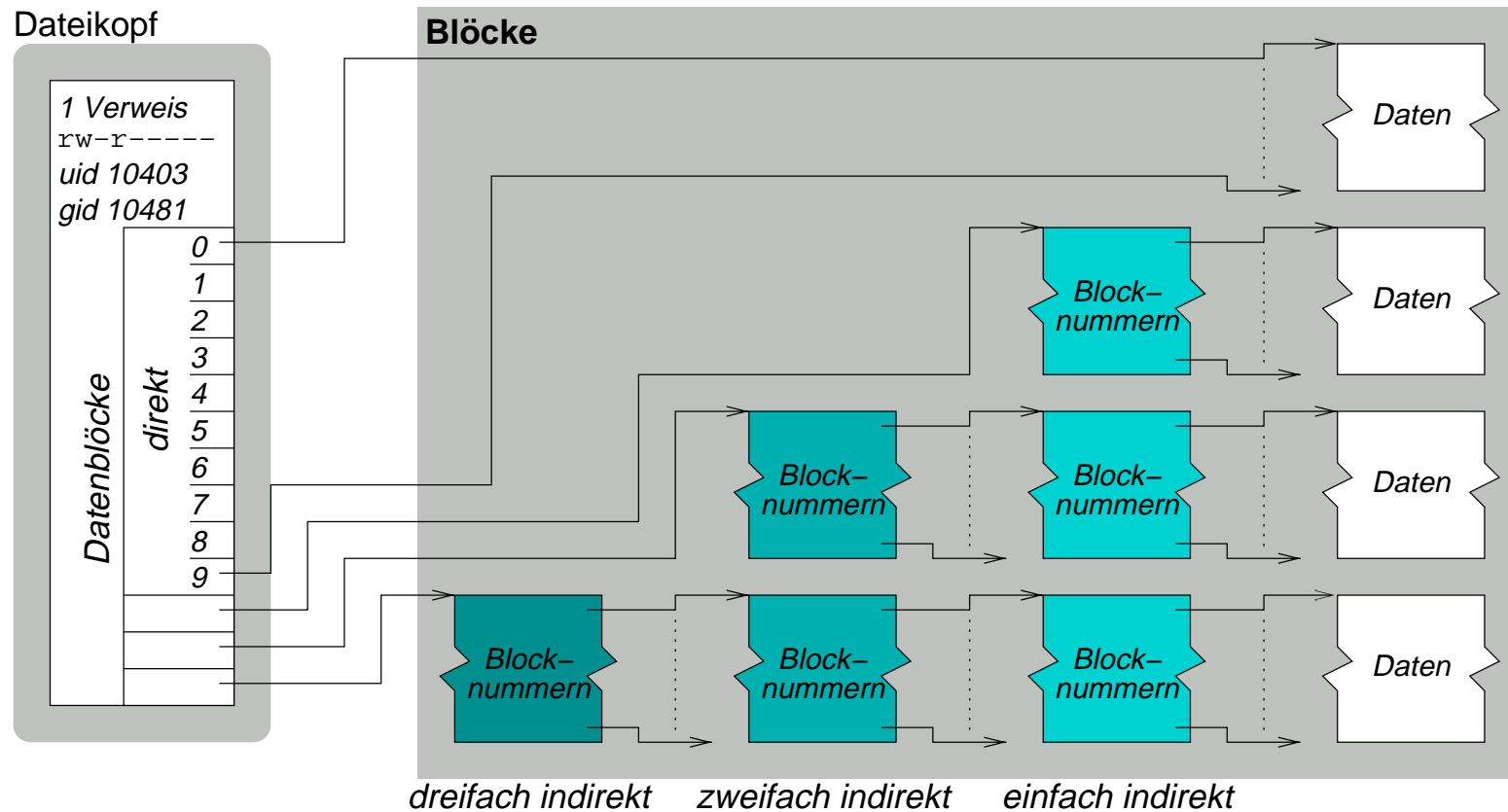
- Dateiinhalte sind permanent gespeichert in Form von Blöcken fester Größe
 - typische Blockgrößen: 512 Bytes (B) – 8 Kilobytes (KB) [14]
 - in Abhängigkeit von $\left\{ \begin{array}{l} \text{den physikalischen Parametern des Massenspeichers} \\ \text{der durchschnittlichen Dateilänge} \\ \text{dem Anwendungsprofil} \end{array} \right.$
- der Dateikopf führt Buch über die „Nummern“ der Datenblöcke der Datei
 - die Blocknummer entspricht einer physikalischen Adresse im Massenspeicher
 - die max. Dateilänge hängt ab von Anzahl/Wertebereich der Blocknummern
- Art und Weise der Blocknummernverwaltung ist abhängig von Zugriffsmustern

Zugriffsmuster

- Blockzugriffe seitens der Anwendung verlaufen sequentiell oder wahlfrei:
 - sequentieller Zugriff** Dateizugriffe erfolgen nach wohlgeordnetem Muster
 - die Datei wird von vorne nach hinten gelesen/beschrieben
 - geeignet für sequentiell organisierte Massenspeicher (z. B. Magnetbänder)
 - wahlfreier Zugriff** Dateizugriffe erfolgen nach beliebigen Muster
 - ein wohlgeordnetes Zugriffsmuster ist nicht erkennbar
 - geeignet für wahlfrei organisierte Massenspeicher (z. B. Festplatte)
- entsprechend verfolgt die Blockverwaltung verschiedene Optimierungskriterien
 - Kompromisslösungen bilden jedoch das „Tagesgeschäft“ (☞ UNIX)

Datenblockadressierung

$$\text{sizeof}(\text{file}) \leq \text{sizeof}(\text{block}) \times (10 + N + N^2 + N^3)$$



Namensraum — *Namespace*

- bildet einen *Kontext*, in dem jeder Name eine eindeutige Bedeutung besitzt

„Java“ bedeutet im Kontext	{	„Geographie“	eine Insel
		„Genussmittel“	ein Heissgetränk
		„Informatik“	eine Programmiersprache
„C“ bedeutet im Kontext	{	„Sprache“	einen Buchstaben
		„Musik“	eine Note
		„Informatik“	eine Programmiersprache

- die Struktur/Organisation ist flach oder (meist) hierarchisch ausgelegt

Aufbau von Namensräumen

flache Struktur definiert nur einen einzigen Kontext

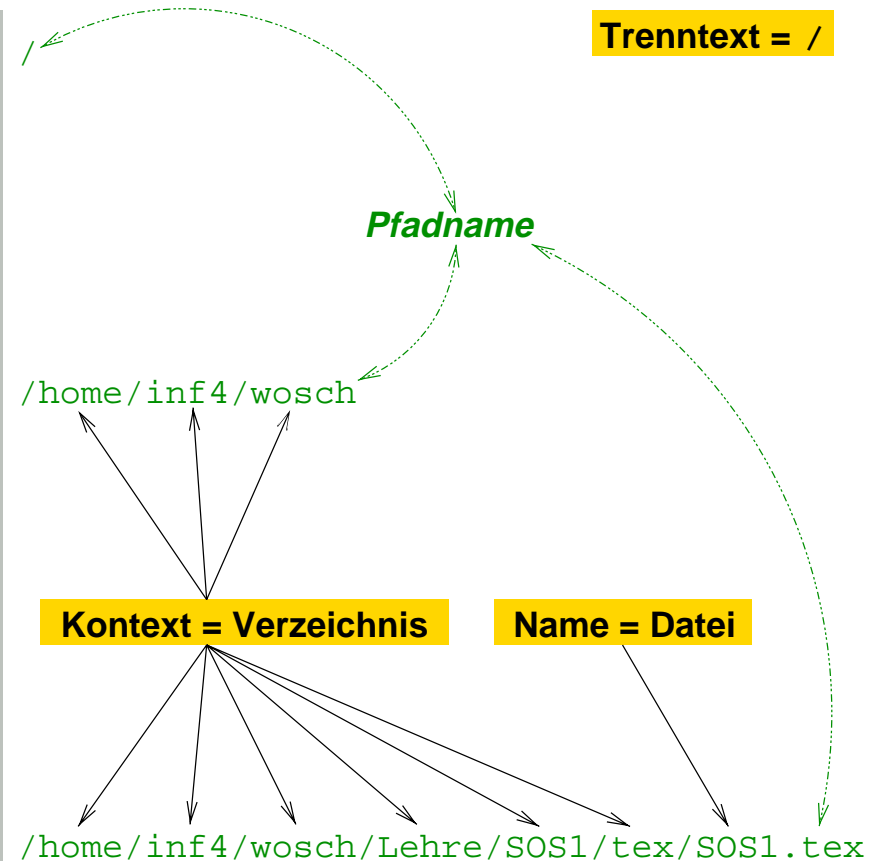
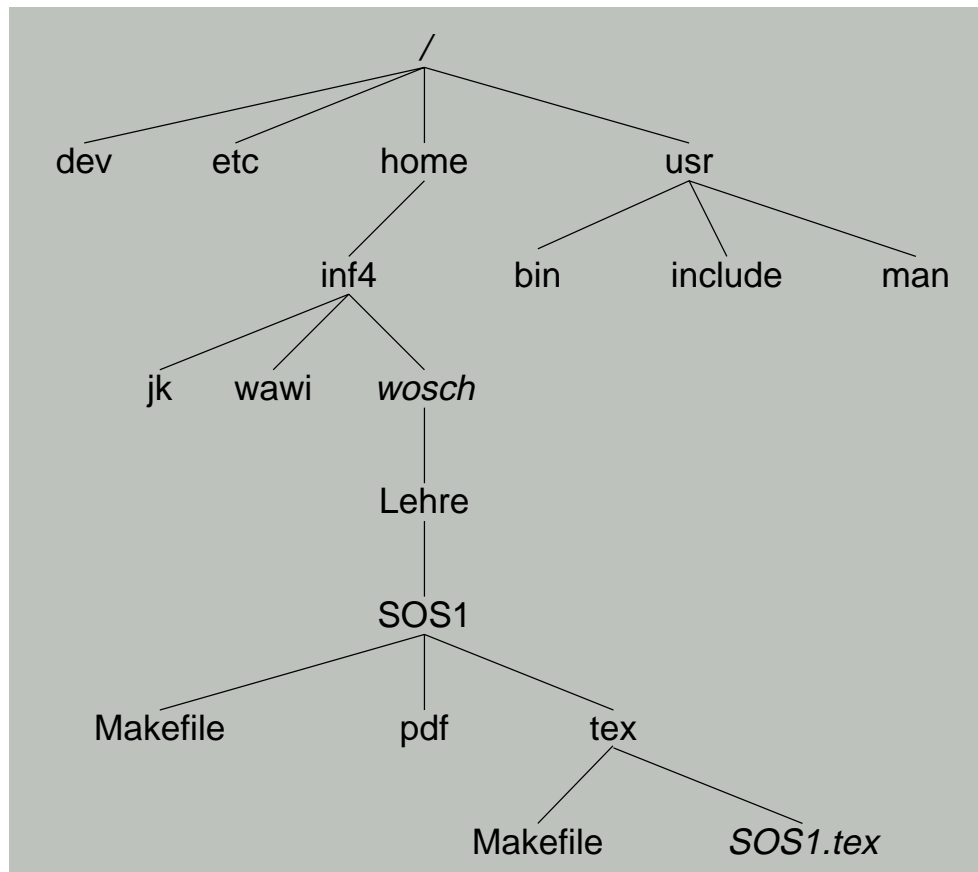
- Eindeutigkeit muss mit der Namenswahl selbst gewährleistet werden

hierarchische Struktur definiert mehrere Kontexte

- Eindeutigkeit wird über einen *Kontextnamen* als Präfix erreicht
- als *Separatoren* werden meist Sonderzeichen („Trenntext“) verwendet:

/	Schrägstrich (<i>slash</i>)	UNIX	☞ foo/bar
\	zurückgelehnter Schrägstrich (<i>backslash</i>)	Windows	☞ foo\bar

Hierarchischer Namensraum — Dateibaum (*file tree*)



Sonderverzeichnisse

Wurzelverzeichnis (*root directory*)

- die Wurzel (bei UNIX bezeichnet mit '/') im Dateibaum
- wird vom System gesetzt (**privilegierte Operation** ➡ `chroot(2)`)

Arbeitsverzeichnis (*working directory*)

- die gegenwärtige Position eines Programms/Prozesses im Dateibaum
- ändert sich beim „Durchklettern“ (➡ `chdir(2)`) des Dateibaums

Heimatverzeichnis (*home directory*)


- das initiale Arbeitsverzeichnis eines Benutzers/Prozesses
- wird vom System gesetzt bei Sitzungsbeginn (*login session* ➡ `login(1)`)

„Spitznamen“ für Verzeichnisse

. (*dot*) das aktuelle *Arbeitsverzeichnis*

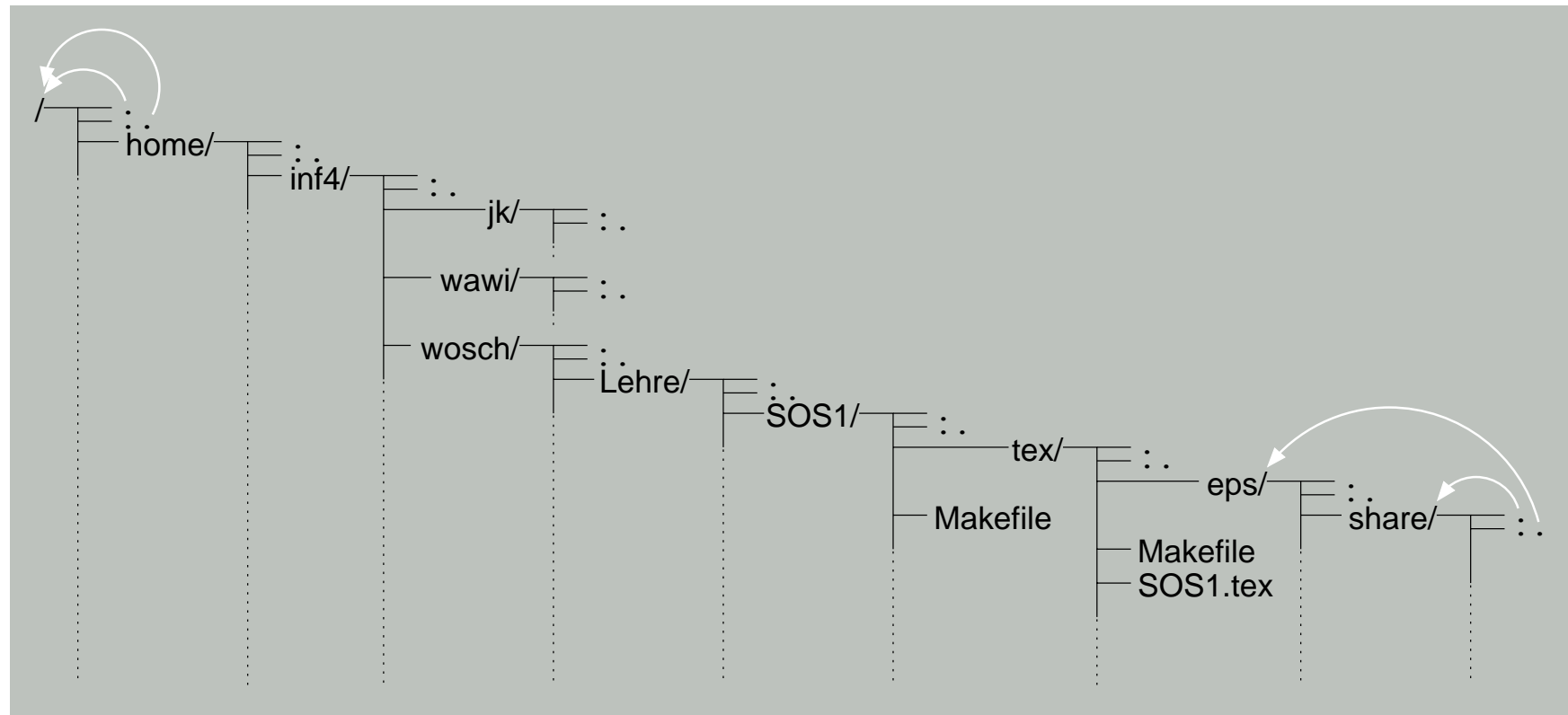
- zur relativen Adressierung des jeweiligen Arbeitsverzeichnisses²⁸
- 1. Eintrag in jedem Verzeichnis (UNIX  `mkdir(2)`)

.. (*dot dot*) das aktuelle *Elternverzeichnis*

- zur relativen Adressierung des jeweils übergeordneten Verzeichnisses
 - entspricht '.', falls es kein Elternverzeichnis gibt (Wurzelverzeichnis)
- 2. Eintrag in jedem Verzeichnis (UNIX  `mkdir(2)`)

²⁸Ermöglicht beispielsweise die Implementierung von `ls(1)`, um die Einträge eines Arbeitsverzeichnisses auflisten zu können, ohne dessen wirklichen Namen kennen zu müssen. Ferner kann durch den Präfix „./“ zum Namen einer ausführbaren Datei der Aufruf lokaler Kommandos/Programme durchgesetzt werden.

Dynamische Datenstruktur „Dateibaum“



„Navigation“ im Namensraum

- im hierarchischen Namensraum beschreiben „Pfade“ die Wege zu Datei(nam)en
 - ein **Pfadname** (*pathname*) ist ein vollständiger Dateiname




- formaler Aufbau eines (UNIX) Datei- bzw. Pfadnamens in EBNF[15]:

```
pathname      = resolver | [resolver], {name, resolver}, name;  
resolver      = {separator}-;  
separator     = "/";  
name          = {character}-;  
character     = character set - separator;  
character set = ASCII;
```

- Beispiele: /, ., .., foo, foo/bar, /foo, bar/, ./bar/.., ../foo/./bar//

Pfadname

relativer \sim vom **Arbeitsverzeichnis** (*working directory*) ausgehend,
z. B. von `/home/inf4/wosch` aus:

 `Lehre/S0S1/S0S1.tex`
 `./Lehre/S0S1/S0S1.tex`
 `../wosch/Lehre/S0S1/S0S1.tex`

oder von `/home/inf4/jk` aus:

 `../wosch/Lehre/S0S1/S0S1.tex`

absoluter \sim vom **Wurzelverzeichnis** (*root directory*) ausgehend:

 `/home/inf4/wosch/Lehre/S0S1/S0S1.tex`

Dateisystem

- Komplex von Datenstrukturen zur Verwaltung von Dateien und Verzeichnissen

Dateisystemkopf (*super block*) zur Beschreibung des Dateisystems

Dateikopftabelle (*inode table*) zur Beschreibung von Dateien/Verzeichnisse

Datenblöcke zur Speicherung der Inhalte der Dateien/Verzeichnisse

- eine (heterogene) dynamische Datenstruktur beschränkten Ausmaßes
 - der Dateisystemkopf legt die jeweiligen „Grenzwerte“ fest
- wird auf genau eine **Partition**²⁹ (z. B. einer Festplatte) abgebildet

²⁹Partitionierung: die logische Unterteilung des Hintergrundspeichers in einen Satz zusammenhängender Sektoren.

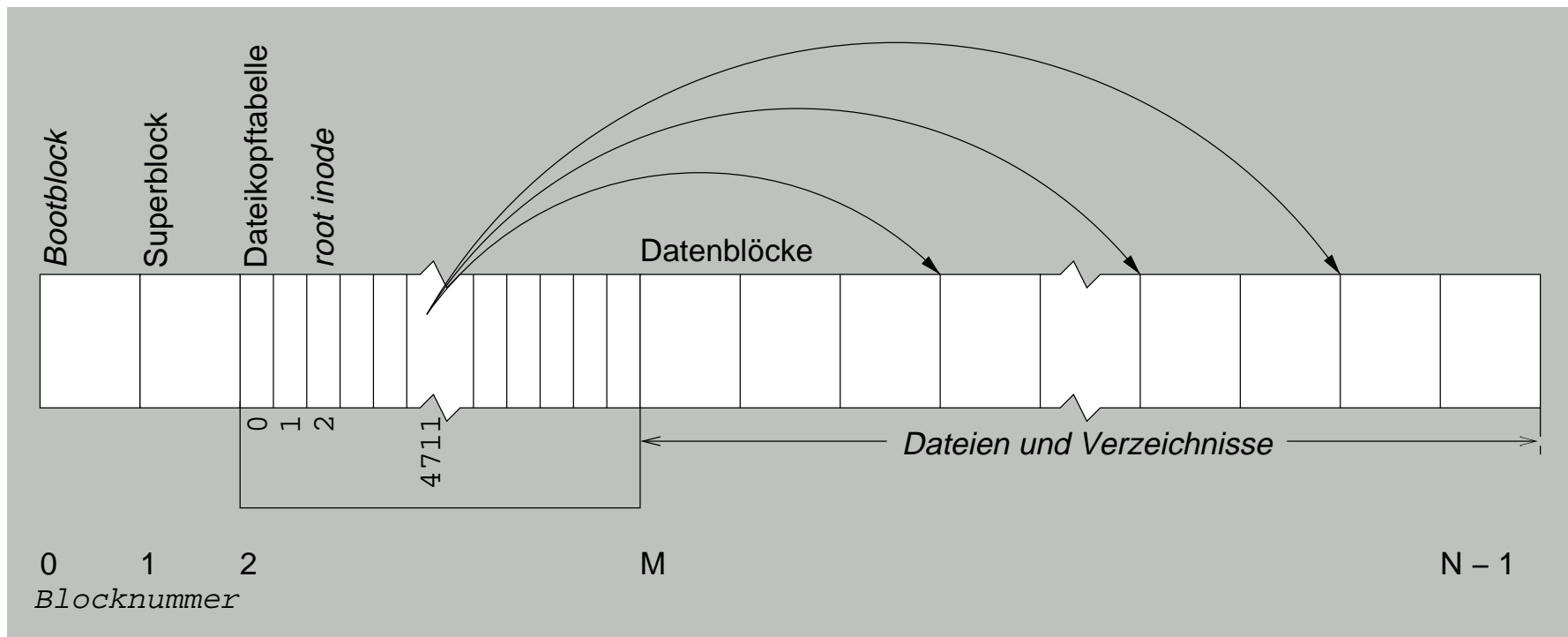
Dateisystemkopf

Superblock (*UNIX File System*, S5FS) logische Blocknummer 1

- enthält Informationen, die das Dateisystem beschreiben und dimensionieren:
 - Name des Wurzelverzeichnisses des Dateisystems
 - Anzahl der dem Dateisystem zugeordneten Blöcke
 - Größe der Dateikopftabelle (d. h., die Anzahl von *inodes*)
 - Anzahl und Liste freier Dateiköpfe bzw. Datenblöcke
- zum „Montieren“ von Dateisystemen benötigte (statische) Systemparameter
 - dynamisches Einbinden von Dateisystemen in (ein) bereits bestehende(s)

☞ Zerstörung des Superblocks z. B. durch einen Absturz (*crash*) der Festplatte hat „unangenehme“ Auswirkungen auf das zugehörige Dateisystem

Dateisystemformat — Partition



Namensauflösung

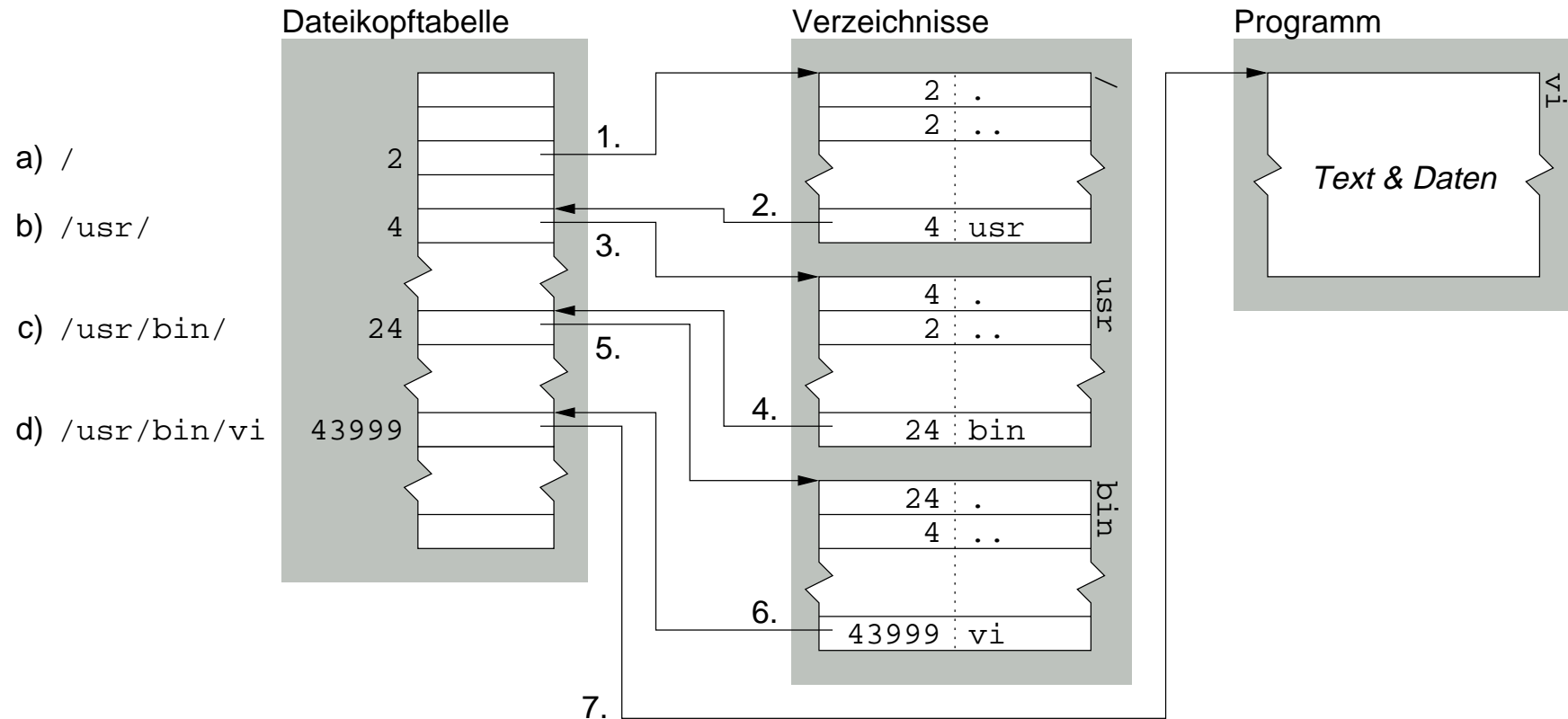
Abbildung (symb. Adresse \leadsto num. Adresse) `creat(2)/link(2)`

- ☞ Bindungsfunktion: den Dateinamen $\left\{ \begin{array}{l} \text{mit einem Dateikopf verknüpfen} \\ \text{in ein Dateiverzeichnis eintragen} \end{array} \right.$
- ☞ einen Pfadnamen mit einem Dateikopf (*inode*) assoziieren

Umsetzung (symb. Adresse \leadsto num. Adresse) `open(2)`

- ☞ Auflösungsfunktion: einen Pfadnamen interpretieren
- ☞ Dateiverzeichnisse nach Namenseinträgen durchsuchen
 - ✓ schrittweise für jeden einzelnen Verzeichnisnamen im Pfad
 - ✓ schließlich für den Dateinamen

Auflösung von /usr/bin/vi



S5FS „*Considered Harmful*“

- Dateiköpfe liegen vorne auf der Platte, Dateien (inkl. Verzeichnisse) hinten
 - zeitaufwändige Namensauflösung, einhergehend mit „kostspieliger“ E/A
 - nicht jede Dateiänderung ist durchgängig in Bezug auf Hintergrundspeicher
- mehrere „*single point of failure*“³⁰ wirken Zuverlässigkeit/Robustheit entgegen
 - der Dateikopf (*inode*) — enthält alle Attribute einer Datei
 - die Dateikopftabelle (*inode table*) — enthält alle Dateiköpfe
 - der Superblock — enthält die Parameter des gesamten Dateisystems
- Stromausfall oder „*head crash*“ können sich äußerst unangenehm auswirken

³⁰Eine beliebige Komponente eines Systems, die im Fehlerfall den Ausfall des Gesamtsystems nach sich zieht.

Dateisystemmontage

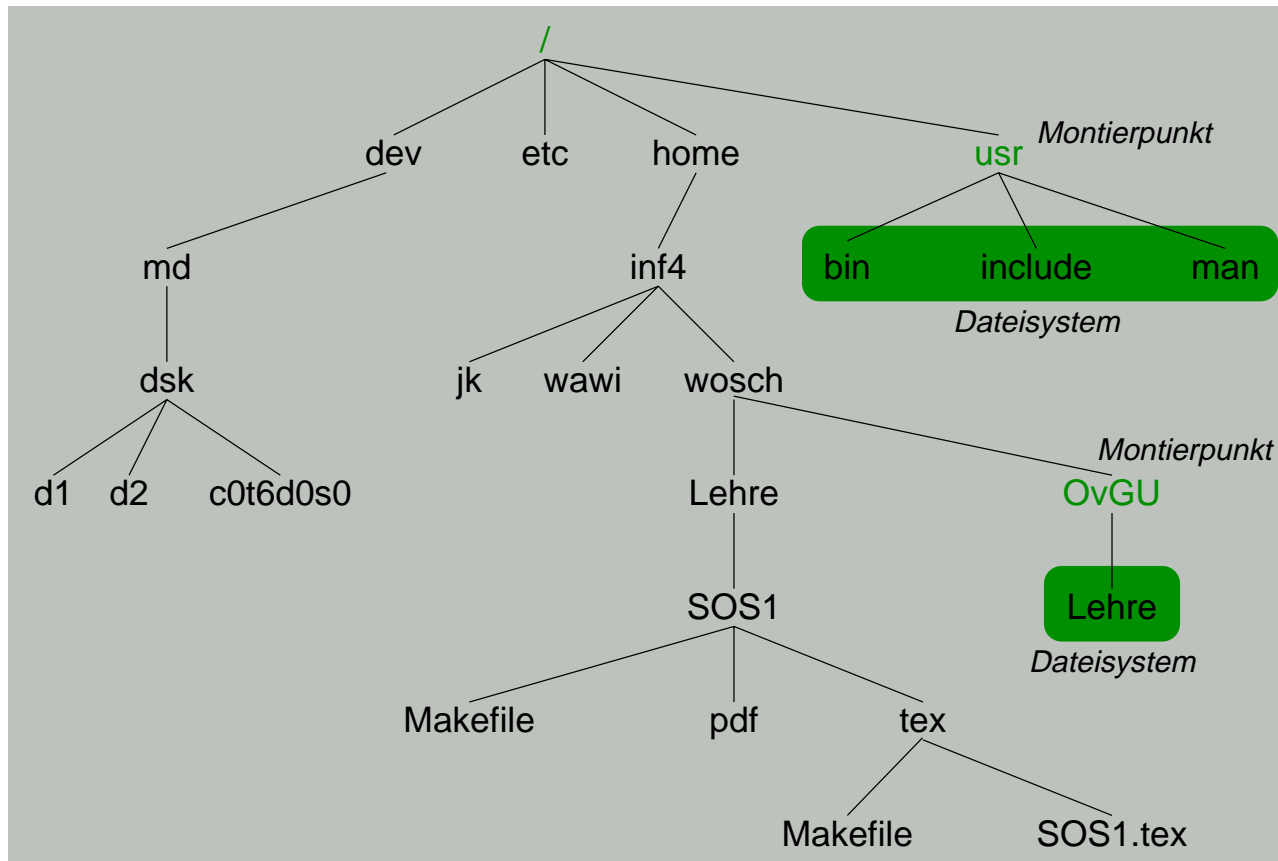
Montierpunkt (*mount point*) eine Stelle im „Wirtdateisystem“, an der die Einbindung eines anderen Dateisystems möglich ist

☞ ein Verzeichnis wird mit der Wurzel eines anderen Dateisystems überlagert

- in Bezug auf das Wirtdateisystem sind montierbare Dateisysteme . . .
 - von gleicher/verschiedener Struktur (z. B. S5FS, UFS, FFS, EXT2, NTFS)
 - auf demselben/einem anderen Gerät (ggf. einem anderen Rechner) resident
- sie bilden eine eigene *Partition*, definiert über eigene Systemparameter

Dateisystemmontage

~wosch/OvGU, /usr



/dev/md/dsk/d1

/dev/md/dsk/d2

mount /dev/md/dsk/d2 /usr

/dev/md/dsk/c0t6d0s0

Dateisystemadressraum

- zwei Sorten von Adressen kommen im (S5FS) Dateisystem zur Verwendung:
 - Blockadressen** (Blocknummern) verweisen auf die für die Speicherung der Informationen zur Verfügung stehenden Datenblöcke
 - ☞ Dateisysteme sind auf *blockorientierte Massenspeicher* ausgerichtet
 - Dateikopfadressen** (Indexknotennummern, *inode numbers*) identifizieren die Strukturen zur Verwaltung von Dateien/Verzeichnissen
- Adressen sind nur gültig (bzw. eindeutig) innerhalb ihres Adressraums
 - Dateiköpfe montierter Dateisysteme sind nicht direkt adressierbar
 - nur über Pfadnamen können Dateien/Verzeichnisse global erreicht werden
- Dateikopfadressen sind kontextabhängig zu betrachten und zu behandeln

Dateikopferenzen

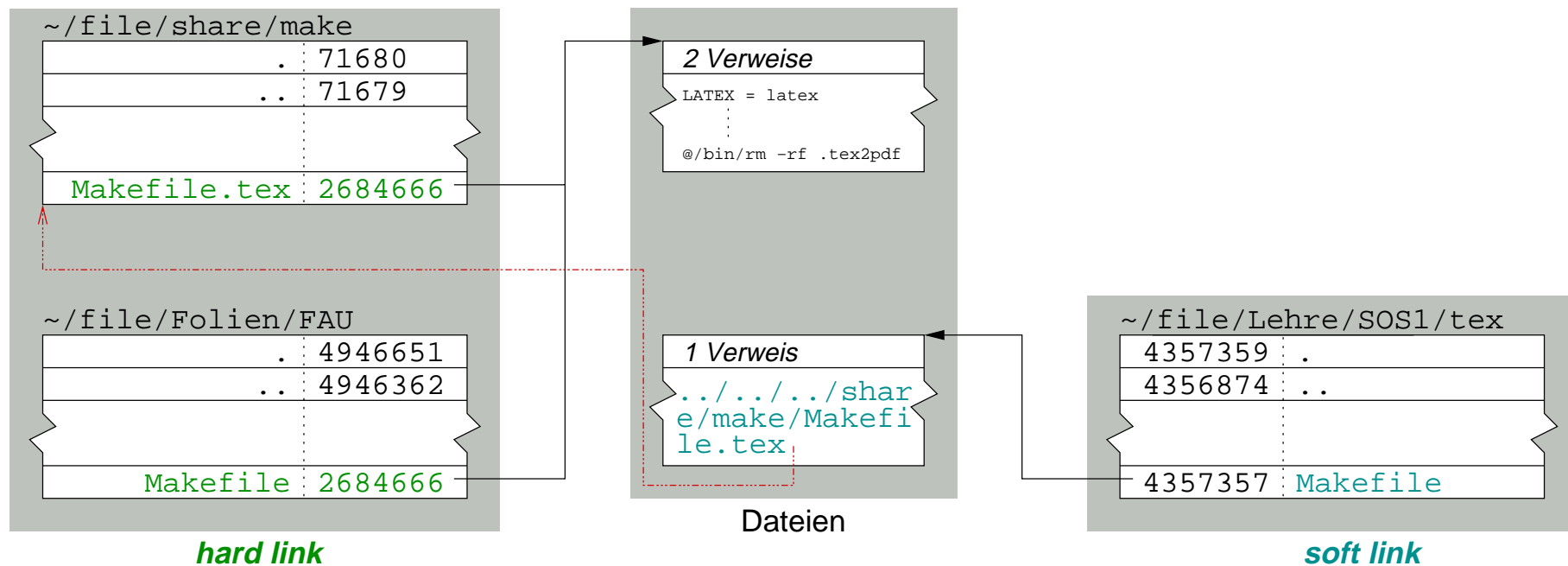
harte **Referenz** (*hard link, link*)

- eine Verknüpfung zwischen Dateiname und Dateikopf (*inode*)
 - ist nur lokal (innerhalb eines Dateisystems) definiert
 - als Paar gespeichert im Verzeichnis (Namenseintrag)
- eine „direkte Adresse“ auf eine Datei

symbolische **Referenz** (*symbolic link, soft link*)

- eine Verknüpfung zwischen Dateiname und Pfadname
 - kann global (d. h, Dateisystem übergreifend) definiert sein
 - gespeichert im Verzeichnis oder als (vom System definierte) Datei
- eine „indirekte Adresse“ auf eine Datei

„Hard Link“ vs. „Soft Link“



Symbolische Referenz „*Considered Harmful*“

```
wosch@lorien 1$ mkdir -p Laptop/faui43w; cd Laptop
wosch@lorien 2$ ln -s faui43w lorien
wosch@lorien 3$ ls -l
total 8
drwxr-xr-x  2 wosch  wosch  68 29 Apr 13:01 faui43w
lrwxr-xr-x  1 wosch  wosch   7 29 Apr 13:02 lorien -> faui43w
wosch@lorien 4$ cd lorien
wosch@lorien 5$ cd ..; rmdir faui43w; cd lorien
-bash: cd: lorien: No such file or directory
wosch@lorien 6$ ls -l
total 8
lrwxr-xr-x  1 wosch  wosch   7 29 Apr 13:02 lorien -> faui43w
wosch@lorien 7$ mkdir faui43w; cd lorien
wosch@lorien 8$ ln -s "gibt es nicht" jippdit
```

`fd = open (path, flags, mode)` erzeugt einen Deskriptor für die mit dem Pfadnamen bezeichnete Datei (öffnet die Datei)³¹

- der *Dateizeiger* (*file pointer*) wird auf den Anfang (0) der Datei gesetzt

`ok = close (fd)` invalidiert einen Dateideskriptor (schließt die Datei)

`newfd = dup (oldfd)` dupliziert einen Dateideskriptor (oldfd)

³¹Der mode Parameter ist erforderlich beim Erzeugen einer Datei, d. h., wenn der Status „O_CREAT“ vorgibt.

`ok = link (path1, path2)` erzeugt eine „harte Referenz“ (*hard link*)

- der Referenzzähler im (direkt adressierten) *Inode* wird um 1 erhöht
 - path1 ist der „ursprüngliche“ Pfadname dieser Datei
 - path2 ist der „alternative“ Pfadname dazu
- ursprünglicher/alternativer Pfadname sind (danach) ununterscheidbar

`ok = unlink (path)` löscht eine „harte Referenz“ (*hard link*)

- der Referenzzähler im (direkt adressierten) *Inode* wird um 1 erniedrigt
- ist der Zählerwert 0, wird die Datei endgültig gelöscht

`ok = symlink (path1, path2)` erzeugt eine „weiche Referenz“ (*soft link*)

- der Referenzzähler im (indirekt adressierten) *Inode* bleibt unverändert
 - vielmehr wird ein neuer *Inode* (für *path2*) angelegt
 - der Referenzzähler dieses *Inodes* wird mit 1 initialisiert
- auf eine (bekannte) Datei wird eine *symbolische Referenz* angelegt
 - löschen einer solchen Referenz, löscht nicht die Datei
 - löschen der Datei, erzeugt eine „verwaiste symbolische Referenz“
- erlaubt das Setzen einer „indirekten Adresse“ auf *alle* Arten von Dateien
 - wie z. B. Verzeichnisse oder etwa Einträge anderer Dateisysteme
 - eine solche Adresse ist zusätzlich „hart“ referenzierbar (`link(2)`)

Datei

UNIX Systemaufrufe (6)

`nr = read (fd, buf, nbytes)` Datei \Rightarrow Puffer

- der Dateizeiger wird um die Anzahl der gelesenen Zeichen weiter gesetzt

`nw = write (fd, buf, nbytes)` Puffer \Rightarrow Datei

- der Dateizeiger wird um die Anzahl der geschriebenen Zeichen weiter gesetzt

`rwp = lseek (fd, offset, whence)` positioniert den Dateizeiger

`ok = fcntl (fd, cmd, arg)` verwaltet eine (geöffnete) Datei

- der optionale Parameter `arg` hängt ab vom Kontrollkommando `cmd`
 - duplizieren von Dateideskriptoren, manipulieren von Dateiattributen
 - freiwilliges Sperrverfahren (*advisory record locking*) zur Koordination
- die geöffnete Datei wird durch einen Dateideskriptor (`fd`) identifiziert

`ok = ioctl (fd, request, argp)` verwaltet ein Gerät oder einen Strom

- die abgesetzten Kommandos lösen gerätespezifische Steuerfunktionen aus
- der Gerätetreiber interpretiert den Befehl samt optionale Parameter `argp`

Prozess . . .

. . . wird durch ein Programm kontrolliert und benötigt zur Ausführung dieses Programms einen Prozessor.

Habermann, *Introduction to Operating System Design*

. . . P ist ein Tripel (S, f, s) , wobei S einen Zustandsraum, f eine Aktionsfunktion und $s \subset S$ die Anfangszustände des Prozesses P bezeichnen. Ein Prozess erzeugt Abläufe, die durch die Aktionsfunktion generiert werden können.

Horning/Randell, *Process Structuring*

. . . ist das Aktivitätszentrum innerhalb einer Folge von Elementaroperationen. Damit wird ein Prozess zu einer abstrakten Einheit, die sich durch die Instruktionen eines abstrakten Programms bewegt, wenn dieses auf einem Rechner ausgeführt wird.

Dennis/van Horn, *Programming Semantics for Multiprogrammed Computations*

. . . ist ein Programm in Ausführung.

unbekannte Referenz, „Mundart“

Prozess \neq Programm (1)

- ein Prozess kann die Ausführung mehrerer Programme zur Folge haben:

Systemaufruf *Anwendungsprogramm* ruft ein *Betriebssystemprogramm* auf
☞ der Prozess ist der *Aktivitätsträger* von mindestens zwei Programmen

- umgekehrt kann ein Programm von mehreren Prozessen ausgeführt werden

nicht-sequentielles Programm im Falle von Uniprozessorsystemen

☞ „präemptive Programmverarbeitung“³²

☞ Faden (*thread*), Unterbrechung (*interrupt*)

paralleles Programm im Falle von Multiprozessorsystemen

³²Die Ausführung einer Aufgabe kann jederzeit von einer höheren Instanz unterbrochen werden.

Prozess \neq Programm (2)

☞ Wissen über das gegenwärtig ausgeführte Programm sagt nicht viel aus über die zu dem Zeitpunkt im System stattfindende Aktivität.

- Welche bzw. wieviel $\left\{ \begin{array}{l} \text{Fäden führen das Programm zur Zeit aus} \\ \vdots \\ \text{Programmunterbrechungen sind zur Zeit aktiv} \end{array} \right\} ?$

☞ Im Betriebssystemkontext ist das Konzept „Prozess“ daher nützlicher als das Konzept „Programm“, um Vorgänge zu beschreiben und zu verwalten.

Prozess \neq Prozessinkarnation

Prozess ist ein **abstraktes Gebilde**

- ein „Programm in Ausführung“[~], ein asynchroner Programmablauf[~]
- ein „Ablauf“[~], der eine Verwaltungseinheit[~] ist

?

Prozessinkarnation ist ein **konkretes Gebilde**

- die „physische Instanz“ des abstrakten Gebildes „Prozess“
 - gebunden an (Software-) Betriebsmittel, verbunden mit einer *Identität*
- die *Verwaltungseinheit*, die einen Prozess beschreibt und repräsentiert

☞ im weitesten Sinne synonyme Begriffe, jedoch nicht in allen Fällen

Prozessmodelle

schwergewichtiger Prozess (*heavyweight process*, „klassischer“ UNIX Prozess)

- Prozessinkarnation und Benutzeradressraum bilden eine Einheit
- Prozesswechsel bedeutet zwei Adressraumwechsel $AR_x \Rightarrow BS \Rightarrow AR_y$

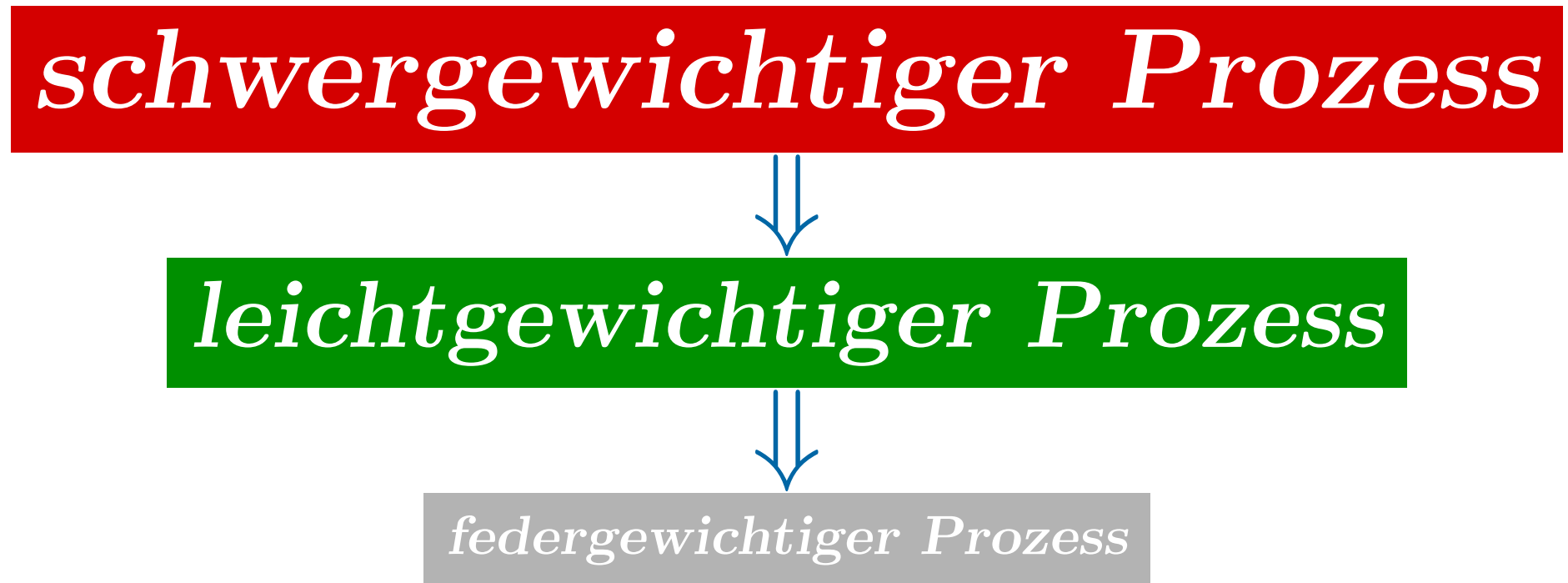
leichtgewichtiger Prozess (*lightweight process*, *kernel-level thread*)

- Prozessinkarnation und Adressraum sind voneinander entkoppelt
- Prozesswechsel bedeutet einen Adressraumwechsel $AR_x \Rightarrow BS \Rightarrow AR_x$

federgewichtiger Prozess (*featherweight process*, *user-level thread*)

- Prozessinkarnationen und Adressraum bilden eine Einheit
- Prozesswechsel entspricht einem *Koroutinenwechsel* (X Kap. 6)

Prozessbenutzthierarchie



Einplanung

Auch „Ablaufplanung“, ist erforderlich, um die um den Prozessor (allgemein: die Betriebsmittel) konkurrierenden Prozesse geordnet ablaufen lassen zu können.

Scheduling stellt sich allgemein zwei grundsätzlichen Fragestellungen:

1. Zu welchem *Zeitpunkt* sollen Prozesse ins System eingespeist werden?
2. In welcher *Reihenfolge* sollen Prozesse ablaufen?

Ein *Scheduling-Algorithmus* verfolgt das Ziel, den von einem Rechnersystem zu leistenden Arbeitsplan so aufzustellen (und zu aktualisieren), dass ein gewisses Maß an Benutzerzufriedenheit maximiert wird.

Ablaufplanung — *Process Scheduling*

Prozessen das Betriebsmittel $\left\{ \begin{array}{l} \text{Prozessor} \\ \text{Speicher}^{\dagger} \\ \text{Gerät}^{\ddagger} \end{array} \right\}$ zuteilen

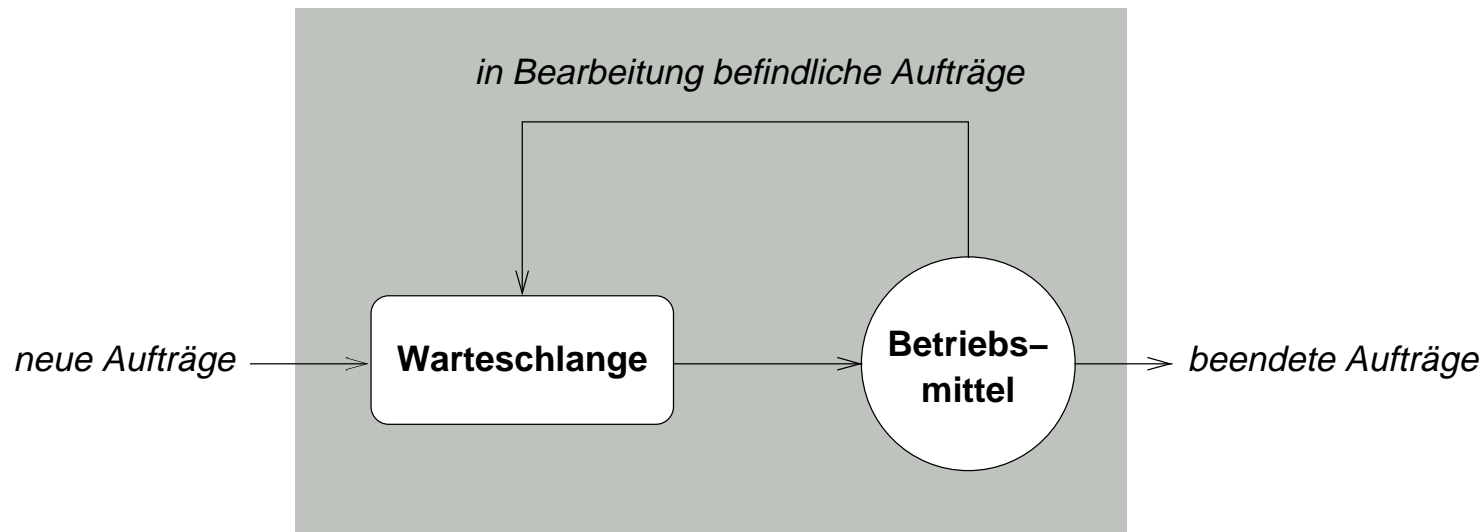
[†] Das Hardware-Betriebsmittel ($\{\text{Vorder,Hinter}\}$ grund-) *Speicher* steht (insbesondere im Falle von RAM) auch stellvertretend für (wiederverwendbare/konsumierbare) Software-Betriebsmittel wie z. B. Puffer, Nachrichten, Signale.

[‡] *Gerät* steht stellvertretend für Drucker, Netzwerk, Platte, . . .

Ablaufplan — *Process Schedule*

- *Fahrplan* zur Belegung der {Hard,Soft}ware-Betriebsmittel durch Prozesse
 - geordnet nach Ankunft, Zeit, Termin, Dringlichkeit, Gewicht, . . .
 - die Ordnung ist eine Funktion der Scheduling-Strategie (bzw. -Algorithmen)
- technisch (zumeist) realisiert auf Basis dynamischer Datenstrukturen
 - eine oder mehrere Queues bzw. Schlangen: *Warteschlangen*
 - die Elemente der Datenstruktur sind die *Prozessdeskriptoren*
- die gewählte Scheduling-Strategie bestimmt u. a. die Rechnerbetriebsart

Scheduling-Modell



Ein einzelner Scheduling-Algorithmus charakterisiert sich durch die Reihenfolge von Prozessen in der Warteschlange und die Bedingungen, unter denen die Prozesse der Warteschlange zugeführt werden.

Warteschlangentheorie

- Betriebssysteme durch die „theoretische/mathematische Brille“ gesehen:
 - ➡ R. W. Conway, L. W. Maxwell, L. W. Millner. *Theory of Scheduling*.
 - ➡ E. G. Coffmann, P. J. Denning. *Operating System Theory*.
 - ➡ L. Kleinrock. *Queuing Theory*.
- die Verfahren stehen und fallen mit den Vorgaben der jeweiligen *Zieldomäne*
 - eine „Eier-legende Wollmilchsau“ des Scheduling wird es nie geben
 - Kompromisslösungen sind geläufig, aber nicht in allen Fällen tragfähig
- *Scheduling* ist als „Querschnittsbelang“ von Betriebssystemen zu verstehen
 - diese Belange zu behandeln, ist eine der (praktischen) Herausforderungen

UNIX Prozesse sind schwergewichtig: Prozess und Adressraum bilden eine Einheit.

```
int foo;  
int hal = 42;  
  
int main () {  
    for (;;)   
        printf("Die Antwort auf alle Fragen lautet %d\n", hal + foo);  
}
```

Wie ist der Adressraum bzw. Speicher des ausführenden Prozesses organisiert?

Speichermodell

UNIX Prozess (2)

```
wosch@fauai40 40$ gcc -O6 -static -o hal hal.c; ./hal
```

Die Antwort auf alle Fragen lautet 42 ... ^Z

```
wosch@fauai40 41$ ps
```

PID	TTY	TIME	CMD
28426	pts/4	0:00	hal
205	pts/4	0:00	ps
25965	pts/4	0:00	tcsh-6.0

```
wosch@fauai40 42$ pmap -x 28426
```

28426: ./hal

Address	Kbytes	RSS	Anon	Locked	Mode	Mapped File	
00010000	216	216	-	-	r-x--	hal	⇐ <i>Textsegment</i>
00054000	16	16	8	-	rwX--	hal	} ⇐ <i>Datensegment</i>
00058000	8	8	8	-	rwX--	[heap]	
FFBFE000	8	8	8	-	rw---	[stack]	⇐ <i>Stapelsegment</i>

total Kb	248	248	24	-	

```

        .file      "hal.c"
        .global   hal
        .section   ".data"
        .align    4
        .type     hal,#object
        .size     hal,4
hal:
        .uaword   42
        .common   foo,4,4
        .section   ".rodata"
        .align    8
.LLC0:
        .asciz    "Die Antwort auf
                  alle Fragen lautet %d\n"
        .section   ".text"
        .align    4
        .global   main
        .type     main,#function
        .proc     04
main:
        !#PROLOGUE# 0
        save      %sp, -112, %sp
        !#PROLOGUE# 1
        sethi     %hi(hal), %l2
        sethi     %hi(foo), %l1
        sethi     %hi(.LLC0), %l0
        ld        [%l2+%lo(hal)], %g1
.LL5:
        or        %l0, %lo(.LLC0), %o0
        ld        [%l1+%lo(foo)], %o3
        call      printf, 0
        add       %g1, %o3, %o1
        b         .LL5
        ld        [%l2+%lo(hal)], %g1
.LLfe1:
        .size     main,.LLfe1-main
        .ident    "GCC: (GNU) 3.0.4"

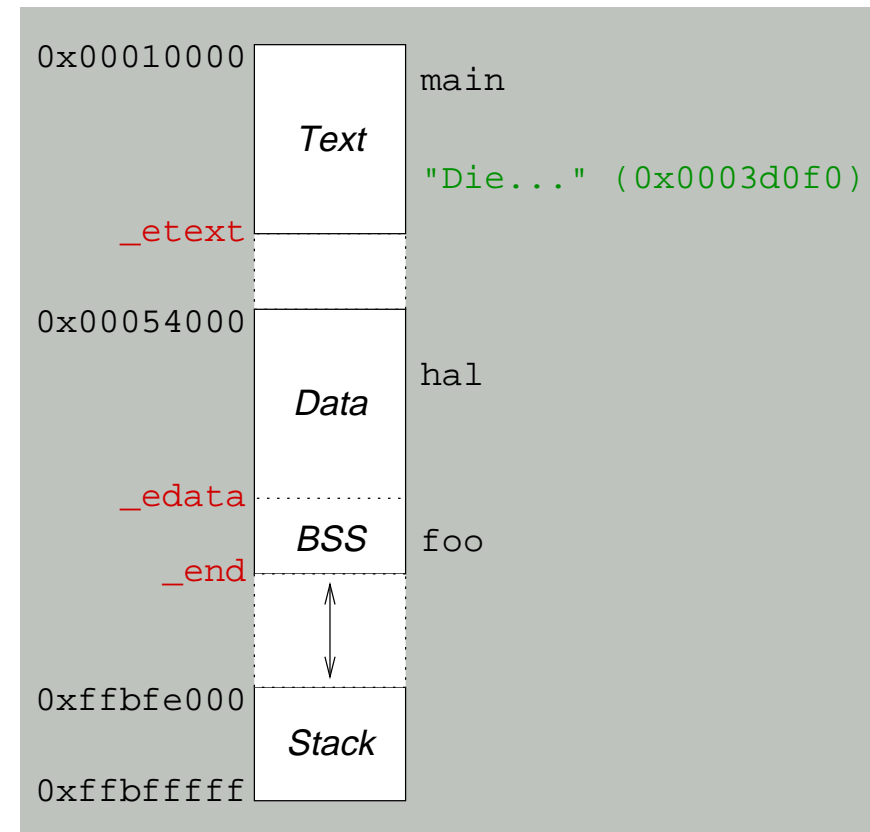
```

gcc -O6 -S hal.c

Speichermodell

```
wosch@fau40 43> nm -p -g hal
:
0000066112 T main      ➡ 0x00010240
0000352140 D hal      ➡ 0x00055f8c
0000360336 B foo      ➡ 0x00057f90
:
0000286461 D _etext    ➡ 0x00045efd
0000358433 D _edata    ➡ 0x00057821
0000361444 D _end      ➡ 0x000583e4
:
```

UNIX Prozess (4)



Vom Binder definierte Symbole:

`extern etext` ist die erste Adresse nach dem Programmtext

`extern edata` ist die erste Adresse nach dem initialisierten Datenbereich

`extern end` ist die erste Adresse nach dem uninitialisierten Datenbereich³³

- mit Ausführungsbeginn entspricht `end` der „Bruchstelle“ des Programms
 - kann zur Ausführungszeit mit `brk(2)/sbrk(2)` verschoben werden
- `sbrk((intptr_t*)0)` liefert den für den Prozess aktuell gültigen Wert

³³Dem „*Block Started by Symbol* (BSS)“ Segment, das vom Lader mit 0 vorinitialisiert wird.

Prozess

UNIX Systemaufrufe (8)

`pid = fork ()` erzeugt einen Kindprozess (exakte Kopie des Elternprozesses)

- zwei Prozesse kehren zurück: Rückgabewert $\left\{ \begin{array}{l} = 0 \\ > 0 \end{array} \right\} \Rightarrow$ ~~Elternprozess~~ Kind
 \Rightarrow ~~Kindprozess~~ Eltern

`pid = wait (status)` wartet auf Stillstand/Termination eines Kindprozesses

- ein terminierter Kindprozess („Zombie“) wird erst jetzt entsorgt
- der Aufrufer erfährt (über `status`) die Stillstands-/Terminationsursache

`exit (status)` terminiert den aufrufenden Prozess

- der Prozess wird zum „Zombie“ und vermerkt den Grund der Termination

`nv = nice (incr)` ändert die Priorität des aufrufenden Prozesses

- der angegebene Wert wird zum „*nice value*“ (nv) des Prozesses addiert
 - eine positive Zahl: je höher ihr Wert, desto niedriger die Priorität
- für die vom Benutzer einstellbare Priorität gilt: $0 \leq nv \leq (2 * NZERO) - 1$

`pid = getpid ()` liefert die Identifikation des aufrufenden Prozesses

`pid = getppid ()` liefert die Identifikation des Elternprozesses

`ok = execv (path, argv)` führt eine Datei aus

- das Speicherabbild des laufenden Prozesses wird ersetzt durch ein Programm
 - die angegebene „ausführbare Datei“ legt das neue Speicherabbild fest
 - diese Datei wird ausgeführt (CPU) oder interpretiert (z. B. von `sh(1)`)
- das auszuführende Programm erhält einen Argumentenvektor (`argv[]`)
 - die Größe des Vektors wird bestimmt und als `argc` weitergegeben

`ok = execve (path, argv, envp)` dito

- das auszuführende Programm erhält zusätzlich Umgebungsparameter

Zusammenfassung

- Betriebssysteme bieten eine „Hand voll“ nützlicher Abstraktionen an
 - Adressräume, Speicher, Dateien, Prozesse
- von zentraler Bedeutung ist die Abbildung von Namen auf Adressen
 - symbolische \rightsquigarrow numerische \rightsquigarrow logische \rightsquigarrow virtuelle \rightsquigarrow physikalische Adresse
 - ein Konzept, das im Kontext von Rechnernetzen seine Fortführung findet³⁴
- je nach Betriebssystemart sind die Abstraktionen unterschiedlich ausgelegt
 - manche Abstraktionen müssen auch überhaupt nicht angeboten werden

³⁴So sind z. B. Email-Adresse und URL nichts anderes als symbolische Adressen, die von „Namensdiensten“ auf korrespondierende Rechneradressen (Internet-Hosts) abgebildet werden.