

# Systemnahe Software (Allgemeine Informatik IV)

Prof. Dr. Franz Schweiggert und Dr. Matthias Grabert

1. August 2001

**Fakultät Mathematik u. Wirtschaftswissenschaften**  
**Abteilung Angewandte Informationsverarbeitung**  
**Vorlesungsbegleiter (gültig ab SS 2000)**



Anmerkungen:

- Auf eine systematische Unterscheidung zwischen *BSD Unix* und *System V Unix (Linux, Solaris)* wird hier verzichtet.
- Die enthaltenen Beispiel-Programme wurden zum großen Teil unter Linux entwickelt und sind weitestgehend unter Solaris getestet (für konstruktive Hinweise sind die Autoren dankbar).
- Die Beispiele sollen jeweils gewisse Aspekte verdeutlichen und erheben nicht den Anspruch von Robustheit und Zuverlässigkeit. Man kann alles anders und besser machen.
- Details zu den behandelten / verwendeten System Calls sollten jeweils im Manual bzw. den entsprechenden Header-Files nachgelesen werden.
- Dieses Skript ersetzt kein Lehrbuch und erhebt auch nicht den Anspruch eines zu sein. Die Erklärungen der Vorlesung sind zum Verständnis unverzichtbar.
- In das Skript sind absichtlich einige Schreibfehler eingebaut worden, damit Sie diese finden können. Wer am meisten Fehler entdeckt hat, erhält zwei Zusatzpunkte für die Übungen.

**Literatur: (jeweils neueste Ausgabe beachten!)**

- [Bach86] M. J. Bach: The Design of the UNIX Operating System. Prentice Hall 1986
- [Comer91] D. E. Comer: Internetworking with TCP/IP. Volume I; Principles, Protocols, and Architecture. Prentice Hall 1991
- [Comer91] D. E. Comer, David L. Stevens: Internetworking with TCP/IP. Volume II; Design, Implementation and Internals. Prentice Hall 1991
- [Etter99] M. Etter: Zur Qualitätsverbesserung in der Entwicklung und Pflege verteilter Anwendungen - Konzeption und Realisierung einer Bibliothek auf der Basis von TCP/IP. Dissertation an der Universität Ulm, Abt. Angewandte Informationsverarbeitung, 1999.
- [Rochkind85] M. Rochkind: UNIX-Programmierung für Fortgeschrittene. Hanser Verlag 1985
- [Santifaller] M. Santifaller: TCP/IP and ONC/NFS. Internetworking in a UNIX Environment. Addison Wesley 1994
- [Stevens92] W. R. Stevens: Advanced Programming in the UNIX Environment. Addison-Wesley 1992
- [Stevens95] W. R. Stevens: TCP/IP Illustrated, Volume 2 - The Implementation. Addison-Wesley 1995
- [Stevens96] W.R. Stevens: TCP/IP Illustrated, Volume 3 - TCP for Transactions, HTTP, NNTP, and the UNIX Domain Protocols. Addison-Wesley 1998
- [Stevens98] W.R. Stevens: UNIX Network Programming, Volume 1 - Networking APIs: Sockets and XTI. Prentice Hall 1998.
- [Stevens99] W.R. Stevens: UNIX Network Programming, Volume 2 - Interprocess Communications. Prentice Hall 1999.
- [Tanenbaum87] A. S. Tanenbaum: Operating Systems - Design and Implementation. Prentice Hall 1987

# Inhaltsverzeichnis

<b>1</b>	<b>Unix Process Subsystem</b>	<b>1</b>
1.1	"Theorie" der Prozesse . . . . .	1
1.2	Prozess-Baum . . . . .	3
1.3	Prozessmanagement . . . . .	3
1.4	Synchrone und asynchrone Prozesse . . . . .	4
1.5	CPU-intensive Prozesse . . . . .	5
1.6	Prozesszustände und Zustandsübergänge . . . . .	6
1.7	Regions . . . . .	7
1.8	Grundlegende System Call's . . . . .	9
1.8.1	System Call "fork" . . . . .	9
1.8.2	System Call "exit" . . . . .	12
1.8.3	System Call "exec" . . . . .	13
1.8.4	System Call "wait" . . . . .	16
1.9	Bootstrapping . . . . .	18
1.10	Der init-Prozess . . . . .	19
1.11	Signale . . . . .	21
1.11.1	Grundlegendes . . . . .	21
1.11.2	Reaktion auf Signale: signal() . . . . .	22
1.11.3	Reaktion auf Signale: sigaction() . . . . .	22
1.11.4	Anwendungsbeispiele . . . . .	28
1.12	Implementierung einer Midi-Shell . . . . .	35
<b>2</b>	<b>Inter-Prozess-Kommunikation (IPC)</b>	<b>39</b>
2.1	Einführung . . . . .	39
2.2	IPC - Client-Server Beispiel . . . . .	40
2.3	System Call dup . . . . .	41
2.4	Unnamed Pipes . . . . .	41
2.5	Client-Server mit "Unnamed Pipes" . . . . .	44
2.6	Standard I/O Bibliotheksfunktion . . . . .	49
2.7	Pipes in der Shell . . . . .	50
2.8	Nicht behandelte IPC-Mechanismen . . . . .	57
<b>3</b>	<b>Netzwerk-Kommunikation</b>	<b>59</b>
3.1	Übersicht . . . . .	59
3.2	Lokale Netze . . . . .	62
3.3	LAN-Standards . . . . .	63
3.4	Ethernet . . . . .	63
3.5	Internetworking - Concept & Architectural Model . . . . .	65
3.6	Transport-Protokolle . . . . .	72
3.6.1	Ports . . . . .	72
3.6.2	UDP . . . . .	73

<b>4</b>	<b>Berkeley Sockets</b>	<b>77</b>
4.1	Grundlagen	77
4.2	Ein erstes Beispiel: Timeserver an Port 11011	78
4.3	Die Socket-Abstraktion	79
4.4	Die Socket-Programmierschnittstelle	80
4.4.1	Vorbemerkungen	80
4.4.2	Überblick/Einordnung	81
4.4.3	Erzeugung eines Socket	85
4.4.4	Benennung eines Socket	85
4.4.5	Aubau einer Kommunikationsverbindung	87
4.4.6	Client-Beispiel: Timeclient für Port 11011	91
4.4.7	Überblick: Gebrauch von TCP-Ports	92
4.4.8	Der Datentransfer	97
4.4.9	Terminierung einer Kommunikationsverbindung	99
4.4.10	Verbindungslose Kommunikation	99
4.4.11	Feststellen gebundener Adresse	101
4.4.12	Socket-Funktionen im Überblick	101
4.5	Konstruktion von Adressen	103
4.5.1	Socket-Adressen	103
4.5.2	Socket-Adressen der UNIX-Domäne	105
4.5.3	Socket-Adressen in der Internet-Domäne	106
4.5.4	Byte-Ordnung	107
4.5.5	Spezifikation von Internet-Adressen	108
4.5.6	Hostnamen	109
4.5.7	Lokale Hostnamen und IP-Adressen	112
4.5.8	Portnummern und Dienste	114
4.5.9	Protokoll- und Netzwerkinformationen	116
4.5.10	Zusammenfassung der Netzwerkinformationen	117
4.5.11	IPv6	118
<b>5</b>	<b>Netzwerk-Programmierung</b>	<b>119</b>
5.1	Client/Server	119
5.1.1	Vorbemerkungen	119
5.1.2	concurrent server	119
5.1.3	iterative server	120
5.2	echo-Server und echo-Client	121
5.3	Erste Implementierungen	121
5.3.1	Headerfiles	121
5.3.2	TCP-Verbindung - Concurrent Server	122
5.3.3	UDP-Verbindung - Iterative Server	127
5.3.4	TCP-Verbindung in der UNIX Domain	131
5.3.5	Modifikation der ersten Implementierung	134
5.3.6	Anmerkungen	135
5.4	Verbesserte Implementierungen	138
5.4.1	Zeilenorientierter Echo-Server	138
5.4.2	Ein „stream“-basierter Echo-Server	143

# Kapitel 1

## Unix Process Subsystem

### 1.1 "Theorie" der Prozesse

- Die Ausführung eines Programms heißt **Prozess**.
- Der **Kontext** eines Prozesses ist seine Ausführungsumgebung. Dazu gehören: Anweisungen, Daten, Stack und Heap - Register - Betriebssystem-Ressourcen (offene I/O-Verbindungen, IPC-Elemente, ...)  
(IPC: *InterProcess Communication*)
- Ausführbare Programme sind z.B.:
  - vom Compiler generierter Objektcode ("a.out")
  - Datei mit einer Folge von Shell-Kommando's → **chmod+x**)
- Das Betriebssystem verwaltet eine endliche Menge von Prozessen und versucht, die vorhandenen Ressourcen (Speicher, Rechenzeit, I/O-Operationen) fair auf die einzelnen Prozesse zu verteilen.  
Ein Prozess folgt bei der Ausführung einer genau festgelegten Folge von Anweisungen, die in sich abgeschlossen sind. Schutzmechanismen des Betriebssystems und der Hardware verhindern, dass ein Prozess auf Anweisungen oder Daten außerhalb seines Adreßraums zugreift. Die gesamte Kommunikation mit anderen Prozessen (*IPC*) oder mit dem Kernel muß über System Calls erfolgen.
- (Fast) alle Prozesse entstehen aus einem bereits existierenden Prozess durch Aufruf von **int fork()**.
- Jeder Prozess besitzt eine Prozessnummer (**PID**), als eindeutige Identifikation, die einen Index in die Prozesstabelle darstellt. In dieser Tabelle verwaltet das Betriebssystem die wichtigsten Daten aller Prozesse.  
Mit der Funktion **int getpid()** kann ein Prozess seine eigene *PID* abfragen:.
- (Fast) jeder Prozess hat einen "Erzeuger" (siehe oben: *fork()*), der sog. **parent process**  
Dessen *PID* kann ein Prozess mit der Funktion **int getppid()** erhalten.
- Jeder Prozess gehört zu einer **Prozessgruppe**, die ebenfalls durch eine kleine positive Integerzahl identifiziert wird (*process group ID*).  
Ist die *process ID* eines Prozesses gleich der *process group ID*, so wird dieser Prozess als **process group leader** bezeichnet.

Mit **kill** (Shell-Kommando und auch System Call) können Signale an alle Prozesse einer *process group* gesandt werden.

Bestimmung der *process group ID*: Funktion **int getpgrp()**;  
in BSD4.3: *int getpgrp(int pid)*; ist das Argument *pid* 0, liefert diese Funktion die *process group ID* des aktuellen Prozesses, ansonsten die der Prozessgruppe des über *pid* angegebenen Prozesses.

Die Zuordnung zu einer Prozessgruppe kann geändert werden:  
unter System V:

**int setpgrp();**

damit wird der aufrufende Prozess zum *process group leader*; als *process group ID* wird die *process ID* geliefert.

unter 4.3.BSD:

*int setpgrp(int pid, int pgrp);*

ist *pid* 0, so ist der aktuelle Prozess gemeint, ansonsten muß der angesprochene Prozess dieselbe effektive *user ID* wie der aktuelle Prozess haben oder der aktuelle Prozess muß Superuser Privilegien haben.

### Beispiel:

```
swg@hypatia$ cat pid.c
#include <stdio.h>
#include <unistd.h>

int main() {
    printf("pid = %ld, ppid = %ld, groupPID = %ld\n",
        getpid(), getppid(), getpgrp() );
    exit(0);
}

swg@hypatia$ gcc pid.c
swg@hypatia$ a.out
pid = 229, ppid = 212, groupPID = 229

swg@hypatia$ echo $$ # PID der aktuellen Shell
212
swg@hypatia$
```

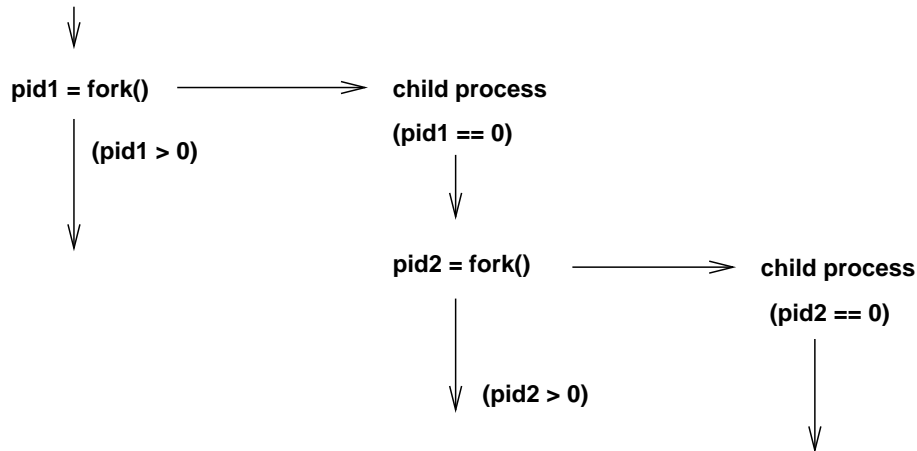
- Jeder Prozess kann Mitglied einer **terminal group** sein, die selbst wieder über eine positive ganze Zahl (*terminal group ID*) identifiziert wird. Diese ist die *process ID* des *process group leader*, der das Terminal "geöffnet" hat - typischerweise die *login shell*. Dieser Prozess wird **control process** für dieses Terminal genannt. Jedes Terminal hat nur einen *control process*.

Die *terminal group ID* identifiziert das Kontroll-Terminal für jede Prozessgruppe. Das Kontroll-Terminal wird benutzt, um Signale abzuarbeiten (von der Tastatur aus erzeugt, oder wenn die *login shell* terminiert). Wenn der Prozess, der gleichzeitig Kontroll-Prozess eines Terminals und Prozessgruppenführer ist (typischerweise die *login shell*), ein *exit* zum Beenden aufruft, so wird das **hangup** Signal an alle Prozesse dieser Gruppe gesandt. Das Kontrollterminal wird automatisch referenziert durch */dev/tty*. Ein **Dämon-Prozess** hingegen ist ein Hintergrundprozess ohne kontrollierendes Terminal. (Z.B. *httpd*)

## 1.2 Prozess-Baum

- Hierarchie und Vererbung

parent process



**Vererbung:** Prozesse vererben Teile ihrer Umgebung an ihre Kindprozesse:

- die Standard-Dateiverbindungen
- den aktuellen Katalog
- Umgebungsvariablen

Auch die Regeln für die Weitergabe von Informationen folgen streng der hierarchischen Vererbungslehre. Elternprozesse können ihren Kindern Teile ihrer Umgebung als Kopie mitgeben. Die umgekehrte Richtung ist nicht möglich. Prozesse können Daten untereinander nur über **IPC-Mechanismen** (*interprocess communication*) austauschen.

Eine Ausnahme bildet der **Beendigungsstatus** (kleine ganze Zahl), den ein Prozess bei seiner Beendigung an seinen Erzeugerprozeß zurückgeben kann (**exit()** → **wait()**). Der Beendigungsstatus 0 signalisiert einen erfolgreichen Ablauf des Prozesses (UNIX-Konvention); ein von Null verschiedener Wert ist das Zeichen für eine nicht erfolgreiche Beendigung und wird in der Regel mit einer Fehlermeldung auf der Standard-Fehlerausgabe *stderr* verbunden (Exit-Status > 0: nicht notwendig Fehlersituation, s. z.B. *egrep()*).

## 1.3 Prozessmanagement

- Kommando **ps** - Ausgabe der Prozesstabelle

Das Kommando *ps* gibt alle wesentlichen Informationen über die aktuelle Prozesshierarchie (→ Prozesstabelle) aus.

Je nach Option (und Betriebssystemversion) besteht die Ausgabe von *ps* aus mehr oder weniger Informationen zu den einzelnen Prozessen. Die Ausgabe ist, wie bei *ls*, in tabellarischer Form.

```
swg@byron$ ps
  PID TTY STAT TIME COMMAND
   205 p1  S    0:01 bash
   207 p2  S    0:00 bash
   214 p3  S    0:00 bash
  1382 p3  S    0:00 script
  1385 p4  R    0:00 ps
swg@byron$
```

#### Bedeutung:

PID	die eindeutige Prozessnummer (process id)
TTY	Name des Terminals (Windows), an dem der Prozess gestartet wurde
STAT	Prozeßstatus - R für Running, S für Sleeping T für Stopped, Z für <b>Zombie</b> (tot, aber nicht vom Vater beerdigt) (siehe man ps)
TIME	verbrauchte Rechenzeit
COMMAND	der zum Prozess gehörende Kommandoname

- Kommando **top** - Darstellung der CPU-intensiven Prozesse → **man top**
- Kommando **kill** - Prozesse abbrechen

Mit dem Kommando **kill nummer** kann man den Prozess mit der Prozessnummer *nummer* gewaltsam beenden (mit *kill* werden an einen Prozess Signale verschickt; siehe *signal.h*, *man kill* oder *kill -l*).

Auf diesem Weg lassen sich Hintergrundprozesse abbrechen, falls sie blockiert sind oder ungewöhnlich viel CPU -Zeit verbrauchen (Endlosschleife?).

Ein blockiertes Terminal läßt sich durch ein *kill*-Kommando für die betreffende Login-Shell, von einem anderen Terminal aus, im Normalfall wieder in einen benutzbaren Zustand zurückbringen. Die Login-Shell ist in der Spalte *COMMAND* als *-sh* oder schlicht als *-* aufgeführt.

Als gewöhnlicher Benutzer darf man nur seine eigenen Prozesse mit *kill* beenden. Als Super-User darf man alle Prozesse beenden; von dieser Möglichkeit sollte man aber nur dann Gebrauch machen, wenn man sich über die Konsequenzen im klaren ist...

## 1.4 Synchron und asynchrone Prozesse

- Vordergrundverarbeitung

In der Regel verarbeitet die Shell Kommandos *synchron*. Sie startet ein Kommando (→ erzeugt Prozess) und erst nach dessen Beendigung meldet sie sich mit dem *Prompt* (beginnt sie mit dem nächsten Kommando). Der Erzeugerprozeß wartet auf das Ende seines Kindprozesses.

- Hintergrundverarbeitung

Kommandos, deren Abarbeitung länger dauert (etwa in der Statistik oder Numerik), blockieren bei synchroner Abarbeitung den interaktiven Betrieb. Schließt ein **&**-Zeichen das Kommando ab, wartet die Shell (Erzeugerprozeß) nicht auf das Ende des gerade gestarteten Prozesses. Sie ist statt dessen sofort bereit, das nächste Kommando entgegenzunehmen.



- Fußangeln

Hintergrundkommandos, die von der Standardeingabe lesen wollen, erhalten das Signal **SIGTTIN** - der Prozeß wird angehalten, existiert aber weiterhin! Damit wird verhindert, dass sich ein Vordergrundprozeß (meistens die Shell selbst) und mehrere aktive Hintergrundprozesse um die Eingabe vom Terminal streiten.

```
swg@byron$ cat bg.c
# include <stdio.h>

int main() {
    int n, zahl;
    printf("Zahl: ");
    n = scanf("%d", &zahl);
    if(n==0) {
        printf("Nichts gelesen!\n");
        exit(2);
    } else {
        printf("Gelesen: %d\n", zahl);
        exit(1);
    }
}

swg@byron$ gcc -Wall bg.c
swg@byron$ a.out &
[1] 1026
Zahl:
[1]+  Stopped                  a.out
swg@byron$ ps | grep a.out
 1026  p4 T    0:00 a.out
262 swg@byron$ kill -18 1026 # SIGCONT

[1]+  Stopped                  a.out
swg@byron$ ps | grep a.out
 1026  p4 T    0:00 a.out
swg@byron$ kill -9 1026
[1]+  Killed                   a.out
swg@byron$
```

Mit dem Kommando **fg** kann ein via SIGTTIN angehaltener Prozeß in den Vordergrund geholt werden.

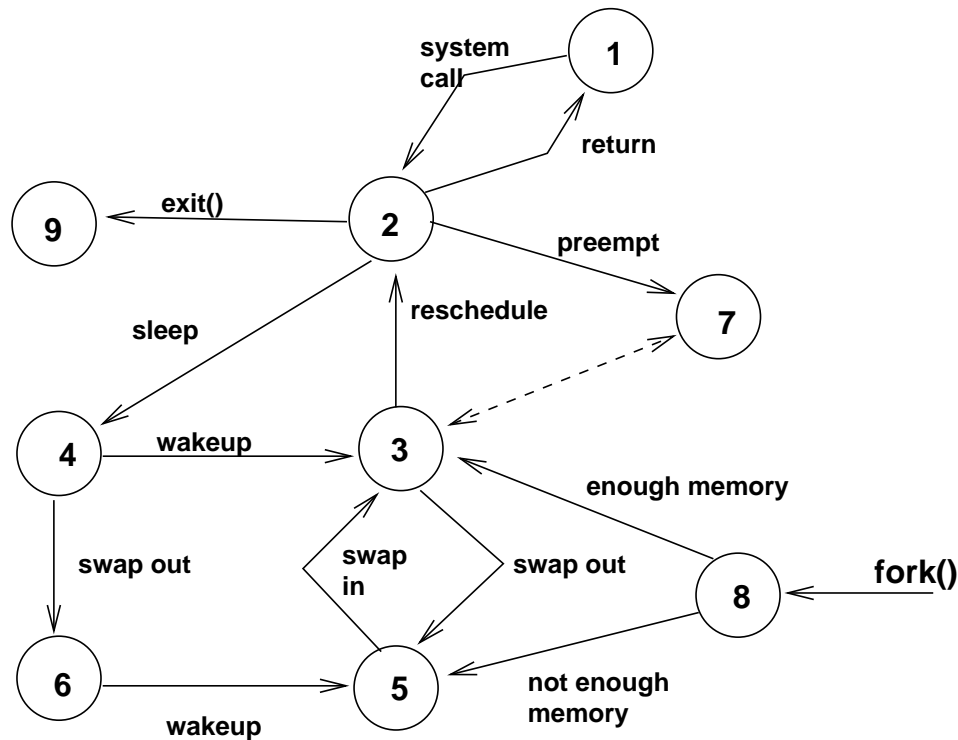
Ein beliebiges Hintergrundkommando terminiert automatisch mit der Terminalsitzung, aus der es gestartet wurde; dies läßt sich für Programme mit extrem langer Laufzeit verhindern:

**nohup kommando &**

## 1.5 CPU-intensive Prozesse

Mit dem Kommando **nice** kann die Priorität eines Prozesses verändert (herabgesetzt) werden!

## 1.6 Prozesszustände und Zustandsübergänge



- 1 Ausführung im **User Mode** (Prozess arbeitet in seinem Adreßraum)
  - 2 Ausführung im **Kernel Mode** (Prozess arbeitet im Kernel-Adreßraum)
  - 3 Prozess wartet auf Zuteilung der CPU durch Scheduler
  - 4 Prozess "schläft" im Hauptspeicher, wartet auf bestimmtes Ereignis, das ihn "aufweckt"
  - 5 Prozess ist "ready to run", aber ausgelagert - muß vom "Swapper" (Prozess 0) in Hauptspeicher geholt werden, bevor der Kern ihn zuteilen kann
  - 6 Prozess "schläft" im Sekundärspeicher (ausgelagert)
  - 7 Prozess geht vom *Kernel Mode* in *User Mode*, aber der Kern suspendiert ihn und vollzieht einen Kontext-Switch (anderer Prozess kommt zur Ausführung (bis auf Kleinigkeiten identisch mit 3))
  - 8 Prozess neu erzeugt (Startzustand für jeden Prozess - ausgenommen Prozess 0)
  - 9 Prozess hat **exit** -Aufruf ausgeführt, existiert nicht weiter, hinterläßt aber Datensatz (Exit-Code, Statistik-Angaben) für Vater-Prozess - Endzustand
- Übergänge  
System Calls und Kernel-Serviceroutinen, sowie äußere Ereignisse (Interrupts durch Uhr, HW, SW, ...) können Zustandswechsel bei den Prozessen bewirken.

Ein *Context Switch* zwischen Prozessen ist nur beim Zustandswechsel des aktiven Prozesses von Zustand (2) *Running in Kernel Mode* nach (4) *Asleep* möglich. Beim *Context Switch* wählt der Scheduler den Prozess mit der größten Priorität (Zur Erinnerung: Kommando **nice**!) aus der Ready-Liste (alle Prozesse im Zustand (3) *Ready to Run*) aus und befördert ihn in den Zustand (2). Der Scheduler teilt ihm die CPU zu und bringt ihn damit zur Ausführung.

## 1.7 Regions

Der *Kernel* unterteilt den **virtuellen** Adreßraum eines Prozesses in logische (abstrakte) Bereiche (**regions**) zusammenhängender Adressen; diese werden als verschiedene Objekte behandelt, die als Ganzes geschützt (*protected*) bzw. geteilt (*shared*) werden können

Text, Daten, Stack bilden typischerweise separate Bereiche

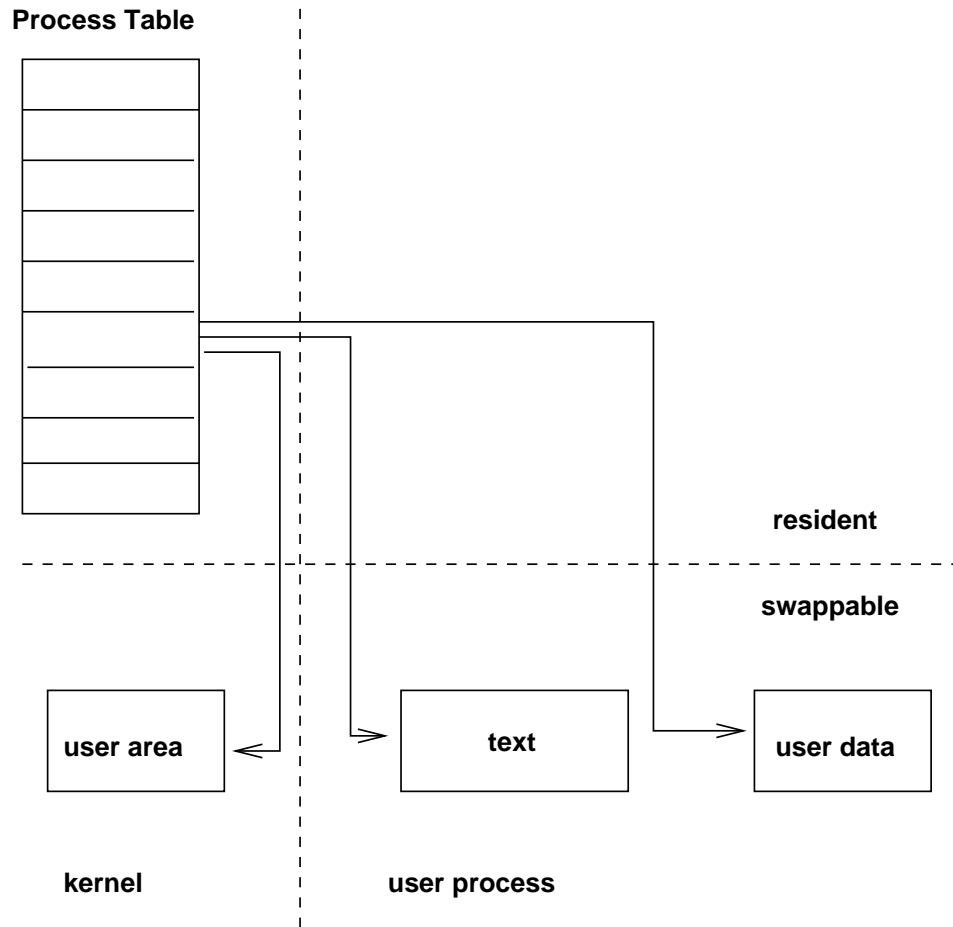
Der *Kernel* unterhält eine Tabelle mit Verweisen auf die aktiven Bereiche (u.a. Information, wo diese im physischen Speicher liegen)

Jeder Prozess unterhält eine private *per process region table* (z.B. in der **u area**)

- Context

Der *UNIX*-Jargon unterteilt den *Context* in drei disjunkte Bereiche: *Text*, *User Data* und *System Data*. Davon hält der Kernel die für ihn wichtigen Daten (*System-Daten*) an zwei definierten Plätzen in seinem Adreßraum - in der **Prozesstabelle** und in der **"User Area"** (*per process data region*)

Text und User-Daten Bereich bilden den *User Adreßraum*.



- Context Switch

Ein "*Context Switch*" findet statt, wenn der Scheduler beschließt, dass ein anderer Prozess zur Ausführung kommen soll. Der Kernel hält den laufenden Prozess an, sichert dessen *Context* und lädt den *Context* des nächsten Prozesses.

Bei jedem Wechsel sichert der Kernel soviel Informationen, dass er später zu dem unterbrochenen Prozess zurückkehren und dessen Ausführung fortsetzen kann. Für den Prozess bleibt die Unterbrechung völlig transparent.

Ein "*Context Switch*" kann nur stattfinden, während der Prozess im Kernel Mode arbeitet. Ein Wechsel des "*Execution Mode*" ist kein "*Context Switch*".

#### Einträge in der Prozesstabelle:

- Prozesszustand (siehe obiges Beispiel *bg.c*: Zustand T)
- Information, wo der Prozess und die *u area* im Haupt-/Sekundärspeicher liegen sowie über Speicherbedarf (☞ *context switch*)
- "Verwandtschaftsverhältnisse"

- *scheduling* Parameter
- diverse "Uhren" ( $\rightarrow$  *time*)
- u.a.m.

Mehr dazu: **man ps** !

### Einträge in der "u area"

- Reale und effektive UserID
- Vektor zur Signalverarbeitung
- Hinweis auf *login terminal*
- Fehlereintrag für Fehler bei einem Systemaufruf
- Rückgabewert von Systemaufrufen
- Aktueller Katalog, Aktuelle Wurzel
- *User File Descriptor Table*
- u.a.m.

## 1.8 Grundlegende System Call's

### 1.8.1 System Call "fork"

#### Syntax

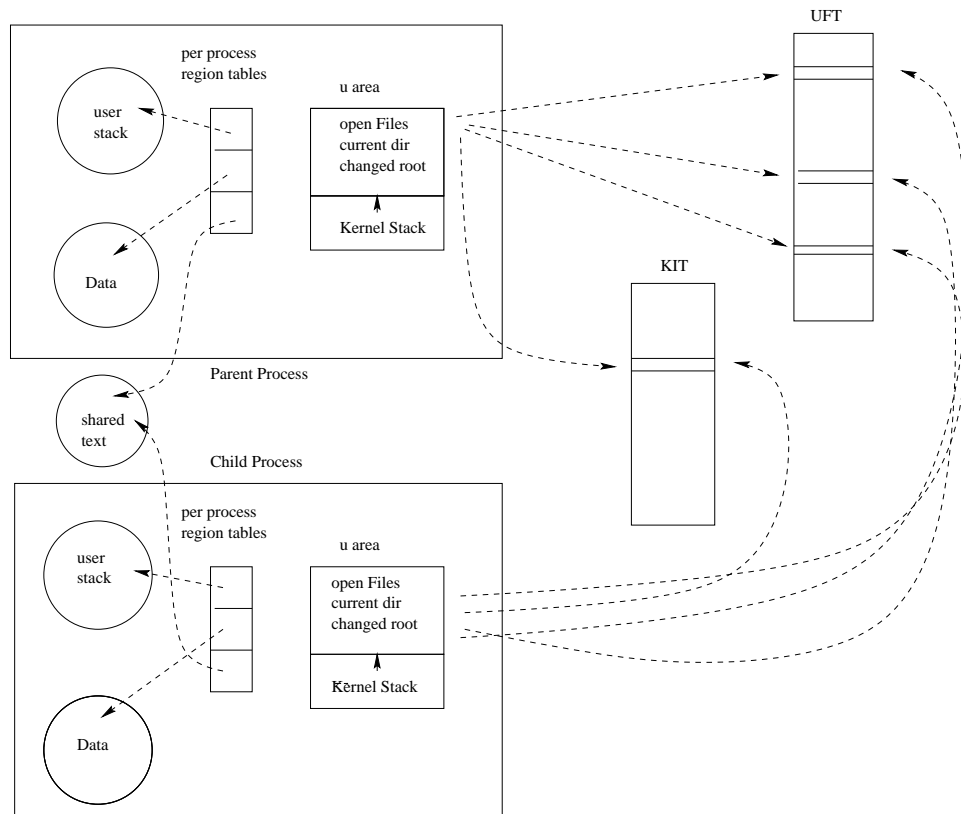
```
int fork () /* create new processes */
/* returns PID and 0 on success or -1 on error */
```

#### Beschreibung

Der System Call *fork* erzeugt einen neuen Prozess, indem er den Context des ausführenden Prozesses dupliziert.

- Der neu erzeugte Prozess wird als Kind-Prozess (**child process**), der aufrufende als Vater-Prozess (**parent process**) bezeichnet.
- *fork* wird einmal aufgerufen und kehrt im Erfolgsfall zweimal zurück: der Aufrufer erhält als Rückgabewert die PID des erzeugten Prozesses, der Kindprozeß erhält 0; im Fehlerfall (z.B. bereits zuviele Prozesse erzeugt) kehrt er einmal mit -1 zurück.

#### Beispiel



```

/* fork1.c */

#include <stdio.h>
#include <unistd.h>

int main()
{
    int pid;

    if ( (pid = fork()) > 0 ) {
        printf("Parent: created process %d, my pid is %ld\n",
            pid, getpid());
    }
    else if ( pid == 0 ) {
        printf("Child: after fork, my pid is %ld\n", getpid());
        printf("Child: my parent is %ld\n", getppid());
    }
    else
        perror( "fork() - can't fork a child" );

    printf("fertig\n");
    exit( 0 );
}

```

Ausgabe:

```
Parent: created process 319, my pid is 318
Child: after fork, my pid is 319
Child: my parent is 1
```

Der Vater-Prozess terminiert hier vor dem Sohn, der Init-Prozess "erbt" den "verlorenen Sohn". (Das Geschehen ist **nicht** deterministisch!)  
Oder so (Erzeuger "schläft" ein bißchen, bevor er terminiert):

```
/* fork2.c */

#include <stdio.h>
#include <unistd.h>

int main()
{
    int pid;
    switch ( pid = fork() ) {
        case -1:
            perror("fork() - can't fork a child");
            break;
        case 0:
            printf("Child: after fork, my pid is %ld\n", getpid());
            printf("Child: my parent is %ld\n", getppid());
            break;
        default:
            sleep(5);
            printf("Parent: created process %d, my pid is %ld\n",
                  pid, getpid());
            break;
    }
    printf("Who am I?\n");
    exit( 0 );
}
```

Ausgabe:

```
Child: after fork, my pid is 347
Child: my parent is 346
Who am I?
Parent: created process 347, my pid is 346
Who am I?
```

**Vererbt werden:**

- *real user ID*
- *real group ID*

- *effective user ID*
- *effective group ID*
- *process group ID*
- *terminal group ID*
- *root directory*
- *current working directory*
- *signal handling settings*
- *file mode creation mask*

#### Unterschiede:

- *process ID*
- *parent process ID*
- eigene File Deskriptoren (Kopie)
- Zeit bis zu einem ggf. gesetzten *alarm* ist beim Kind auf 0 gesetzt

#### Gebrauch:

1. Ein Prozess erzeugt von sich selbst eine Kopie, so dass diese "die eine oder andere" Operationen ausführt (→ Server).
2. Ein Prozess will ein anderes Programm zur Ausführung bringen - der Kind-Prozess führt via *exec* ein neues Programm aus (überlagert sich selbst) (→Shell)

### 1.8.2 System Call "exit"

#### Syntax

```
void exit( int status ) /* terminate process */
                        /* does NOT return */
                        /* 0 <= status < 256 */
```

#### Beschreibung

- Mit dem System Call *exit* beendet ein Prozess aktiv seine Existenz. Von diesem System Call gibt es keine Rückkehr in den User Mode.
- Unterscheide: System Call *exit* und C-Bibliotheks-Funktion *exit*  
C-Funktion leert erst alle Puffer und ruft dann den System Call auf.
- Die Bibliotheks-Funktion **\_exit** ruft unmittelbar den System Call, Puffer werden also nicht mehr geleert.
- Der Kernel gibt den User Adreßraum (Text und User-Daten) des Prozesses frei, sowie auch einen Teil des System-Daten Bereichs (User Area). Übrig bleibt nur der Eintrag in der Prozesstabelle, der den Beendigungsstatus und die Markierung des Prozesses als **Zombie** enthält, bis ihn der init-Prozess erbt und abräumt.



### 1.8.3 System Call “exec”

#### Syntax

```

int execl( char *path,          /* path of program file */
           char *arg0,          /* 1st argument (cmd name) */
           char *arg1, ...,    /* 2nd, ... argument */
           char *argn,         /* last argument */
           (char *) 0
        );

int execlp( char *filename,      /* name of program file */
           char *arg0, char *arg1, ... char *argn,
           (char *) 0
        );

int execl( char *path,          /* path of program file */
           char *arg0, char *arg1, ... char *argn,
           (char *) 0,
           char **envp) /* pointer to environment */
        );

int execv( char *path,          /* path of program file */
           char *argv[]; /* pointer to array of argument */
        );

int execvp( char *filename,      /* name of program file */
           char *argv[]; /* pointer to array of argument */
        );

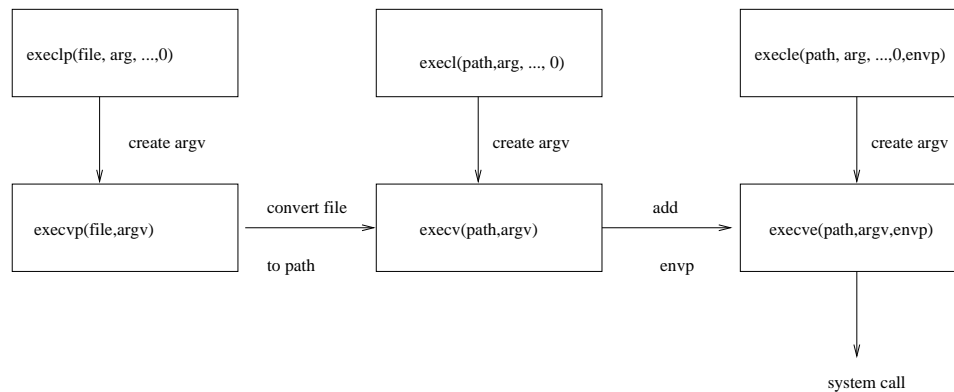
int execve( char *path,          /* path of program file */
           char *argv[], /* pointer to array of argument */
           char *envp[], /* pointer to environment */
        );
/* all return with -1 on error only */

```

#### Beschreibung

Die **exec** - System Calls überlagern im Context des ausführenden Prozesses den Text und den User-Daten Bereich mit dem Inhalt der Image-Datei. Anschließend bringen sie die neuen Instruktionen zur Ausführung.

## Zusammenhang



1. Die drei Funktionen in der oberen Reihe enthalten jedes Kommando-Argument als separaten Parameter, der NULL-Zeiger ( (char \*) 0 - !!!) schließt die variable Anzahl ab (kein *argc*!). Die drei Funktionen in der unteren Reihe fassen die Kommando-Argumente in einen Parameter **argv** zusammen, das Ende wird entsprechend wieder durch den NULL-Zeiger definiert.
2. Die zwei Funktionen in der linken Spalte definieren die Programm-Datei nur durch den Datei-Namen; diese wird über die Einträge in der Umgebungsvariable **PATH** gesucht (konvertiert in vollen Pfadnamen). Falls diese nicht gesetzt ist, wird als *default*-Suchpfad *"/bin:/usr/bin"* genommen. Enthält das Argument *path* einen *slash*, so wird die Variable *PATH* nicht verwendet.
3. Bei den vier Funktion in den beiden linken Spalten wird das Environment über die externe Variable *environ* an das neue Programm übergeben. Die beiden Funktionen in der rechten Spalte spezifizieren explizit eine Environment-Liste (muß ebenfalls mit einem NULL-Zeiger abgeschlossen sein).

### Vererbung bei *exec()*:

- *process ID*
- *parent process ID*
- *process group ID*
- *terminal group ID*
- *time left until an alarm clock signal*
- *root directory*
- *current working directory*
- *file mode creation mask*
- *real user ID*

- *real group ID*
- *file locks*

**Mögliche Änderungen mit *exec()*:**

*set user ID* / *set group ID* - *bit* der neuen Datei gesetzt

- *effective user ID* auf *user ID* des Besitzers der Programm-Datei
- *effective group ID* auf *group ID* des Besitzers der Programm-Datei

**Signale:**

- Terminieren bleibt Terminieren
- Ignorieren bleibt Ignorieren
- Speziell abgefangene Signale werden wegen des Überlagers auf Terminieren gesetzt

**Beispiel:**

```
/* ----- execl.c ----- */

#include <stdio.h>
#include <unistd.h>

void exec_test()          /* execl version */ {
    printf( "The quick brown fox jumped over " );
    fflush( stdout );
    execl( "/bin/echo", "echo", "the", "lazy", "dogs.",
          (char *) 0 );
    perror("execl - can't exec /bin/echo" );
}

int main() {
    int pid;

    if( (pid = fork() ) > 0)
        printf("PP: Parent-PID: %ld / Child-PID: %d\n",
              getpid(), pid);
    else if (pid == 0) {
        printf("CP: Child-PID: %ld\n", getpid() );
        exec_test();
    }
    else
        perror("fork() - can't fork");
    exit(0);
}
```

**Ausgabe:**

```
PP: Parent-PID: 662 / Child-PID: 663
CP: Child-PID: 663
The quick brown fox jumped over the lazy dogs.
```

Läßt man in obigem Beispiel den Funktionsaufruf **fflush()** weg, so passiert folgendes:

```
/* ----- exec2.c ----- */

#include <stdio.h>
#include <unistd.h>
#include <sys/types.h>

void exec_test()          /* execl version */ {
    printf( "The quick brown fox jumped over " );
    execl( "/bin/echo", "echo", "the", "lazy", "dogs.",
           (char *) 0 );
    perror("execl - can't exec /bin/echo" );
}

int main() {
    int pid;

    if( (pid = fork() ) > 0)
        printf("PP: Parent-PID: %ld / Child-PID: %d\n",
               getpid(), pid);
    else if (pid == 0) {
        printf("CP: Child-PID: %ld\n", getpid() );
        exec_test();
    }
    else
        perror("fork() - can't fork");
    exit(0);
}
```

Ausgabe:

```
PP: Parent-PID: 671 / Child-PID: 672
the lazy dogs.
```

*printf()* erledigt die Ausgabe gepuffert (in Blöcken zu 512 Bytes), wenn in eine Datei oder in eine Pipe geschrieben wird. Bei der Ausgabe ans Terminal (s.o.) kann die Ausgabe auch dann bruchstückhaft (oder gar nicht) erscheinen, wenn Zeilenpufferung stattfindet. Normalerweise führt dies zu keinen Problemen; der letzte Puffer wird automatisch entleert, wenn der Prozess endet. In obigem Beispiel war der Prozess noch nicht beendet, als der **exec** -Aufruf erfolgte - der Benutzerdatensegment liegende Puffer wurde überlagert, bevor er entleert werden konnte.

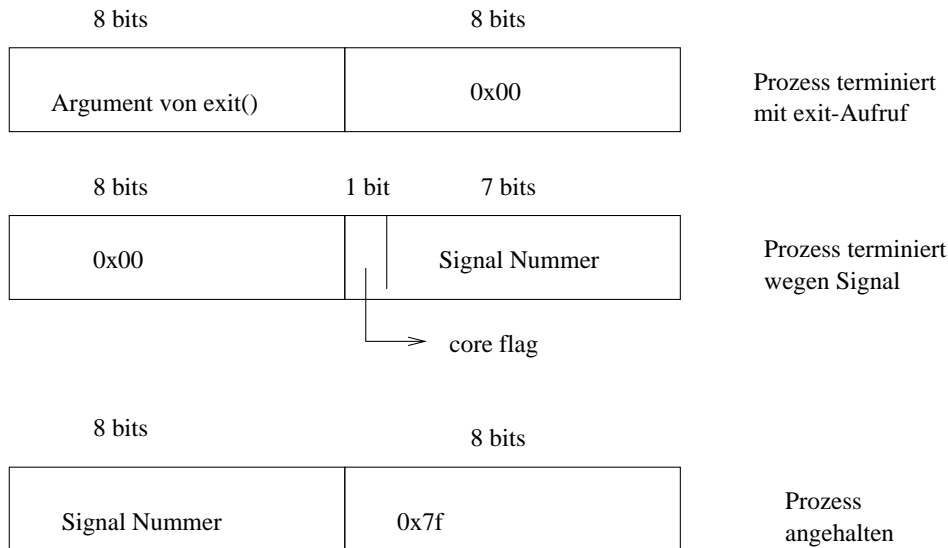
#### 1.8.4 System Call “wait”

##### Syntax

```
int wait( int * pstatus );
/* returns PID of child or -1 on error */
```

##### Beschreibung

- Falls der Aufrufer von *wait()* keinen Kind-Prozess erzeugt hat, kehrt *wait()* mit -1 zurück.
- Ansonsten blockiert *wait()* (wird vom Kernel suspendiert), bis einer der erzeugten Prozesse terminiert. In diesem Fall liefert *wait()* die PID des terminierten Kind-Prozesses. Über das Argument von *wait()* erhält der Prozess Informationen über den terminierten Kind-Prozess:



- Signal-Nummern sind größer 0!

Aufgaben des Kernel, wenn ein Prozess **exit()** aufruft:

- Wartet der Erzeuger-Prozess mit **wait()**, so wird er von der Termination des Kind-Prozesses benachrichtigt (*wait()* kehrt zurück). Im Argument von *wait()* wird der *exit*-Status (das Argument in *exit()*) geliefert, als Rückgabewert die PID des Kind-Prozesses.
- Wartet der Erzeuger-Prozess nicht, so wird der Kind-Prozess als **Zombie**-Prozess markiert. Der Kernel gibt dessen Ressourcen frei, bewahrt aber den *exit*-Status (Beendigungsstatus) in der Prozesstabelle auf.
- Sind *process ID*, *process group ID*, *terminal group ID* des terminierenden Prozesses alle gleich, so wird das Signal **SIGHUP** an jeden Prozess mit gleicher *process group ID* gesandt.

Wenn der Erzeuger vor dem Kind-Prozess terminiert:

- Die PPID der Kind-Prozesse wird auf 1 gesetzt, der *init*-Prozess "erbt" die "Waisen".  
Der *init*-Prozess terminiert nie, und wenn doch, so wäre es im Multiuser-Betrieb nicht mehr möglich, dass sich weitere Benutzer anmelden. In 4.3BSD würde es zu einem automatischen *reboot* kommen.

Wenn ein Prozess mit **exit()** terminiert, so wird dem Erzeuger das Signal **SIGCLD** (**SIGCHLD**) zugeleitet (default-Reaktion: ignorieren). In System V

ist es möglich, dass ein Prozess verhindert, dass seine Kind-Prozesse zu *Zombies* werden; dazu hat er nur den Aufruf

```
signal(SIGCLD, SIG_IGN)
```

abzusetzen; der *exit*-Status der Kindprozesse wird bei deren Termination nicht weiter aufbewahrt.

## 1.9 Bootstrapping

- Henne oder Ei?

Alle Aktionen in einem *UNIX* System erfolgen durch **Prozesse**. Für alle Prozesse besteht eine Erzeuger-Kindprozeß-Beziehung.

Bleibt die Frage, woher kommt der erste Prozess, die Wurzel dieses Prozessbaums?

- **Bootstrapping**

Das Starten eines Betriebssystems heißt *Bootstrapping*.

Für *UNIX* fallen darunter alle Aktionen vom Stromeinschalten bis zum Erreichen des stabilen Zustands in dem Prozess 1 läuft. Anschließend können alle weiteren Aktionen mit Prozessen und nach den durch die System Calls festgelegten Regeln erfolgen.

- Phase 0 Hardware

Das Einschalten der Stromversorgung bewirkt verschiedene Aktionen, die sehr von der Hardware und der Architektur des Rechners abhängen. Alle haben aber das gleiche Ziel: Selbsttest und Grundinitialisierung der einzelnen Komponenten. Die Hardware befindet sich danach in einem definierten Startzustand.

- Phase 1: First Stage Boot

Die Hardware verfügt über einige "festeingebaute" Routinen, die in der Lage sind, einen Boot-Block von der Platte zu lesen und zur Ausführung zu bringen.

Der Boot-Block enthält ein **loader** - Programm. Dieses Programm muß mit der Console kommunizieren können und bereits über das *UNIX* File System Bescheid wissen. Falls der Platz im Boot-Block für ein Programm mit diesen Fähigkeiten nicht ausreicht, muß das minimale *Boot Loader* Programm ein weiteres, größeres *Loader* Programm laden und ausführen können.

- Phase 2: Second Stage Boot

Aufgabe des **loader** - Programms ist es, den Kernel, meist die Datei */unix*, von der Platte in den Hauptspeicher zu laden und zu starten.

Anschließend beginnt *UNIX* zu "leben".

Loader Programme sind meist trickreiche Assemblerprogramme. Im Boot-Block können nur wenige Anweisungen untergebracht werden, die müssen aber komplexe, Hardware-nahe Aufgaben erledigen.

- Phase 3: Standalone

Sobald der *UNIX* Kernel gestartet wurde, übernimmt er die gesamte Kontrolle über den Rechner. Der Kernel kann jetzt alles weitere aus eigener Kraft erledigen (*Standalone*).

In seiner Startphase

- ☞ setzt er die Interruptvektoren im *low mem*,
- ☞ initialisiert die Speicherverwaltungs-Hardware,
- ☞ baut seine Tabellen (Prozess-, Open File-, Inode-, usw.) auf
- ☞ führt eine *mount*-ähnliche Operation für das *root* File System aus...

Nun fehlt noch etwas "Magie", um den ersten Prozess zu erschaffen. Die Prozessverwaltung generiert in ihrer Startphase den **Prozess 0** "von Hand". Dieser Prozess, meist **swapper** genannt, besteht nur aus dem System-Daten Bereich: Slot 0 in der Prozesstabelle plus *per process data region*. Er besitzt keinen Text oder User-Daten Bereich. Dafür existiert er während der gesamten Systemlaufzeit und ist für das *Scheduling* zuständig. Er benötigt hierfür nur Instruktionen und Daten aus dem Kernel-Adreßraum.

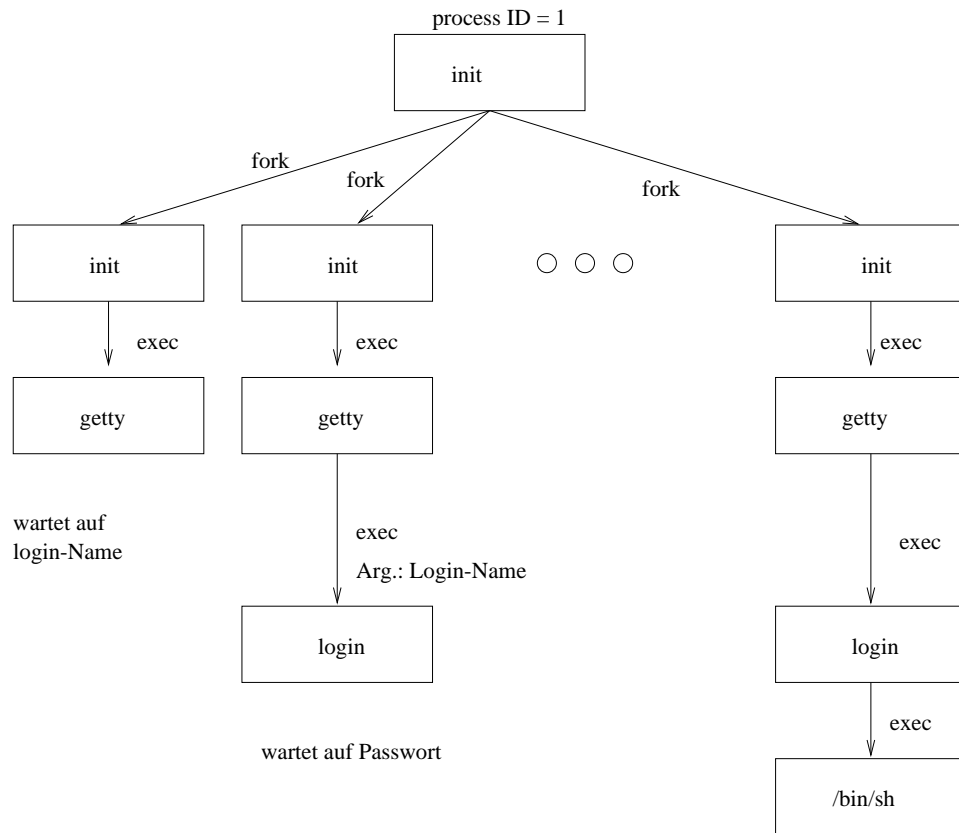
Da Prozess 0 eigentlich kein echter Prozess ist, erschafft der Kernel auch **Prozess 1** manuell. Soweit als möglich benutzt oder imitiert der Kernel hier bereits den *fork* Mechanismus, um den Prozess 1 vom Prozess 0 abzuspalten. Prozess 1 erhält jedoch einen ganz regulären Context und kann anschließend vom Scheduler als normaler Prozess zur Ausführung gebracht werden. Sein "hand crafted" Text Bereich enthält einzig die Anweisung

```
exec1( '/etc/init'', 'init'', 0 )
```

Nach dem *exec* System Call läuft im Context des Prozesses 1 das Programm */etc/init* und Prozess 1 heißt nun **init - Prozess**.

## 1.10 Der init-Prozess

- *user Id* 0
- Operationen:
  1. 4.3BSD  
*init* führt das Shell-Skript */etc/rc* aus. Dabei werden u.a. einige Daemon-Prozesse gestartet. *init* entnimmt der Datei */etc/ttys*, welche Terminals für den Multiuser-Betrieb aktiviert werden müssen.
  2. System V  
*init* liest die Datei */etc/inittab*, in der spezifiziert ist, was zu tun ist. Zum normalen Multiuser-Betrieb wird die Datei */etc/rc* ausgeführt; dieses Programm startet die meisten Daemon-Prozesse. Nach dessen Beendigung werden wie in */etc/inittab* definiert die angeschlossenen Terminals aktiviert. (mehr dazu in einem *Administrator Reference Manual*).



- *getty* setzt die Übertragungsgeschwindigkeit des Terminals, gibt irgendeine Begrüßungsformel aus und wartet auf Eingabe eines login-Namens.
- Login-Name eingegeben → *exec(/bin/login)*
- *login* sucht den eingegebenen Login-Namen in der Paßwortdatei und fordert die Eingabe eines Paßwortes.  
 Alle bis hier ausgeführten Programme (*init*, *getty*, *login*) laufen als Prozesse mit *user ID* und *effective user ID* 0 (*superuser*) - mit dem System Call *exec* ändert sich die *process ID* nicht.  
 Danach wird das *current working directory* auf den entsprechenden Eintrag für das *login directory* aus der Paßwortdatei gesetzt.  
*Group ID* und *user ID* (in dieser Reihenfolge) werden via *setgid* und *setuid* wie in der Paßwortdatei definiert gesetzt.  
 Über *exec* wird das in der Paßwortdatei spezifizierte Programm gestartet (falls keine Angabe: */bin/sh*).

Falls der sich anmeldende Benutzer nicht der Superuser ist (*login-Name* meist *root*):

- *setgid* und *setuid* reduzieren Prozessprivilegien
- beide System Calls sind dem Superuser vorbehalten (daher diese Reihenfolge)



## 1.11 Signale

### 1.11.1 Grundlegendes

- Nachricht an einen Prozess bzgl. eines eingetretenen Ereignisses
- auch als "Software-Interrupt bezeichnet
- tritt i.a. asynchron auf

#### Signale werden geschickt

- von einem Prozess an einen anderen Prozess (oder auch an sich selbst)
- vom Kernel an einen Prozess

#### Signale werden identifiziert

- intern über Integers (**Signalnummer**)
- "extern" über Namen → **signal.h**

#### Reaktion auf Signale:

- *default* → **signal.h**
- kann mit dem System Call **signal()** und der Bchereifunktion **sigaction()** eingestellt werden (s.u.)

#### Bedingungen, die Signale erzeugen:

1. System Call **kill**: damit kann ein Prozess sich oder einem anderen Prozess ein Signal senden (führt nicht notwendig zur Termination)  
**int kill(int pid, int sig);**  
Der sendende Prozess muß entweder ein superuser-Prozess sein oder der sendende und empfangende Prozess müssen dieselbe effektive userId haben.
  - Ist das pid-Argument 0, so geht das Signal an alle Prozesse in der Prozess-Gruppe des Senders
  - Ist das pid-Argument -1 und der Sender nicht der Superuser, so geht das Signal an alle Prozesse, deren *real user ID* gleich der effektive user ID des Senders ist
  - Ist das pid-Argument -1 und der Sender der Superuser, so geht das Signal an alle Prozesse ausgenommen die Superuser-Prozesse (i.a. PID's 0 oder 1)
  - Ist das pid-Argument negativ, aber ungleich -1, so geht das Signal an alle Prozesse, deren Prozess-Gruppen-Nummer gleich dem Absolutbetrag von pid ist
  - Ist das sig-Argument 0, so wird nur eine Fehler-Prüfung durchgeführt, aber kein Signal geschickt (z.B. zur Prüfung der Gültigkeit des pid-Arguments)
2. Kommando **kill**: nutzt den System Call **kill()** (s.o. und **man**)

3. Tastatureingaben, z.B.:  
**ctrl-c** (oder **delete**) beendet einen im Vordergrund laufenden Prozess (→ **SIGINT**), genauer alle Prozesse in der Kontrollgruppe dieses Terminals  
 - vom Kernel verschickt  
**quit**-Zeichen (meist **ctrl-\**) erzeugt **SIGQUIT**
4. Hardware-Bedingungen, z.B.:  
 Gleitpunkt-Arithmetik-Fehler (*Floating Point Exception*) (**SIGFPE**)  
 Adreßraum-Verletzungen (*Segmentation Violation*) (**SIGSEGV**)
5. Software-Bedingungen, z.B.:  
 Schreiben in eine Pipe, an der kein Prozess zum Lesen "hängt" (**SIGPIPE**).

### 1.11.2 Reaktion auf Signale: **signal()**

1. Ein Prozess kann eine Funktion definieren, die bei Eintreten eines bestimmten Signals ausgeführt werden soll (**signal handler**).
2. Signale (außer **SIGKILL** und **SIGSTOP**) können ignoriert werden.

**System Call** *signal*:

```
#include <signal.h>

int (*signal (int sig, void (*func) (int))) (int);
/* ANSI C signal handling */
```

Also:

- *signal* ist eine Funktion, die einen Zeiger auf eine Funktion zurückliefert, die selbst eine *int* zurückliefert (die bisherige Reaktion auf das Signal oder bei Fehler **SIG\_ERR**)
- Das erste Argument ist die Signalnummer (Makro aus *signal.h*), das zweite (*func*) ist Zeiger auf eine Funktion, die *void* liefert (die neue Reaktion auf das Signal)

Wegen der detaillierteren Reaktionsmöglichkeit und einiger anderen Fußangeln ist die Funktion **sigaction()** vorzuziehen!

### 1.11.3 Reaktion auf Signale: **sigaction()**

NAME

sigaction - detailed signal management

SYNOPSIS

```
#include <signal.h>

int sigaction(int sig, const struct sigaction *newact,
              struct sigaction *oldact);
```

DESCRIPTION

The **sigaction()** function allows the calling process to examine or specify the action to be taken on delivery of a specific signal. (See **signal(5)** for an explanation of gen-

eral signal concepts.)

The `sig` argument specifies the signal and can be assigned any of the signals specified in `signal(5)` except `SIGKILL` and `SIGSTOP`.

If the argument `newact` is not `NULL`, it points to a structure specifying the new action to be taken when delivering `sig`. If the argument `oldact` is not `NULL`, it points to a structure where the action previously associated with `sig` is to be stored on return from `sigaction()`.

The `sigaction` structure includes the following members:

```
void      (*sa_handler)();
void      (*sa_sigaction)(int, siginfo_t *, void *);
sigset_t  sa_mask;
int       sa_flags;
```

The `sa_handler` member identifies the action to be associated with the specified signal, if the `SA_SIGINFO` flag (see below) is cleared in the `sa_flags` field of the `sigaction` structure. It may take any of the values specified in `sig-`

`nal(5)` or that of a user specified signal handler. If the `SA_SIGINFO` flag is set in the `sa_flags` field, the `sa_sigaction` field specifies a signal-catching function.

The `sa_mask` member specifies a set of signals to be blocked while the signal handler is active. On entry to the signal handler, that set of signals is added to the set of signals already being blocked when the signal is delivered. In addition, the signal that caused the handler to be executed will also be blocked, unless the `SA_NODEFER` flag has been specified. `SIGSTOP` and `SIGKILL` cannot be blocked (the system silently enforces this restriction).

The `sa_flags` member specifies a set of flags used to modify the delivery of the signal. It is formed by a logical OR of any of the following values (more: see man page):

<code>SA_RESTART</code>	If set and the signal is caught, certain functions that are interrupted by the execution of this signal's handler are transparently restarted by the system; namely, <code>read(2)</code> or <code>write(2)</code> on slow devices like terminals, <code>ioctl(2)</code> , <code>fcntl(2)</code> , <code>wait(2)</code> , and <code>waitid(2)</code> . Otherwise, that function returns an <code>EINTR</code> error.
<code>SA_NOCLDWAIT</code>	If set and <code>sig</code> equals <code>SIGCHLD</code> , the system will not create zombie

processes when children of the calling process exit. If the calling process subsequently issues a `wait(2)`, it blocks until all of the calling process's child processes terminate, and then returns -1 with `errno` set to `ECHILD`.

#### RETURN VALUES

On success, `sigaction()` returns 0. On failure, it returns -1 and sets `errno` to indicate the error. If `sigaction()` fails, no new signal handler is installed.

#### ERRORS

see man page

#### SEE ALSO

`kill(1)`, `intro(2)`, `exit(2)`, `fcntl(2)`, `ioctl(2)`, `kill(2)`, `pause(2)`, `read(2)`, `sigaltstack(2)`, `sigprocmask(2)`, `sig-`  
`send(2)`, `sigsuspend(2)`, `wait(2)`, `waitid(2)`, `write(2)`,  
`signal(3C)`, `sigsetops(3C)`, `thr_create(3T)`, `attributes(5)`,  
`siginfo(5)`, `signal(5)`, `ucontext(5)`

#### Beispiel:

```
/* ----- sign.h ----- */
/* Signal Handler */

#include <signal.h>

typedef void (*Sigfunc)(int);

#ifndef SIGN_H
#define SIGN_H
/* avoid multiple includes */

Sigfunc ignoresig(int);
/* interrupt und quit ignorieren */
Sigfunc entrysig(int);
/* interrupt und quit restaurieren */
#endif

/* ----- sign.c ----- */

# include "sign.h"
# include <stdio.h>

void myignore(int sig){
    printf("signal handler: SIGNAL = %d<\n",sig);
```

```

    return;
}

struct sigaction newact, oldact;

Sigfunc ignoresig(int sig) {
    static int first = 1;
    newact.sa_handler = myignore;
    if (first) {
        first = 0;
        if (sigemptyset(&newact.sa_mask) < 0)
            return SIG_ERR;
        /* Durch diese Initialisierung der Komponente sa_mask mit
         * der leeren Menge wird bewirkt, dass waehrend der Aus-

         * fuehrung der installierten Signalbehandlungsfunktion
         * mit Ausnahme des im Argument sig angegebenen Signals
         * keine weiteren Signale von der Zustellung durch den
         * Systemkern zurueckgehalten werden
         */

        newact.sa_flags = 0;
        newact.sa_flags |= SA_RESTART;
        if (sigaction(sig, &newact, &oldact) < 0)
            return SIG_ERR;
        else
            return oldact.sa_handler;
    } else {
        if (sigaction(sig, &newact, NULL) < 0)
            return SIG_ERR;
        else
            return NULL;
    }
}

Sigfunc entrysig(int sig) {
    if (sigaction(sig, &oldact, NULL) < 0 )
        return SIG_ERR;
    else
        return NULL;
}

/* ----- main.c ----- */
# include <stdio.h>
# include <unistd.h>
# include "sign.h"

int main() {
    int sleep_time;
    /*install reaction*/
    if ( (ignoresig(SIGINT) == SIG_ERR) ||
         (ignoresig(SIGQUIT) == SIG_ERR) ) {

```

```

        perror("ignoresig");
        exit(1);
    }
    printf("Break?\n");
    while(1) { /*loop forever*/
        sleep_time = 5;
        do {
            printf("go asleep for %d sec\n", sleep_time);
            sleep_time = sleep(sleep_time);
        } while(sleep_time != 0);
        ;
        printf("And now?\n");
        /*restore reaction */
        if ( (entrysig(SIGINT) == SIG_ERR) ||
            (entrysig(SIGQUIT) == SIG_ERR) ) {
            perror("entrysig");
            exit(2);
        }
    }
    exit(0);
}

```

### Ausführung:

```

hypatia$ show
Break?
go asleep for 5 sec
#ctrl-c
signal handler: SIGNAL = >2<
go asleep for 4 sec
#ctrl-\
Quit
hypatia$ show
Break?
go asleep for 5 sec
# ctrl-c
signal handler: SIGNAL = >2<
go asleep for 4 sec
#ctrl-c
signal handler: SIGNAL = >2<
go asleep for 2 sec
And now?
go asleep for 5 sec
#ctrl-c

hypatia$

```

### Modifikation wie oben:

```

/* ----- main1.c ----- */
# include <stdio.h>
# include <unistd.h>

```

```
# include "sign.h"

int main() {
    int sleep_time;
    if (ignore_sig(SIGINT) == SIG_ERR) /*set new reaction*/ {
        perror("ignore_sig");
        exit(1);
    }
    printf("Break ?\n");
    while(1) { /*do forever*/
        sleep_time = 5;
        do {
            printf("go asleep for %d s\n", sleep_time);
            sleep_time = sleep(sleep_time);
            printf("sleep_time: %d\n", sleep_time);
        } while( (sleep_time != 0) );

        printf("and now?\n");
    }
    exit(0);
}
```

**Ausführung:**

```
hypatia$ showl
Break ?
go asleep for 5 s
# ctrl-c
signal handler: SIGNAL = >2<
sleep_time: 4
go asleep for 4 s
# ctrl-c
signal handler: SIGNAL = >2<
sleep_time: 1
go asleep for 1 s
sleep_time: 0
and now?
go asleep for 5 s
# ctrl-c
signal handler: SIGNAL = >2<
sleep_time: 5
go asleep for 5 s
# ctrl-c
signal handler: SIGNAL = >2<
sleep_time: 3
go asleep for 3 s
# ctrl-\
Quit
hypatia$
```

### 1.11.4 Anwendungsbeispiele

Im folgenden soll ein Prozess Kind-Prozesse erzeugen, die vor dem Erzeuger-Prozess terminieren und auf die der Erzeuger auch nicht wartet - er wird allerdings über ein Signal **SIGCHLD** / **SIGCLD** benachrichtigt. Die Kind-Prozesse werden zu **Zombie-Prozessen**.

```
/* ----- sign.h ----- */

#include <signal.h>

typedef void (*Sigfunc)(int);

#ifdef SIGN_H
#define SIGN_H
/* avoid multiple includes */

Sigfunc ignoresig(int);
Sigfunc entrsig(int);
#endif

/* ----- sign.c ----- */

# include "sign.h"
# include <stdio.h>

void myignore(int sig){
    printf("signal handler: SIGNAL = >%d<\n",sig);
    return;
}

struct sigaction newact, oldact;

Sigfunc ignoresig(int sig) {
    static int first;
    newact.sa_handler = myignore;
    first = 1;
    if (first) {
        first = 0;
        if (sigemptyset(&newact.sa_mask) < 0)
            return SIG_ERR;

        newact.sa_flags = 0;
        newact.sa_flags |= SA_RESTART;

        if (sigaction(sig, &newact, &oldact) < 0)
            return SIG_ERR;
        else
            return oldact.sa_handler;
    } else {
        if (sigaction(sig, &newact, NULL) < 0)
            return SIG_ERR;
        else
```



```

        return NULL;
    }
}

Sigfunc entrysig(int sig) {
    if (sigaction(sig, &oldact, NULL) < 0 )
        return SIG_ERR;
    else
        return NULL;
}

/* ----- main.c ----- */
# include <stdio.h>
# include <stdlib.h>
# include <unistd.h>
# include <signal.h>
# include "sign.h"

int main() {
    char cmdstr[128];
    int pid[3], i = 3;
    if (ignoresig(SIGCHLD) == SIG_ERR) {
        perror("ignoresig");
        exit(1);
    }

    printf("Create child processes:\n");
    while(i>0) {
        switch(pid[3-i]=fork()) {
            case -1: perror("fork");
                    exit(1);
            case 0 : printf("I'm child %d with PID = %ld!\n", 4-i, getpid());
                    sleep(5);
                    exit(0);
            default:
                    sleep(1);
                    if (kill(pid[3-i],2) < 0 ) {
                        perror("kill");
                        exit(1) ;
                    }
        }
        i--;
    }
    sleep(5);
    printf("Process table (extract child processes above): \n");
    sprintf(cmdstr,"ps -lp %d,%d,%d", pid[0], pid[1], pid[2]);
    system(cmdstr);
    sleep(1);
    printf("Parent: going to exit!\n");
    exit(0);
}

```

**Ausführung:**

```

hypatia$ show
Create children:
I'm child 1 with PID = 18542!
I'm child 2 with PID = 18543!
signal handler: SIGNAL = >17<
I'm child 3 with PID = 18544!
signal handler: SIGNAL = >17<
signal handler: SIGNAL = >17<
Process table (extract):
18541  p2 S    0:00 show
18542  p2 Z    0:00 (show <zombie>)
18543  p2 Z    0:00 (show <zombie>)
18544  p2 Z    0:00 (show <zombie>)
18545  p2 S    0:00 sh -c sleep 1; ps | grep show
18548  p2 R    0:00 sh -c sleep 1; ps | grep show
signal handler: SIGNAL = >17<
Parent: going to exit!
hypatia$

```

**Vermeidung der Zombies**

```

/* ----- signal.h ----- */

#include <signal.h>

typedef void (*Sigfunc)(int);

#ifdef SIGN_H
#define SIGN_H
/* avoid multiple includes */

Sigfunc ignoresig(int);
Sigfunc entrysig(int);
#endif

/* ----- signal.c ----- */

# include "sign.h"
# include <stdio.h>
# include <sys/wait.h>

void myignore(int sig){
    int status;
    printf("signal handler: SIGNAL = >%d<\n",sig);
    waitpid(0, &status, WNOHANG);
    return;
}

struct sigaction newact, oldact;

```

```

Sigfunc ignoresig(int sig) {
    static int first;
    newact.sa_handler = myignore;
    first = 1;
    if (first) {
        first = 0;
        if (sigemptyset(&newact.sa_mask) < 0)
            return SIG_ERR;

        newact.sa_flags = 0;
        newact.sa_flags |= SA_RESTART;

        if (sigaction(sig, &newact, &oldact) < 0)
            return SIG_ERR;
        else
            return oldact.sa_handler;
    } else {
        if (sigaction(sig, &newact, NULL) < 0)
            return SIG_ERR;
        else
            return NULL;
    }
}

Sigfunc entriysig(int sig) {
    if (sigaction(sig, &oldact, NULL) < 0 )
        return SIG_ERR;
    else
        return NULL;
}

/* ----- main1.c ----- */
# include <stdio.h>
# include <stdlib.h>
# include <unistd.h>
# include "sign1.h"

int main() {
    char cmdstr[128];
    int pid[3], i = 3;
    if (ignoresig(SIGCLD) == SIG_ERR) {
        perror("ignoresig");
        exit(1);
    }

    printf("Create child processes:\n");
    while(i>0) {
        switch(pid[3-i]=fork()) {
            case -1: perror("fork");
                    exit(1);
            case 0 : printf("I'm child %d with PID = %ld!\n", 4-i, getpid());

```

```

        sleep(5);
        exit(0);
    default: sleep(1);
        kill(pid[3-i],9);
    }
    i--;
}
sleep(5);
printf("Process table (extract child processes above): \n");
sprintf(cmdstr, "ps -lp %d,%d,%d", pid[0], pid[1],pid[2]);
system(cmdstr);
sleep(1);
printf("Parent: going to exit!\n");
exit(0);
}

```

### Ausführung:

```

hypatia$ showl
Create children:
I'm child 1 with PID = 334!
I'm child 2 with PID = 335!
signal handler: SIGNAL = >17<
I'm child 3 with PID = 336!
signal handler: SIGNAL = >17<
signal handler: SIGNAL = >17<
Process table (extract):
  333  p2 S    0:00 showl
  337  p2 S    0:00 sh -c sleep 1; ps | grep show
  340  p2 S    0:00 grep show
signal handler: SIGNAL = >17<
Parent: going to exit!
hypatia$

```

### Weitere Anwendung: Signale zur "Kommunikation" - SIGUSR1 und SIGUSR2

```

/* ----- sign.h ----- */

#include <signal.h>

typedef void (*Sigfunc)(int);

#ifndef SIGN_H
#define SIGN_H
/* avoid multiple includes */

Sigfunc ignoresig(int);
Sigfunc entrysig(int);
#endif

```

```

/* ----- sign.c ----- */

#include "sign.h"
#include <stdio.h>
#include <siginfo.h>

struct sigaction newact, oldact;
siginfo_t info;

void my_handler(int sig){
    printf("signal handler: SIGNAL = >%d<\n",sig);
    return;
}

Sigfunc ignoresig(int sig) {
    static int first = 1;
    newact.sa_handler = my_handler;
    if (first) {
        first = 0;
        if (sigemptyset(&newact.sa_mask) < 0)
            return SIG_ERR;

        newact.sa_flags = 0;
        newact.sa_flags |= SA_RESTART;
        /*SA_SIGINFO -> man sigaction */

        if (sigaction(sig, &newact, &oldact) < 0)
            return SIG_ERR;
        else
            return oldact.sa_handler;
    } else {
        if (sigaction(sig, &newact, NULL) < 0)
            return SIG_ERR;
        else
            return NULL;
    }
}

Sigfunc entrysig(int sig) {
    if (sigaction(sig, &oldact, NULL) < 0 )
        return SIG_ERR;
    else
        return NULL;
}

/* ----- main.c ----- */
#include <stdio.h>
#include <stdlib.h>
#include <unistd.h>
#include <signal.h>
#include "sign.h"

```

```

int main() {
    int pid, ppid;
    if (ignoresig(SIGUSR1) == SIG_ERR) {
        perror("ignoresig");
        exit(1);
    }

    switch(pid=fork()) {
        case -1: perror("fork");
                exit(1);
        case 0 :
            printf("I'm child with PID = %ld!\n", getpid());
            ppid=getppid();
            sleep(1);
            if (kill(ppid,SIGUSR1) < 0 ) {
                perror("kill");
                exit(1);
            }
            sleep(10);
            exit(0);
        default: printf("I'm parent with PID = %ld\n", getpid());
                 sleep(1);
    }
    sleep(1);
    printf("Parent: going to exit!\n");
    exit(0);
}

```

### Erläuterungen zu einigen Signalen:

Generell mit Vorsicht zu betrachten, da system-abhängig System V  $\longleftrightarrow$  BSD) und oft nicht sehr zuverlässig zu behandeln!

**SIGALRM** Ein Prozess kann mit `unsigned int alarm(unsigned int sec)`; eine Alarmuhr stellen. Der Parameter *sec* bestimmt, nach wievielen Sekunden (*real time*) der Kernel dem Prozess das *SIGALRM* Signal sendet (*software timeout*); ist der Wert 0, werden alle früher gestellten Alarmuhren *gecancelt*. Der Rückgabewert ist die verbleibende Zeit des letzten Aufrufs (falls einer) von *alarm*.

Die Funktion `unsigned int sleep(unsigned int sec)`; (Prozess für *sec* Sekunden "schlafen" legen) nutzt *SIGALRM*.

**SIGCLD** Wird dem Erzeuger-Prozess gesandt, wenn ein von ihm erzeugter Prozess terminiert. In 4.3BSD besagt dieses Signal, dass sich der Status des Kind-Prozesses geändert hat (Termination, aber auch angehalten durch Signale).

**SIGHUP** Beim Schließen eines Terminals gesandt an alle Prozesse, für die dieses Terminal das Kontroll-Terminal ist. Terminiert ein **process group leader**, so wird dieses Signal an alle Prozesse dieser Gruppe gesandt.

**SIGPIPE** Erzeugt, wenn ein Prozess versucht in eine FIFO-Datei oder Pipe zu schreiben, an der kein Prozess zum Lesen hängt.

**SIGSTOP** Zum gezielten Anhalten eines Prozesses, kann nicht ignoriert oder abgefangen werden; Fortsetzen des Prozesses mit Signal **SIGCONT**.

**SIGUSR1 SIGUSR2** Zur "Kommunikation" zwischen Prozessen; die einzige Information für den Empfänger ist der Signaltyp, die Prozess ID des Senders kann der Empfänger nicht feststellen.

### Zuverlässige Signale

In früheren UNIX-Versionen konnten Laufzeitbedingungen entstehen, die dazu führten, dass Signale verloren gingen, ein Ereignis konnte eintreten und ein Signal hervorrufen, das ein Prozess nicht registrierte. In System VR3 und 4.3BSD wurden Erweiterungen - leider nicht kompatibel - durchgeführt.

- Signal-Handler bleiben installiert, nachdem ein Signal eintrifft. In früheren Versionen wurde eine prozeß-definierte Aktion auf **SIG\_DFL** zurückgesetzt, bevor der Signal-Handler aufgerufen wurde; dies bedeutet, dass das erneute Eintreffen desselben Signals verloren geht, bevor der Prozess in der Lage ist, den `signal()` Aufruf erneut abzusetzen.
- Ein Prozess muß Signale abfangen können - aber sie sollen nicht einfach verschwinden (*discarded*); wenn das gewünscht ist, dann explizit über die **SIG\_IGN** Aktion. Statt dessen sollte das Signal aufbewahrt werden und dann ausgeliefert werden, wenn der Prozess bereit ist. Dazu gibt es in 4.3BSD das Konzept des *blocking*, in System V des *holding* von Signalen.
- Während ein Signal einem Prozess zugestellt wird, wird dieses "gehalten" (blockiert); d.h., dass wenn ein Signal ein zweites Mal erzeugt wird, während der Prozess noch das erste Auftreten bearbeitet, der Signal Handler kein zweites Mal aufgerufen wird. Statt dessen wird das zweite Auftreten aufbewahrt; nur wenn der Signal Handler von der Behandlung des ersten Auftretens normal zurückkehrt, wird er ein zweites Mal zur Behandlung des zweiten Auftretens aufgerufen.

Mehr zu Signalen: /Stevens90/

## 1.12 Implementierung einer Midi-Shell

### Leistungsmerkmale:

- Ausführung einer einfachen Kommandozeile, die von *stdin* eingegeben wird. Unterstützt werden IO-Umlenkung in Dateien sowie Hintergrundkommandos.
- Möglichkeiten, Kommandos mit *ctrl-c* bzw. *ctrl-(Backslash)* abzubrechen (dazu muss die *minishell* gegen diese Signale immun gemacht werden)

```
/*
 * Midishell via System
 * basierend auf Loesung von Blatt 3 AI4 SS01
 * Vater nimmt Kommando entgegen und Sohn fuehrt dieses via execvp() aus
 * Signale SIGINT und SIGQUIT werden ggf. abgefangen und an den Sohn gesandt
 * Eingabeumlenkung "<" und ">" und Hintergrundkommandos "&" wer-
den unterstuetzt
 * (c) mg sai 2001
 */
#include <stdio.h>
```

```

#include <strings.h>
#include <signal.h>
#include <wait.h>
#include <fcntl.h>
#include <unistd.h>

static int pid; // pid des Sohnes
// das ist uebrigens auch ein Kommentar :-))

/*
 * wird bei Eintreffen von SIGINT oder SIGQUIT vom Vater hier abgefangen
 * und an den Sohn gesandt
 */

void sendtoson(int sig)
{
    if (!pid) // kein Sohn aktiv ...
        return;

    printf("tinish gets signal %d ... forwarding to %d\n", sig, pid);
    kill(pid, sig); // forward signal to son
}

int main(int argc, char **argv)
{
    char buf[200]; // Kommandozeile
    char *tok; // braucht strtok ... argv
    char *args[20]; // Argumentvektor fuer Kommando
    char *bg; // gesetzt, falls "&" in Kommandozeile
    char infile[100]; // hier steht ggf. die Ausgabedatei
    char outfile[100]; // hier steht ggf. die Eingabedatei
    char *out; // gesetzt, falls > kam
    char *in; // gesetzt, falls < kam
    int i = 0;
    int stat; // fuer waitpid
    struct sigaction new;

    new.sa_handler = sendtoson;
    new.sa_flags = SA_RESTART; // Einlesen nicht unterbrechen
    sigaction(SIGINT, &new, NULL); // SIGINT und SIGQUIT abfangen
    sigaction(SIGQUIT, &new, NULL);
    printf("tinish > ");
    while (fgets(buf, sizeof(buf)-1, stdin) && strcmp(buf, "exit\n")) {
        buf[strlen(buf)-1] = 0;
        if (*buf==0) { // leere Kommandozeile
            printf("tinish > ");
            continue;
        }
        if (bg = strchr(buf, '&')) // Hintergrundstart
            *bg = ' '; // & loeschen

        in = strchr(buf, '<'); // Eingabeumlenkung ???
        if (out = strchr(buf, '>')) { // Ausgabeumlenkung ???
            *out = 0; // > loeschen
            sscanf(out+1, "%s", outfile); // Dateinamen lesen

```



```

}
if (in) {
*in = 0; // < loeschen
sscanf(in+1, "%s", infile); // Dateinamen lesen
}
tok = strtok(buf, " "); // Kommandozeile bauen
i = 0;
while (tok) {
args[i++] = strdup(tok);
tok = strtok(NULL, " ");
}
args[i] = NULL;
switch (pid = fork()) {
case -1: /* error */
break;
case 0: /* Sohn */
// Signale zuruecksetzen
new.sa_handler = SIG_DFL;
sigaction(SIGINT, &new, NULL);
sigaction(SIGQUIT, &new, NULL);
if (out) { // Ausgabeumlenkung ??
close(1); // Stdout schliessen
// open nimmt 1 als Filedes-
// kriptor (kleinster freier)
if (open(outfile,
O_CREAT|O_TRUNC|O_WRONLY, 0644)==-1)
perror(outfile),
exit(1);
}
if (in) {
close(0); // Stdin schliessen
// open nimmt 0 als Filedes-
// kriptor (kleinster freier)
if (open(infile, O_RDONLY)==-1)
perror(infile),
exit(1);
}
execvp(args[0], args); // Kommando starten !
printf("tinish: cmd >%s< not found\n", args[0]);
exit(0);
default:
if (!bg) // nicht im Hintergrund: warten
waitpid(pid, &stat, 0);
}
printf("tinish > ");
pid = 0; // braucht Fktn. sendtoson
}
return 0;
}

```

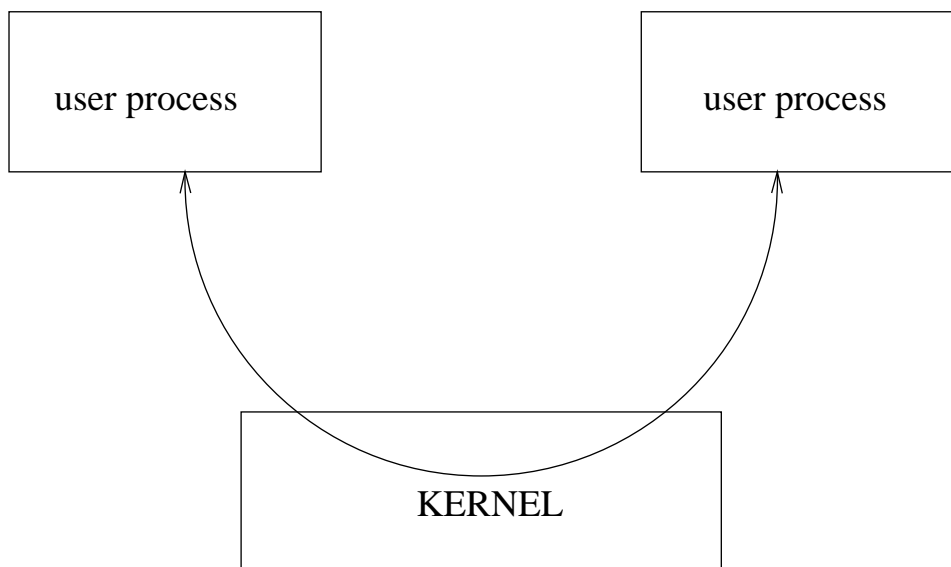


## Kapitel 2

# Inter-Prozess-Kommunikation (IPC)

### 2.1 Einführung

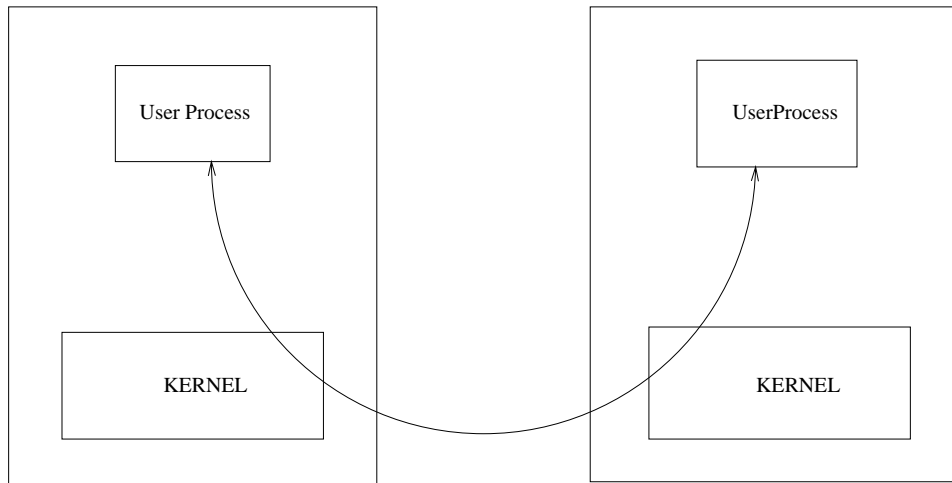
Jeder *UNIX*-Prozess besitzt seinen eigenen Context. Innerhalb eines Prozesses können die verschiedenen Moduln über Parameter und Rückgabewerte bei Funktionsaufrufen oder über globale Variablen Daten austauschen. Wollen jedoch zwei eigenständige Prozesse Daten miteinander austauschen, so kann dies nur über den Kernel via System Calls erfolgen. Denn der Kernel verhindert unkontrollierte Übergriffe eines Prozesses in den Adreßraum eines anderen Prozesses.



Bei diesem Konzept müssen beide Prozesse explizit der Kommunikation zustimmen. Der *UNIX*-Kernel bietet mit seinen *Interprocess Communication Facilities* nur die Möglichkeit zur Kommunikation an.

## Netzwerk-Kommunikation

Das UNIX-IPC-Konzept lässt sich orthogonal erweitern von der lokalen Kommunikation zwischen Prozessen innerhalb eines Systems auf Netzwerk-Kommunikation zwischen Prozessen, die auf verschiedenen Systemen laufen. Vor allem die Entwicklungsarbeiten der University of California at Berkeley brachte hier einige neue Ansätze zur Interprocess Communication in *UNIX* ein.



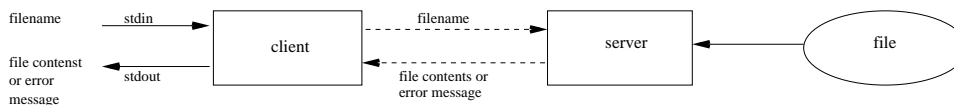
Für die Prozesse kann es völlig transparent sein, wo der jeweilige Partnerprozeß abläuft.

## 2.2 IPC - Client-Server Beispiel

Der *Client* liest einen Dateinamen von *stdin* ein und schreibt ihn in den *IPC*-Kanal. Anschliessend wartet er auf die Reaktion des Servers.

Der *Server* liest einen Dateinamen von dem *IPC*-Kanal und versucht die Datei zu öffnen. Gelingt es dem Server, die Datei zu öffnen, kopiert er ihren Inhalt in den *IPC*-Kanal. Läßt sich die Datei nicht öffnen, schickt der Server eine Fehlermeldung über den *IPC*-Kanal.

Der Client wartet auf Daten am *IPC*-Kanal, er liest sie von dort und schreibt sie nach *stdout*. Konnte der Server die Datei öffnen, zeigt der Client so den Dateiinhalt an, sonst kopiert der Client die Fehlermeldung durch.



Die beiden gestrichelten Pfeile zwischen dem Server und dem Client entsprechen dem jeweiligen "Interprocess Communication" Kanal.

## 2.3 System Call dup

Im Zusammenhang mit **unnamed pipes** ist der Systemaufruf **dup** nützlich:

```
int dup(int fd)    /*duplicate file descriptor*/
/* returns new file descriptor or -1 on error */
```

**dup** verdoppelt einen bestehenden Filedeskriptor und liefert als Resultat einen neuen Filedeskriptor (mit der **kleinsten** verfügbaren Nummer), der mit der gleichen Datei oder der gleichen Pipe verbunden ist. Beide File Deskriptoren haben denselben Positionszeiger. Damit kann z.B. ein Filedeskriptor mit der Nummer **0** erhalten werden, falls dieser vorher geschlossen wurde. Falls so mit **exec** ein Programm ausgeführt wird, das von Filedeskriptor **0** liest, kann es so dazu gebracht werden, aus einer Pipe zu lesen. Ähnlich kann so der Filedeskriptor **1** "manipuliert" werden.

## 2.4 Unnamed Pipes

- System Call **pipe()**

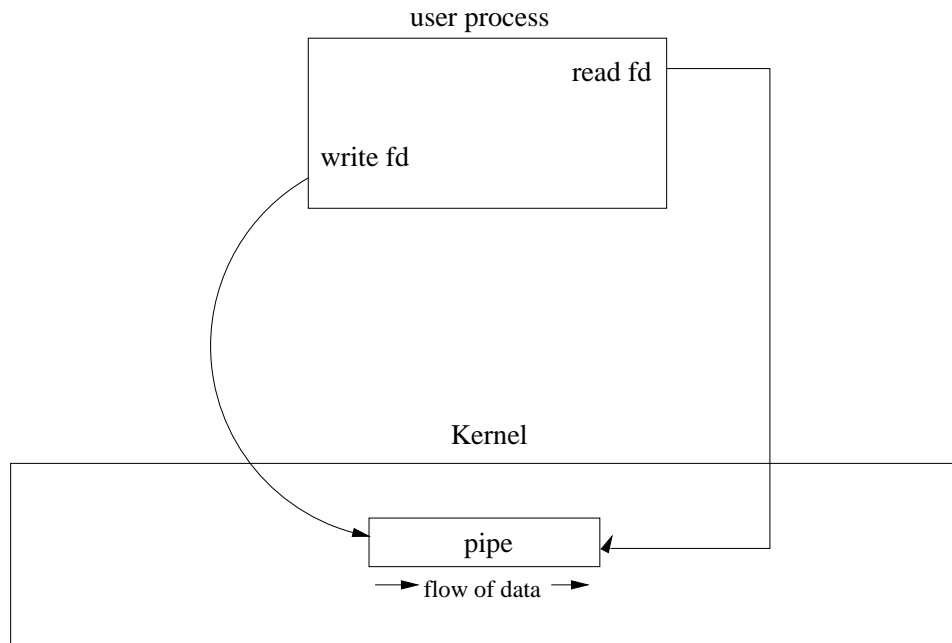
```
int pipe( int pipefd[2] )  /* create a pipe */

/* pipefd[2]: file descriptors */
/* returns 0 on success or -1 on error */
```

- Beschreibung

Pipes sind der älteste *IPC*-Mechanismus. Seit Mitte der 70er Jahre existieren sie auf allen Versionen und Arten von *UNIX*.

Eine Pipe besteht aus einem *unidirektionalen* Datenkanal. Zwei File Deskriptoren repräsentieren die Pipe im User Prozess. Der System Call *pipe()* kreiert die Pipe, er liefert die beiden Enden als File Deskriptoren über sein Vektorargument an den Prozess. Dabei ist **pipefd[1]** das "Ende" zum Schreiben, **pipefd[0]** das "Ende" zum Lesen.



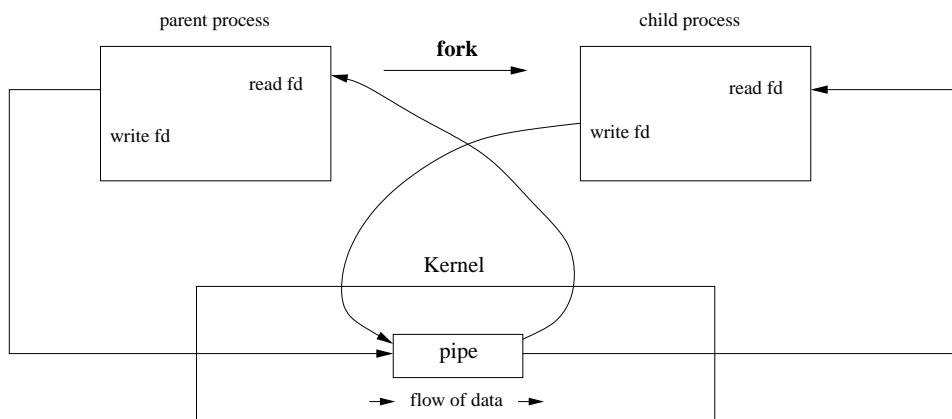
- Deadlock

In dieser Konstellation läßt sich die Pipe nur als “Zwischenspeicher” für Daten außerhalb des User Adreßraums benutzen.

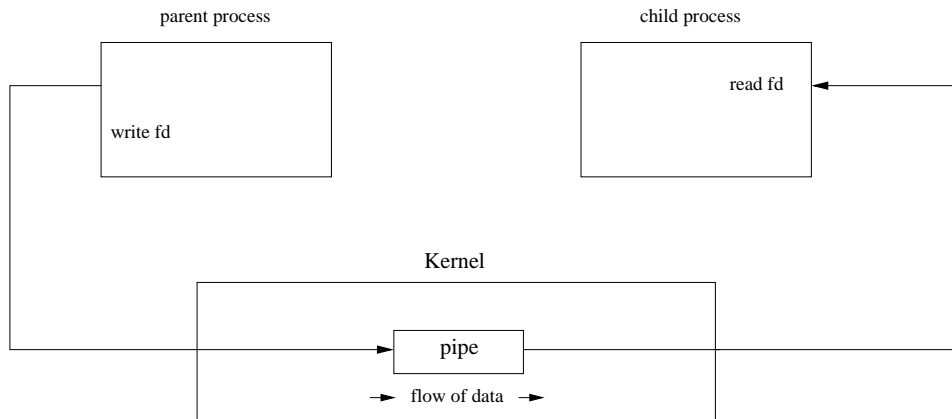
**VORSICHT!** Der Kernel synchronisiert Prozesse, die in Pipes schreiben oder aus Pipes lesen – *Producer / Consumer* Modell. Sollte hier in dem Ein-Prozess-Beispiel ein `read()` oder `write` System Call blockieren, entsteht ein *Deadlock*, denn der blockierte Prozess kann den befreienden, komplementären `write()` oder `read()` System Call nicht absetzen.

- Kommunikation zwischen verwandten Prozessen

Durch einen `fork()` System Call entstehen zwei eigenständige, aber verwandte Prozesse, die insbesondere die gleichen I/O-Verbindungen besitzen.



Schließt nun ein Prozess sein Lese-Ende und der andere Prozess sein Schreib-Ende, so entsteht ein unidirektionaler Kommunikationspfad zwischen den beiden Prozessen.



Wiederholen die Prozesse die `fork`, `pipe()` und `close` System Calls, entstehen längere Pipelines. Dieses Datenverarbeitungs-Prinzip ist untrennbar mit *UNIX* verbunden.

- Beispiel: Pipe zwischen zwei Prozessen

Das Programm kreiert eine *pipe* und einen zweiten Prozess, dem diese beiden *pipe*-Deskriptoren vererbt werden. Der Erzeugerprozeß schreibt Text in die Pipe und wartet auf das Ableben des Kindprozesses. Der Kindprozeß liest Text aus der Pipe, schreibt ihn nach *stdout* und beendet seine Ausführung. Folgende Schritte sind der Reihe nach auszuführen:

- ☞ Parent führt `pipe()` aus
- ☞ Parent führt `fork()` aus
- ☞ Child schließt sein Schreib-Ende der Pipe und wartet an seinem Lese-Ende der Pipe.
- ☞ Parent schließt sein Lese-Ende der Pipe, schreibt Text in sein Schreib-Ende der Pipe, und führt `wait()` für sein Kind aus.
- ☞ Child liest von seinem Lese-Ende, gibt gelesenen Text aus und terminiert.
- ☞ Parent hat auf Child gewartet und kann jetzt auch terminieren.

- Realisierung

```
/*----- ./pipe1/pipe.c -----*/

#include <stdio.h>
#include <unistd.h>
#include <sys/wait.h>

#define RD_FD 0
```

```

#define WR_FD    1

int main() {
int childpid, gone, pipefd[2];

if ( pipe( pipefd ) < 0 ) {
    perror("pipe" );
    exit(1);
}

switch( childpid = fork() ){
    case -1:
        perror("fork");
        exit(1);
    case 0: /* child: */ { /* <-- */
        char buf[ 128 ];    int nread;

        /* close WRITE end of pipe */
        close( pipefd[ WR_FD ] );

        /* read from pipe and copy to stdout */
        nread = read( pipefd[ RD_FD ], buf, sizeof(buf));

        printf( "Child: read '%.*s', going to exit\n",
                                   nread, buf );

        break;
    }

    default: /* parent: */
        /* close READ end of pipe */
        close( pipefd[ RD_FD ] );

        /* write something into pipe */
        write( pipefd[ WR_FD ], "hello world", 11 );

        printf("parent: wrote 'hello world' into pipe\n");

        do /* wait for child */ {
            if ( (gone = wait( (int *) 0 )) < 0 ) {
                /*no interest for exit_status*/
                perror("wait");
                exit(2);
            }
        } while ( gone != childpid );

    }
exit(0);
}

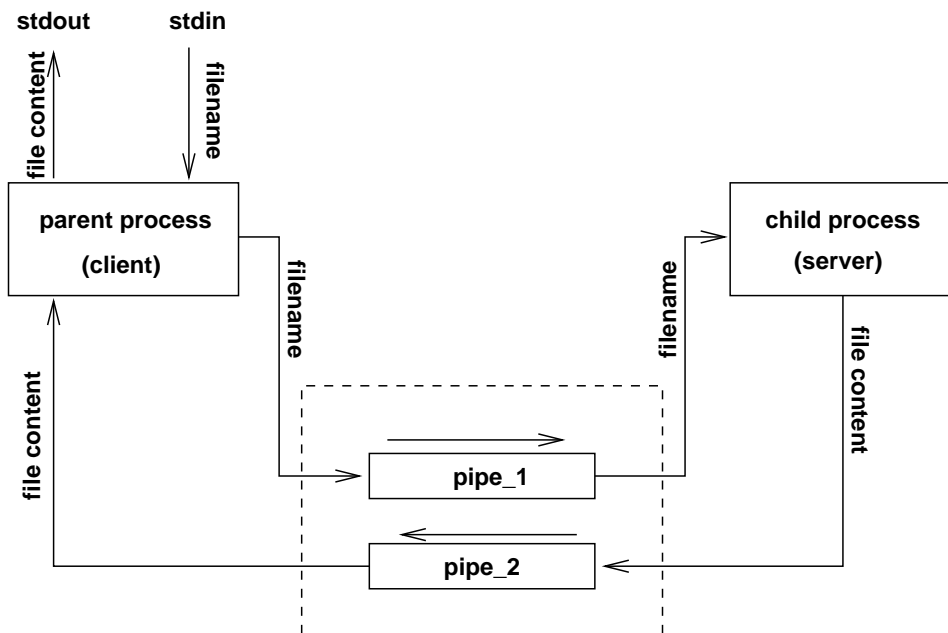
```

## 2.5 Client-Server mit “Unnamed Pipes”

- Bidirektional



Durch eine Pipe fließen Daten nur in genau eine Richtung. Zur Realisierung unseres Client-Server Beispiels benötigen wir aber einen *bidirektionalen* Kommunikationskanal. Wir müssen dazu zwei Pipes kreieren und eine Pipe für jede Richtung konfigurieren.



- Vorgehen
  - ☞ Pipe1 und Pipe2 kreieren
  - ☞ `fork()` ausführen
  - ☞ Linker Prozess (Erzeuger) schließt
    - \* Lese-Ende von Pipe1 und
    - \* Schreib-Ende von Pipe2
  - ☞ Rechter Prozess (Kind) schließt
    - \* Schreib-Ende von Pipe1 und
    - \* Lese-Ende von Pipe2

### Realisation mit zwei Pipes

- Hauptprogramm:

```
/*
 * main.c: two (unnamed) pipes to realize the IPC-channel
 * two pipes between two processes
 * main - uses two functions for client / server functionality
 */

#include <stdio.h>
#include <unistd.h>
#include <sys/wait.h>
```

```

#include "client.h"
#include "server.h"

#define RD_FD  0
#define WR_FD  1

int main() {
int childpid, gone, pipe_1[2], pipe_2[2];

if ( pipe( pipe_1 ) < 0 || pipe( pipe_2 ) < 0 ) {
    perror("pipe(): can't creat pipes");
    exit(1);
}

if ( (childpid = fork()) > 0 ) {
    /* ----- client <- parent ----- */
    printf( "Client/Parent: my pid is %ld\n", getpid() );

    close( pipe_1[ RD_FD ] ); /* close READ  end of pipe1 */
    close( pipe_2[ WR_FD ] ); /* close WRITE end of pipe2 */

    /*----- */
    client( pipe_2[ RD_FD ], pipe_1[ WR_FD ] );
    /*----- */

    /* wait for child */
    do {
        if ( (gone = wait( (int *) 0 )) < 0 ) {
            perror("wait");
            exit(2);
        }
    } while ( gone != childpid );

    printf( "Cli/Par: server/child %d terminated\n", gone );
    printf( "Cli/Par: going to exit\n" );

} else if ( childpid == 0 ) {
    /* ----- server/child ----- */
    printf( "S/Ch: after fork, my pid is %ld\n", getpid() );

    close( pipe_1[ WR_FD ] ); /* close WRITE end of pipe1 */
    close( pipe_2[ RD_FD ] ); /* close READ  end of pipe2 */

    /*----- */
    server( pipe_1[ RD_FD ], pipe_2[ WR_FD ] );
    /*----- */

    printf( "Server/Child: going to exit\n" );

} else {
    perror("fork" );
    exit(3);
}

```

```

    exit( 0 );
}

```

- Realisation des Client-Teils

```

/*
 *  client.c: realize the client part
 *  read line from stdin,
 *  write this text to IPC-channel,
 *  copy text from IPC-channel to stdout
 */

# define CLIENT_H
# include "client.h"

# include <stdio.h>
# include <unistd.h>

# define BUFSIZE 256

void client( int readfd, int writefd ) {
    char buf[ BUFSIZ ];
    int  n;

    /*
     * read filename from stdin,
     * write it to IPC-channel
     */
    printf("give filename: ");

    if ( fgets(buf, BUFSIZ, stdin) == (char *) 0 ) {
        fprintf(stderr, "Client: filename read error" );
        exit(1);
    }

    n = strlen( buf );
    if ( buf[ n-1 ] == '\n' )
        n--; /* zap NL */

    if ( write( writefd, buf, n ) != n ) {
        fprintf(stderr, "write(): Client: can't write to IPC-
channel" );
        exit(2);
    }

    /*
     * read data from IPC-channel
     * write it to stdout
     */

    while ( (n = read( readfd, buf, BUFSIZ )) > 0 )
        if ( write( 1, buf, n ) != n ) /* fd 1 == stdout */ {

```

```

        fprintf(stderr, "write(): Client: can't write to stdout" );
        exit(3);
    }

    if ( n < 0 ) {
        fprintf(stderr, "read(): Client: can't read from IPC-
channel" );
        exit(4);
    }
}

```

- Realisation des Server-Teils

```

/*
 * server.c: realize the server part
 *   read filename from IPC-channel,
 *   open this file,
 *   copy data from file to IPC-channel.
 */

# define SERVER_H
# include "server.h"

#include <stdio.h>
#include <unistd.h>
#include <fcntl.h>

#define BUFSIZE 512

void server( int readfd, int writefd ) {
    char buf[ BUFSIZ ];
    int n, fd;

    /* read filename from IPC-channel */

    if ( (n = read( readfd, buf, BUFSIZ )) < 0 ) {
        fprintf(stderr, "Server: filename read error" );
        exit(1);
    }
    buf[ n ] = '\0';

    /* try to open file */

    if ( (fd = open( buf, O_RDONLY )) < 0 ) {
        /* Format and send error mesg to client */
        char errmesg[ BUFSIZ ];

        (void) sprintf( errmesg, "Server: can't open infile (%.s)\n",
            BUFSIZ/2, buf );
        n = strlen( errmesg );
        if ( write( writefd, errmesg, n ) != n ) {
            fprintf(stderr, "S: can't write errmesg to IPC-channel" );

```

```

        exit(2);
    }
    return;
}

/*
 * read data from file and
 * write to IPC-channel
 */

while ( (n = read( fd, buf, BUFSIZ )) > 0 )
    if ( write( writefd, buf, n ) != n ) {
        fprintf(stderr, "Server: can't write to IPC-channel" );
        exit(3);
    }
if ( n < 0 ) {
    fprintf(stderr, "Server: can't read file" );
    exit(3);
}
}

```

## 2.6 Standard I/O Bibliotheksfunktion

### `popen()` und `pclose()`

```

#include <stdio.h>

FILE *popen( char * cmd, char * mode )
/* create a pipe to a cmd:
   cmd:  cmd to be executed
   mode: read from or write to pipe/cmd
   returns file pointer on success or NULL on error
*/

int pclose( FILE * fp )
/* close pipe with cmd */
/* returns exit status of cmd or -1 on error */

```

#### • Beschreibung

Die Funktion `popen()` aus der Standard I/O Bibliothek kreiert eine **unidirektionale (!)** Pipe und einen neuen Prozess, der von der Pipe liest oder in die Pipe schreibt. In dem neuen Prozess startet eine Shell und führt die mitgegebene Kommandozeile **cmd** (unter Berücksichtigung von **PATH**) aus. Die Pipe wird abhängig vom Argument **mode** (entweder "r" oder "w") so konfiguriert, dass sie *stdout* oder *stdin* des erzeugten Kommandos mit dem aufrufenden Prozess verbindet. Der aufrufende Prozess erhält sein Pipe-Ende als "File Pointer" von `popen`.

`pclose` schließt eine mit `popen` geöffnete I/O-Verbindung ab (NICHT: `fclose()`). Die Funktion blockiert bis das Kommando terminiert und liefert den Exit-Status des Kommandos zurück.

- Beispiel:

```

/* ----- popen.c: popen() ----- */

#include <stdio.h>

#define BUFFER_SIZE 1024

int main() {

    char buf[ BUFFER_SIZE ];
    FILE * fp;

    if ( (fp = popen( "/bin/pwd", "r" )) == (FILE *) 0 ) {
        perror("popen()");
        exit(1);
    }

    if ( fgets( buf, BUFFER_SIZE, fp ) == (char *) 0 ) {
        perror("fgetsr");
        exit(2);
    }

    printf( "The current working directory is:\n" );
    printf( "\t%s", buf );    /* pwd inserts newline */

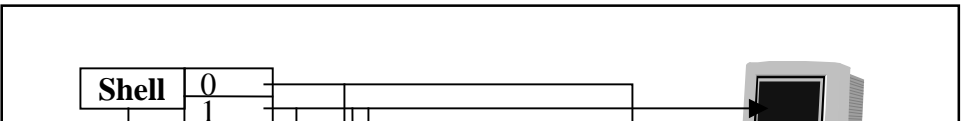
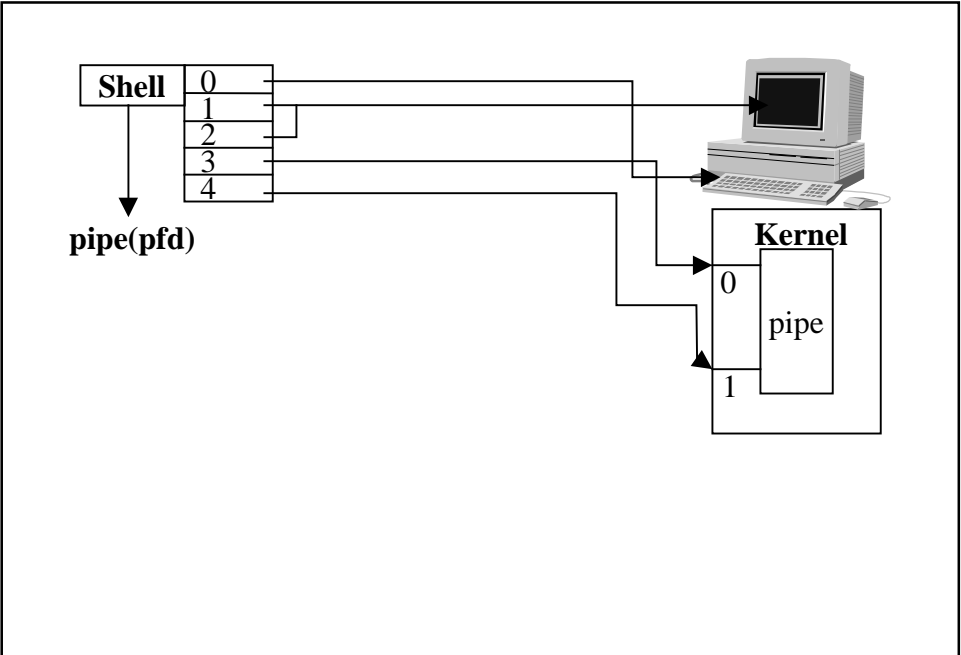
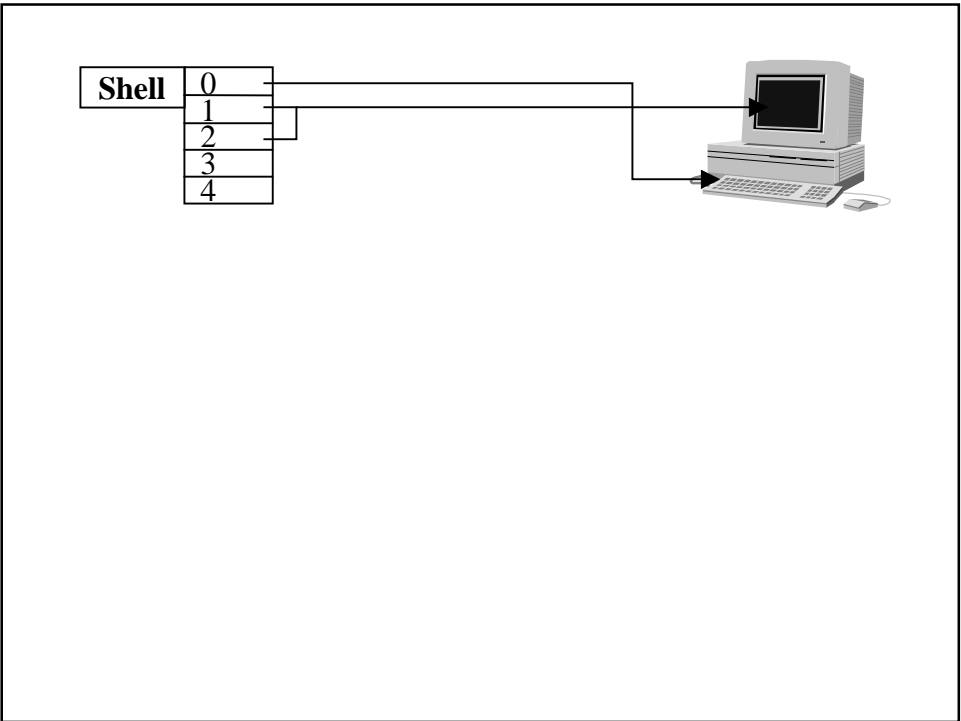
    pclose( fp );
    exit(0);
}

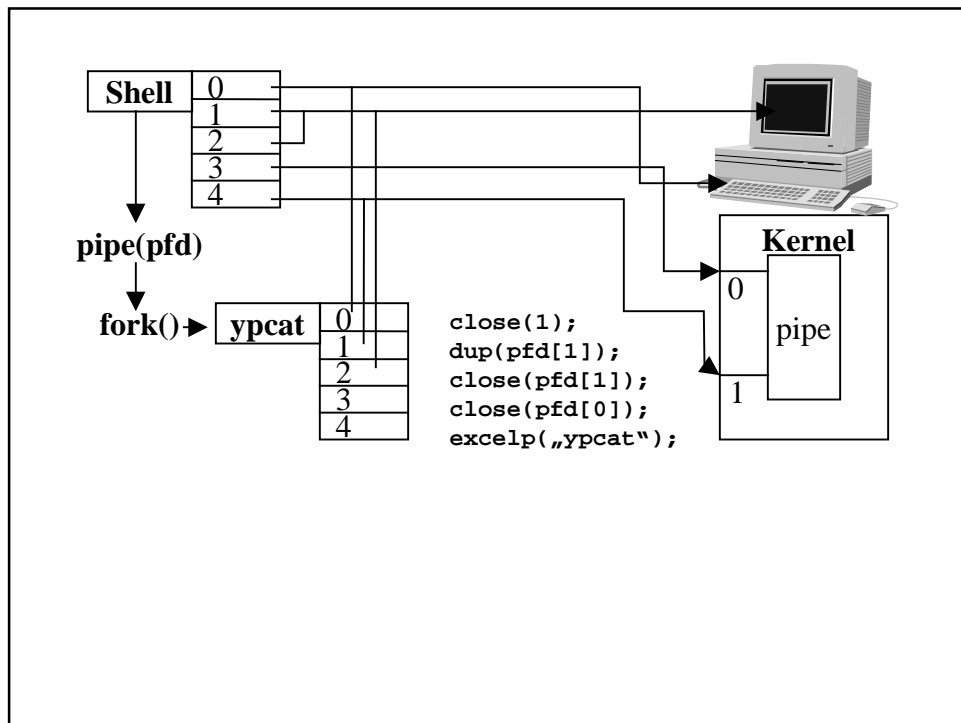
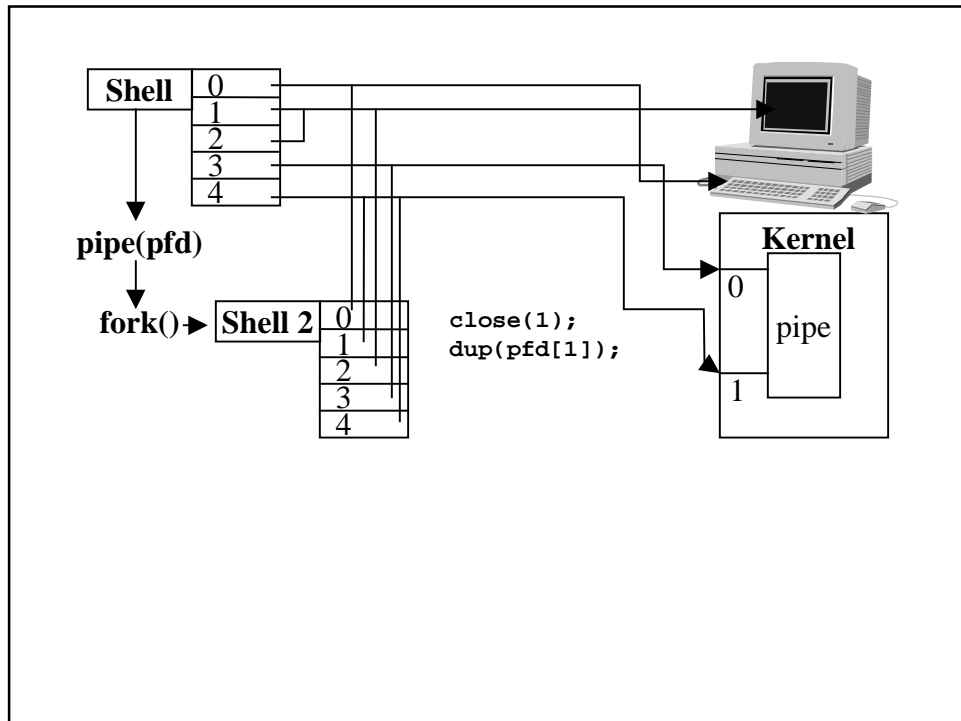
```

## 2.7 Pipes in der Shell

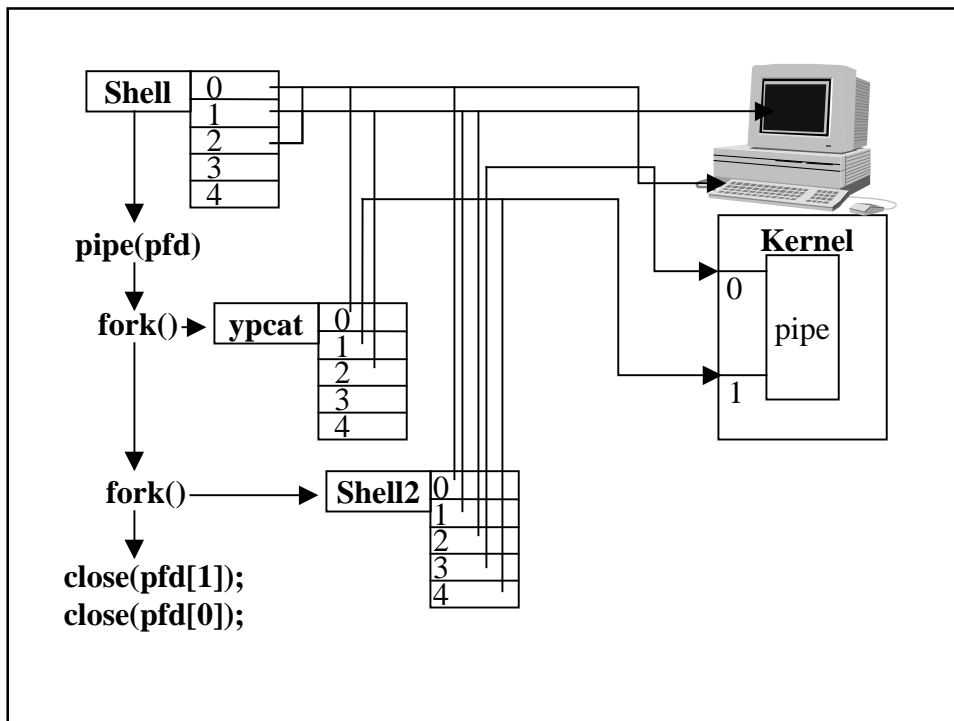
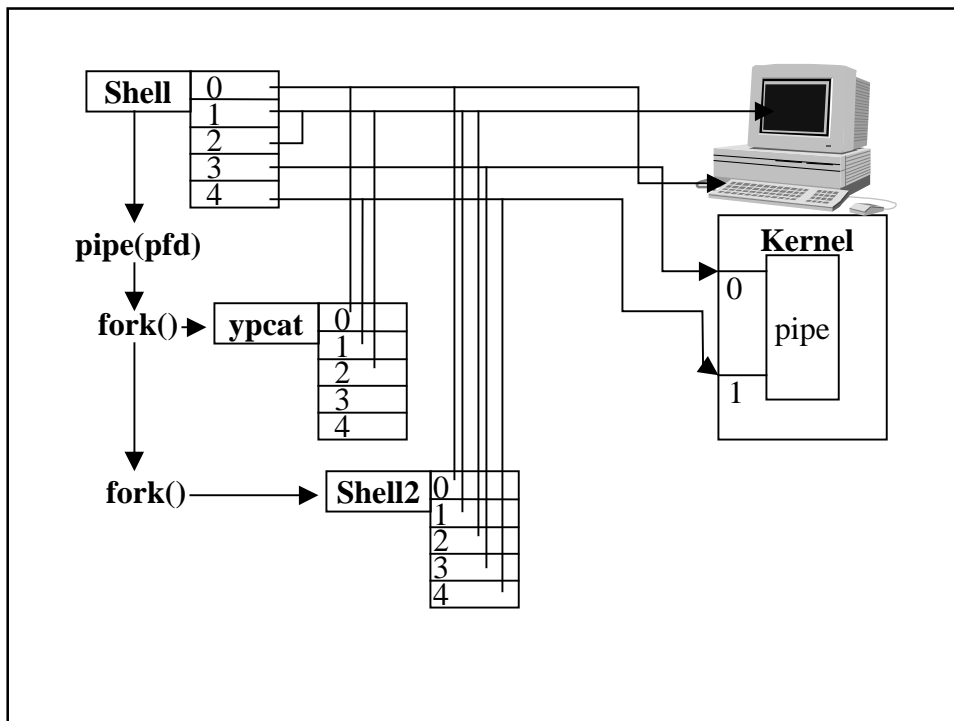
In der Midi-Shellimplementierung fehlen noch die Pipes. Das folgende Beispiel zeigt den Ablauf in der Shell für die Kommandozeile

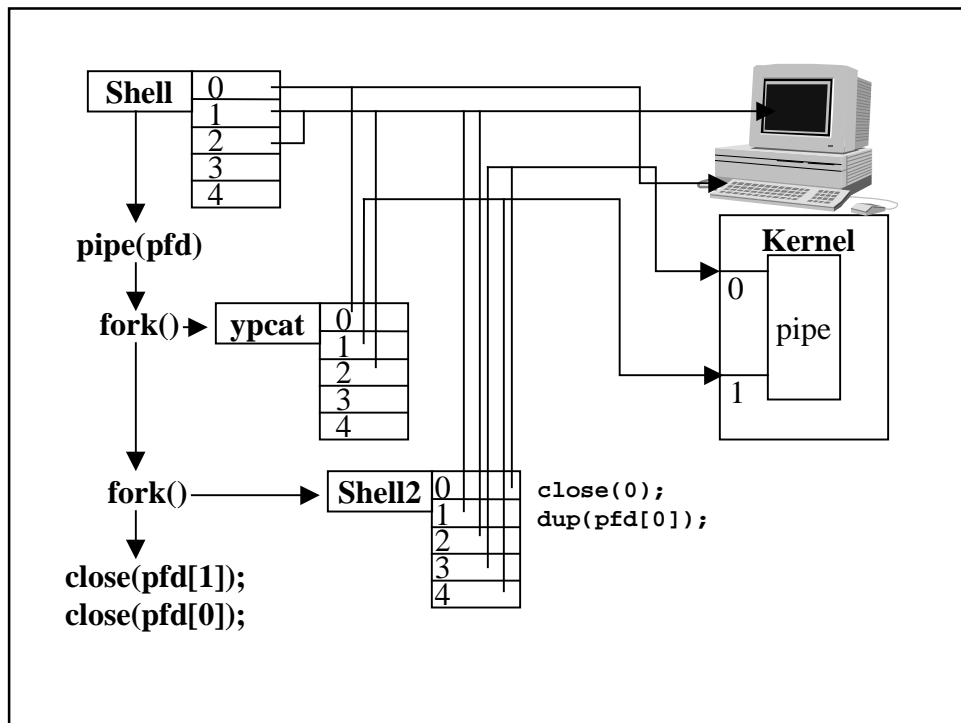
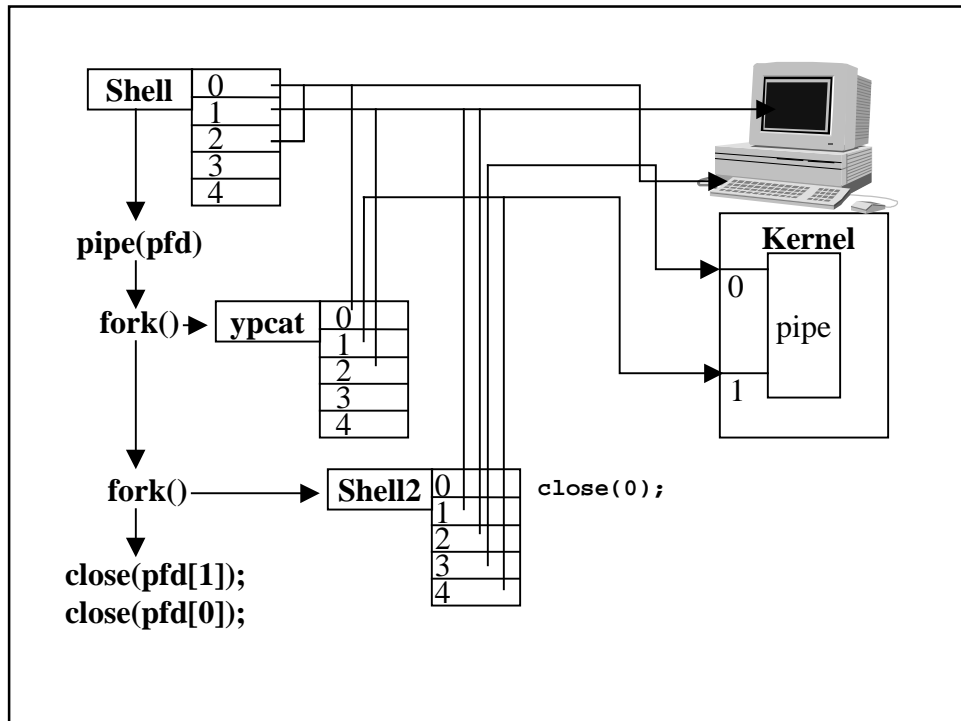
`ypcat passwd|grep grabert:`

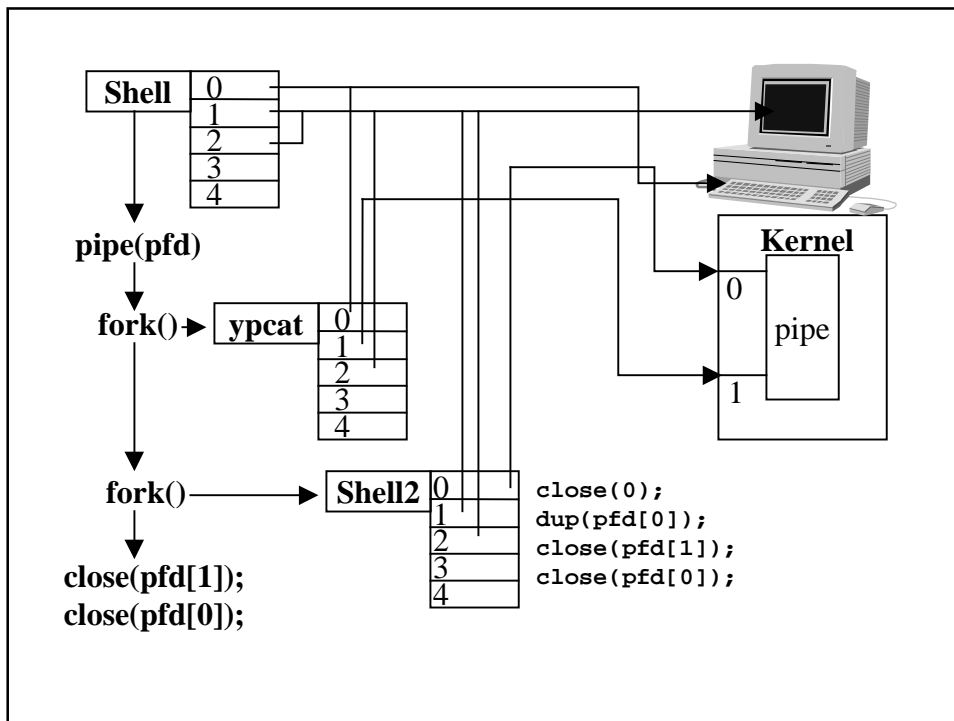
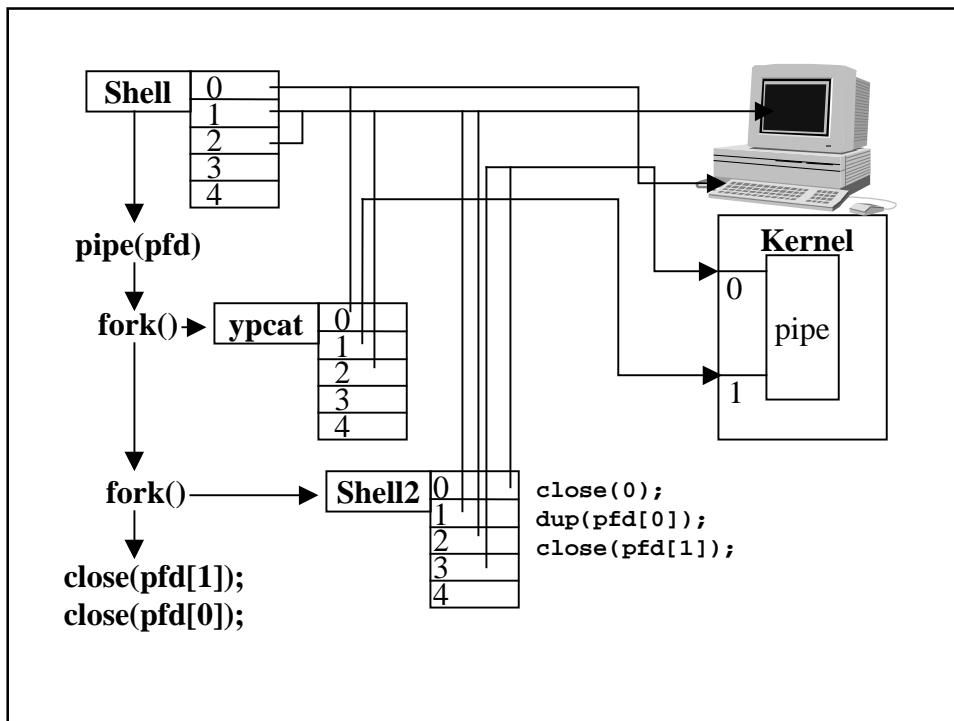


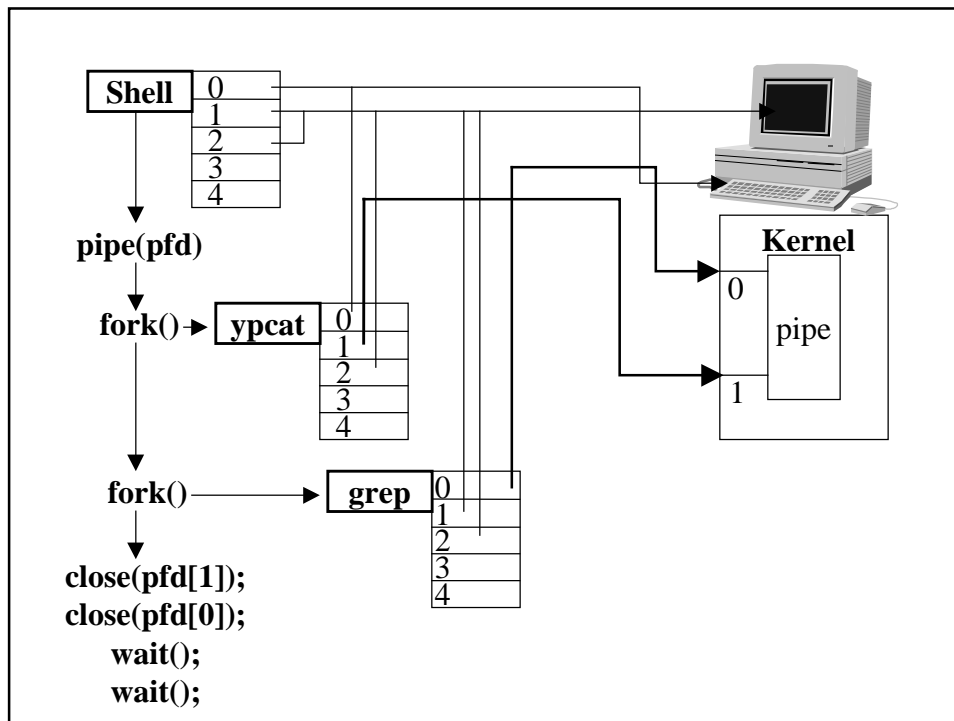
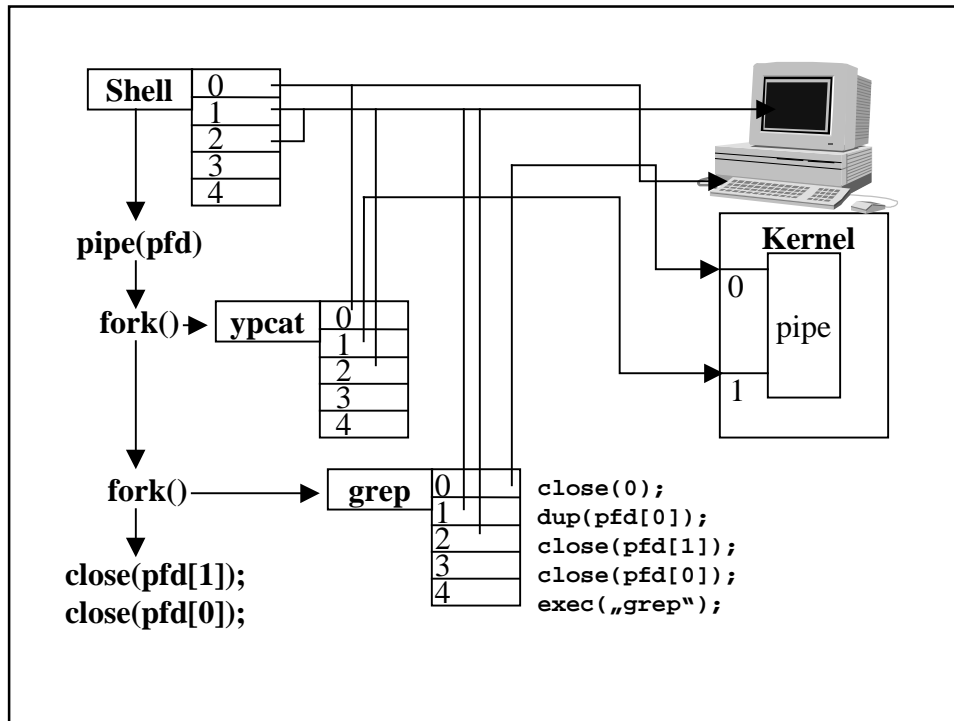










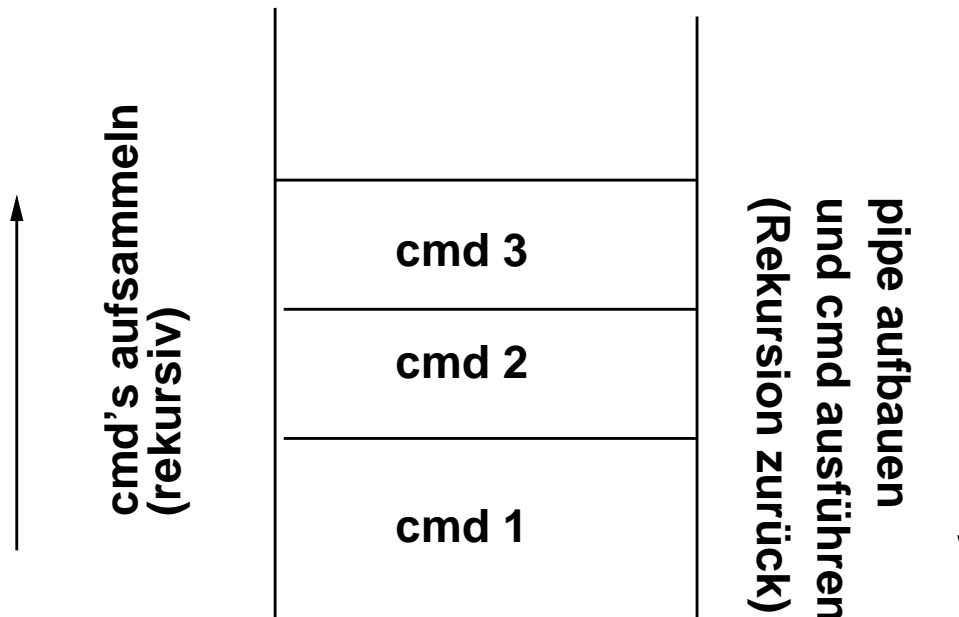


In der 'wirklichen Shell' mit beliebig vielen aneinandergereihten Pipes stellt sich die Frage, wie soll man die Kommandozeile (rekursiv) aufbauen (von vorne oder von hinten her ??)

beides ist möglich - hier im Beispiel von 'hinten' her! → Funktion `command`

```
[< infile] cmd1 [arg1, ...] | cmd2 [arg1, ...] |  
    cmd3 [arg1, ...] [> outfile | >> outfile]
```

## Stack via Rekursion



## 2.8 Nicht behandelte IPC-Mechanismen

- **named pipes (FIFO-Dateien)**

ähnlich den *unnamed pipes*, aber Eintrag in *directory*

Erzeugung:

```
int mknod(char * pathname, int mode, int dev);
```

**/etc/mknod** *pathname* **p**

für FIFO's entfällt das Argument *dev*; siehe *man mknod*

- **Message Queues**

siehe `<sys/ipc.h>`, `<sys/msg.h>`

- **Semaphore**

Mechanismus zur Synchronisation, weniger zum Datenaustausch

siehe `<sys/ipc.h>`, `<sys/sem.h>`

- Shared Memory

siehe `<sys/ipc.h>`, `<sys/shm.h>`

# Kapitel 3

## Netzwerk-Kommunikation

### 3.1 Übersicht

- räumlich verteiltes System von Rechnern, Steuereinheiten und Peripheriegeräten, verbunden mit Datenübertragungseinrichtungen
- aktive / passive Komponenten

#### Ausbreitung

- globales Netz (**GAN** *global area network*)  
Internet, EUNet, VNET (IBM), u.a.
- Weitverkehrsnetz (**WAN** *wide area network*)  
DATEV-Netz, DFN (Deutsches Forschungsnetz)
- lokale Netze (**LAN** *local area network*)  
innerhalb eines Unternehmens; i.a. mehrere, die selbst wieder vernetzt sind
  - ☞ **Front-end-Lan** (Netz innerhalb einer Abteilung / Instituts)
  - ☞ **Backbone-LAN** (Verbindung von Front-end-LAN's)

#### Netzwerktopologie

- physikalische / logische Verbindung der Rechner im Netz (Stern, Ring, Bus)

#### Protokolle

- Regeln (Vereinbarungen), nach denen Kommunikationspartner (Rechner) eine Verbindung aufbauen, die Kommunikation durchführen und die Verbindung wieder abbauen

Grundlegendes herstellerunabhängiges Konzept:

### DIN/ISO-OSI-Referenzmodell

- **OSI - Open Systems Interconnection**

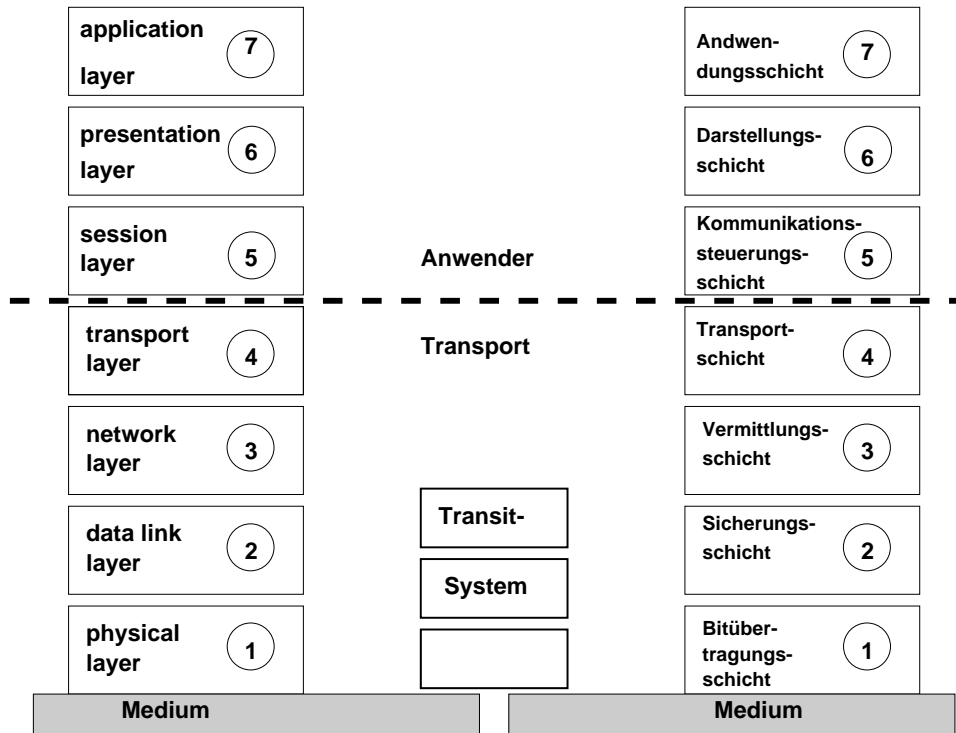


Abbildung 3.1: ISO-OSI-Referenzmodell

- 1. Bitübertragungsschicht**  
Regelung aller physikalisch-technischer Eigenschaften der Übertragungsmedien zwischen den verschiedenen End-/Transitsystemen  
Darstellung von Bits via Spannungen, Stecker, ...
- 2. Sicherungsschicht**  
Sicherung der Schicht **1** gegen auf den Übertragungstrecken auftretenden Übertragungsfehler (elektromagnetische Einflüsse)  
z.B. Prüfziffern, *parity bits*
- 3. Vermittlungsschicht**  
Adressierung der Zielssysteme über das (die) Transitsystem(e) hinweg sowie Wegsteuerung der Nachrichten durch das Netz  
Flußkontrolle zwischen End- und Transitsystemen (Überlastung von Übertragungswegen und Rechnern / Transitsystemen, faire Verteilung der Bandbreite)
- 4. Transportschicht**  
Stellt die mithilfe der Schichten **1 2 3** hergestellten Endsystemverbindungen für die Anwender zur Verfügung  
z.B. Abbildung logischer Rechnernamen auf Netzadressen



5. **Kommunikationssteuerungsschicht**  
Bereitsstellung von Sprachmitteln zur Steuerung der Kommunikationsbeziehung (*session*)  
Aufbau, Wiederaufnahme nach Unterbrechung, Abbau
6. **Datendarstellungsschicht**  
Vereinbarungen bzgl. Datenstrukturen für Datentransfer
7. **Anwendungsschicht**  
Berücksichtigung inhaltsbezogener Aspekte (Semantik)

### Quasistandard

- **TCP/IP**  
*Transmission Control Protocol / Internet Protocol*

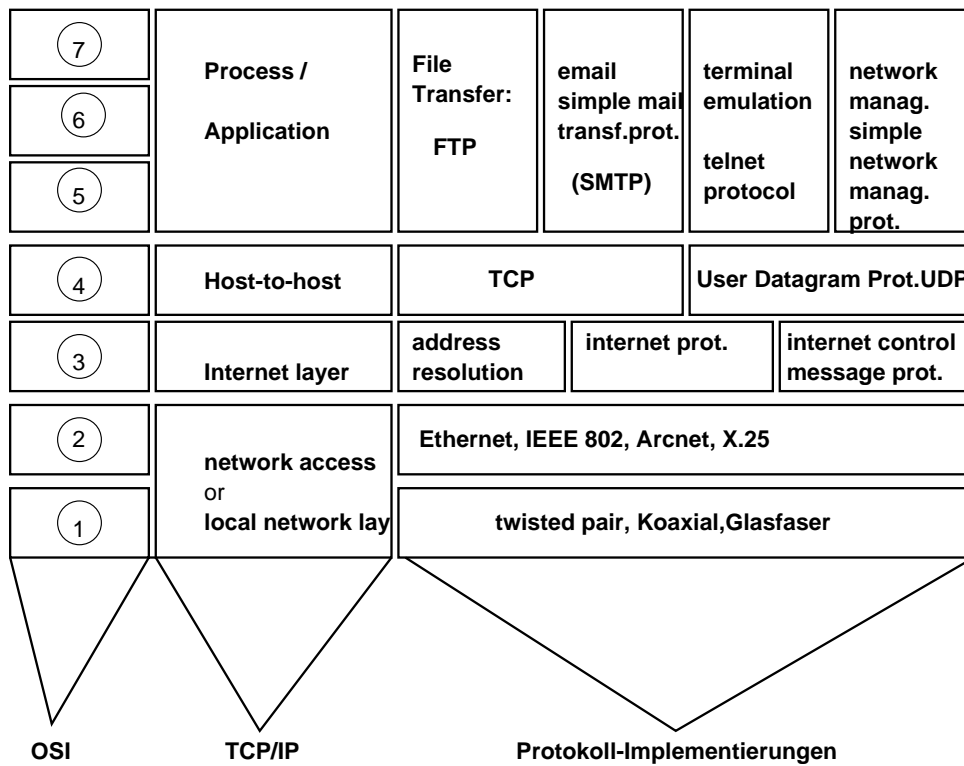


Abbildung 3.2: TCP/IP

### Kopplung von Netzen

- Kopplungseinheiten zur Verbindung von Netzen (*internetworking units*)  
Adreßumwandlung, Wegewahl (*routing*), Flußkontrolle, Fragmentierung und Wiederezusammenfügung von Datenpaketen, Zugangskontrolle, Netzwerkmanagement
  - **Repeater**  
Verstärker; Empfangen, Verstärken, Weitersenden der Signale auf Bitübertragungsschicht; die zu verbindenden Netze müssen identisch sein; Verbindung von Netzsegmenten
  - **Bridge**  
Verbindung von Netzen mit unterschiedlichen Übertragungsmedien, aber mit gleichem Schichtaufbau; operiert auf Sicherungsschicht
  - **Router**  
operiert auf Vermittlungsschicht
  - **Gateway**  
Verknüpfung von Netzen, die in Schicht 3 (und aufwärts) unterschiedliche Struktur aufweisen
  - **Hub**  
Eine Art Multiplexer, der das Eingangssignal sternförmig an die angeschlossenen Geräte weiterleitet (keine eigene Adresse); reduziert den Verkabelungsaufwand

## 3.2 Lokale Netze

### Zugangsverfahren

→ wer darf wann senden ?

- a) über strenge Vorschrift wird festgelegt, wer wann senden darf
- b) jeder sendet wie er will, bis Fehler auftritt, der dann korrigiert wird

ad a) Token-Verfahren

ein besonderes Bitmuster (*Token*) "kursiert" im Netz  
senden darf der, der es besitzt  
das Token wird an das Ende der Sendung angefügt

ad b) CSMA/CD-Verfahren

*carrier sense multiple access with collision detection*

1. viele beteiligte Sender (*multiple access*)
2. vor dem Senden in den Kanal "horchen" (*carrier sense*)  
wenn frei, senden, sonst warten
3. während des Sendens den Kanal prüfen, ob andere senden, um Kollisionen zu erkennen (*collision detection*)  
wenn Kollision, müssen alle Sender abbrechen; jeder wartet eine zufällig gewählte Zeitspanne und wiederholt Sendevorgang - Sender mit der kürzesten Zeitspanne "gewinnt"

### 3.3 LAN-Standards

#### IEEE 802.3 CSMA/CD - Bussystem (ISO 8802-2)

→ ETHERNET

relativ kostengünstig, weit verbreitet;

ursprünglich war *Ethernet* als Standard Anfang der 80er Jahre von Xerox, DEC

und Intel erarbeitet worden; ISO-8802-3 ist eine Erweiterung hiervon

ursprüngliche Version: Übertragungsrate bis 10MBit/s, inzwischen sind 1GBit/s realisiert und 10GBit/s geplant

#### Kabeltypen:

- "Thick Ethernet" (Yellow Cable)
- "Thin Ethernet" (RG-58)
- Koaxialkabel
- verdrehte Kupferkabel (*twisted pair*)

#### IEEE 802.5 Tokenring mit Token-Zugangsverfahren (ISO 8802-5)

- verdrehte Kupferkabel (1-4MBit/s)
- Koaxialkabel (4-40MBit/s)

#### ISO-8802-4 Tokenbus

logischer Ring (Token-Zugang) auf physikalischem Bus

Anwendung: MAP-Protokolle (*manufacturing automation protocol*), zur Vernetzung von Robotern, CNC-Maschinen

#### ISO-8802-8 FDDI *fiber distributed data interface*

für Nahverkehrsnetze mit hoher Übertragungsgeschwindigkeit (100 MBit/s), meist als *Backbone* eingesetzt

#### IEEE-802.11 WLAN *wireless lan*

drahtlose Datenübertragung via Funkverbindung; Reichweite: ca. 30-50m (je nach Umgebung); Übertragungsrate: 11MBit/s (durch Kanalbündelung auch 22MBit/s); Sterntopologie

### 3.4 Ethernet

(LAN)

- ca. 1970 XEROX PARC
- Standardisiert 1978 von XEROX, INTEL, DEC
- Kabel: koaxial
  - ☞ jeweils am Ende des Kabels ein Widerstand zur Vermeidung von Reflexionen elektrischer Signale
  - ☞ Kabel rein passiv
    - Anbindung an Ethernet → "T"-Stück

- elektronische Komponenten:
  - **Transceiver** (Übertragen/Empfangen in/von Ethernet)
  - **host interface** (Rechner-Bus)

□ Eigenschaften:

- alle teilen sich einen Kanal (BUS)
- **broadcast** - alle Transceiver hören alles, host interface stellt fest, ob Nachricht für diesen Rechner gedacht ist
- **best-effort delivery** - "Bemühensklausel": ob Sendung wohl ankommt?
- CSMA/CD-Zugangsverfahren

□ Varianten:

- twisted pair ethernet → "Telefondrähte"
- thin-wire ethernet → "Kabelfernsehen"  
geringere Leistung, geringerer Preis, auch für interface
- Breitbandtechnik → multiplexing (weniger Kabel)

□ Adressierung

host interface bildet Filter, alle Pakete werden dahin weitergeleitet, nur die dem Transceiver entsprechenden Pakete (Hardware-Adressen) werden an Rechner weitergeleitet

Adresse eines Rechners: 48 Bit Integer, vom Hersteller auf Interface festgelegt, von IEEE gemanaged

- Adreßtypen:
  - physische Adresse einer Schnittstelle
  - network broadcast address (alle Bits auf 1, "an alle")
  - multicast broadcast (Teilmengen broadcast)  
BS initialisiert Schnittstelle: welche Adressen sollen erkannt werden?

**Frame Format:**

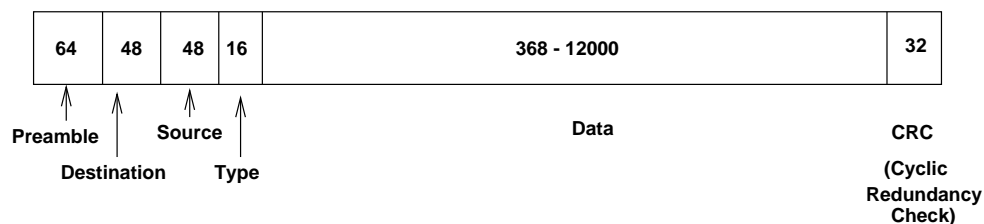


Abbildung 3.3: Ethernet Frame Format

- Preamble: zur Synchronisation der Knoten, alternierende 0/1-Folge
- Type: welches Protokoll (für BS)? → self-identifying

## 3.5 Internetworking - Concept & Architectural Model

- bislang:  
ein Netz (ein physikalisches Netz) - z.B. Ethernet, Token Ring, Adressen von Hosts waren physikalische Adressen
- jetzt:  
 "Netz über Netzen über Netzen über ... über physikalischen Netzen"
  - ☞ notwendig:  
 Abstraktion von zugrundeliegenden physikalischen Netzen  
Ansatz 1: spezielle (Applikations-) Programme, die aus der Heterogenität der physikalischen Netze / Hardware eine softwaremäßige Homogenität herstellen  
Ansatz 2: Verbergen von Details im Betriebssystem des jeweiligen Rechners, layered-system Architecture

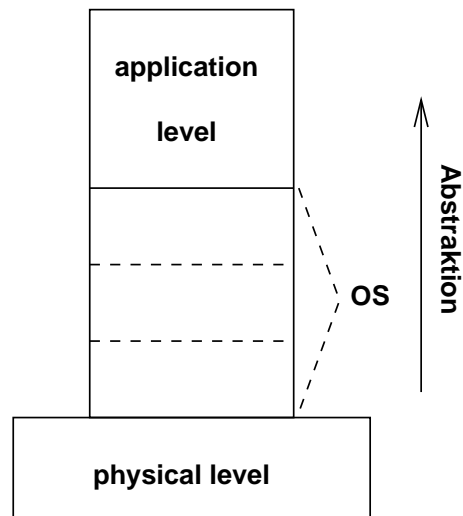


Abbildung 3.4: Layered System Architecture

### Internet Architecture

□ *In a TCP/IP internet, computers called gateways provide all interconnections among physical networks*

Frage: Muß ein Gateway alle in allen Netzen erreichbare Rechner kennen (ein Super-Computer)?

Antwort: Nein, *gateways route packets based on destination network, not on destination host* (alle Netze müssen i.P. bekannt sein, → lernende Gateways, → Mini-Computer)

### Benutzer-Sicht:

ein Internet als ein großes, virtuelles Netzwerk

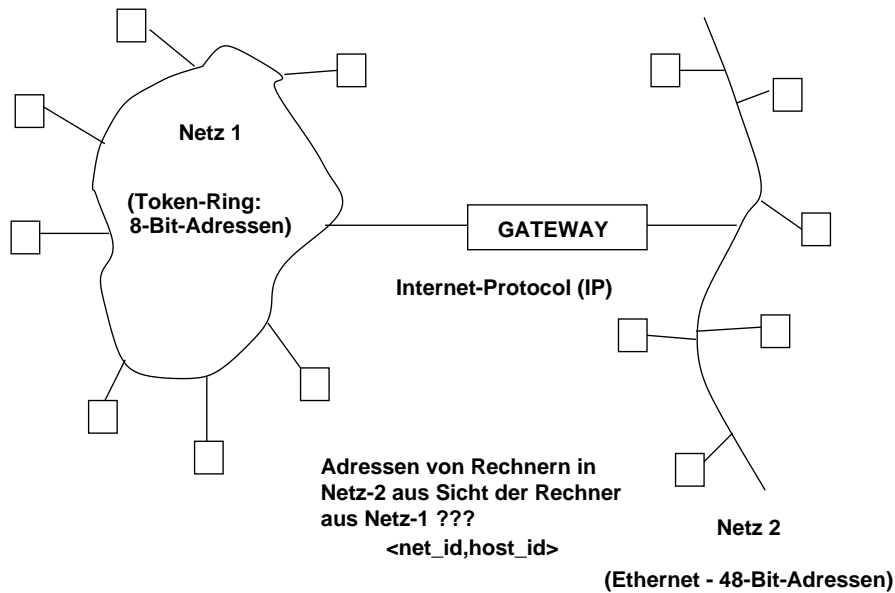


Abbildung 3.5: Internet Architecture

□ The **TCP/IP** internet protocols treat all networks equally (independant of their delay and throughput characteristics, maximum packet size or geographic scale). A local area network like *Ethernet*, a wide area network like *NSFNET* backbone, or a point-to-point link between two machines each count as one network

#### Internet-Adressen

- globale Identifikation von hosts im (in ihrem) Netz

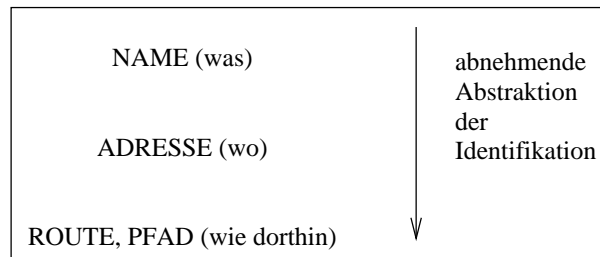


Abbildung 3.6: Abstraktion der Identifikation

- **IP-Adresse (logische Adresse):**  
Integer, die die Route unterstützen (32-Bit): (**netId,hostId**)

**A:** wenige Netze mit sehr vielen Hosts (mehr als  $2^{16}$  Hosts)

**B:** mehr Netze mit ziemlich vielen Hosts (zwischen  $2^8$  und  $2^{16}$  Hosts)

**C:** viele Netze mit wenigen Hosts (weniger als  $2^8$  Hosts)

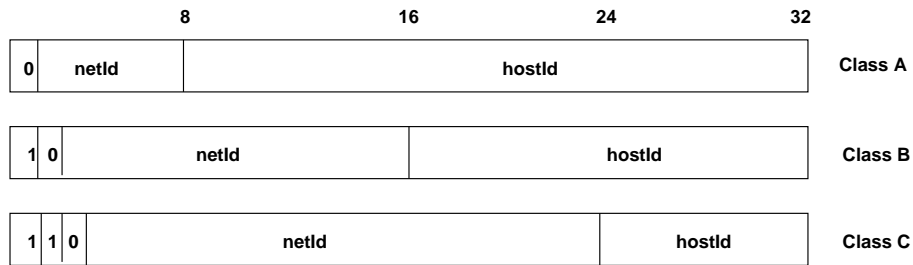


Abbildung 3.7: 32-Bit IP-Adresse

NB: IP-Adresse spezifiziert die Verbindung eines Hosts zu einem Netz!

NB: HostId == 0: spezifiziert Netz selbst

NB: HostId == 1: Broadcast an alle Host's des Netzes (falls dieses Netz Broadcast kennt)

lesbare Adressen ( **punktiertes Dezimalformat**): (*dotted decimal notation*)

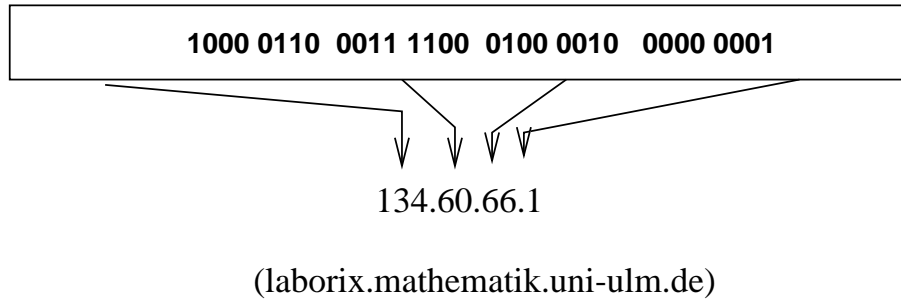


Abbildung 3.8: IP-Adresse: dotted decimal form

#### Adressenvergabe:

- NIC (*Network Information Center*) vergibt NetId
- Verantwortung für HostId delegiert an jeweilige Organisation
  - Problem: lowest-Byte (Bit)  $\leftrightarrow$  highest-Byte (Bit)

Protocol network standard **byte order** : most significant byte send first!

Internet Adresse  $\rightarrow$  physikalische Adresse: **ARP** - **A**ddress **R**esolution **P**rotocol  
 IP-Adresse als "physikalische" Adresse in einem virtuellen Netzwerk

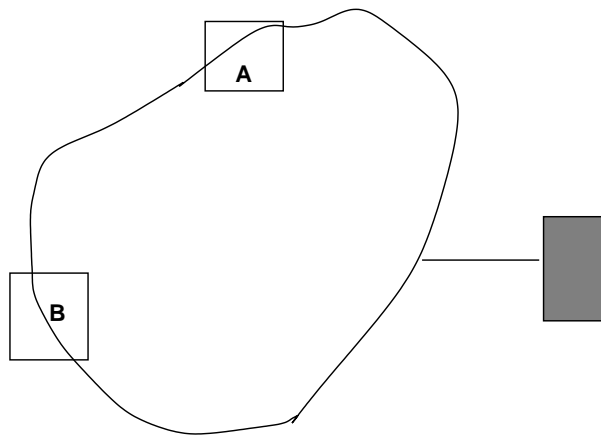


Abbildung 3.9: Adress-Auflösung

$I_A, I_B$  Internet-Adressen (netId, hostId)

$P_A, P_B$  physikalische Adressen (z.B. Ethernet)

Address resolution problem:

$I_A \rightarrow P_A$  ??? (low level software)



Ethernet	48 Bit
proNet	8 Bit
IP-hostId	8 .. 24 Bit

- proNet:  $P_{\text{Adr}} \subseteq I_{\text{Adr}}$  (kein Problem)
- X.25: Tabelle  $\subseteq I_{\text{Adr}} \times P_{\text{Adr}}$  + hashing (Zuordnungstabelle)
- Bei Ethernet:
  - A will an B senden, kennt aber nur  $I_B$  und nicht  $P_B$ 
    1. A sendet via *broadcast* an alle mit  $I_B$ , B kennt (hoffentlich) seine eigene Internet-Adresse  
 → **broadcasting ist sehr teuer**, insbesondere auch weil ein IP-Paket u.U. in sehr viele Pakete auf dem Ethernet zerlegt wird (später: Fragmentation)
    2. A sendet *broadcast*: "Rechner mit Internet-Adresse  $I_B$  möge sich melden mit  $P_B$ "  
 B meldet sich mit  $(I_B, P_B)$  bei A  
 A merkt sich für die Zukunft  $(I_B, P_B)$

#### TCP/IP

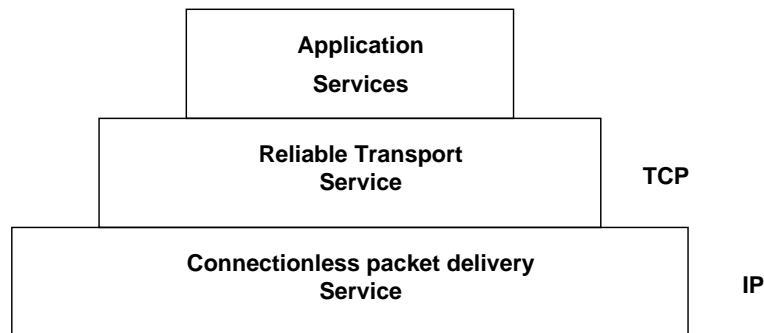


Abbildung 3.10: TCP/IP-Schichtenmodell

1. Concept of **Unreliable Delivery**  
 Auslieferung von Paketen ist nicht garantiert (*may be lost, duplicated, delayed, delivered out of order* → *service will not detect nor inform sender or receiver*)
2. **Connectionless**  
 Jedes Paket wird losgelöst von den anderen behandelt (evt. auch anders gerouted)
3. **best-effort delivery**  
 "bemühe mich um Auslieferung"

#### IP Internet Protocol:

- Regeln, die den *unreliable, connectionless, best-effort delivery* - Mechanismus definieren
  1. *basic unit of data transfer* durch ein TCP/IP-Internet
  2. enthält routing-Funktion (Auswahl des Pfades)
  3. Regeln, wie Host's und Gateway's Pakete verarbeiten sollen, wie und wann Fehlermeldungen produziert werden sollen, Bedingungen für das "Wegwerfen" von Paketen

### Internet Datagram - *basic transfer unit*

- grober Aufbau:

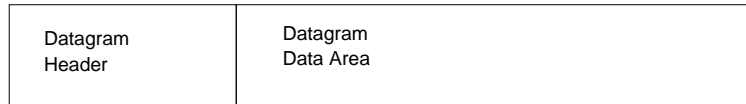


Abbildung 3.11: IP Datagramm - grober Aufbau

- genauer:

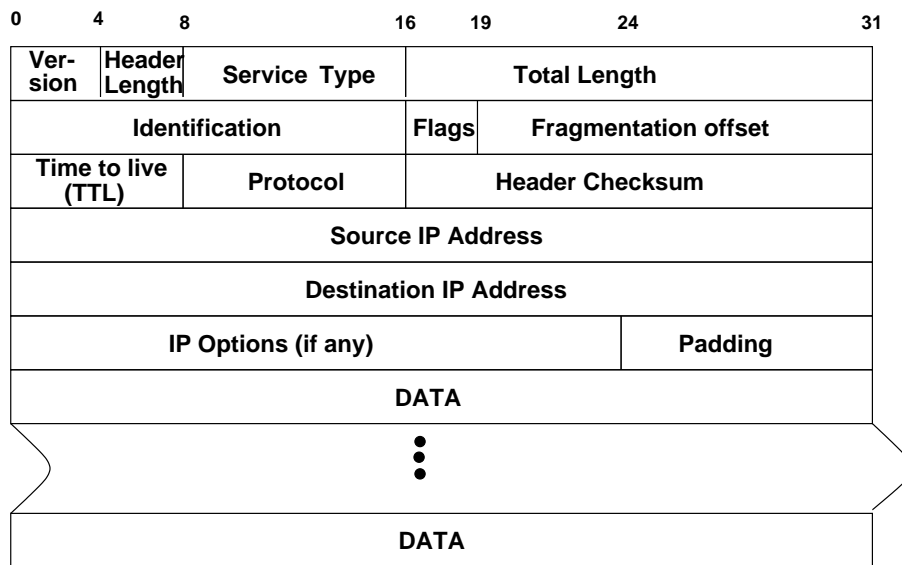


Abbildung 3.12: IP Datagram - im Detail

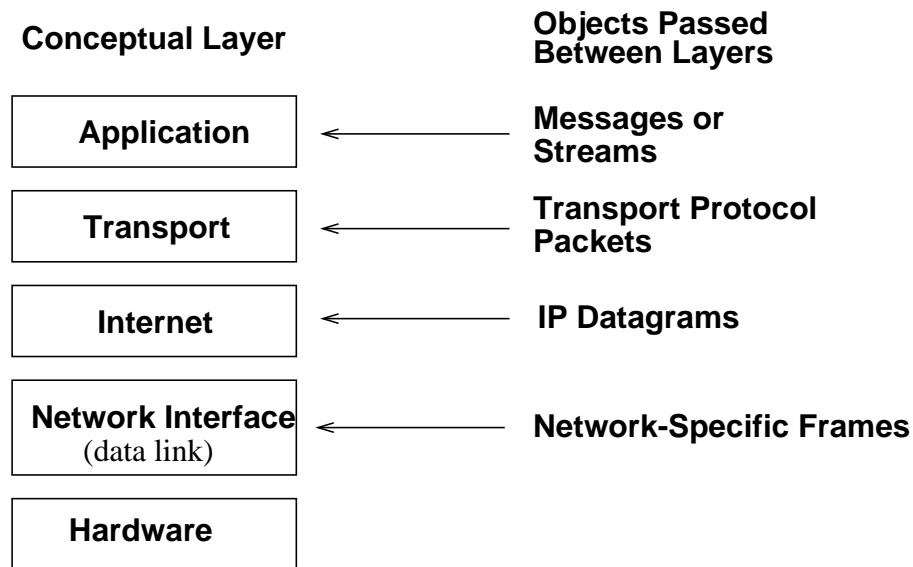
**TCP/IP Internet Layering Model**

Abbildung 3.13: TCP/IP Layering Model

**Application Layer**

Auf der obersten Ebene starten Benutzer Anwendungsprogramme, die auf im Netz verfügbare Dienste zugreifen wollen. Ein Anwendungsprogramm interagiert mit dem / den Transport-Protokoll(en), um Daten zu senden / empfangen.

Jedes Anwendungsprogramm wählt die benötigte Transport-Art, z.B. eine Folge individueller Botschaften (*messages*) oder einfach eine Folge von Bytes. Das Anwendungsprogramm übergibt diese Daten in der verlangten Form an die Transport-Ebene zur Auslieferung.

**Transport Layer**

Die Hauptaufgabe dieser Schicht besteht darin, die Verbindung zur Kommunikation zwischen zwei Anwendungsprogrammen bereitzustellen (*end-to-end-communication*).

- Regulieren des Informationsflusses
- Bereitstellen eines zuverlässigen Transports
- Sicherstellen, dass Daten korrekt und in Folge ankommen
- Warten auf Empfangsbestätigung des Empfängers
- erneutes Senden verlorengegangener Pakete
- Aufteilen des Datenstroms in kleine Stücke (*packets*)

- Übergeben jeden Paketes mit Zieladresse für die nächste Schicht
- i.a. arbeiten viele Applikation mit der Transport-Schicht (Senden, Empfangen)
- jedem Paket wird zusätzliche Information beigefügt (welche Appl.?)
  - Prüfsumme

### ***Internet Layer***

- Erstellen der IP Datagrams
- (weiter-) senden der Datagrams
- "auspacken" der Datagrams (auf Zielmaschine)
- Übergabe an das richtige Transport-Protokoll (Zielmaschine)
- ICMP

### ***Network Interface Layer***

- übernimmt Datagrams und schickt sie auf spezifischem Netz weiter

## **3.6 Transport-Protokolle**

### **3.6.1 Ports**

- Internet Protokolle adressieren *Host's*
- Adressierung von Applikationen (letztliches Ziel) innerhalb des Zielrechners?

Unterstellt seien *multiprocess*-Systeme als Zielrechner (z.B. UNIX-Rechner).

- Prozess als letztlisches Ziel?  
Prozesse werden dynamisch erzeugt und terminiert, Sender kann Prozesse auf anderen Maschinen nicht identifizieren  
Dienste (sprich Funktionen, Leistungen) auf anderen Rechnern sind das Ziel, unabhängig von welchem Prozess (welchen Prozessen) diese realisiert sind!
- Jede Maschine hat eine Menge sog. Protokoll- **Port's** (abstrakter Endpunkt), identifiziert durch positive *Integer*. Das jeweilige Betriebssystem bietet Mechanismus an, mit dem Prozesse Ports spezifizieren und nutzen können.  
Die meisten Betriebssysteme unterstützen synchronen Zugriff auf Ports. Wenn dabei eine Prozess Daten von einem Port lesen will, noch keine Daten da sind, wird er solange blockiert, bis (genügend) Daten eingetroffen sind.  
Ports sind typischerweise gepuffert.
- Ziel-Adresse: (IP Adresse + Port-Nummer)  
Jede Meldung enthält neben *destination port* auch *source port* (z.B. zum Antworten).

### 3.6.2 UDP

#### Format der UDP-Messages:

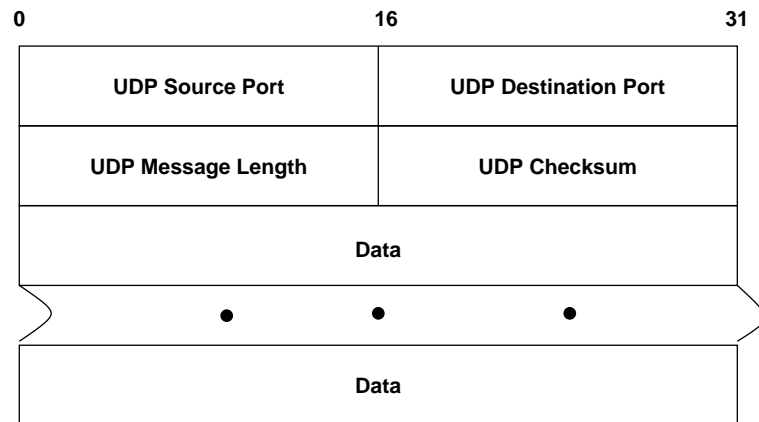


Abbildung 3.14: UDP - Format

□ *The User Datagram Protocol (UDP) provides unreliable connectionless delivery service using IP to transport messages between machines. It adds ability to distinguish among multiple destinations within a given host computer*

- keine Empfangsbestätigungen (*Acknowledgement*)
- Kein Ordnen ankommender Meldungen
- keinerlei Feedback
  - *UDP messages* können verloren gehen, dupliziert werden, ungeordnet ankommen, schneller ankommen als verarbeitet werden!
  - Anwendungsprogramme, die auf UDP aufbauen, müssen selbst für "Zuverlässigkeit" sorgen.
- *Source Port* ist optional, wenn nicht genutzt, so NULL
- *LENGTH* Anzahl Bytes (*octets*) im UDP Datagram inkl. Header.
- *CECKSUM* ebenfalls optional; wenn keine Prüfsumme berechnet, so NULL.  
**NB:** IP berechnet keine Prüfsumme über den Datenteil!  
 Berechnung der Prüfsumme (wie bei IP): Daten werden in 16-Bits-Einheiten aufgeteilt, Summe über 1-er-Komplement wird gebildet, davon wieder 1-er-Komplement gebildet; ist Prüfsumme identisch Null, so wird davon das 1-er-Komplement gebildet (alles auf 1); unproblematisch, da es bei 1-er-Komplement zwei Darstellungen des Zahl Null gibt: alle Bits auf 0 oder alle auf 1!

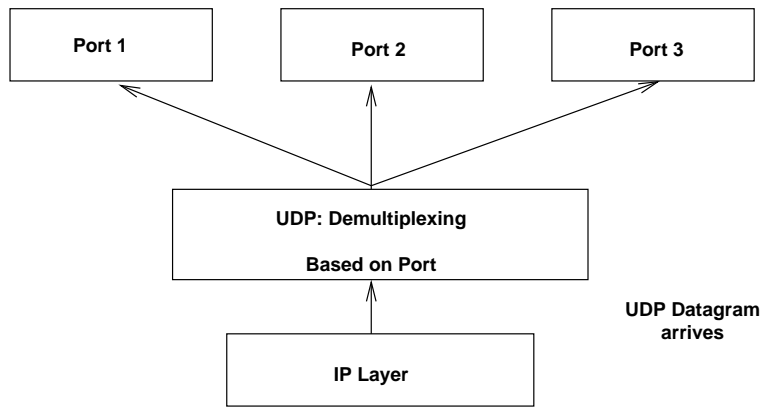


Abbildung 3.15: UDP-Demultiplexing

**Port-Nummern:**

- einige zentral vergeben (*well-known port numbers, universal assignment*)
- *dynamic binding approach*  
Port ist nicht global bekannt; wenn ein Programm einen Port benötigt, wird von der Netzwerk-Software einer bereitgestellt. Um eine Port-Nummer auf einem anderen Rechner zu erfahren, wird eine entsprechende Anfrage gestellt und beantwortet.

Decimal	Keyword	UNIX Keyword	Description
0	-	-	Reserved
7	ECHO	echo	Echo
9	DISCARD	discard	Discard
11	USERS	systat	Active Users
13	DAYTIME	daytime	Daytime
37	TIME	time	Time
42	NAMESERVER	name	Host Name Server
43	NICNAME	whois	Who Is?
53	DOMAIN	nameserver	Domain Name Server
67	BOOTPS	bootps	Bootstrap Protocol Server
68	BOOTPC	bootpc	Bootstrap Protocol Client
69	TFTP	tftp	Trivial File Transfer
111	SUNRPC	sunrpc	Sun Microsystems RPC
123	NTP	ntp	Network Time Protocol
161	-	snmp	SNMP net protocol
.			
.			
.			

**Beispiel einer TCP/IP Sitzung an Port 80 (Webserver):**

```
thales$ telnet www.mathematik.uni-ulm.de 80
GET /sai/beutten/ HTTP/1.0          <<< Diese Zeile wird eingetippt
                                     <<< zweites Return!
HTTP/1.1 200 OK                     <<< ab hier: Antwort des Servers
Date: Wed, 01 Aug 2001 09:35:30 GMT
Server: Apache/1.3.9 (Unix)
Last-Modified: Thu, 07 Jun 2001 11:31:17 GMT
ETag: "6d8ca-398-3blf6605"
Accept-Ranges: bytes
Content-Length: 920
Connection: close
Content-Type: text/html

<html><HEAD>
<TITLE>
Universit&auml;t Ulm, SAI, Gerhard Beuttenm&uuml;ller
</TITLE>
<LINK REV="made" HREF="mailto:beutten@mathematik.uni-ulm.de">
</HEAD>
<BODY>
...
</BODY>
</HTML>
thales$
```





# Kapitel 4

## Berkeley Sockets

### 4.1 Grundlagen

- **API** (*application program interface*) zu Kommunikations-Protokollen
- In UNIX: *Berkeley sockets* und *System V Transport Layer Interface (TLI)*
- Beide Schnittstellen wurden für C entwickelt
- Die Implementierungen der Internet-Protokolle und der Socket-Schnittstelle wurden erstmals 1982 in der Version 4.1c der *Berkeley Software Distributions (BSD)* an der Universität von Kalifornien in Berkeley integriert.
- Mit der Version 4.2BSD war 1983 das erste Unix-System mit Netzwerkunterstützung verfügbar.

Die beiden zentralen Abstraktionen der Socket-Schnittstelle sind der *Socket* und die **Kommunikationsdomäne**. Aus Sicht der Anwendung repräsentiert ein Socket einen **Kommunikationsendpunkt** eines Benutzerprozesses innerhalb der gewählten Kommunikationsdomäne. Die Domäne spezifiziert die Protokollfamilie und definiert allgemeine Eigenschaften und Vereinbarungen wie z.B. die Adressierung des Kommunikationsendpunktes.

Durch die Abstraktion Socket und Kommunikationsdomäne wird eine Trennung zwischen den Funktionen der Socket-Schnittstelle und den im System implementierten Kommunikationsprotokollen erreicht. Dem Anwender steht somit eine einheitliche Schnittstelle zur Verfügung, die verschiedene Netzwerkprotokolle und lokale Interprozesskommunikationsprotokolle gleichermaßen unterstützt. Die Implementierung ist vollständig im Betriebssystem integriert.

- Der Zugriff einer Anwendung auf die Funktionen der Socketschicht (*socket layer functions*) im Betriebssystemkern erfolgt über die Systemaufrufe der Socket-Schnittstelle. Hier findet die Umsetzung der protokollunabhängigen Operationen in die protokollspezifischen Implementierungen der darunterliegenden Protokollschicht (*protocol layer*) statt.
- Die Auswahl des konkreten Kommunikationsprotokolls wird bei der Erzeugung eines Socket über die Spezifikation der Domäne festgelegt.
- In der Protokollschicht (*protocol layer*) sind die Implementierungen der vom System unterstützten Protokollfamilien zusammengefaßt.

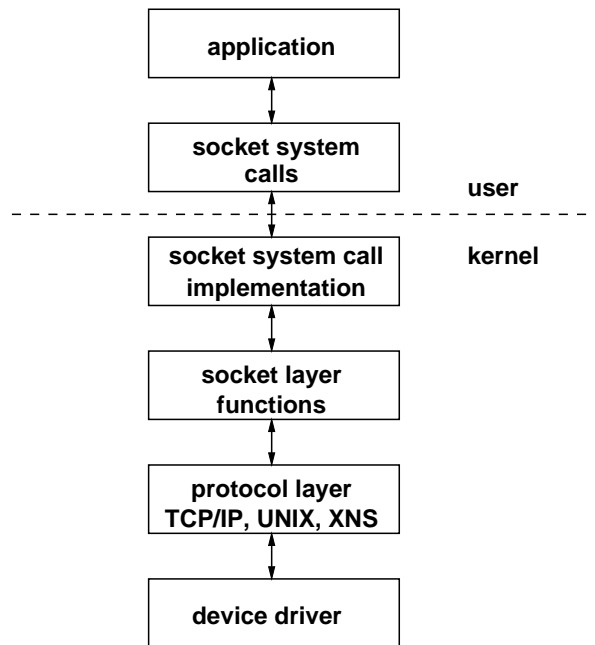


Abbildung 4.1: Vereinf. Modell der Implementierung von Sockets unter BSD

- Die darunterliegende Datenübertragungsschicht enthält die Implementierungen der Gerätetreiber (*device driver*) zur Datenübertragung über verschiedene Medien oder auch systeminterne Mechanismen.

## 4.2 Ein erstes Beispiel: Timeserver an Port 11011

```

thales$ cat timserv.c
#include <netdb.h>
#include <netinet/in.h>
#include <stdio.h>
#include <strings.h>
#include <sys/socket.h>
#include <sys/time.h>
#include <unistd.h>

#define TPORT 11011

int main()
{
    struct sockaddr_in addr, client_addr;
    size_t client_add_len = sizeof(client_addr);
    int sfd, fd;
    int optval = 1;

    bzero(&addr, sizeof(addr)); // mit 0en füllen
    addr.sin_family = AF_INET; // TCP/IP-Verbindung
    addr.sin_port = htons(TPORT); // Port eintragen, Network Byte Order

    if ((sfd = socket(AF_INET, SOCK_STREAM, 0)) < 0)
  
```

```

        perror("socket"),    // Socket erzeugen
        exit(1);
if (setsockopt(sfd,SOL_SOCKET, SO_REUSEADDR, &optval, sizeof(optval))<0)
    perror("setsockopt"),    // rasche Wiederzuweisung ermöglichen
    exit(1);
if (bind(sfd, (struct sockaddr *) &addr, sizeof(addr)) <0)
    perror("bind"),          // Socket mit Port verknuepfen
    exit(1);

if (listen(sfd, SOMAXCONN) <0)
    perror("listen"),        // Maximalanzahl Verbindungen einstellen
    exit(1);
        // accept nimmt *eine* Verbindung an
        // blockiert bis einer sich verbindet
while ((fd = accept(sfd, &client_addr, &client_add_len)) >=0) {
    time_t clock;
    char *tbuf;

    time(&clock);            // aktuelle Uhrzeit holen
    tbuf = ctime(&clock);    // in String schreiben
    write(fd, tbuf, strlen(tbuf)); // in Clientverbindung schreiben
    close(fd);
}
exit(0);                    // seldomly reached
}

```

```

thales$ timserv&
thales$ telnet thales 11011
Wed Aug  1 12:00:35 2001
thales$

```

### 4.3 Die Socket-Abstraktion

- Der Socket repräsentiert als zentrale Abstraktion einen Kommunikationsendpunkt eines Benutzerprozesses innerhalb der gewählten Kommunikationsdomäne.
- Letztere wiederum ist eine weitere Abstraktion, die allgemeine Kommunikationseigenschaften der vom System bereitgestellten Mechanismen für die Prozesskommunikation mit Sockets zusammenfaßt. Dazu zählen beispielsweise der zu verwendende Namensraum und die damit verbundene Art der Adreßspezifikation bei der Benennung eines konkreten Socket.
- Die Kommunikation zwischen je zwei Sockets verläuft normalerweise immer innerhalb derselben Domäne.

Die in der Regel von allen Unix-Systemen unterstützten und heute hauptsächlich verwendeten Kommunikationsdomänen sind die **Unix-Domäne** für die Prozesskommunikation auf dem lokalen System und die **Internet-Domäne** für die Prozesskommunikation nach den Internet-Standard-Protokollen. Die innerhalb der Internet-Domäne miteinander kommunizierenden Prozesse können

sehr wohl auch auf demselben lokalen System laufen; hierzu bietet jedoch die Verwendung der Unix-Domäne entscheidende Vorteile durch erweiterte protokoll-spezifische Eigenschaften und durch bessere Performance.

#### Socket-Typen:

- In der Socket-Abstraktion entsprechen die Socket-Typen den für die Anwendung jeweils sichtbaren Kommunikationsverfahren (**verbindungslos - verbindungsorientiert**).
- Zusätzlich definiert ein Socket-Typ spezielle Eigenschaften eines Protokolls wie z.B. fehlerfreie Datenübertragung, Flusskontrolle, Einhaltung von Datensatzgrenzen).
- Letztlich wird über den Socket-Typ die **Kommunikationssemantik** definiert.
- Wie bei den Domänen gilt: Prozesskommunikation nur zwischen Sockets vom selben Typ!

#### Die wichtigsten Socket-Typen unter Unix:

**STREAM:** Ein **Stream Socket** stellt ein Kommunikationsverfahren bereit, das einen bidirektionalen, kontrollierten und verlässlichen Datenfluss garantiert, d.h. alle transferierten Datenpakete kommen in derselben Reihenfolge vollständig und ohne Duplikate beim Empfänger an (vgl. in der Semantik zu bidirektionale Unix-Pipes). In BSD-Unix sind Pipes auf diese Weise implementiert.

**DGRAM:** Ein **Datagramm Socket** stellt ebenfalls ein bidirektionales Kommunikationsverfahren bereit, aber ohne Flusskontrolle und ohne Garantie, dass gesendete Pakete den Empfänger erreichen, dass die Reihenfolge erhalten bleibt oder dass Duplikate ankommen. Sie sind zudem im Datenvolumen beschränkt. Datensatzgrenzen bleiben beim Datentransfer allerdings erhalten.

**RAW:** Ein **Raw Socket** stellt eine allgemeine Schnittstelle zu den meist der Transportschicht zugrundeliegenden Kommunikationsprotokollen bereit, welche die Socket-Abstraktion unterstützen. Dieser Socket-Typ ist normalerweise Datagramm-orientiert, obwohl die exakte Charakteristik von der Kommunikationssemantik des konkreten Protokolls abhängt. In der Internet-Protokollfamilie kann mit Raw Sockets beispielsweise direkt das *Internet Protocol (IP)*, das **ICMP** (*Internet Control Message Protocol*) oder das **IGMP** (*Internet Group Management Protocol*) verwendet werden.

## 4.4 Die Socket-Programmierschnittstelle

### 4.4.1 Vorbemerkungen

Ein Ziel bei der Entwicklung der Socket-Schnittstelle war, die bestehenden Systemaufrufe des Unix-I/O-Systems weitestgehend auch für die Netzwerkkommunikation zu nutzen (orthogonale Erweiterung). Aufgrund wesentlicher Unterschiede in der Semantik von File-I/O und Netzwerk-I/O konnte die Abstraktion *Everything is a file* nicht befriedigend erweitert werden. Als Kompromiss wurden deshalb insgesamt 17 neue Systemaufrufe für die Kommunikation mit Sockets bereitgestellt.

In BSD-basierten Systemen sind Sockets vollständig im Betriebssystem realisiert und somit erfolgt der Zugriff auf die Socket-Schnittstelle vollständig über System Calls. In SVR4-basierten Systemen sind Sockets auf der Basis des Streams-Subsystems implementiert. Die Socket-Funktionen sind dabei entweder als Bibliotheks-Funktionen unter Verwendung der Streams-Systemaufrufe oder auch im Systemkern mit einer Systemaufrufchnittstelle realisiert. Für die Anwendungsprogrammierung ist dies in der Regel irrelevant.

Sockets werden in gleicher Weise wie Dateien über Deskriptoren realisiert. Viele Systemaufrufe des I/O-Subsystems (wie z.B. *read()* oder *write()*) können so auch auf Sockets angewandt werden.

#### **4.4.2 Überblick/Einordnung**

**Berkeley  
Sockets  
API**



**ISO/OSI**

4

**TCP**  
reliable

**UDP =**  
unreliable

**IP + PORT +**  
**Checksum**

3

IP (unreliable) / ARP

2

Ethernet-Verkabelung/Protokoll

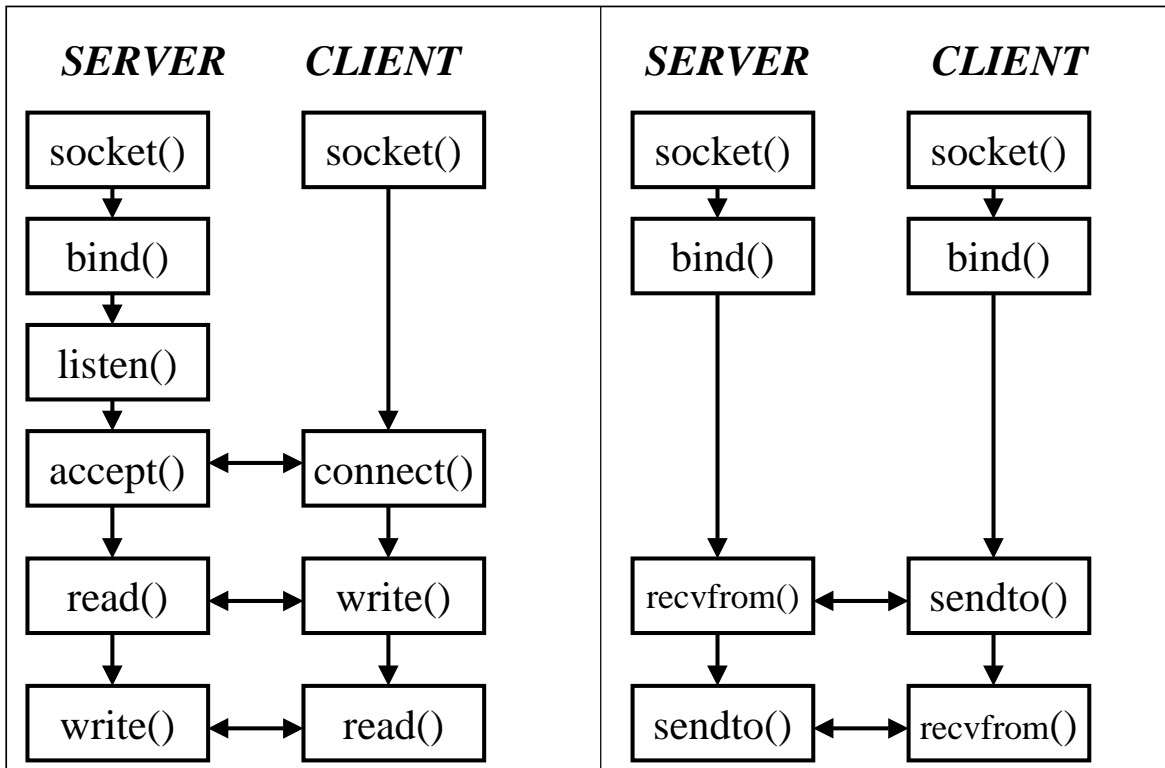
1

Twisted Pair, Glasfaser, Koaxial

**(TCP/UDP)-IP = 5-Tupel :**

(SendAddr, SendPort, Protokoll,  
RecvAddr, RecvPort)

## API-Aufrufe: TCP



## InternetDomain, UnixDomain

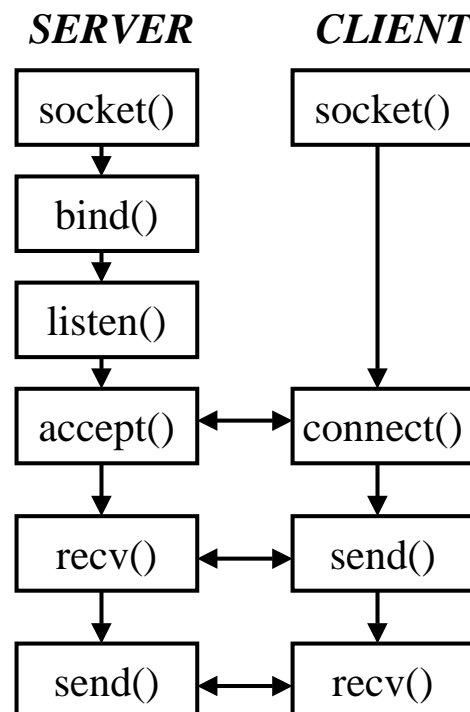
- InternetDomain = (bidirektionale)  
Kommunikation zwischen Rechnern
- UnixDomain = (bidirektionale)  
Kommunikation auf EINEM (Unix-)Rechner

# UNIXDomain

3er Tupel:

- (Protocol, local-pathname, foreign-pathname)
- Kommunikation *innerhalb* eines Dateisystems!

## UNIXDomain





### 4.4.3 Erzeugung eines Socket

Die Erzeugung einer Socket-Instanz erfolgt mit der Funktion **socket()**:

```
# include <sys/types.h>
# include <sys/socket.h>

int socket(
    int domain,
    int type,
    int protocol
)
```

- Die zu verwendende Domäne wird über den Parameter *domain* angegeben.
- Die zu verwendende Kommunikationssemantik wird über den Parameter *type* angegeben.
- Das konkrete Kommunikationsprotokoll kann über den Parameter *protocol* angegeben werden. Wird hier **0** angegeben (also nichts spezifiziert), so selektiert das System eine geeignetes Protokoll passend zu den ersten beiden Parametern.
- Rückgabewert ist ein Deskriptor, der den erzeugten Socket referenziert und in allen weiteren darauf operierenden Funktionen benutzt wird.
- Für den Parameter *domain* sind in *sys/socket.h* Konstanten definiert: **AF\_UNIX** (auch als **PF\_UNIX**) für die Unix-Domäne, **AF\_INET** (auch als **PF\_INET**) für die Internet-Domäne. **AF** steht für *address family*, **PF** für *protocol family*.
- Für den Parameter *type* sind in *sys/socket.h* ebenfalls Konstanten definiert: **SOCK\_STREAM** für den Socket-Typ STREAM, **SOCK\_DGRAM** für DGRAM und **SOCK\_RAW** für RAW.
- Der Rückgabewert **-1** signalisiert einen Fehler:  
**EPROTONOSUPPORT**, falls das spezifizierte Protokoll nicht unterstützt wird  
**ENOPROTOPTYPE**, falls der Socket-Typ innerhalb der gewählten Domäne unzulässig ist oder nicht unterstützt wird;  
 weitere Fehler können aufgrund systeminterner Ressourcenknappheit oder mangelnden Zugriffsrechten entstehen.

#### Beispiel:

```
int sd = socket(PF_INET, SOCK_STREAM, 0);
```

Damit wird ein Stream-Socket der Internet-Domäne erzeugt, welches das voreingestellte Transportprotokoll (hier **TCP**) als das dem Socket zugrundeliegende Kommunikationsprotokoll verwendet.

### 4.4.4 Benennung eines Socket

Mit **socket()** werden die internen, den Socket repräsentierenden Datenstrukturen allokiert und initialisiert. In diesem Zustand kann der Socket als **unbenannter Socket** bezeichnet werden. Solange der Socket noch mit keiner

Adresse verknüpft ist, kann der Socket noch nicht von fremden Prozessen angesprochen werden und somit auch keine Kommunikation stattfinden.

Die Benennung eines Socket erfolgt durch Zuweisung einer Adresse an den Socket:

- **Socket-Adresse in der Internet-Domäne:**

```
<protocol,local-address,local-port,foreign-address,foreign-port>
```

Ein Halb-Tupel `<protocol,address,port>` definiert dabei jeweils einen Kommunikationsendpunkt, *foreign* bezieht sich auf den zukünftigen Kommunikationspartner.

- **Socket-Adressen in der Unix-Domäne:**

Hier werden die Kommunikationsverbindungen über Pfadnamen identifiziert, also über Tupel der Form

```
<protocol,local-pathname,foreign-pathname>
```

Mit der Funktion **bind()** wird ein Kommunikationsendpunkt (ein Halbtupel, s.o.) festgelegt:

```
# include <sys/types.h>
# include <sys/socket.h>

int bind(
    int          sd,          /*socket descriptor*/
    struct sockadr * address,  /*address*/
    int          addresslen /*length of address*/
);
```

Aufgrund der meist unterschiedlichen Adressformate der einzelnen Domänen sind Socket-Adressen als eine Folge von Bytes variabler Länge mittels einer generischen Datenstruktur anzugeben. Alle Funktionen der Socket-Schnittstelle, die Adressen verwenden, referenzieren diese nur über diese generische Datenstruktur, die in **sys/socket.h** wie folgt definiert ist:

```
struct sockadr {
    u_char  sa_len;          /*total address length */
    u_char  sa_family;       /*address family */
    char    sa_data[14];     /*protocol-specific address*/
};
```

- Die Komponente *sa\_len* gibt die gesamte Länge der Socket-Adresse in Bytes an.
- *sa\_family* bezeichnet die Socket-Adressfamilie, die dem Adressformat der verwendeten Kommunikationsdomäne entspricht.
- *sa\_data* enthält die ersten 14 Bytes der Adresse selbst. Damit wird die Länge der tatsächlichen Adresse nicht beschränkt!

Anm.: Dieser Punkt ist sehr implementierungsspezifisch und wird mit dem Übergang auf IPv6 geändert werden müssen.

**Beispiel für die Initialisierung und Zuweisung einer Adresse in der Unix-Domäne** (Pfadname */tmp/my\_socket*):

```
# include <sys/un.h>
# include <string.h>

int sd;
struct sockaddr_un sun_addr;    /*socket address
                                *defined in <sys/un.h>
                                */

/* Create a socket in the UNIX domain: */

sd = socket(PF_UNIX,SOCK_STREAM,0);

/* Initialize the socket address: */
(void) memset(&sun_addr,0,sizeof(sun_addr));
(void) strcpy(sun_addr.sun_path,"/tmp/my_socket");
sun_addr.sun_family = AF_UNIX;
sun_addr.sun_len      = (u_char) ( sizeof(sun_addr.sun_len) +
                                sizeof(sun_addr.sun_family) +
                                sizeof(sun_addr.sun_path) + 1);

/* bind the address to the socket: */
bind(sd,(struct sockaddr *)&sin_addr, sizeof(sun_addr));
```

**Benennung einer Adresse in der Internet-Domäne:**

- Hier ist aus Rechneradresse und Portnummer eine Netzwerkadresse zu konstruieren
- Dazu gibt es eine ganze Reihe noch zu besprechender Bibliotheksfunktionen

**Prinzip:**

```
# include <netinet/in.h>

int sd;
struct sockaddr_in sin_addr; /*defined in <netinet.h>*/

/* create a socket: */
sd = socket(PF_INET,SOCK_STREAM,0);

/* Initialize the address: ... */

/* Bind the name to the socket: */
bind(sd, (struct sockaddr *)&sin_addr,sizeof(sin_addr));
```

#### 4.4.5 Aufbau einer Kommunikationsverbindung

Der Aufbau einer Kommunikationsverbindung zwischen zwei nicht notwendigerweise verschiedenen Prozessen verläuft in der Regel asymmetrisch, wobei ein Prozess als **Client**, der andere als **Server** bezeichnet wird. Der Server

stellt normalerweise einen Dienst zur Verfügung, wartet also auf die Inanspruchnahme dieses Dienstes durch einen anderen Prozess, den Client; entsprechend werden die Kommunikationsendpunkte als passiver bzw. aktiver Kommunikationsendpunkt bezeichnet.

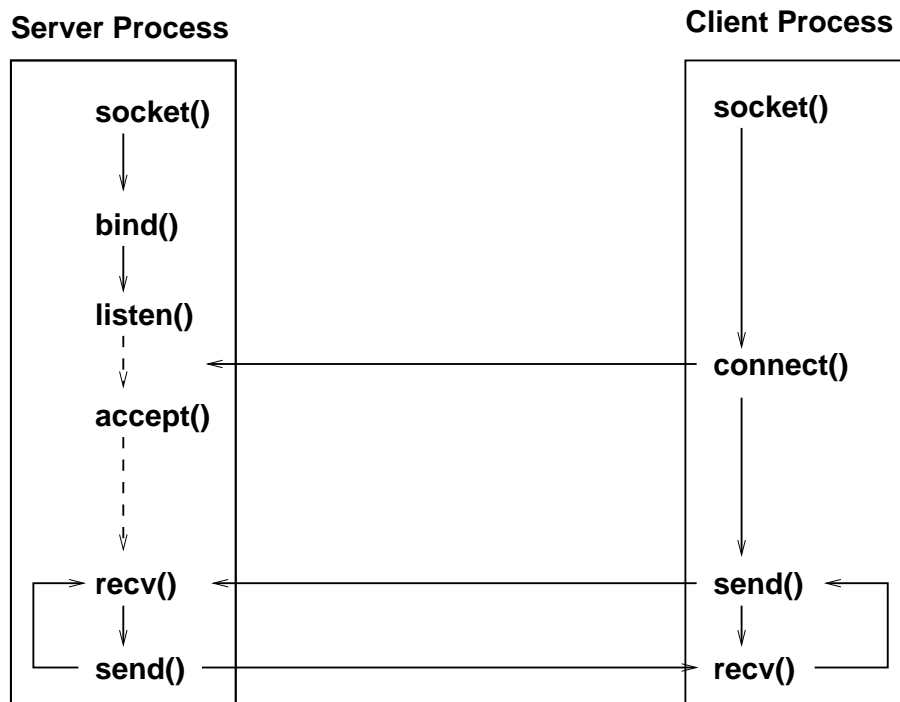


Abbildung 4.2: Verbindungsorientierte Client-Server-Kommunikation

Die Benennung des Server-Socket spezifiziert dabei die eine Hälfte einer möglichen Kommunikationsverbindung; die Vervollständigung wird durch einen Client initiiert, der aktiv die Verbindung zu einem Server anfordert und dabei die bekannte Adresse des Servers spezifiziert und seine eigene (implizit) mitliefert. Dazu dient die Funktion **connect()**:

```
# include <sys/types.h>
# include <sys/socket.h>

int connect(
    int sd, /*client's socket descriptor*/
    struct sockaddr * address, /* server's address*/
    int adresslen
);
```

Das zweite und dritte Argument ist analog zu den entsprechenden Parametern von *bind()* auf Server-Seite zu verstehen. Der Server hat kein explizites *bind()* zur Adresszuweisung an seinen Socket durchzuführen, da es automatisch vom Betriebssystem vorgenommen wird, sofern der Socket zum Zeitpunkt der Verbindungsaufnahme noch unbenannt ist.

Anm.: Die Socket-Adresse des Kommunikationspartners wird in der Socket-Terminologie auch als **Peer-Adresse** bezeichnet.

**Beispiel:**

```
# include <netinet/in.h>

int sd;
struct sockaddr_in sin_addr; /* for address of the server*/

sd = socket(PF_INET, SOCK_STREAM, 0);

/* Initialize the socket address of the server (see below)*/

/*connect to the specified server:*/
connect(sd, (struct sockaddr *)&sin_addr, sizeof(sin_addr));
```

Der Client wird durch die Ausführung von *connect()* solange blockiert, bis entweder die Kommunikationsverbindung vom Server vervollständigt wurde (somit erfolgreich hergestellt wurde) oder bis ein Fehler auftritt.

Bevor der Server Verbindungsanforderungen entgegen nehmen kann, muss sein erzeugter (*socket()*) und benannter (*bind()*) Socket als passiver Kommunikationsendpunkt gekennzeichnet werden (als Bereitschaft, Verbindungen zu akzeptieren) - dazu dient die Funktion **listen()**:

```
# include <sys/socket.h>

int listen(
    int sd,
    int backlog
);
```

- Der Parameter **backlog** spezifiziert die Länge einer Warteschlange des passiven Socket, in der Verbindungsanforderungen von Clients solange gehalten werden, bis sie vom Server explizit akzeptiert werden.

Dies erfolgt durch die Funktion **accept()**:

```
# include <sys/types.h>
# include <sys/socket.h>

int accept(
    int sd,
    struct sockaddr * address,
    int * addresslen
);
```

- Der Parameter *sd* ist der Deskriptor des benannten, passiven Socket.
- Die Argumente *address* und *addresslen* sind Wert-/Resultat-Parameter: sie müssen mit einer Variablen der domänen-spezifischen Socket-Adresse bzw. deren Länge initialisiert werden. Nach erfolgreichem Funktionsaufruf enthalten sie die Socket-Adresse des Client bzw. deren tatsächliche Länge (also die *Peer*-Adresse). Ist der Server an der *Peer*-Adresse nicht interessiert, so ist im Parameter *addresslen* der Nullzeiger anzugeben, der Parameter *address* bleibt dabei unberücksichtigt.
- *accept()* blockiert, bis eine Verbindungsanforderung eines Clients in der Warteschlange des passiven Sockets zur Verfügung steht.

- Als Resultat liefert *accept()* - sofern keine Fehler aufgetreten ist - einen Socket-Deskriptor zurück, der einen neuen Socket referenziert; dieser repräsentiert den Kommunikationsendpunkt für die nun fertige neue Verbindung.

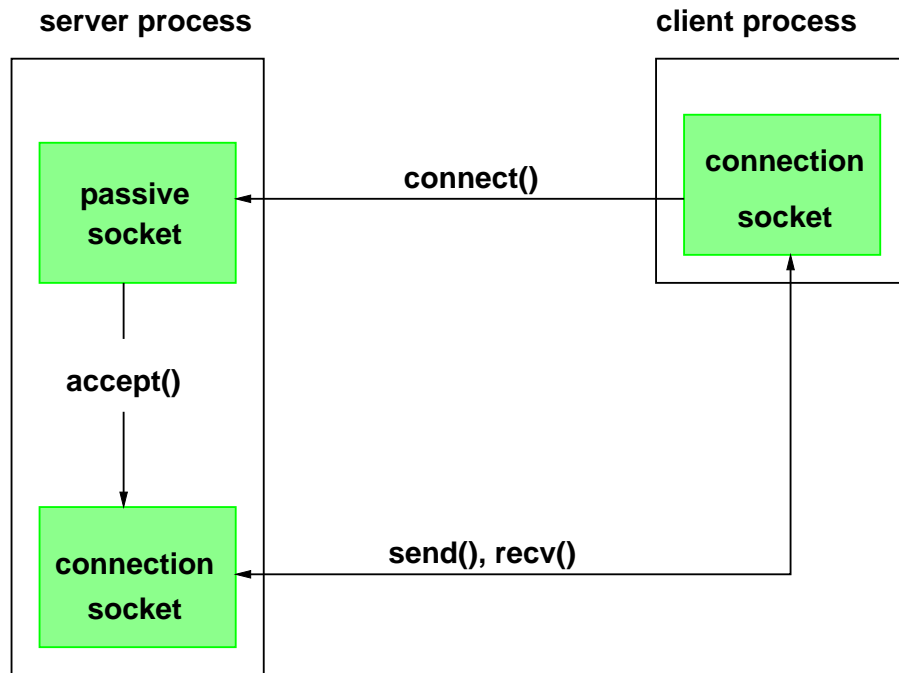


Abbildung 4.3: Aufbau einer verbind.-orient. Client/Server-Kommunikation

- Akzeptieren einer Kommunikationsverbindung in der UNIX-Domäne:

```

int          sd,          /*server socket descriptor*/
              conn_sd;    /*connection socket de-
sc. */
struct sockaddr_un client_addr; /*client socket address */
int            client_len; /*length of client address*/

/*
 * sd = socket(); bind(sd,...); listen(sd,...);
 */

client_len = sizeof(client_addr);
conn_sd    = accept(sd, (struct sockaddr *) &client_addr,
                    &client_len);
  
```

- Akzeptieren einer Kommunikationsverbindung in der Internet-Domäne:

```

int          sd,          /*server socket descriptor*/
              conn_sd;    /*connection socket de-
sc. */
struct sockaddr_in client_addr; /*client socket address */
  
```

```

int                client_len;    /*length of client address*/

/*
 * sd = socket(); bind(sd,...); listen(sd,...);
 */

client_len = sizeof(client_addr);
conn_sd     = accept(sd, (struct sockaddr *) &client_addr,
                    &client_len);

```

#### 4.4.6 Client-Beispiel: Timeclient für Port 11011

```

#include <netdb.h>
#include <netinet/in.h>
#include <stdio.h>
#include <strings.h>
#include <sys/socket.h>
#include <unistd.h>

#define TPORT      11011

int main(int argc, char **argv)
{
    struct sockaddr_in addr;
    struct hostent *hp;
    int fd;
    int nbytes;
    char buf[BUFSIZ];

    if (argc!=2)
        printf("usage: %s hostname\n", argv[0]),
        exit(1);

    if (!(hp = gethostbyname(argv[1])))    // IP-Adresse des Ser-
vers holen
        fprintf(stderr, "unknown host: %s\n", argv[1]),
        exit(1);

    bzero(&addr, sizeof(addr));    // mit 0en fuellen
    addr.sin_family = AF_INET;    // TCP/IP-Verbindung
    addr.sin_port = htons(TPORT); // Port eintragen, Network Byte Order
    // Serveradresse eintragen
    bcopy(hp->h_addr, &addr.sin_addr, hp->h_length);

    if ((fd = socket(AF_INET, SOCK_STREAM, 0)) < 0)
        perror("socket"),    // Socket erzeugen
        exit(1);

    if (connect(fd, &addr, sizeof(addr)) < 0)
        perror("connect"),    // an Socketport verbinden
        exit(1);

    // vom Socket lesen
    while ((nbytes = read(fd, buf, sizeof(buf))) > 0)
        printf("%.s", nbytes, buf);

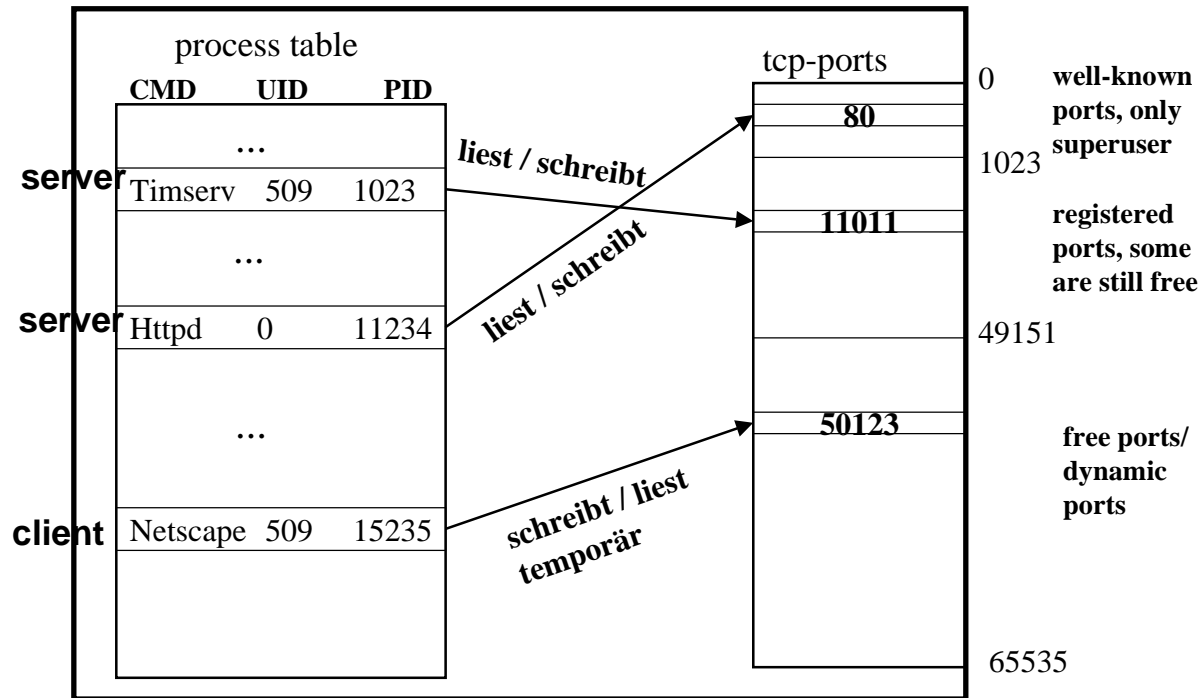
```

```
    close(fd);  
    exit(0);  
}
```

#### 4.4.7 Überblick: Gebrauch von TCP-Ports



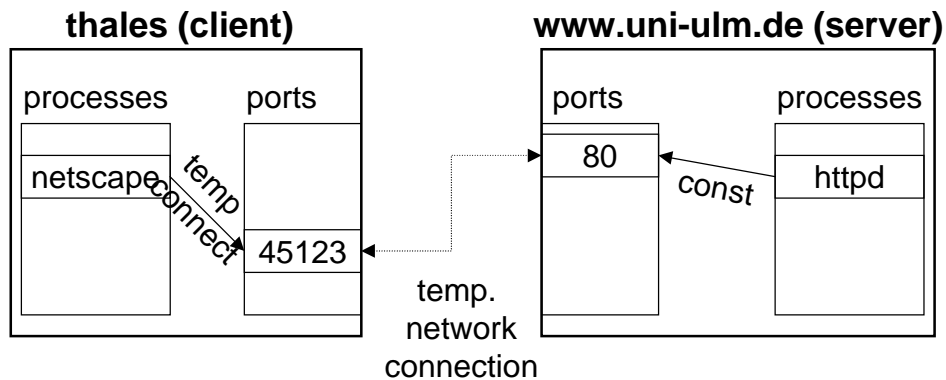
## Unix operating system



- Wie erfahre ich, **welches** Programm unter Unix aktuell einen bestimmten Port belegt?  
Kommando *lsof* = list of open files (and ports)  
(z.B.: `lsof | grep TCP | grep portno`)

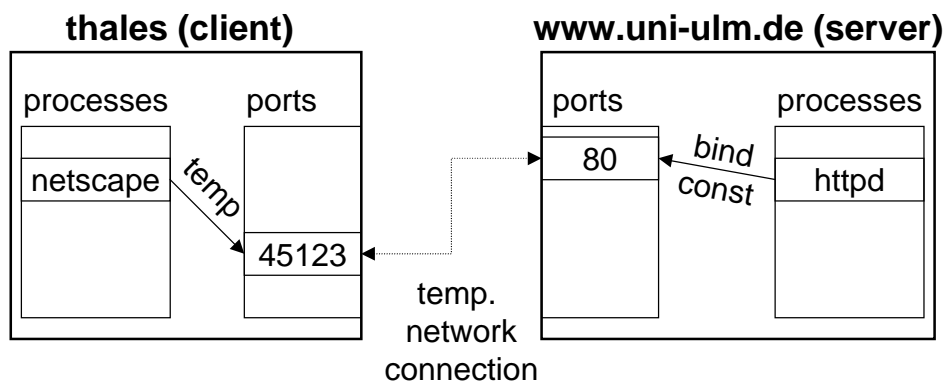
## netscape auf thales liest webseite <http://www.uni-ulm.de>

```
netscape: gethostbyname(„www.uni-ulm.de“) // get IP-Addr. server
fd = socket(AF_INET, SOCKET_STREAM, 0) // TCP-IP socket
sin_port = htons(80); // well-known port
sin_family = AF_INET // Internetdomain
connect(fd, [to host www.uni-ulm.de; port 80]);
write(fd, „GET / HTTP/1.0\n\n“, 16); read(fd, buf, sizeof(buf));
```



## httpd (Webserverprog) auf [www.uni-ulm.de](http://www.uni-ulm.de) erhält Anfrage von Netscape auf [thales.mathematik.uni-ulm.de](http://thales.mathematik.uni-ulm.de)

```
httpd: sin_port = 80; // Dienst an Port 80
sin_family = AF_INET // Internetdomain
fd = socket(AF_INET, SOCKET_STREAM, 0) // TCP-IP socket
sin_port = htons(80); // well-known port
sin_family = AF_INET // Internetdomain
bind(fd, [Internetdomain; port 80]); // socket „anbinden“
listen(fd, SOMAXCON); // max. Anfragen Queue
nfd = accept(fd, ....); // Anfrage von thales via neuem fd annehmen
read(nfd, buf, sizeof(buf)); write(nfd, „<html><body>...“, xyz);
```



# struct sockaddr

- **zentraler Informationsträger für bind, accept und connect ist die Struktur sockaddr**
  - bei Protokollart AF\_INET: struct sockaddr\_in
  - bei Protokollart AF\_UNIX: struct sockaddr\_un
- **wichtigster Inhalt bei (sin\_family=)AF\_INET:**
  - sin\_port = Portnummer
  - sin\_addr = IP-Adresse Server bzw. Client (bei accept)
- **wer braucht was?**
  - bind: braucht Infos **für** Server (Protokoll+Port)
  - accept: schreibt Infos **über** Client (Protokoll+IP+Port)
  - connect: braucht Infos **über** Server (Protokoll+IP+Port)

Wie erfahre ich, **welcher Rechner** von **welchem Port** aus sich an meinen Server angedockt hat?

```
accept(fd, (struct sockaddr*)&c_addr, ...)
```

accept schreibt Infos über den Client in die sockaddr-Struktur c\_addr

```
printf(„port: %d\n“, ntohs(c_addr.sin_port));
```

ntohs = converts network byte order to host byte order

```
printf(„rechner: %s\n“, inet_ntoa(c_addr.sin_addr));
```

inet\_ntoa = konvertiert IP-Adresse in String (dotted-decimal form z.B.134.60.66.21)



### 4.4.8 Der Datentransfer

Nachdem zwischen Client und Server eine Kommunikationsverbindung aufgebaut ist, kann ein Datentransfer stattfinden. Dazu können in gewohnter Weise die System Calls des UNIX-I/O-Systems benutzt werden: **read()**, **readv()** (*read from multiple buffers*), **write()**, **writev()** (*write into multiple buffers*). Die Socket-Schnittstelle stellt 3 weitere Funktionspaare zur Verfügung, die

die Semantik der UNIX-I/O-Funktionen um Socket- und protokollspezifische Eigenschaften und Mechanismen erweitern.

- **send(), recv()**

```
# include <sys/types.h>
# include <sys/socket.h>
```

```
ssize_t send ( int sd, void * buf, size_t len, int flags);
ssize_t recv ( int sd, void * buf, size_t len, int flags);
```

Ist bei beiden Funktion für *flags* der Wert 0 angegeben, so sind identisch zu *write()* und *read()*. Die Spezifikation bestimmter Methoden des Datentransfers oder das Versenden / Empfangen von Kontrollinformationen kann durch entsprechende *flag*-Werte aktiviert werden:

**MSG\_OOB** Die spezifizierten Daten sollen mit hoher Priorität gesendet bzw. empfangen werden. Solche Daten werden im Kontext von Sockets als **out-of-band**-Daten (OOB-Daten) bezeichnet. Sie werden im Vergleich zu normalen Daten auf einem logisch unabhängigen Kanal gesendet. Bei OOB-Daten kann zudem der übliche Pufferungsmechanismus umgangen werden. Dieser Mechanismus unterliegt allerdings Beschränkungen und ist nur für wenige, verbindungsorientierte Protokolle realisiert.

**MSG\_PEEK** Auf der Seite des Empfängers wird hiermit spezifiziert, dass gepufferte Daten nur gelesen, aber nicht "konsumiert" werden sollen. Der nächste Lesezugriff liefert dieselben Daten noch einmal zutücl.

**MSG\_WAITALL** Der Lesezugriff soll solange blockieren, bis die angeforderte Datenmenge insgesamt zur Verfügung steht.

**MSG\_DONTWAIT** Ausführung der Operation im nicht-blockierenden Modus.

**MSG\_DONTROUTE** Ausgehende Datenpakete werden ohne Berücksichtigung einer Routing-Tabelle nur in das lokal angeschlossene Netzwerk gesendet. Dies ist nur für spezielle Diagnoseprogramme interessant.

**MSG\_EOR** Die Flagge **end-of-record** markiert das logische Ende eines Datensatzes; die Daten werden dabei mit zusätzlicher Kontrollinformation versendet. Sie kann aber nur verwendet werden, wenn das zugrundeliegende Kommunikationsprotokoll das Konzept der Datensatz-Übermittlung unterstützt.

**MSG\_EOF** Hiermit wird das Ende der Datenübertragung markiert und als Kontrollinformation zusammen mit den angegebenen Daten übertragen.

Die für den Datentransfer bereitgestellten Flaggen sind zumeist auf spezielle Kommunikationsprotokolle beschränkt und auch nur definiert, wenn die diese Protokolle auf dem System implementiert sind. Die protokollunabhängigen

Flaggen sind **MSG\_PEEK**, **MSG\_WAITALL** und **MSG\_DONTWAIT**; die beiden letzteren stehen nur in neueren Implementierungen zur Verfügung!

- **sendto(), recvfrom()**

```
# include <sys/types.h>
# include <sys/socket.h>

ssize_t sendto ( int sd, void * buf, size_t len, int flags,
                 struct sockaddr * address, int addresslen);

ssize_t recvfrom ( int sd, void * buf, size_t len, int flags,
                  struct sockaddr * address, int * addresslen);
```

Diese Funktionen stehen für den Datentransfer mit verbindungslosen Kommunikationsverfahren zur Verfügung. Dabei ist die Angabe der Sender- und Empfänger-Adresse in jedem Datenpaket erforderlich, da keine virtuelle Verbindung zwischen den Partnern besteht.

Die Angabe der Adressen (Parameter *address* und *addresslen*) sind wie bei den Funktionen *connect()* bzw. *accept()* anzugeben. Läßt man diese Parameter weg, so entsprechen diese Funktionen den obigen Funktionen *send()* und *recv()*.

- **sendmsg(), recvmsg()**

```
# include <sys/types.h>
# include <sys/socket.h>

ssize_t sendmsg ( int sd, struct msghdr * msg, int flags );

ssize_t recvmsg ( int sd, struct msghdr * msg, int flags );
```

Diese beiden Funktionen stellen die allgemeinste Schnittstelle dar und erweitern die Funktionalität der obigen Funktionen. Mehr dazu siehe z.B. im Manual!

#### 4.4.9 Terminierung einer Kommunikationsverbindung

Das Schliessen und die damit verbundene Freigabe der systeminternen Ressourcen erfolgt mit der Funktion **close()** auf den entsprechenden Socket-Deskriptor. Bei Prozess-Termination werden die Socket-Verbindungen in gleicher Weise wie die Datei- oder terminal-Verbindungen geschlossen.

In verbindungsorientierten Kommunikationsprotokollen, die ja einen verlässlichen Datentransfer garantieren, versucht das System beim Schliessen eines Socket für eine gewisse Zeit evt. noch ausstehende, zwischengepufferte Daten zu transferieren. Dies lässt sich durch entsprechende Operationen ändern.

Eine Verbindung zwischen zwei Sockets ist **voll-duplex**; dies bedeutet, dass Sende- und Empfangskanal logisch voneinander unabhängig sind. Mit der Funktion **shutdown()** lässt sich die Verbindung auch nur in einer Richtung terminieren. Dies wird typischerweise bei verbindungsorientierten Protokollen vom Sender dazu verwendet, dem Empfänger das Ende der Eingabe anzuzeigen. Die Empfängerseite bleibt für den Datentransfer weiterhin geöffnet. Das Schliessen der Empfängerseite bewirkt, dass noch nicht konsumierte, zwischengepufferte Daten wie auch alle zukünftig noch eintreffenden Daten verworfen werden.

- **shutdown()**

```
# include <sys/socket.h>
```

```
int shutdown ( int sd, int how );
```

Der Wert von *how*:

- 0** Leseseite (Empfängerseite) wird geschlossen
- 1** Schreibseite (Senderseite) wird geschlossen
- 2** Beide Seiten werden geschlossen

#### 4.4.10 Verbindungslose Kommunikation

In diesem Fall läuft die Kommunikation nach einem symmetrischen Modell, auch wenn ggf. einer der Prozesse die Funktion des Servers, der andere die des Clients einnehmen kann (aber es findet kein Verbindungsaufbau statt). Die Adressen der Kommunikationspartner werden also nicht über eine virtuelle Kommunikationsverbindung festgelegt; in jedem Datenpaket muss stattdessen die Empfängeradresse beigefügt werden. Die Datenpakete werden bei verbindungsloser Kommunikation oft auch als **Datagramme** und die Kommunikationsendpunkte als **Datagramm Sockets** (Socket-Typ: **SOCK\_DGRAM**) bezeichnet.

Soll der Socket mit einer bestimmten lokalen Adresse benannt werden, so muss die Zuweisung der Adresse über die Socket-Funktion **bind()** vor dem ersten Datentransfer stattfinden; ansonsten erfolgt die Benennung des lokalen Socket beim Senden des ersten Datagramms implizit durch das Betriebssystem.

- Versenden von Daten mit gleichzeitiger Spezifikation der Empfängeradresse: Funktionen **sendto()** und **sendmsg()**
- Empfang von Daten mit gleichzeitiger Gewinnung der Absenderadresse: Funktionen **recvfrom()** und **recvmsg()**

Ein Datentransfer zwischen zwei Datagramm-Sockets kann nur dann stattfinden, wenn beide Kommunikationsendpunkte explizit über die Funktion **bind()** oder implizit über die Funktionen **sendto()** oder **sendmsg()** benannt sind; ansonsten werden die gesendeten Datagramme verworfen!

Da verbindungslose Kommunikation auch unzuverlässige Datenzustellung bedeutet, sollte man sich in Client/Server-Anwendungen die gesendeten Datagramme vom Empfänger bestätigen lassen, sofern auf eine Anfrage keine Daten vom Empfänger erwartet werden. Ist ein kontrollierter und zuverlässiger Datentransport nötig, so müssen dies die Anwendungen in diesem selbst regeln!

- Prinzip der Kommunikation zwischen zwei Datagramm Sockets:

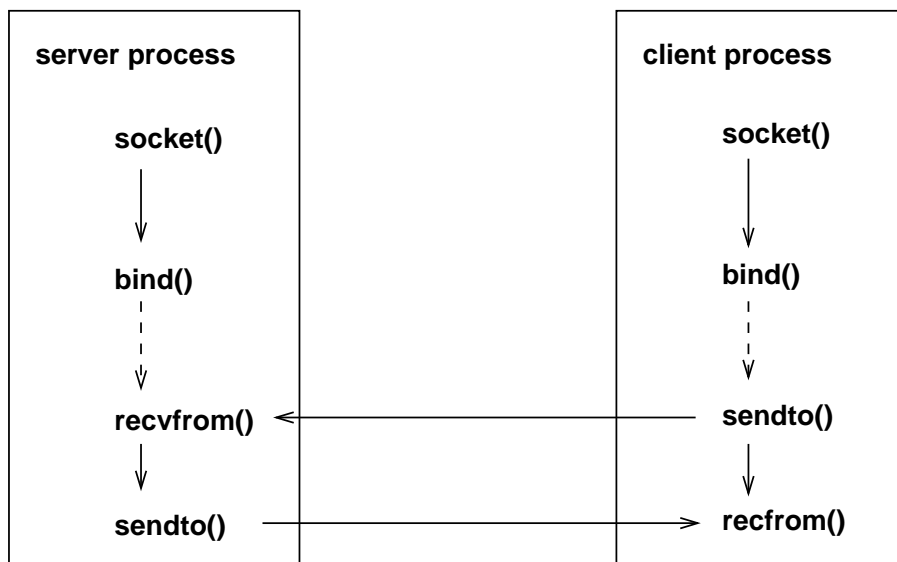


Abbildung 4.4: Aufbau einer verb.-losen Client/Server-Kommunikation

- Der Server-Prozess benennt den erzeugten Datagramm Socket mit einer nach aussen hin bekannten Adresse.
- Die explizite Benennung des Client-Socket mit `bind()` ist **optional** und i.d.R. nicht erforderlich, da nur der Server die Adresse des Kommunikationspartners als Empfängeradresse für die zu sendende Rückantwort benötigt!
- Wichtig ist in diesem Zusammenhang nur, dass die Kommunikationsverbindung innerhalb der gewählten Domäne eindeutig ist; bei der impliziten Benennung eines Socket durch das Betriebssystem ist dies gewährleistet!

Auf Datagramm Sockets kann die Funktion **connect()** angewandt werden; damit wird jedoch keine Kommunikationsverbindung aufgebaut, sondern die dabei spezifizierte Adresse dem Socket als Empfängeradresse zugewiesen, an die



alle folgenden Datagramme zu senden sind. Ausserdem werden dann an diesem Socket nur Datagramme empfangen, deren Adresse mit der in *connect()* angegebenen Adresse übereinstimmt. Damit erübrigt sich die Identifizierung der empfangenen Datagramme. Da hiermit sowohl die Sender- wie Empfängeradresse festgelegt sind, können zum Datentransfer auch die Funktionen **send()** und **recv()** bzw. die UNIX-I/O-Systemaufrufe verwendet werden. In den Funktionen **sendto()** und **sendmsg()** sollten die Socket-Adressen aus Gründen der Portabilität unspezifiziert bleiben. Die Empfängeradresse kann durch einen erneuten Aufruf von *connect()* jederzeit geändert werden sowie eine Beziehung durch Angabe einer ungültigen Adresse gelöscht werden.

#### 4.4.11 Feststellen gebundener Adresse

Manchmal ist es notwendig, die lokal oder entfernt gebundene Socket-Adresse allein anhand des Socket Deskriptors festzustellen; dazu dienen die Funktionen **getsockname()** und **getpeername()**, die als Resultat die lokale bzw. entfernt gebundene Adresse zurückliefern:

```
# include <sys/socket.h>
int getsockname ( int sd,
                  struct sockaddr * address, int * addresslen
                );

int getpeername ( int sd,
                  struct sockaddr * address, int * addresslen
                );
```

Die Parameter *sd*, *address*, *addresslen* sind dabei in gleicher Weise wie in der Funktion *accept()* zu spezifizieren!

Die Funktion **getsockname()** ist insbesondere dann nützlich, wenn die lokale Socket-Adresse vom System zugewiesen wurde. Das Feststellen der entfernt gebundenen Adresse mit der Funktion **getpeername()** ist dann notwendig, wenn ein Prozess diese Adresse benötigt und allein den Socket Deskriptor einer bereits akzeptierten Kommunikationsverbindung erhält und somit keinen Zugriff auf die Peer-Adresse besitzt. Dies ist beispielsweise bei durch den **Internet Superserver inetd** gestarteten Server-Prozessen der Fall.

#### 4.4.12 Socket-Funktionen im Überblick

- **Aufbau**

- **socket()**:  
Erzeugen eines unbenannten Socket
- **socketpair()**:  
Erzeugen eines Paares von miteinander verbundenen Sockets, siehe Manual
- **bind()**:  
Zuweisung einer lokalen Adresse an einen unbenannten Socket

- **Server**

- **listen()**:  
Einen Socket auf Verbindungsanforderungen vorbereiten

- **accept():**  
Eine Verbindungsanforderung akzeptieren
- **Client**
  - **connect():**  
Eine Verbindung zu einem Socket anfordern
- **Empfangen**
  - **read():**  
Daten einlesen
  - **readv():**  
Daten in mehrere Puffer einlesen
  - **recv():**  
Daten einlesen und Angabe von Optionen
  - **recvfrom():**  
Daten einlesen, optional Senderadresse empfangen, Angabe von Optionen
  - **recvmsg():**  
Daten in mehrere Puffer einlesen, optional Senderadresse und Kontrollinformationen empfangen, Angabe von Optionen
- **Senden**
  - **write():**  
Daten senden
  - **writev():**  
Daten aus mehreren Puffern senden
  - **send():**  
Daten senden und Angabe von Optionen
  - **sendto():**  
Daten an die spezifizierte Empfängeradresse senden, Angabe von Optionen
  - **sendmsg():**  
Daten aus mehreren Puffern senden, und Kontrollinformationen an den spezifizierten Empfänger senden, Angabe von Optionen
- **Ereignisse**
  - **select():**  
Multiplexen und auf I/O-Bedingungen warten, siehe Manual
- **Terminieren**
  - **shutdown():**  
Eine Verbindung in eine oder beide Richtungen terminieren
  - **close():**  
Eine Verbindung terminieren und Socket schliessen
- **Administration**
  - **getsockname():**  
Feststellen der lokal gebundenen Socket-Adresse

- **getpeername()**:  
Feststellen der entfernt gebundenen Socket-Adresse
- **setsockopt()**:  
Ändern von Socket- und Protokoll-Optionen, siehe Manual
- **getsockopt()**:  
Auslesen von Socket- und Protokoll-Optionen, siehe Manual
- **fcntl()**:  
Ändern der I/O-Semantik, siehe Manual
- **ioctl()**:  
Verschiedene Socketoperationen, siehe Manual

## 4.5 Konstruktion von Adressen

Adressen für die Lokalisierung eines Dienstes auf einem nicht notwendig entfernten System sind protokoll-spezifisch und setzen sich in der Internet-Domäne aus einer **Rechneradresse** und einer den Dienst identifizierenden **Portnummer** zusammen. Anwendungen spezifizieren den anzufordernden Dienst i.r.G. über einen Namen statt über Rechneradresse plus Portnummer, so z.B. den WWW-Server auf dem Rechner *www.mathematik.uni-ulm.de*, dessen bereitgestellter Dienst mit *http* bezeichnet wird. Namen lassen sich schliesslich leichter merken als numerische Adressen, man erreicht dadurch auch eine gewisse Unabhängigkeit (Änderung von Adresse und Portnummer unter Beibehaltung des Namens).

Für die Konvertierung von Namen in Adressen bzw. Portnummern, für die Konstruktion und Manipulation von Netzwerkadressen sowie für die Lokalisierung eines Rechners gibt es eine Reihe von Bibliotheksfunktionen, die allerdings nicht Bestandteil der Socket-Schnittstelle sind.

### 4.5.1 Socket-Adressen

- Alle Funktionen der Socket-Schnittstelle, die auf Socket-Adressen operieren, verwenden eine generische Socket-Adressstruktur.
- Damit wird die Unabhängigkeit der Socket-Schnittstelle von den konkreten Implementierungen der bereitgestellten Kommunikationsprotokolle erreicht.
- Die Interpretation der Socket-Adressen erfolgt in den protokollspezifischen Funktionen, die bei der Erzeugung einer Socket-Instanz durch die Domäne festgelegt sind.
- Deklaration der Socket-Adressstruktur in auf 4.3BSD basierenden Betriebssystemen:

```
struct sockaddr {
    u_short sa_family;    /* address family          */
    char     sa_data[14]; /* protocol-specific address */
}
```

- Die Komponente **sa\_family** enthält das Adressformat.

- Die Komponente **sa\_data** maximal die ersten 14 Bytes der protokollspezifischen Adresse. Dies ist ein Implementierungsdetail der Socket-Schnittstelle und beschränkt nicht die Länge der protokollspezifischen Adresse.

Die Implementierung von Netzwerkprotokollen stellt mit Blick auf die Performance viele Anforderungen an die Speicherverwaltung des Betriebssystems, so z.B.

- die Handhabung von Datenpuffern unterschiedlicher Länge,
- das einfache Hinzufügen / Entfernen von Kopfdaten oder
- die Datenübergabe zwischen den eigenverantwortlichen Funktionen der verschiedenen Netzwerkschichten.

Auf BSD-Systemen ist für die Kommunikation mit Sockets eine spezielle Speicherverwaltung implementiert, die sogenannten **memory buffers**; dabei wird versucht, Kopieroperationen der Datenpakete zu minimieren. SVR6 basierende Systeme verwenden i.d.R. die Mechanismen des Streams-Subsystems.

Die Schnittstellen zwischen den Schichten des Socket-Modells sind zwar wohldefiniert, die Grenzen in der Implementierung sind eher fließend, da die einer Schicht zugrundeliegenden Datenstrukturen oft auch den Funktionen der darüberliegenden Schicht zugänglich sind. So sind das Adressformat und die ersten 14 Bytes der protokollspezifischen in der Socket-Schicht bekannt, sie sind aber auch der Anwendungsschicht bekannt (Abhängigkeit der Anwendung von den konkreten Kommunikationsprotokollen!). Bezüglich der Kompatibilität und Portabilität entsteht eine weitere Abhängigkeit dadurch, dass die generische Socket-Adresse eine Datenstruktur des Betriebssystemkerns ist und dass sich diese Struktur ab der Version *4.3BSD-Reno* (und aller darauf aufbauenden Systeme) wie folgt geändert hat:

```
struct sockaddr {
    u_char  sa_len;           /* total address length      */
    u_char  sa_family;       /* address family            */
    char    sa_data[14];     /* protocol specific address */
}
```

Die Komponente **sa\_len** enthält die Länge der protokollspezifischen Adresse in Bytes; die Gesamtlänge der Datenstruktur ist unverändert 16 Bytes! Die Hinzunahme dieser Komponente ist für die systeminterne Implementierung notwendig, damit Adressen variabler Länge protokollunabhängig behandelt werden können. Die Anwendung hat davon keinen echten Vorteil, da die Längeninformation Argument der entsprechenden Socket-Funktionen ist, somit bereits verfügbar ist.

Die sehr systemnahe Implementierung von Netzwerkanwendungen mittels der Socket-Schnittstelle hat einige Nachteile, die sich insbesondere auf die Portabilität auswirken. Darauf soll im Folgenden jedoch nicht weiter eingegangen werden (Fragen und Lösungen diesbezüglich wurden in der Dissertation von M. Etter behandelt).



### 4.5.3 Socket-Adressen in der Internet-Domäne

Hier ist die Socket-Adressstruktur in **netinet/in.h** wie folgt definiert:

```
struct in_addr {
    u_long s_addr; /* 32 bit netid/hostid */
};

struct sockaddr_in {
    u_char    sin_len; /*total address length (16 bytes)*/
    u_char    sin_family; /*address family AF_INET */
    u_short   sin_port; /*16 bit port number */
    struct in_addr sin_addr; /*IP address */
    char      sin_zero[8]; /*unused */
};
```

- Die Datenstruktur **in\_addr** enthält nur die Komponente **s\_addr**, in der eine 32 Bit lange Adresse des Internetprotokolls IP in Netzwerkbyteordnung gespeichert wird.
- Die Komponente **sin\_len** der Datenstruktur **sockaddr\_in** ist immer **sizeof(struct sockaddr\_in) = 16** Bytes.
- Die Adressfamilie in Komponente **sin\_family** ist immer **AF\_INET**.
- Die Komponente **sin\_port** ist eine 16 Bit lange Port-Nummer in Netzwerkbyteordnung, die zusammen mit der Internet-Adresse **sin\_addr** einen Kommunikationsendpunkt eindeutig definiert.
- Die Komponente **sin\_zero** ist aus Gründen der Portabilität mit Null-Bytes zu initialisieren und dient lediglich zur Ausweitung der Datenstruktur auf die Länge der generischen Socket-Adressstruktur **sockaddr**:

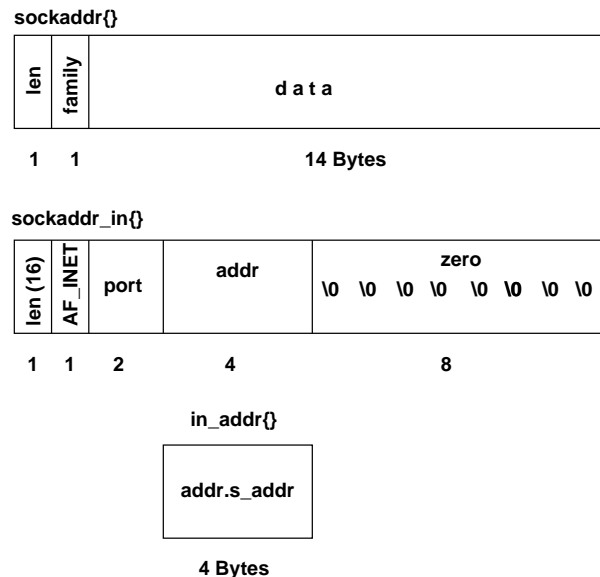


Abbildung 4.5: Organisation der Adress-Struktur **sockaddr\_in**

Für die implizite Selektierung einer geeigneten lokalen IP-Adresse oder auch Portnummer durch das System existieren jeweils eine ausgezeichnete IP-Adresse und eine Portnummer mit Sonderbedeutung: die **IP-Adresse 0 (INADDR\_ANY)** bzw. **0.0.0.0** in punktiertem Dezimalformat sowie die **Portnummer 0**. Bei der Spezifikation der Komponenten **sin\_addr.s\_addr** und **sin\_port** mit diesen sog. *Wildcards* wird die lokale IP-Adresse nach einem erfolgreichen Verbindungsaufbau entsprechend der ein- / ausgehenden Netzwerkschnittstelle automatisch vom System festgelegt und bei der Benennung eines Socket eine frei Portnummer aus dem Bereich der **kurzlebigen Portnummern** gewählt.

#### 4.5.4 Byte-Ordnung

Die Anordnung von Bytes bei Mehr-Byte-Größen erfolgt nicht bei allen Computersystemen in der gleichen Reihenfolge. Für die Speicherung einer 2-Byte-Größe gibt es zwei Möglichkeiten:

- das niederwertige Byte liegt an der Startadresse (**Little-Endian-Anordnung**)
- das höherwertige Byte liegt an der Startadresse (**Big-Endian-Anordnung**)

Bei 4-Byte-Größen können zusätzlich noch die 2-Byte-Größen unterschiedlich angeordnet sein!

Für den Austausch protokollspezifischer Daten werden in den Internet-Protokollen nur 2- und 4-Byte-Integerwerte im Big-Endian-Format verwendet (**network byte order**), die Bit-Ordnung selbst ist ebenfalls in diesem Format!

Die Implementierungen von Netzwerkprotokollen sind somit auf jedem System dafür verantwortlich, dass protokollspezifische Daten in Netzwerkbyteordnung transferiert werden, d.h. die Daten sind von der Byteordnung des Computersystems (*host*) in die Netzwerkbyteordnung zu transferieren, sofern sich die Anordnungen unterscheiden. Zur Konvertierung von 2-Byte-Größen (*short*) und 4-Byte-Größen (*long*) gibt es folgende Funktionen (Makros):

```
u_short htons(u_short hostshort); /*host-to-network short*/
```

```
u_long  htonl(u_long  hostlong); /*host-to-network long*/
```

```
u_short ntohs(u_short netshort); /*network-to-host short*/
```

```
u_long  ntohl(u_long  netlong); /*network-to-host long*/
```

Die protokollspezifischen Mehrbytegrößen, die bei Internet-Protokollen zu spezifizieren sind, sind die Internet-Adresse und die Portnummer in *sockaddr\_in* (Komponenten *sin\_addr.s\_addr* und *sin\_port*).

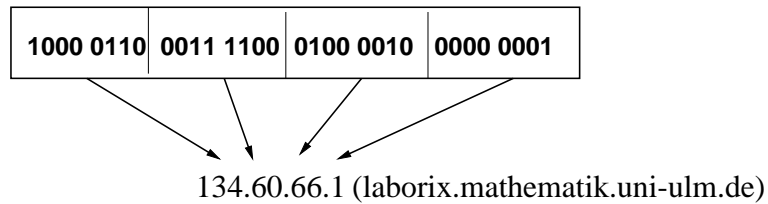


Abbildung 4.6: dotted-decimal notation

### 4.5.5 Spezifikation von Internet-Adressen

Die 32-Bit-IP-Adressen werden meist in der sog. **punktiertes Dezimalformat (dotted decimal notation)** angegeben; diese entsteht dadurch, dass byte-weise Dezimalzahlen gebildet werden, die durch Punkt getrennt sind.

Jeder **Host** in einem IP-Netzwerk ist über eine IP-Adresse eindeutig identifiziert. Ist ein Host an mehrere IP-Netze angeschlossen (**multi-homed host**), so muss dieser für jedes angeschlossene Netz eine IP-Adresse besitzen. Über einen Alias-Mechanismus können einem Host auch mehrere IP-Adressen zugeordnet werden.

Zur Manipulation und Konvertierung von IP-Adressen gibt es wieder einige Bibliotheksfunktionen, die im Headerfile **arpa/inet.h** definiert sind.

- Umwandlung eines als String in *dotted-decimal notation* vorliegende IP-Adresse in eine 32-Bit-IP-Adresse: Funktion **inet\_addr()**

```
# include <arpa/inet.h>

unsigned long inet_addr ( char * ipaddr);
```

Rückgabewerte ist eine 32-Bit-IP-Adresse oder im Fehlerfall die Konstante **INADDR\_NONE** (0xffffffff), die allerdings einer gültigen **Broadcast-Adresse** entspricht. Zu beachten ist auch, dass der Rückgabewert *unsigned long* und nicht *struct in\_addr*. Dies behebt die folgende Funktion:

- **inet\_aton()**

```
# include <sys/socket.h>
# include <netinet/in.h>
# include <arpa/inet.h>

int inet_aton (
    char          * ipaddr;    /*dotted decimal notation */
    struct in_addr * in_addr;   /*result: 32-bit-IP-address*/
);
```

Rückgabewert ist 1 im Erfolgsfall, 0 sonst!

- Konvertierung einer 32-Bit-IP-Adresse in *dotted-decimal notation*: Funktion **inet\_ntoa()**



```
# include <sys/socket.h>
# include <netinet/in.h>
# include <arpa/inet.h>

char * inet_ntoa ( struct in_addr inaddr);
```

- Beispiel für eine Anwendung:

```
struct sockaddr_in addr;

if(inet_aton("134.60.66.5", &addr.sin_addr))
    (void) printf("IP-Address: %s\n",
        inet_ntoa(addr.sin_addr));
```

### 4.5.6 Hostnamen

Die Beziehung zwischen Hostnamen und IP-Adressen werden in der Internet-Domäne jeweils in der Datenstruktur **hostent** repräsentiert, die als Resultat der Funktionen **gethostbyname()** und **gethostbyaddr()** geliefert wird. Diese Funktionen zur **Adress-Resolution** werden als **Resolver** bezeichnet.

Die Struktur *hostent* (in **netdb.h** definiert):

```
struct hostent {
    char *    h_name;           /*official name of host           */
    char **   h_aliases;       /*alias list                      */
    int       h_addrtype;      /*host address type (address family)*/
    int       h_length;        /*length of address              */
    char **   h_addr_list;     /*address list from name server   */
}
# define h_addr h_addr_list[0]
                               /*address for backward compatibility*/
```

Diese Datenstruktur beschreibt den offiziellen Namen des Rechners in der Komponente **h\_name**, eine Liste seiner öffentlichen Alias-Namen in **h\_aliases**, den Adress-Typ in **h\_addrtype**, die Länge einer Adresse in **h\_length** und eine Liste der IP-Adressen (in Netzwerkbyteordnung) in **h\_addr\_list**. Falls es sich bei dem Rechner um einen *multi-homed host* handelt oder Alias-Adressen definiert sind, so enthält die Liste der IP-Adressen entsprechend viele Elemente. Aus Kompatibilitätsgründen verweist die Makrodefinition **h\_addr** auf das erste Element dieser Liste.

Die Funktionen **gethostbyname()** und **gethostbyaddr()**:

```
# include <netdb.h>
struct hostent * gethostbyname (char * name);
struct hostent * gethostbyaddr ( char * addr,
                                int    len,
                                int    type
                                );
```

In der Funktion **gethostbyaddr()** ist die protokollspezifische Adresse in Netzwerkbyteordnung, deren Länge und der Adress-Typ anzugeben. In der Internet-Domäne sind als Parameter eine IP-Adresse, deren Länge, zu erhalten als **sizeof(struct in\_addr)**, und die Konstante **AF\_INET** anzugeben.

Die Spezifikation von **Hostnamen** erfolgt entweder über einfache Namen wie beispielsweise **thales** oder über absolute Namen wie **thales.mathematik.uni-ulm.de.**. Ein absoluter Namen wird auch als **Fully Qualified Domain Name (FQDN)** bezeichnet; diese müssen mit einem Punkt enden, der die Wurzel des hierarchisch geordneten Namensraums bezeichnet. Relative Hostnamen werden abhängig von der administrativen Konfiguration des Systems mit Hilfe der auf dem lokalen System voreingestellten Namensdomäne vervollständigt. In Benutzeranwendungen wird der abschliessende Punkt bei absoluten Namen meist weggelassen.

#### Beispiel:

```
/* hostent.c: print hostent */

# include <stdio.h>
# include <netdb.h>
# include <sys/socket.h>
# include <netinet/in.h>
# include <arpa/inet.h>

void print_hostent( char * host ) {
    struct hostent * hp;

    (void) printf("%s:\n", host);
    if ( (hp = gethostbyname(host)) ) {
        char ** ptr;

        (void) printf("Offizieller Host-Name: %s\n", hp->h_name);

        for(ptr = hp->h_aliases; ptr && *ptr; ptr++)
            (void) printf("    alias: %s\n", *ptr);

        if( hp->h_addrtype == AF_INET)
            for(ptr = hp->h_addr_list; ptr && *ptr; ptr++)
                (void) printf("    Adresse: %s\n",
                              inet_ntoa(*(struct in_addr *) *ptr));
    } else
        (void) printf("----> Kein Eintrag!\n");

    (void) printf("-----\n");
}

int main(int argc, char ** argv) {
    int i;
    for(i = 1 ; i < argc; i++)
        print_hostent(argv[i]);
    exit(0);
}
```

```

thales$ gcc -o hostent -Wall -lxnet hostent.c
thales$ hostent thales turing
thales:
Offizieller Host-Name: thales
  alias: ftp
  alias: www
  alias: pop
  alias: glueck
  alias: adi
  Adresse: 134.60.66.5
-----
turing:
Offizieller Host-Name: turing
  alias: loghost
  alias: mailhost
  Adresse: 134.60.166.1
-----
thales$

```

Die Beziehungen zwischen **Hostnamen** und **Internet-Adressen** werden in der verteilten Datenbank des **Domain Name System (DNS)** verwaltet; die darin enthaltenen Informationen sind über sogenannte **Nameserver** zugänglich. Alternativ dazu werden Informationen über Hostnamen und Internet-Adressen auf dem lokalen System in der Datei **/etc/hosts** oder über den **Network Information Service (NIS)** bereitgestellt. Die unterschiedlichen Möglichkeiten der **Adress-Resolution** zeigt die folgende Abbildung:

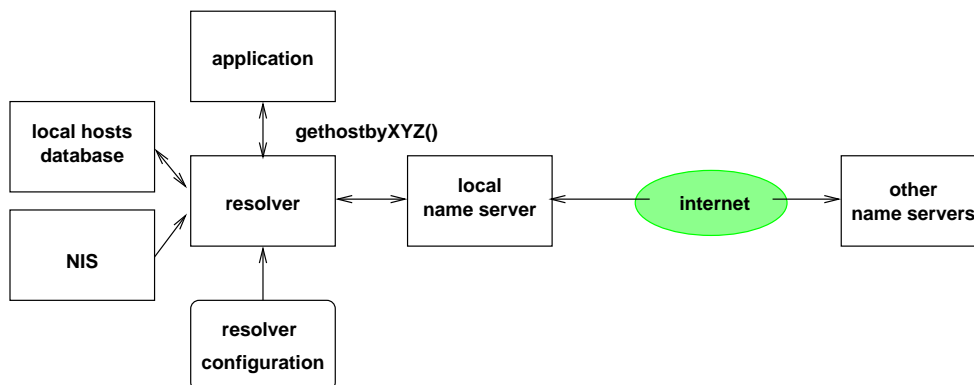


Abbildung 4.7: Methoden der Adress-Resolution

Die auf eine Anfrage gelieferten Informationen variieren aufgrund der verschiedenen Zugangsverfahren und auch unterschiedlichen Organisationen der Datenbanken. Die Zugangsverfahren hängen auch von der Implementierung wie der administrativen Konfiguration des Resolvers ab. Die Resolver-Funktionen liefern auf jeden Fall den offiziellen Hostnamen und eine IP-Adresse in der Struktur **hostent** zurück. Bei Verwendung lokaler Mechanismen liefern einige Systeme jedoch nur einfache und keine absoluten Hostnamen zurück. Die wesentlichen Unterschiede zeigen sich bei Aliasnamen und Aliasadressen. Wird

beispielsweise die Host-Tabelle oder NIS verwendet, erhält man genau eine Adresse und alle Aliasnamen. Werden die Informationen über Nameserver angefordert, so erhält man eventuell Aliasadressen und höchstens einen Aliasnamen, sofern es sich bei dem in der Anfrage spezifizierten Hostnamen um einen Aliasnamen handelt; dazu noch einmal obiges Programm:

```
thales$ hostent www #using NIS
www:
Offizieller Host-Name: thales
  alias: ftp
  alias: www
  alias: pop
  alias: glueck
  alias: adi
  Adresse: 134.60.66.5
-----
thales$ hostent www.mathematik #using DNS
www.mathematik:
Offizieller Host-Name: thales.mathematik.uni-ulm.de
  Adresse: 134.60.66.5
-----
```

#### 4.5.7 Lokale Hostnamen und IP-Adressen

Besteht bereits eine Verbindung, so können der lokale **Hostname** und die lokale **IP-Adresse** mit Hilfe der Socket-Funktion **gethostname()** und der Resolver-Funktion **gethostbyaddr()** ermittelt werden.

- Die Funktion **gethostname()**:

```
# include <unistd.h>

int gethostname(char * name, size_t namelen);
```

##### Beispiel:

```
hypatia$ cat gethost.c
/* gethost.c: Hostnamen bestimmen */

# include <stdio.h>
# include <unistd.h>

void main() {

    char name[20];
    if ( gethostname(name,20) == 0 )
        (void) printf("Hostname: %s\n", name);
}
hypatia$ gcc -Wall -o gethost gethost.c
hypatia$ gethost
Hostname: hypatia
hypatia$
```

Die maximale Länge eines Hostnamens ist auf den meisten Systemen über die Konstante **MAXHOSTNAMELEN** im Headerfile **sys/param.h** festgelegt.

- Das Kommando **uname**:

```
hypatia$ uname -a
Linux hypatia 2.2.13 #5 Mon Apr 3 12:46:02 MEST 2000 i586 unknown
hypatia$
```

- Die Funktion **uname()**:

Dazu ist im Headerfile **sys/utsname.h** folgende Datenstruktur definiert:

```
struct utsname {
    char    sysname[SYS_NMLN]; /*operating system name          */
    char    nodename[SYS_NMLN]; /*node name (host name)    */
    char    release[SYS_NMLN]; /*operating system release level*/
    char    version[SYS_NMLN]; /*operating system version level*/
    char    machine[SYS_NMLN]; /*hardware type            */
}
```

Die Funktion selbst ist

```
int uname( struct utsname * buf);
```

#### Beispiel:

```
hypatia$ cat uname.c
/* uname.c: get host info */

# include <stdio.h>
# include <sys/utsname.h>

void main(){

    struct utsname buf;
    if( uname(&buf) == 0 ) {
        (void) printf("Host: %s\nOS: %s\nRelease: %s\n",
                      buf.nodename, buf.sysname, buf.release);
        (void) printf("Version: %s\nHardware: %s\n",
                      buf.version, buf.machine);
    }
}

hypatia$ gcc -Wall -o my_uname uname.c
hypatia$ my_uname
Host: hypatia
OS: Linux
Release: 2.2.13
Version: #5 Mon Apr 3 12:46:02 MEST 2000
Hardware: i586
hypatia$
```

### 4.5.8 Portnummern und Dienste

In der Internet-Protokollfamilie werden in den Protokollen der **Transportschicht (TCP und UDP)** 16 Bit lange **Portnummern** für die Identifizierung eines Dienstes bereitgestellt. Diese 65536 Portnummern werden von der **IANA** (*Internet Assigned Numbers Authority*) in drei Bereiche eingeteilt: **well-known ports** (0–1023), **registered ports** (1024–49151) und **dynamic and/or private ports** (49152 – 65535). Der aktuelle Stand ist in der Datei

- <ftp://ftp.isi.edu/in-notes/iana/assignments/port-numbers>

verfügbar.

*Well-known ports* identifizieren bekannte Internet-Dienste. So wird in allen TCP/IP-Implementierungen dem **Telnet-Server telnetd** die TCP-Portnummer 23 und dem **TFTP-Server tftpd** (*Trivial File Transfer Protocol*) die UDP-Portnummer 69 zugewiesen, sofern dieses Anwendungsprotokoll unterstützt wird und die entsprechenden Netzwerkanwendungen auf dem System bereitgestellt sind. Auf UNIX-Systemen werden die Portnummern aus dem Bereich 1 – 1023 als **reservierte Portnummern** bezeichnet und können nur von Prozessen mit Superuser-Privilegien zur Benennung von Sockets verwendet werden. Die sog. *well-known ports* belegen hier die Portnummern 1 – 511 und die Portnummern 512 – 1023 sind für Client-Anwendungen mit Supreuser-Privilegien reserviert, die eine reservierte Portnummer als Bestandteil der Client/Server-Authentifizierung benötigen. Beispiele dafür sind **rlogin** und **rsh**.

Von der IANA nicht verwaltet werden Dienste, die registrierte Portnummern verwendet (lediglich als Konvention aufgelistet). So sind z.B. die Portnummern 6000–6063 für einen **X Window Server** für beide Protokolle (allerdings derzeit nur TCP verwendet) registriert. Nach Konvention binden X-Server für passive Sockets die Portnummern 6000 +  $x$ , wobei  $x$  die Nummer des Displays angibt. Wengleich die registrierten Portnummern frei verfügbar sind zur Benennung eines Socket beliebig genutzt werden können, ist dies keinesfalls zu empfehlen. Über die dynamischen und privaten Portnummern, häufig auch als **kurzlebige Portnummern** (*ephemeral ports*) bezeichnet, wird von der IANA nichts ausgesagt. Sie sind für eine implizite Benennung eines Socket durch das Betriebssystem reserviert. Die meisten UNIX-Systeme binden heute noch die Nummern 1024 – 5000 für kurzlebige Portnummern; damit können maximal 3977 Sockets (typischerweise für Client-Anwendungen) zu einem Zeitpunkt implizit benannt sein. Solaris-Betriebssysteme verwenden kurzlebige Portnummern aus dem Bereich 32768 – 65535.

Die Beziehungen zwischen Portnummern und den offiziellen Namen der Dienste sowie deren Aliasnamen werden in der Datei **/etc/services** oder einer entsprechenden NIS-Tabelle definiert. Dazu ist im Headerfile **netdb.h** folgende Struktur definiert:

```
struct servent {
    char *    s_name;      /*official service name          */
    char **   s_aliases;   /*alias list                      */
    int       s_port;      /*port number (network byte order)*/
    char *    s_proto;     /*protocol to use                 */
}
```

Die Zugriffsfunktionen:

- **getservbyname()**

```
# include <netdb.h>

struct servent * getservbyname(
    char * name,
    char * protocol
);
```

- **getservbyport()**

```
# include <netdb.h>

struct servent * getservbyport(
    int port;
    char * protocol
);
```

Beispiel:

```
thales$ cat getserv.c
/* getserv.c: get service */

# include <stdio.h>
# include <netdb.h>
# include <netinet/in.h>

void print_servent( char * name, char * protocol) {
    struct servent * sp;

    if( (sp = getservbyname(name,protocol)) ) {
        char ** ptr;

        (void) printf("Offizieller Name des Dienstes: %s\n",
            sp->s_name);

        for( ptr = sp->s_aliases; ptr && *ptr; ptr++)
            (void) printf("    Alias: %s\n", *ptr);

        (void) printf("    Port-#: %d\n",
            ntohs( (u_short) sp->s_port));
        (void) printf("    Protokoll: %s\n", sp->s_proto);
    } else
        (void) printf("Kein Eintrag!\n");
}

int main(int argc, char ** argv) {
    if( argc != 3 ) {
        fprintf(stderr, "Usage: %s service protocol\n", argv[0]);
        exit(1);
    } else
```

```

        print_servent(argv[1], argv[2]);
    exit(0);
}
thales$ gcc -Wall -o getserv -lxnet getserv.c
thales$ getserv telnet tcp
Offizieller Name des Dienstes: telnet
    Port-#: 23
    Protokoll: tcp
thales$

```

### 4.5.9 Protokoll- und Netzwerkinformationen

Die Beziehungen zwischen Protokollnamen und Protokollnummern sowie zwischen Netzwerknamen und Netzwerkadressen werden in zwei weiteren Tabellen (auf dem lokalen System in den Dateien **/etc/protocols** und **/etc/networks**) sowie ggf. in entsprechenden NIS-Tabellen gehalten. Der Zugriff darauf erfolgt vergleichbar zu Hostnamen und Diensten über zwei in **netdb.h** definierte Datenstrukturen und entsprechenden Zugriffsfunktionen.

- Die Datenstruktur **struct protoent**:

```

struct protoent {
    char *   p_name;        /*official protocol name */
    char **  p_aliases;     /*alias list                */
    char *   p_proto;       /*protocol number           */
};

```

- Die Funktionen **getprotobyname()** und **getprotobynumber()**:

```

#include <netdb.h>

struct protoent * getprotobyname( char * name );

struct protoent * getprotobynumber( int number );

```

- Die Datenstruktur **struct netent**:

```

struct netent {
    char *   n_name;        /*official net name*/
    char **  n_aliases;     /*alias list                */
    int      n_addrtype;    /*net type                  */
    u_long   n_net;        /*net number                */
};

```

- Die Zugriffsfunktionen **getnetbyname()** und **getnetbyaddr()**:

```

#include <netdb.h>

struct netent * getnetbyname( char * name );

struct netent * getnetbyaddr( u_long  addr,
                             int type );

```



Das Internet-Protokoll verwendet eine 8 Bit lange Protokollnummer zur Identifizierung der nächst höheren Protokollschicht, an die das transportierte IP-Datagramm weiterzuleiten ist. Die für das Internet-Protokoll definierten Protokollnummern werden von der **IANA** verwaltet und sind in RFC1700 enthalten; die aktuelle Version findet sich unter

- <ftp://ftp.isi.edu/in-notes/iana/assignments/protocol-numbers>

#### 4.5.10 Zusammenfassung der Netzwerkinformationen

Die Erzeugung und Benennung eines Kommunikationsendpunktes erfordert

- die Selektion eines Kommunikationsprotokolls und
- die Konstruktion einer protokollspezifischen Adresse, die in der Internet-domäne aus
  - aus einer Netzwerkadresse,
  - einer Rechneradresse und
  - einer den Dienst spezifizierenden Portnummer

zusammengesetzt ist.

Für eine unabhängige Spezifikation dieser Informationen über Namen statt über protokollspezifische Adressen und Nummern stehen den Anwendungen vier Tabellen zur Verfügung, die die Beziehungen zwischen Namen, deren Aliasnamen und den protokollspezifischen Informationen enthalten.

Information	Tabelle	Datenstruktur	Funktionen
Host	/etc/hosts	hostent	gethostbyname(), gethostbyaddr()
Dienst	/etc/services	servent	getservbyname(), getservbyport()
Protokoll	/etc/protocols	protoent	getprotobyname(), getprotobynumber()
Netzwerk	/etc/networks	netent	getnetbyname(), getnetbyaddr()

Die in der Spalte **Tabelle** angegebenen Informationen werden oftmals zentral mit Hilfe des **Network Information Service (NIS)** verwaltet, damit Änderungen und Erweiterungen nicht auf jedem System lokal zu aktualisieren sind. Die Schnittstellen der Zugriffsfunktionen sind unabhängig von der Lokalität und Organisation der Tabellen.

Zusätzlich existieren für jede der vier Tabellen je drei weitere Funktionen, die das sequentielle Auslesen der gesamten Tabelle unabhängig von deren Lokalität und Organisation ermöglichen.

In den Folgenden Funktionsnennungen sind die Buchstaben **XXX** jeweils zu ersetzen durch einen der Namen in der Spalte *Datenstruktur*:

- **struct XXX \* getXXX(void);**

Diese Funktionen öffnen ggf, die entsprechende Tabelle, lesen jeweils den

nächsten Eintrag und liefern als Resultat einen Zeiger auf die entsprechend initialisierte Datenstruktur **XXX** zurück. Ist das Ende der Tabelle erreicht, liefern sie den Null-Zeiger.

- **void setXXX(int stayopen);**

Diese vier Funktionen öffnen jeweils die entsprechende Tabelle und setzen deren internen Positionszeiger auf den Anfang zurück. Ist das Argument *stayopen* ungleich 0, so können auch die in der jeweiligen Spalte *Funktionen* angegebenen Zugriffsfunktionen mit entsprechenden Selektionskriterien zum sequentiellen Auslesen der Tabelle benutzt werden. Die Verbindung zu einer vt. zugrundeliegenden Netzwerkdatenbank (NIS) wird durch die Ausführung der Zugriffsfunktionen in diesem Fall nicht geschlossen.

- **void endXXX(void);**

Diese Funktionen schliessen die zugehörige Tabelle und beenden eine evt. bestehende NIS-Verbindung;

Das Auslesen von Host- / Netzwerkinformationen, die mit Hilfe des **Domain Name System (DNS)** verwaltet werden, ist nicht möglich. Hier liefern die Zugriffsfunktionen den Inhalt der lokalen bzw. der via NIS verwalteten Tabellen zurück. Die Funktion *gethostent()* kann also **nicht zum Auslesen der gesamten Hostinformationen des Internet benutzt werden!**

#### 4.5.11 IPv6

Das Internet-Protocol der nächsten Generation **IP next generation** ist die neue version des aktuellen Internet-Protokolls in der Version 4 (**IPv4**). Die offizielle Bezeichnung des neuen Protokolls ist **IPv6** (z.Zt. noch als *Draft Standard*).

- Anstatt der bislang 32 Bit langen Adressen werden in IPv6 128 Bit lange Adressen benutzt (vergrößerter, besser strukturierter Adressraum)
- Limitierungen bzgl. der Wegwahl von IP-Paketen (*routing*) und der Konfiguration von Netzwerken werden beseitigt.
- Verbesserte Sicherheitsmechanismen wie Verschlüsselung und Authentifizierung werden integriert.
- Vorgesehen sind Mechanismen für eine prioritätsgesteuerte Datenflussskontrolle von Echtzeitanwendungen (z.B. Übermittlung von multimediale Daten in Echtzeit).

##### Literatur zu IPv6 - z.B.:

- Gilligan, R.E.; Nordmark, E.: Transition Mechanisms for IPv6 Hosts and Routers. RFC 1933, Sun Microsystems, Inc., April 1996
- Huitema, C.: IPv6 - The New Internet Protocol. Prentice Hall, Inc., Englewood Cliffs, 1996

# Kapitel 5

## Netzwerk-Programmierung

### 5.1 Client/Server

#### 5.1.1 Vorbemerkungen

In diesem Abschnitt sollen einige der zuletzt dargestellten Konzepte und Funktionen an einer einfachen Client/Server-Implementierung exemplarisch dargestellt werden. Für die Implementierung des Servers gibt es im Prinzip zwei Möglichkeiten:

- Der Server arbeitet die ankommenden Anforderungen sukzessive ab (**iterative server**).
- Der Server arbeitet die ankommenden Anforderungen parallel ab (**concurrent server**).

#### 5.1.2 concurrent server

Server *forkt*, neuer Prozess bearbeitet Anforderung

```
/* Concurrent Server: conc_srv.c
 * simple error handling
 */

int sockfd, newsockfd;

if ( (sockfd = socket( ... ) ) < 0 ) {
    perror("socket error");
    exit(1);
}
if ( bind(sockfd, ...) < 0 ) {
    perror("bind error");
    exit(1);
}
if ( listen(sockfd, 5) < 0 ) {
    perror("listen error");
    exit(1);
}

while (1) {
    newsockfd = accept(sockfd, ...); /* blockiert */
```

```

    if ( newsockfd < 0 ) {
        perror("accept error");
        exit(1);
    }
    switch (pid = fork() ) {
        case -1:
            perror("fork error");
            exit(1);
        case 0:
            close(sockfd);
            /*verarbeite Anforderung:*/
            doit(newsockfd);
            exit(0);
        default:
            close(newsockfd);
            break;
    }
}

```

### 5.1.3 iterative server

```

/* Iterative Server: iter_srv.c
 * simple error handling
 */

int sockfd, newsockfd;

if ( (sockfd = socket( ... ) ) < 0 ) {
    perror("socket error");
    exit(1);
}
if ( bind(sockfd, ...) < 0 ) {
    perror("bind error");
    exit(1);
}
if ( listen(sockfd, 5) < 0 ) {
    perror("listen error");
    exit(1);
}

while (1) {
    newsockfd = accept(sockfd, ...); /* blockiert */
    if ( newsockfd < 0 ) {
        perror("accept error");
        exit(1);
    }
    /*verarbeite Anforderung:*/
    doit(newsockfd);
    close(newsockfd);
}

```

## 5.2 echo-Server und echo-Client

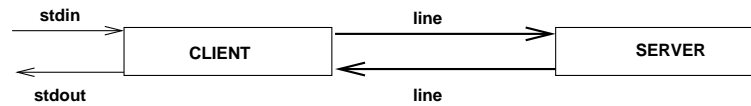


Abbildung 5.1: echo: Client/Server

- Client liest eine Zeile (Folge von Zeichen bis zu einem *newline*!) von *stdin* und übergibt diese an den Server
- Server liest die Zeile über seine Netzwerk-Eingabe und gibt sie über seine Netzwerk-Ausgabe an den Client zurück (*echo*)
- Der Client liest die Zeile von der Socket Schnittstelle, gibt sie an *stdout* aus und terminiert.

## 5.3 Erste Implementierungen

### 5.3.1 Headerfiles

Die wesentlichen Include-Anweisungen für die Inet-Domäne sind im Headerfile **inet.h**, die für die Unix-Domäne in **unix.h** zusammengefasst:

```

/* ----- inet.h -----
 * Definitions for TCP and UDP client / server programs
 */

#include <stdio.h>
#include <unistd.h>
#include <signal.h>
#include <string.h>
#include <errno.h>
#include <sys/wait.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>

#define SERV_UDP_PORT 6000
#define SERV_TCP_PORT 5529
    /* must not conflict with
     * any other TCP server's port
     */
#define SERV_HOST_ADDR "127.0.0.1"
/* "134.60.66.5": it's thales */

/* "127.0.0.1" */
/*local host*/

#define MAXLINE 256
  
```

```

/* unix.h
 *
 * Definitions for UNIX domain stream and datagram
 * client / server programs
 */

#include <stdio.h>
#include <unistd.h>
#include <strings.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <sys/un.h>

#define UNIXSTR_PATH "./s.unixstr"
#define UNIXDG_PATH  "./s.unixdg"

#define MAXLINE 512

```

### 5.3.2 TCP-Verbindung - Concurrent Server

Der **echo-Server** wird als *concurrent server* realisiert. Um **Zombie**-Prozesse zu vermeiden, wird der bereits früher behandelte **signal handler** verwendet:

```

/* ----- sign.h ----- */
/* Signal Handler */

#ifdef SIGN_H
#include <signal.h>
#include <sys/wait.h>
#include <unistd.h>

typedef void (*Sigfunc)(int);

#else

#define SIGN_H
/* avoid multiple includes */
#include <signal.h>

typedef void (*Sigfunc)(int);

Sigfunc ignoresig(int);
/* ignore interrupt and avoid zombies
 * just for midishell (parent)
 */
Sigfunc ignoresig_bg(int);
/* ignore interrupt -
 * just for execution of background commands
 */
Sigfunc entrysig(int);
/* restore reaction on interrupt */

```

```

#endif

/* ----- sign.c ----- */

#define SIGN_H

# include <stdio.h>
# include "sign.h"

void shell_handler(int sig){
    if( (sig == SIGCHLD) || (sig == SIGCLD)) {
        int status;
        waitpid(0, &status, WNOHANG);
    }
    return;
}

struct sigaction newact, oldact;

Sigfunc ignoresig(int sig) {
    static int first = 1;
    newact.sa_handler = shell_handler;
    if (first) {
        first = 0;
        if (sigemptyset(&newact.sa_mask) < 0)
            return SIG_ERR;
        newact.sa_flags = 0;
        newact.sa_flags |= SA_RESTART;
        if (sigaction(sig, &newact, &oldact) < 0)
            return SIG_ERR;
        else
            return oldact.sa_handler;
    } else {
        if (sigaction(sig, &newact, NULL) < 0)
            return SIG_ERR;
        else
            return NULL;
    }
}

struct sigaction newact_bg, oldact_bg;

Sigfunc ignoresig_bg(int sig) {
    newact_bg.sa_handler = SIG_IGN;
    if (sigemptyset(&newact_bg.sa_mask) < 0)
        return SIG_ERR;
    newact_bg.sa_flags = 0;
    newact_bg.sa_flags |= SA_RESTART;
    if (sigaction(sig, &newact_bg, &oldact_bg) < 0)
        return SIG_ERR;
    else

```

```

        return oldact_bg.sa_handler;
    }

```

```

Sigfunc entrysig(int sig)    {
    if (sigaction(sig, &oldact, NULL) < 0 )
        return SIG_ERR;
    else
        return NULL;
}

```

### Der Echo-Server:

```

/* -----main_srv.c -----
 * server using TCP protocol
 *
 * simple error handling
 */

#include "inet.h"
#include "sign.h"

int main() {
    int sockfd,newsockfd,clilen,childpid, n;
    struct sockaddr_in cli_addr, serv_addr;
    char recvline[MAXLINE];

    if( (ignore sig(SIGCHLD) == SIG_ERR) ) {
        perror("Signal");
        exit(1);
    }
    /*
     * open a TCP socket - Internet stream socket
     */
    if ( (sockfd = socket(AF_INET, SOCK_STREAM,0)) < 0) {
        perror("server: can't open stream socket");
        exit(1);
    }

    /*
     * bind our local address so that the client can send us
     */
    /*bzero((char *) &serv_addr, sizeof(serv_addr));*/
    memset((char *) &serv_addr, 0, sizeof(serv_addr));
    serv_addr.sin_family = AF_INET;

    serv_addr.sin_addr.s_addr = htonl(INADDR_ANY);
    /* INADDR_ANY: tells the system that we'll accept a connection
     * on any Internet interface on the system, if it is multihomed
     * Address to accept any incoming messages (-> in.h).
     * INADDR_ANY is defined as ((unsigned long int) 0x00000000)
     */

```



```

serv_addr.sin_port = htons(SERV_TCP_PORT);

if ( bind(sockfd, (struct sockaddr *)&serv_addr,
        sizeof(serv_addr)) < 0) {
    perror("server: can't bind local address");
    exit(1);
}

listen(sockfd,5);

while(1) {
    /*
     * wait for a connection from a client process
     * - concurrent server -
     */
    clilen=sizeof(cli_addr);
    newsockfd = accept(sockfd, (struct sockaddr *) &cli_addr,
                       &clilen);

    if (newsockfd < 0 ) {
        if (errno == EINTR)
            continue; /*try again*/
        perror("server: accept error");
        exit(1);
    }

    if ( (childpid = fork()) < 0) {
        perror("server: fork error");
        close(sockfd);
        close(newsockfd);
        exit(1);
    }
    else if (childpid == 0) {
        close(sockfd);
        if( (n=recv(newsockfd, recvline,MAXLINE,0)) < 0)
            exit(2);
        if( send(newsockfd,recvline,n,0) < n)
            exit(3);
        close(newsockfd);
        exit(0);
    }

    close(newsockfd); /*parent*/
}
}

```

**Der Echo-Client:**

```

/* ----- main_cli.c -----
 * client using TCP protocol
 */

```

```

#include "inet.h"

int main(){
    int sockfd;
    struct sockaddr_in serv_addr;
    char sendline[MAXLINE], recvline[MAXLINE];
    int n;

    /*
     * fill in the structure "serv_addr" with address
     * of server we want to connect with
     */

    /*bzero( (char *) &serv_addr, sizeof(serv_addr));*/
    memset( (char *) &serv_addr, 0, sizeof(serv_addr));
    serv_addr.sin_family = AF_INET;
    serv_addr.sin_addr.s_addr = inet_addr(SERV_HOST_ADDR);
    serv_addr.sin_port = htons(SERV_TCP_PORT);

    /*
     * open a TCP socket - internet stream socket
     */

    if ( (sockfd = socket(AF_INET, SOCK_STREAM,0)) < 0 ) {
        perror("client: can't open stream socket");
        exit(1);
    }

    /*
     * connect to the server
     */

    if (connect(sockfd, (struct sockaddr *) &serv_addr,
        sizeof(serv_addr) ) < 0) {
        perror("client: can't connect to server");
        exit(2);
    }

    printf("%25s", "Client - give input: ");
    if( fgets(sendline,MAXLINE,stdin) != NULL) {
        n = strlen(sendline);
        if (send(sockfd, sendline,n,0) < n)
            exit(3);
        shutdown(sockfd,1);
        if(recv(sockfd,recvline,n,0) < 0)
            exit(4);
        recvline[n] = '\0';
        printf("%25s%s\n","Client - got: ", recvline);
    }

    close(sockfd);
    exit(0);
}

```

**Ausführung:**

```

hypatia$ make -f LinMakefile.srv
gcc -Wall -c main_srv.c
gcc -Wall -c sign.c
gcc -Wall -o server main_srv.o sign.o
hypatia$ make -f LinMakefile.cli
gcc -Wall -c main_cli.c
gcc -Wall -o client main_cli.o
hypatia$ server &
[1] 658
hypatia$ client
    Client - give input: Eine Eingabezeile
          Client - got: Eine Eingabezeile

hypatia$ ps | grep server
    658 pts/0    00:00:00 server
hypatia$ kill 658
[1]+  Terminated                  server
hypatia$

```

**5.3.3 UDP-Verbindung - Iterative Server****Der Echo-Server:**

```

/* -----main_srv.c -----
 * server using UDP protocol
 *
 * simple error handling
 */

#include "inet.h"

int main() {
    int sockfd, clilen, n;
    struct sockaddr_in cli_addr, serv_addr;
    char recvline[MAXLINE];

    /*
     * open a UDP socket - Internet datagramm socket
     */
    if ( (sockfd = socket(AF_INET, SOCK_DGRAM, 0)) < 0) {
        perror("server: can't open datagramm socket");
        exit(1);
    }

    /*
     * bind our local address so that the client can send us
     */
    bzero((char *) &serv_addr, sizeof(serv_addr));
    serv_addr.sin_family = AF_INET;

```

```

serv_addr.sin_addr.s_addr = htonl(INADDR_ANY);
/* INADDR_ANY: tells the system that we'll accept a connection
 * on any Internet interface on the system, if it is multihomed
 * Address to accept any incoming messages (-> in.h).
 * INADDR_ANY is defined as ((unsigned long int) 0x00000000)
 */

serv_addr.sin_port = htons(SERV_UDP_PORT);

if ( bind(sockfd, (struct sockaddr *)&serv_addr,
        sizeof(serv_addr)) < 0) {
    perror("server: can't bind local address");
    exit(1);
}
clilen=sizeof(cli_addr);

while(1) {
    /*
     * wait for a connection from a client process
     * - iterative server -
     */

    if( (n=recvfrom(sockfd, recvline,MAXLINE,0,
        (struct sockaddr *)&cli_addr,&clilen)) < 0) {
        perror("server - recvfrom");
        exit(2);
    }
    if( sendto(sockfd,recvline,n,0, (struct sockaddr *)&cli_addr,
        sizeof(cli_addr)) < n) {
        perror("server - sendto");
        exit(3);
    }
}
}

```

### Der Echo-Client:

```

/* ----- main_cli.c -----
 * client using UDP protocol
 */

#include "inet.h"

int main(int argc, char * argv){
    int sockfd;
    struct sockaddr_in cli_addr, serv_addr;
    char sendline[MAXLINE], recvline[MAXLINE];
    int n;

    /*
     * fill in the structure "serv_addr" with address
     * of server we want to connect with
     */

```

```

bzero( (char *) &serv_addr, sizeof(serv_addr));
serv_addr.sin_family = AF_INET;
serv_addr.sin_addr.s_addr = inet_addr(SERV_HOST_ADDR);
serv_addr.sin_port = htons(SERV_UDP_PORT);

/*
 * open a UDP socket - internet datagramm socket
 */

if ( (sockfd = socket(AF_INET, SOCK_DGRAM,0)) < 0 ) {
    perror("client: can't open datagramm socket");
    exit(1);
}

/*
 * bind any local address for us
 */
bzero( (char *) &cli_addr, sizeof(cli_addr));
cli_addr.sin_family = AF_INET;
cli_addr.sin_addr.s_addr = htonl(INADDR_ANY);
cli_addr.sin_port = htons(0);

if (bind(sockfd, (struct sockaddr *) &cli_addr,
        sizeof(cli_addr)) < 0 ) {
    perror("client: can't bind local address");
    exit(2);
}

printf("%25s", "Client - give input: ");
if( fgets(sendline,MAXLINE,stdin) != NULL) {
    n = strlen(sendline);
    if (sendto(sockfd, sendline,n,0,(struct sockaddr *)&serv_addr,
        sizeof(serv_addr)) < n) {
        perror("client - sendto");
        exit(3);
    }
    shutdown(sockfd,1);
    if(( n=recvfrom(sockfd,recvline,n,0,(struct sockaddr *)0,
        (int *)0)) < 0) {
        perror("client - recvfrom");
        exit(4);
    }
    recvline[n] = '\0';
    printf("%25s%s\n","Client - got: ", recvline);
}

close(sockfd);
exit(0);
}

```

**Ausführung:**

```
hypatia$ make -f LinMakefile.srv
```

```
gcc -Wall -c main_srv.c
gcc -Wall -o server main_srv.o
hypatia$ make -f LinMakefile.cli
gcc -Wall -c main_cli.c
gcc -Wall -o client main_cli.o
hypatia$ server &
[1] 946
hypatia$ client
    Client - give input: Noch 'ne Zeile
        Client - got: Noch 'ne Zeile

hypatia$ ps | grep server
    946 pts/0    00:00:00 server
hypatia$ kill 946
[1]+  Terminated                  server
hypatia$
```

### 5.3.4 TCP-Verbindung in der UNIX Domain

#### Der Echo-Server:

```

/* main_srv.c
 * server using UNIX domain stream protocol
 */

#include "unix.h"
#include "sign.h"

int main() {
    int sockfd,newsockfd,clilen,childpid,servlen,n;
    struct sockaddr_un cli_addr, serv_addr;
    char recvline[MAXLINE];

    if( (ignore sig(SIGCHLD) == SIG_ERR) ) {
        perror("Signal");
        exit(1);
    }

    /*
     * open a UNIX domain stream socket
     */
    if ((sockfd=socket(AF_UNIX,SOCK_STREAM,0)) < 0 ) {
        perror("server: can't open a stream socket");
        exit(2);
    }

    /*
     * bind our local address so that the client can send to us
     */

    /* set all with zeros */
    bzero( (char *) &serv_addr, sizeof(serv_addr));

    serv_addr.sun_family = AF_UNIX;
    strcpy(serv_addr.sun_path, UNIXSTR_PATH);

    /* determine length of address: */
    servlen = strlen(serv_addr.sun_path) + sizeof(serv_addr.sun_family);

    if (bind(sockfd, (struct sockaddr *) &serv_addr, servlen) < 0 ) {
        perror("server: can't bind local address");
        exit(3);
    }

    listen(sockfd,5);

    while(1) {
        /*
         * wait for a connection from a client process
         * - concurrent server -
         */

```

```

    clilen = sizeof(cli_addr);
    newsockfd = accept(sockfd, (struct sockaddr *) &cli_addr,
        &clilen);
    if (newsockfd < 0 ) {
        perror("server: accept error");
        exit(4);
    }

    if ( (childpid = fork() ) < 0 ) {
        perror("server: can't fork");
        exit(5);
    }
    else if (childpid == 0) { /* child */
        close(sockfd);
        if( (n=recv(newsockfd, recvline, MAXLINE,0)) < 0)
            exit(6);
        if( send(newsockfd, recvline,n,0) < n)
            exit(7);
        close(newsockfd);
        exit(0);
    }
    /*parent:*/
    close(newsockfd);
}
}

```

### Der Echo-Client:

```

/* main_cli.c
 *
 * client using Unix domain stream protocol
 */

#include "unix.h"

int main() {
    int sockfd, servlen;
    struct sockaddr_un serv_addr;
    char sendline[MAXLINE], recvline[MAXLINE];
    int n;

    /*
     * Fill in the structure "serv_addr" with the
     * address of the server that we want to sent do
     */
    bzero( (char *) &serv_addr, sizeof(serv_addr) );
    serv_addr.sun_family = AF_UNIX;
    strcpy(serv_addr.sun_path, UNIXSTR_PATH);
    servlen = strlen(serv_addr.sun_path) + sizeof(serv_addr.sun_family);

    /*
     * open an Unix domain stream socket

```



```

    */
    if ( (sockfd = socket(AF_UNIX, SOCK_STREAM, 0) ) < 0) {
        perror("client: can't open stream socket");
        exit(1);
    }

    /*
     * connect to the server
     */
    if (connect(sockfd, (struct sockaddr *) &serv_addr, servlen) < 0) {
        perror("client: can't connect to server");
        exit(1);
    }

    printf("%25s", "Client - give input: ");
    if( fgets(sendline,MAXLINE,stdin) != NULL) {
        n = strlen(sendline);
        if (send(sockfd, sendline,n,0) < n)
            exit(3);
        shutdown(sockfd,1);
        if(recv(sockfd,recvline,n,0) < 0)
            exit(4);
        recvline[n] = '\0';
        printf("%25s%s\n","Client - got: ", recvline);
    }

    close(sockfd);
    exit(0);
}

```

**Ausführung:**

```

hypatia$ make -f LinMakefile.sr v
gcc -Wall -c main_srv.c
gcc -Wall -c sign.c
gcc -Wall -o server main_srv.o sign.o
hypatia$ make -f LinMakefile.cl i
gcc -Wall -c main_cli.c
gcc -Wall -o client main_cli.o
hypatia$ server &
[1] 959
hypatia$ ls -l s.*
srwxr-xr-x  1 swg      users          0 Apr 30 12:12 s.unixstr
hypatia$ client
    Client - give input: one line
    Client - got: one line

hypatia$ ps | grep server
  959 pts/4    00:00:00 server
hypatia$ kill 959
[1]+  Terminated                  server
hypatia$ server &

```

```
[1] 966
hypatia$ server: can't bind local address: Address already in use

[1]+  Exit 3                  server
hypatia$ rm s.unixstr
hypatia$
```

### 5.3.5 Modifikation der ersten Implementierung

Der Client liest Zeile für Zeile von der Standardeingabe, schickt diese sukzessive an den Server, liest sie wieder und gibt sie „markiert“ wieder an die Standardausgabe. Im folgenden wird nur der geänderte Teil dargestellt:

#### Der Client-Teil:

```
while ( fgets(sendline,MAXLINE,stdin) != NULL) {
    n = strlen(sendline);
    if (send(sockfd, sendline,n,0) < n)
        exit(3);
    if(recv(sockfd,recvline,n,0) < 0)
        exit(4);
    recvline[n] = '\0';
    printf(">>> %s", recvline);
}
```

#### Der Server-Teil:

```
else if (childpid == 0) {
    close(sockfd);
    while ( (n=recv(newsockfd, recvline,MAXLINE,0)) > 0) {
        if( send(newsockfd,recvline,n,0) < n)
            exit(3);
    }
    close(newsockfd);
    exit(0);
}
```

#### Ausführung:

```
hypatia$ server &
[1] 1282
hypatia$ client < cli.src
>>> while ( fgets(sendline,MAXLINE,stdin) != NULL) {
>>>     n = strlen(sendline);
>>>     if (send(sockfd, sendline,n,0) < n)
>>>         exit(3);
>>>     if(recv(sockfd,recvline,n,0) < 0)
>>>         exit(4);
>>>     recvline[n] = '\0';
>>>     printf(">>> %s", recvline);
```

```
>>> }
hypatia$
```

### 5.3.6 Anmerkungen

Diese Implementierungen sind wenig robust (triviales Fehlerhandling); sie sollten die prinzipielle Anwendung der Socket Funktionen demonstrieren. Lese- und Schreiboperationen auf *Stream Sockets* können auch weniger als die spezifizierte Anzahl von Bytes als Resultatwert liefern. Dies ist generell bei allen über einen Deskriptor referenzierten Objekten möglich, wenn beispielsweise der zugrundeliegende Systemaufruf durch ein Signal unterbrochen wurde oder der Deskriptor sich im nicht-blockierenden Modus befindet. Dies kann auch dadurch passieren, dass beim Lesen oder Schreiben Ressourcenbeschränkungen verletzt werden oder momentan keine weitere Eingabe zur Verfügung steht. Die exakte Semantik ist von dem konkret referenzierten Objekt abhängig.

Bei *Stream Sockets* ist das Ein- / Ausgabeverhalten vom Erreichen der Socket-Puffergrenzen im Betriebssystemkern abhängig. Die Kommunikationsverbindung ist hier **bidirektional** und **voll-duplex**; jeder Socket besitzt einen **Sendepuffer** und einen **Empfangspuffer**, die beide voneinander unabhängig sind. Die Verlässlichkeit des Datentransfers wie auch die Datenflusskontrolle werden über das TCP-Protokoll und die internen Algorithmen der TCP-Implementierungen mit Hilfe der Sende- und Empfangspuffer der miteinander verbundenen Sockets geregelt.

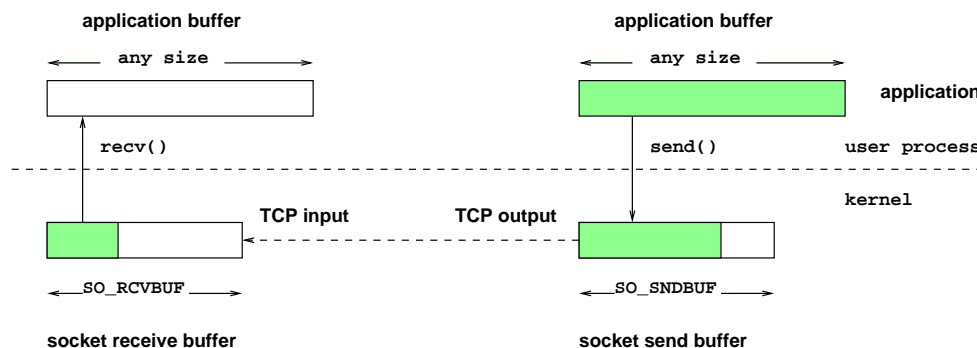


Abbildung 5.2: Socket-Pufferung (vereinfacht)

Die Dimensionen des Sende- und Empfangspuffers sind systemabhängig voreingestellt und lassen sich über die beiden Socket-Optionen **SO\_SNDBUF** und **SO\_RCVBUF** modifizieren. Die Dimension des Puffers in der Anwendung ist frei wählbar.

Die Ausführung der I/O-Funktionen auf Sockets bewirkt letztlich nur, dass die Daten zwischen dem Puffer der Anwendung (im Benutzeradressraum) und dem entsprechenden Socket-Puffer (im Adressraum des Kerns) kopiert werden. Im Fall einer Schreiboperation zeigt die als Resultat gelieferte Anzahl von Bytes nur an, dass diese Anzahl in den Sendepuffer des Socket kopiert

wurden und nicht, dass diese Zahl von Bytes den Empfänger erreicht hat. Entsprechend werden beim Lesen die momentan im Empfangspuffer des Socket gehaltenen Daten in den Anwendungspuffer kopiert. Der tatsächliche Datentransfer über den Kommunikationskanal wird von den Ein- / Ausgabefunktionen der TCP-Implementierung **asynchron** vorgenommen und kann von der Anwendung nicht beeinflusst werden. Das Verhalten der Lese- und Schreibfunktionen ist also von den momentan im Empfangspuffer enthaltenen Daten bzw. dem freien Bereich im Sendepuffer abhängig (ähnlich der Funktionsweise von Pipes).

Lesefunktionen liefern (im Erfolgsfall) immer das Minimum aus der angeforderten Datenmenge und den im Empfangspuffer gehaltenen Bytes zurück – sofern der Systemaufruf nicht durch ein Signal unterbrochen wurde – maximal aber **SO\_RCVBUF** Bytes! Schreibfunktionen versuchen, die spezifizierte Anzahl von Bytes zu senden und blockieren solange, bis die gesamte Datenmenge in den Puffer kopiert wurde. Im Erfolgsfall gilt also, dass die als Resultat gelieferte Anzahl der spezifizierten Anzahl entspricht – sofern kein Signal den Systemaufruf unterbricht! Diese Funktionalität kann auch in den Lesefunktionen der Socket-Schnittstelle durch Spezifikation der Flagge **MSG\_WAITALL** erreicht werden. Die Existenz diese Flagge wie auch die Eigenschaft, dass Schreibfunktionen versuchen, die spezifizierte Datenmenge insgesamt zu versenden, sind allerdings von der jeweiligen Implementierung der Socket-Schnittstelle abhängig!

In vielen Netzerkanwendungen ist es oft notwendig, logisch unvollständige Lese- und Schreiboperationen – die im Fall eines unterbrochenen Systemaufrufs ja auch keinen Fehler darstellen – entsprechend zu behandeln. Dazu können z.B. die folgenden beiden Funktionen verwendet werden:

- **sendn():**

```
/* ipc_send.c */
# define IPC_SEND_H
# include "ipc_send.h"

ssize_t sendn(int fd, char * buf, size_t len) {

    char * ptr = buf;
    size_t nc;
    ssize_t n;

    for(nc = len; nc > 0; ptr += n, nc -=n) {
        send_again:
        if( (n=send(fd,ptr,nc,0)) <= 0)
            if(errno == EINTR) {
                errno = 0;
                goto send_again;
            } else
                return ( len -= nc) ? len : -1;
    }
    return len;
}
```

- **recvn():**

```

/* ipc_recv.c */
# define IPC_RECV_H
# include "ipc_recv.h"

ssize_t recvn(int fd, char * buf, size_t len) {

    char *    ptr = buf;
    size_t    nc;
    ssize_t    n;
    int        flags = 0;

    for(nc = len; nc > 0; ptr += n, nc -= n) {
        recv_again:
        if( (n = recv(fd, ptr, nc, flags)) < 0 )
            if(errno == EINTR) {
                errno = 0;
                goto recv_again;
            } else
                return (len -= nc) ? len : -1;
        else if( n == 0 ) {
            errno = ENOTCONN;
            return (len -= nc) ? len : 0;
        }
        if(buf[n-1] == '\n') return n;
    }
    return len;
}

```

Diese beiden Funktionen implementieren inkrementelles Schreiben und Lesen, die die Anwendung jeweils solange blockieren, bis die spezifizierte Datenmenge insgesamt gesendet resp. empfangen wurde; sie berücksichtigen zudem die Möglichkeit, dass Systemaufrufe durch Signale unterbrochen werden können.

Die Funktion **sendn()** ist unabhängig von der Semantik und Implementierung der Socket-Funktion **send()** realisiert: in einer Schleife werden die noch ausstehenden Daten durch weitere Aufrufe von **send()** gesendet, falls die als Resultat gelieferte Anzahl von Bytes kleiner als die spezifizierte Datenmenge ist. Die Funktion **recvn()** ist ähnlich realisiert; unterstützt die Socket-Funktion **recv()** die Flagge **MSG\_WAITALL**, so wird diese auch zur Performance-Steigerung beim Einlesen der Daten genutzt. Ist inkrementelles Lesen erforderlich, so übernimmt in diesem Fall die Funktion **recv()** diese Aufgabe selbst und erspart der Anwendung weitere Systemaufrufe. Dennoch ist auch in diesem Fall die Schleifenkonstruktion erforderlich, falls **recv()** durch ein Signal unterbrochen wird und einen weiteren Aufruf notwendig macht.

Der Resultatwert beider Funktionen ist im Erfolgsfall die in der Komponente *len* angegebene Anzahl Bytes. Dies gilt auch bei der Spezifikation von 0 Bytes.

Im Fehlerfall liefert **sendn()** als Resultat **-1**, sofern noch keine Daten gesendet wurden, andernfalls die Anzahl der bis zum Auftreten des Fehlers erfolgreich gesendeten Bytes. Entsprechendes Fehlerverhalten gilt auch für **recv()** mit der Ausnahme, dass eine ordnungsmäßige Termination der Verbindung durch den Kommunikationspartner als logischer Fehler behandelt wird und den Resultatwert **0** liefert, sofern noch keine Daten empfangen wurden.

## 5.4 Verbesserte Implementierungen

### 5.4.1 Zeilenorientierter Echo-Server

Wie in den vorigen Beispielen wird ein zeilenorientierter Echo-Server betrachtet; der Client liest also eine Zeile (Folge von Bytes bis *newline*), schickt diese an den Server und dieser schickt sie wieder zurück. Dazu sollen die beiden Funktionen **sendn()** und **recv()** verwendet werden. Dabei ist aber zu beachten, dass Zeilen deutlich kürzer sein können als die entsprechenden Puffer – die Flagge **MSG\_WAITALL** würde hier zu einer Blockade der Kommunikation führen.

Die Dateien insgesamt:

- **sign.h** und **sign.c** wie bisher
- Das Hauptprogramm des Servers:

```
/* -----main_srv.c -----
--
* server using TCP protocol
*
* simple error handling
*/

#include "inet.h"
#include "sign.h"
#include "str_echo.h"

int main(int argc, char **argv) {

    int sockfd,newsockfd,clilen,childpid;
    struct sockaddr_in cli_addr, serv_addr;

    if( (ignoresig(SIGCHLD) == SIG_ERR) ) {
        perror("Signal");
        exit(1);
    }
    /*
    * open a TCP socket - Internet stream socket
    */
    if ( (sockfd = socket(AF_INET, SOCK_STREAM,0)) < 0) {
        perror("server: can't open stream socket");
        exit(1);
    }
}
```

```

/*
 * bind our local address so that the client can send us
 */
bzero((char *) &serv_addr, sizeof(serv_addr));
serv_addr.sin_family = AF_INET;

serv_addr.sin_addr.s_addr = htonl(INADDR_ANY);
/* INADDR_ANY: tells the system that we'll accept a connection
 * on any Internet interface on the system, if it is multihomed
 * Address to accept any incoming messages (-> in.h).
 * INADDR_ANY is defined as ((unsigned long int) 0x00000000)
 */

serv_addr.sin_port = htons(SERV_TCP_PORT);

if ( bind(sockfd, (struct sockaddr *)&serv_addr,
        sizeof(serv_addr)) < 0) {
    perror("server: can't bind local address");
    exit(1);
}

listen(sockfd,5);

while(1) {
    /*
     * wait for a connection from a client process
     * - concurrent server -
     */
    clilen=sizeof(cli_addr);
    newsockfd = accept(sockfd, (struct sockaddr *) &cli_addr,
                       &clilen);

    if (newsockfd < 0 ) {
        if (errno == EINTR)
            continue; /*try again*/
        perror("server: accept error");
        exit(1);
    }

    if ( (childpid = fork()) < 0) {
        perror("server: fork error");
        close(sockfd);
        close(newsockfd);
        exit(1);
    }
    else if (childpid == 0) {
        close(sockfd);
        /*******/
        str_echo(newsockfd);
        /*******/

        exit(0);
    }
}

```

```

        close(newsockfd); /*parent*/
    }
}

```

- Die Funktionalität des Servers:

```

/* str_echo.c
 * function used in connection-oriented servers
 * read a stream socket one line at a time,
 * and write each line back to the sender
 * return when the connection is terminated
 */
# define STR_ECHO_H
# include "str_echo.h"

void str_echo(int sockfd) {
    int n;
    char line[MAXLINE];

    while(1) {
        n = recvn(sockfd, line, MAXLINE);
        if (n == 0)
            return; /*connection terminated*/
        else if ( n < 0 ) {
            perror("str_echo: readline error");
            exit(3);
        }
        /*Ausgabe auf stderr des Servers: */

        if (sendn(sockfd, line, n) != n) {
            perror("str_echo: write error");
            exit(3);
        }
    }
}

```

- Das Hauptprogramm des Client:

```

/* ----- main_cli.c -----
-
 * client using TCP protocol
 */

#include "inet.h"
#include "str_cli.h"
#include "str_echo.h"

int main(){
    int sockfd;
    struct sockaddr_in serv_addr;

```



```

/*
 * fill in the structure "serv_addr" with address
 * of server we want to connect with
 */

bzero( (char *) &serv_addr, sizeof(serv_addr));
serv_addr.sin_family = AF_INET;
serv_addr.sin_addr.s_addr = inet_addr(SERV_HOST_ADDR);
serv_addr.sin_port = htons(SERV_TCP_PORT);

/*
 * open a TCP socket - internet stream socket
 */

if ( (sockfd = socket(AF_INET, SOCK_STREAM,0)) < 0 ) {
    perror("client: can't open stream socket");
    exit(2);
}

/*
 * connect to the server
 */

if (connect(sockfd, (struct sockaddr *) &serv_addr,
    sizeof(serv_addr) ) < 0) {
    perror("client: can't connect to server");
    exit(3);
}

/*****/
str_cli(stdin, sockfd);
/*****/

close(sockfd);
printf("\n");
exit(0);
}

```

- Die Funktionalität des Client:

```

/* str_cli.c
 * function used by connection-oriented clients
 *
 * read the contents of the FILE *fp, write each line to the
 * stream socket (to the server process), then read a li-
ne back
 * from the socket and write it to stdout
 *
 * return to caller when an EOF is encountered on the in-
put file
 */

```

```

# define STR_CLI_H
# include "str_cli.h"

void str_cli(FILE *fp, int sockfd) {
    int n;
    char sendline[MAXLINE], recvline[MAXLINE+1];

    while (fgets(sendline, MAXLINE, fp) != NULL) {
        n = strlen(sendline);
        if (sendn(sockfd, sendline, n) != n) {
            perror("str_cli: write error on socket");
            exit(1);
        }

        /*
         * now read a line from the socket and
         * write it to stdout
         */
        n = recvn(sockfd, recvline, n);
        if (n < 0) {
            perror("str_cli: readline error");
            exit(1);
        }
        recvline[n] = '\0';
        printf(">>> ");
        fputs(recvline, stdout);

    }

    if (ferror(fp)) {
        perror("str_cli: error reading file");
        exit(1);
    }
}

```

- Die Funktion **recvn()** — ohne die Flagge **MSG\_WAITALL**

```

/* ipc_recv.c */
# define IPC_RECV_H
# include "ipc_recv.h"

ssize_t recvn(int fd, char * buf, size_t len) {

    char *    ptr = buf;
    size_t    nc;
    ssize_t    n;
    int        flags = 0;

    for(nc = len; nc > 0; ptr += n, nc -= n) {

```

```

recv_again:
if( (n = recv(fd, ptr, nc, flags)) < 0 )
    if(errno == EINTR) {
        errno = 0;
        goto recv_again;
    } else
        return (len -= nc) ? len : -1;
else if( n == 0 ) {
    errno = ENOTCONN;
    return (len -= nc) ? len : 0;
}
if(buf[n-1] == '\n') return n;
}
return len;
}

```

- Ausführung:

```

hypatia$ make -f LinMakefile.srv
gcc -Wall -c main_srv.c
gcc -Wall -c ipc_send.c
gcc -Wall -c ipc_recv.c
gcc -Wall -c str_echo.c
gcc -Wall -c sign.c
gcc -Wall -o server main_srv.o str_echo.o ipc_send.o ipc_recv.o sign.o
hypatia$ make -f LinMakefile.cli
gcc -Wall -c main_cli.c
gcc -Wall -c str_cli.c
gcc -Wall -o client main_cli.o str_cli.o ipc_recv.o ipc_send.o
hypatia$ server &
[1] 1883
hypatia$ client
eine Zeile
>>> eine Zeile
und noch eine Zeile
>>> und noch eine Zeile
^d
hypatia$

```

### 5.4.2 Ein „stream“-basierter Echo-Server

Der Client liest eine beliebige Folge von Zeichen aus der Standardeingabe, schickt sie an den Server, der schickt sie zurück und der Client kopiert das Ganze an die Standardausgabe. Dazu müssen auf Client-Seite kleinere Änderungen (statt **fgets()** wird **read()** verwendet) vorgenommen werden.

Die Dateien insgesamt:

- **sign.h** und **sign.c** wie bisher
- Das Hauptprogramm des Servers:

```

/* -----main_srv.c -----
--
* server using TCP protocol
*
* simple error handling
*/

#include "inet.h"
#include "sign.h"
#include "str_echo.h"

int main(int argc, char **argv) {

    int sockfd,newsockfd,clilen,childpid;
    struct sockaddr_in cli_addr, serv_addr;

    if( (ignore sig(SIGCHLD) == SIG_ERR) ) {
        perror("Signal");
        exit(1);
    }
    /*
    * open a TCP socket - Internet stream socket
    */
    if ( (sockfd = socket(AF_INET, SOCK_STREAM,0)) < 0) {
        perror("server: can't open stream socket");
        exit(1);
    }

    /*
    * bind our local address so that the client can send us
    */
    bzero((char *) &serv_addr, sizeof(serv_addr));
    serv_addr.sin_family = AF_INET;

    serv_addr.sin_addr.s_addr = htonl(INADDR_ANY);
    /* INADDR_ANY: tells the system that we'll accept a connection
    * on any Internet interface on the system, if it is multihomed
    * Address to accept any incoming messages (-> in.h).
    * INADDR_ANY is defined as ((unsigned long int) 0x00000000)
    */

    serv_addr.sin_port = htons(SERV_TCP_PORT);

    if ( bind(sockfd, (struct sockaddr *)&serv_addr,
        sizeof(serv_addr)) < 0) {
        perror("server: can't bind local address");
        exit(1);
    }

    listen(sockfd,5);

    while(1) {
        /*
        * wait for a connection from a client process

```

```

    * - concurrent server -
    */
    clilen=sizeof(cli_addr);
    newsockfd = accept(sockfd, (struct sockaddr *) &cli_addr,
                       &clilen);

    if (newsockfd < 0 ) {
        if (errno == EINTR)
            continue; /*try again*/
        perror("server: accept error");
        exit(1);
    }

    if ( (childpid = fork()) < 0) {
        perror("server: fork error");
        close(sockfd);
        close(newsockfd);
        exit(1);
    }
    else if (childpid == 0) {
        close(sockfd);

        /*******/
        str_echo(newsockfd);
        /*******/

        exit(0);
    }

    close(newsockfd); /*parent*/
}
}

```

- Die Funktionalität des Servers:

```

/* str_echo.c
 * function used in connection-oriented servers
 * read a stream socket one line at a time,
 * and write each line back to the sender
 * return when the connection is terminated
 */
# define STR_ECHO_H
# include "str_echo.h"

void str_echo(int sockfd) {
    int n;
    char line[MAXLINE];

    while(1) {
        n = recvn(sockfd,line,MAXLINE);
        if (n == 0)
            return; /*connection terminated*/
    }
}

```

```

        else if ( n < 0 ) {
            perror("str_echo: readline error");
            exit(3);
        }
        /*Ausgabe auf stderr des Servers: */

        if (sendn(sockfd, line, n) != n) {
            perror("str_echo: write error");
            exit(3);
        }
    }
}

```

- Das Hauptprogramm des Client:

```

/* ----- main_cli.c -----
-
* client using TCP protocol
*/

#include "inet.h"
#include "str_cli.h"
#include "str_echo.h"

int main(){
    int sockfd;
    struct sockaddr_in serv_addr;

    /*
     * fill in the structure "serv_addr" with address
     * of server we want to connect with
     */

    bzero( (char *) &serv_addr, sizeof(serv_addr));
    serv_addr.sin_family = AF_INET;
    serv_addr.sin_addr.s_addr = inet_addr(SERV_HOST_ADDR);
    serv_addr.sin_port = htons(SERV_TCP_PORT);

    /*
     * open a TCP socket - internet stream socket
     */

    if ( (sockfd = socket(AF_INET, SOCK_STREAM, 0)) < 0 ) {
        perror("client: can't open stream socket");
        exit(2);
    }

    /*
     * connect to the server
     */

    if (connect(sockfd, (struct sockaddr *) &serv_addr,

```

```

        sizeof(serv_addr) ) < 0) {
    perror("client: can't connect to server");
    exit(3);
}

/*****/
str_cli(sockfd);
/*****/

close(sockfd);
printf("\n");
exit(0);
}

```

- Die Funktionalität des Client:

```

/* str_cli.c
 * function used by connection-oriented clients
 *
 * read from stdin (0), write to
 * stream socket (to the server process), then read back
 * from the socket and write to stdout
 *
 * return to caller when an EOF is encountered on the input
 */

# define STR_CLI_H
# include "str_cli.h"

void str_cli(int sockfd) {
    int n;
    char sendline[MAXLINE], recvline[MAXLINE+1];

    for(;;) {
        if( (n = read(0, sendline, MAXLINE)) < 0 ) {
            perror("read");
            exit(1);
        } else if (n == 0)
            return; /*done*/

        if (sendn(sockfd, sendline, n) != n) {
            perror("str_cli: write error on socket");
            exit(2);
        }

        /*
         * now read from the socket and write to stdout
         */
        if( (n = recvn(sockfd, recvline, n)) != n) {
            perror("str_cli: read error on socket");
            exit(3);
        }
    }
}

```

```

        if(write(1,recvline,(size_t) n) < 0) {
            perror("write");
            exit(4);
        }
    }
}

```

- Die Funktion **recvn()** — jetzt mit der Flagge **MSG\_WAITALL**

```

/* ipc_recv.c */
# define IPC_RECV_H
# include "ipc_recv.h"

ssize_t recvn(int fd, char * buf, size_t len) {

    char *    ptr = buf;
    size_t    nc;
    ssize_t    n;
    int        flags = 0;

# ifdef MSG_MSG_WAITALL
    flags |= MSG_WAITALL
# endif

    for(nc = len; nc > 0; ptr += n, nc -= n) {
        recv_again:
        if( (n = recv(fd, ptr, nc, flags)) < 0 )
            if(errno == EINTR) {
                errno = 0;
                goto recv_again;
            } else
                return (len -= nc) ? len : -1;
        else if( n == 0 ) {
            errno = ENOTCONN;
            return (len -= nc) ? len : 0;
        }
        if(buf[n-1] == '\n') return n;
    }
    return len;
}

```

- Ausführung:

```

hypatia$ make -f LinMakefile.srv
gcc -Wall -c main_srv.c
gcc -Wall -c ipc_send.c
gcc -Wall -c ipc_recv.c

```



```

gcc -Wall -c str_echo.c
gcc -Wall -c sign.c
gcc -Wall -o server main_srv.o str_echo.o ipc_send.o ipc_recv.o sign.o
hypatia$ make -f LinMakefile.cli
gcc -Wall -c main_cli.c
gcc -Wall -c str_cli.c
gcc -Wall -o client main_cli.o str_cli.o ipc_recv.o ipc_send.o
hypatia$ server &
[1] 1997
hypatia$ client < ipc_recv.c
/* ipc_recv.c */
# define IPC_RECV_H
# include "ipc_recv.h"

ssize_t recvn(int fd, char * buf, size_t len) {

    char *    ptr = buf;
    size_t    nc;
    ssize_t    n;
    int        flags = 0;

# ifdef MSG_WAITALL
    flags |= MSG_WAITALL
# endif

    for(nc = len; nc > 0; ptr += n, nc -= n) {
        recv_again:
        if( (n = recv(fd, ptr, nc, flags)) < 0 )
            if(errno == EINTR) {
                errno = 0;
                goto recv_again;
            } else
                return (len -= nc) ? len : -1;
        else if( n == 0 ) {
            errno = ENOTCONN;
            return (len -= nc) ? len : 0;
        }
        if(buf[n-1] == '\n') return n;
    }
    return len;
}

hypatia$

```



# Abbildungsverzeichnis

3.1	ISO-OSI-Referenzmodell . . . . .	60
3.2	TCP/IP . . . . .	61
3.3	Ethernet Frame Format . . . . .	64
3.4	Layered System Architecture . . . . .	65
3.5	Internet Architecture . . . . .	66
3.6	Abstraktion der Identifikation . . . . .	66
3.7	32-Bit IP-Adresse . . . . .	67
3.8	IP-Adresse: dotted decimal form . . . . .	68
3.9	Adress-Auflösung . . . . .	68
3.10	TCP/IP-Schichtenmodell . . . . .	69
3.11	IP Datagramm - grober Aufbau . . . . .	70
3.12	IP Datagram - im Detail . . . . .	70
3.13	TCP/IP Layering Model . . . . .	71
3.14	UDP - Format . . . . .	73
3.15	UDP-Demultiplexing . . . . .	74
4.1	Vereinf. Modell der Implementierung von Sockets unter BSD . . . . .	78
4.2	Verbindungsorientierte Client-Server-Kommunikation . . . . .	88
4.3	Aufbau einer verbind.-orient. Client/Server-Kommunikation . . . . .	90
4.4	Aufbau einer verb.-losen Client/Server-Kommunikation . . . . .	100
4.5	Organisation der Adress-Struktur sockaddr_in . . . . .	106
4.6	dotted-decimal notation . . . . .	108
4.7	Methoden der Adress-Resolution . . . . .	111
5.1	echo: Client/Server . . . . .	121
5.2	Socket-Pufferung (vereinfacht) . . . . .	135

# Index

- /etc/hosts, 144
- /etc/networks, 149
- /etc/protocols, 149
- /etc/services, 147
  
- accept(), 128, 135
- address resolution, 108
- Address Resolution Protocol, 108
- Adress-Resolution, 142, 144
- AF\_INET, 124
- AF\_UNIX, 124, 138
- alarm(), 37
- Anwendungsschicht, 101
- API, 121
- argv, 14, 81
- ARP, 108
- asynchron, 169
  
- backlog, 128
- Beendigungsstatus, 4, 17
- best-effort delivery, 104, 109
- bidirektional, 168
- Big-Endian, 140
- bind(), 125, 132–134
- Bitübertragungsschicht, 100
- Bootstrapping, 18
- Bridge, 102
- broadcast, 104
- Broadcast-Adresse, 141
- BSD, 121
- builtins, 73
- byte order, 108
  
- cd, 73
- child process, 9
- Client, 127
- Client-Server, 66
- close(), 65, 132, 135
- concurrent server, 152, 155, 164
- connect(), 127, 133, 135
- Connectionless, 109
- control process, 2
- CSMA/CD-Verfahren, 102
- ctrl-\, 21
- ctrl-c, 21
  
- current working directory, 12, 14
  
- Datagram Encapsulation, 110
- Datagramm, 132
- Datagramm Socket, 123, 132, 133
- Datendarstellungsschicht, 101
- delete, 21
- DIN/ISO-OSI, 100
- DNS, 144, 151
- Domain Name System, 144, 151
- dotted decimal notation, 108, 141
- dup(), 63
- dynamic and/or private ports, 147
  
- echo-Client, 158, 161, 165
- echo-Server, 155, 160, 164
- effective group ID, 12
- effective user ID, 12
- effektive group ID, 15
- effektive user ID, 15
- Empfangspuffer, 168
- end-of-record, 131
- endhostent(), 151
- endnetent(), 151
- endprotoent(), 151
- endservent(), 151
- Environment, 90
- ephemeral ports, 147
- Ethernet, 103
- exec, 13
- execl(), 13
- execle(), 13
- execlp(), 13
- execv(), 13
- execve, 13
- execvp(), 13, 81
- exit(), 4, 12, 17
- exit-Status, 17
- export, 73
  
- fcntl(), 136
- FDDI, 103
- fflush(), 16
- FIFO-Dateien, 97
- file locks, 15

- file mode creation mask, 12, 14
- fork(), 1, 9, 64, 65, 152
- FQDN, 143
- Fragmentierung, 110
- Fully Qualified Domain Name, 143
  
- GAN, 99
- Gateway, 102
- gethostbyaddr(), 142, 145
- gethostbyname(), 142
- gethostent(), 150
- gethostname(), 145
- getnetbyaddr(), 149
- getnetbyname(), 149
- getnetent(), 150
- getpeername(), 134, 136
- getpgrp(), 2
- getpid(), 1
- getppid(), 1
- getprotobyname(), 149
- getprotobynumber(), 149
- getprotoent(), 150
- getservbyname(), 148
- getservbyport(), 148
- getservent(), 150
- getsockname(), 134, 135
- getsockopt(), 136
- gettoken(), 74
- global area network, 99
  
- hangup-Signal, 2
- Hintergrund, 44
- Host, 141
- host interface, 103
- hostent, 142
- Hostname, 145
- Hostnamen, 143, 144
- htonl(), 140
- htons(), 140
- Hub, 102
  
- I/O-Umlenkung, 44, 72, 78
- IANA, 147, 150
- ICMP, 113, 123
- IGMP, 123
- INADDR\_ANY, 140
- INADDR\_NONE, 141
- inet\_addr(), 141
- inet\_aton(), 141
- inet\_ntoa(), 141
- inetd, 134
- init-Prozess, 17, 19
- Inter-Process Communication, 61
- Internet Adresse, 108
- Internet Assigned Numbers Authority, 147
- Internet Datagram, 110
- Internet Domain, 122
- Internet Protocol, 109, 123
- Internet Superserver, 134
- Internet-Adressen, 105, 141, 144
- ioctl(), 136
- IP, 109, 123
- IP next generation, 151
- IP-Adresse, 107, 145
- IP-Adresse 0, 140
- IPC, 4, 61, 62
- IPv4, 151
- IPv6, 151
- iterative server, 152, 153, 160
  
- Kernel Mode, 7
- kill, 2, 4
- kill(), 21
- kill-Kommando, 21
- Kommando-Pipeline, 73
- Kommando-Sequenz, 73
- Kommunikationsdomäne, 121
- Kommunikationsendpunkt, 121
- Kommunikationssemantik, 123
- Kommunikationssteuerungsschicht, 101
- Kontext, 1
- kurzlebige Portnummern, 140, 147
  
- LAN, 99
- lex, 73
- listen(), 128, 134
- Little-Endian, 140
- loader, 18
- local area network, 99
  
- MAX\_PATH, 138
- MAXHOSTNAMELEN, 146
- Maxi-Shell, 72
- maximum transfer unit, 112
- memory buffers, 137
- message queues, 98
- Midi-Shell, 44
- Mini-Shell, 38
- mknod(), 98
- MSG\_WAITALL, 169–171, 175, 181
- MTU, 112
- multi-homed host, 141
  
- named pipes, 97
- Nameserver, 144
- netdb.h, 147, 149

- Netzwerktopologie, 99
- network byte order, 140
- Network Information Service, 144, 150
- nice-Kommando, 6
- NIS, 144, 150
- nohup-Kommando, 6
- ntohl(), 140
- ntohs(), 140
  
- OSI, 100
- out-of-band-data, 130
  
- parent process, 1, 9
- parent process ID, 14
- PATH, 14
- pclose(), 71
- Peer-Adresse, 128
- PF\_INET, 124
- PF\_UNIX, 124
- PID, 1
- pipe(), 63, 65
- pipefd[0], 63
- pipefd[1], 63
- Pipeline, 79
- popen(), 71
- Port, 118
- Port-Nummer, 139
- Port-Nummern, 120
- Portnummer, 136, 147
- Portnummer 0, 140
- process group ID, 12, 14
- process group leader, 1, 37
- process ID, 14
- Protokolle, 99
- Prozess, 1, 18
- Prozessgruppe, 1
- Prozesstabelle, 8
- ps-Kommando, 4
- punktiertes Dezimalformat, 108, 141
  
- quit, 21
  
- Raw Socket, 123
- read(), 64, 130, 135
- readv(), 130, 135
- real group ID, 11, 15
- real user ID, 11, 14
- Rechneradresse, 136
- recv(), 130, 134, 135, 170
- recvfrom(), 131, 132, 135
- recvmsg(), 131, 132, 135
- recvn(), 170, 171
- regions, 7
  
- registered ports, 147
- Repeater, 102
- reservierte Portnummern (UNIX), 147
- Resolver, 142
- rlogin, 147
- root directory, 12, 14
- Router, 102
- routing, 151
- rsh, 147
  
- sa\_data, 137
- sa\_family, 136
- sa\_len, 137
- select(), 135
- Semaphore, 98
- send(), 130, 134, 135, 170
- Sendepuffer, 168
- sendmsg(), 131, 132, 134, 135
- sendn(), 169–171
- sendto(), 131, 132, 134, 135
- Server, 127
- set, 73
- set group ID, 15
- set user Id, 15
- sethostent(), 151
- setnetent(), 151
- setpgrp(), 2
- setprotoent(), 151
- setservent(), 151
- setsockopt(), 136
- shared memory, 98
- shutdown(), 132, 135
- Sicherungsschicht, 100
- SIG\_DFL, 22, 38
- SIG\_ERR, 22
- SIG\_IGN, 22, 38
- sigaction(), 25
- SIGALRM-Signal, 37
- SIGCHLD-Signal, 17, 31
- SIGCLD-Signal, 17, 31, 37
- SIGCONT-Signal, 38
- SIGFPE-Signal, 22
- SIGHUP-Signal, 17, 37
- SIGINT-Signal, 21
- SIGKILL-Signal, 22
- signal handler, 22, 155
- signal handling settings, 12
- signal(), 18, 21, 22, 24, 38
- Signalbehandlung, 95
- Signale, 15, 21
- Signalnummer, 21
- SIGPIPE-Signal, 22, 37
- SIGQUIT-Signal, 21
- SIGSEGV-Signal, 22

SIGSTOP-Signal, 22, 38  
SIGTTIN-Signal, 5  
SIGUSR1-Signal, 35, 38  
SIGUSR2-Signal, 35, 38  
sleep(), 37  
SO\_RCVBUF, 168  
SO\_SNDBUF, 168  
SOCK\_DGRAM, 124, 132  
SOCK\_RAW, 124  
SOCK\_STREAM, 124  
Socket Adresse, 125  
socket(), 124, 134  
Socket-Typen, 122  
socketpair(), 134  
Sockets, 121  
Stream Socket, 123, 168  
strtok(), 39  
struct in\_addr, 139  
struct netent, 149  
struct protoent, 149  
struct sockaddr, 125, 136, 137  
struct sockaddr\_in, 126, 139  
struct sockaddr\_un, 126, 138  
sun\_family, 138  
sun\_len, 138  
sun\_path, 138  
swapper, 19  
sys/param.h, 146  
sys/utsname.h, 146  
  
TCP, 124, 147, 164  
TCP/IP, 101, 105, 109  
Telnet-Server, 147  
telnetd, 147  
terminal group, 2  
terminal group ID, 12, 14  
TFTP-Server, 147  
tftpd, 147  
Token-Verfahren, 102  
top-Kommando, 4  
Transceiver, 103  
Transportschicht, 100, 147  
Trivial File Transfer Protocol, 147  
TTL, 113  
  
u area, 7, 9  
UDP, 147, 160  
uname(), 146  
uname-Kommando, 146  
UNIX domain, 122, 164  
unnamed pipes, 63  
Unreliable Delivery, 109  
User Mode, 6  
  
verbindungslos, 123, 132  
verbindungsorientiert, 123  
Vererbung, 3, 11, 14  
Vermittlungsschicht, 100  
voll-duplex, 132, 168  
  
wait(), 4, 16, 17  
WAN, 99  
well-known ports, 147  
well-known ports (UNIX), 147  
wide area network, 99  
write(), 64, 130, 135  
writev(), 130, 135  
  
X Window Server, 147  
  
Zombie, 12, 17, 18, 31, 33, 155