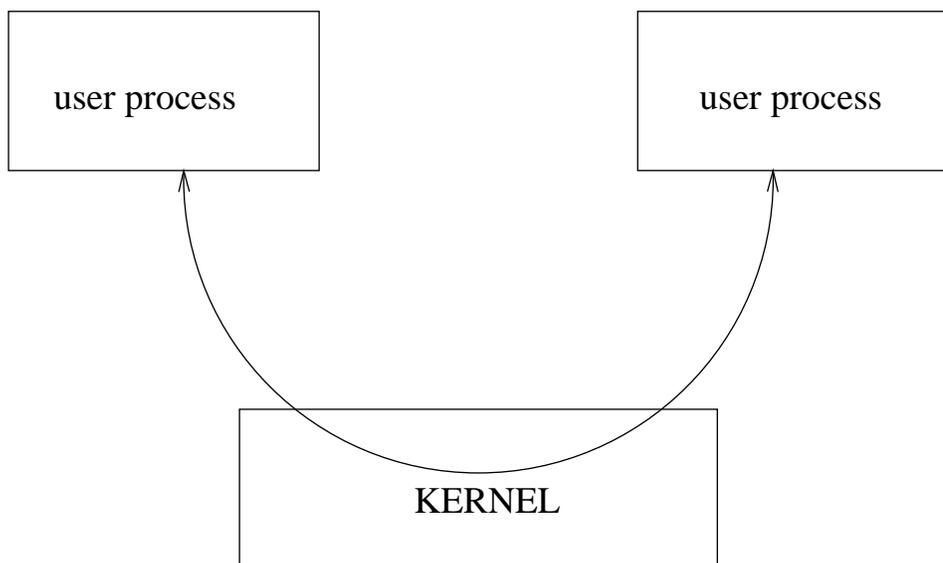


Kapitel 2

Inter-Prozess-Kommunikation (IPC)

2.1 Einführung

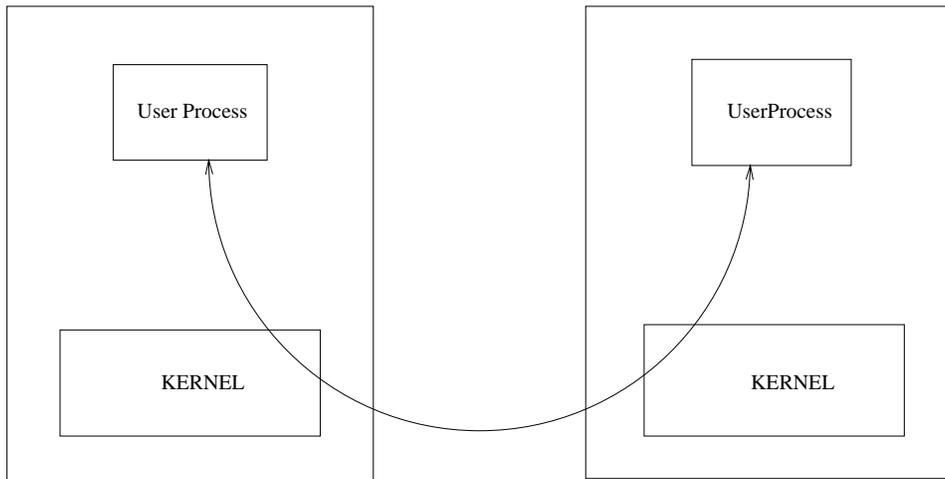
Jeder *UNIX*-Prozess besitzt seinen eigenen Context. Innerhalb eines Prozesses können die verschiedenen Moduln über Parameter und Rückgabewerte bei Funktionsaufrufen oder über globale Variablen Daten austauschen. Wollen jedoch zwei eigenständige Prozesse Daten miteinander austauschen, so kann dies nur 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 System laufen. Vor allem die Entwicklungsarbeiten der University of California at Berkeley brachte hier einige neue Ansätze zur Interprocess Communication in *UNIX* ein.



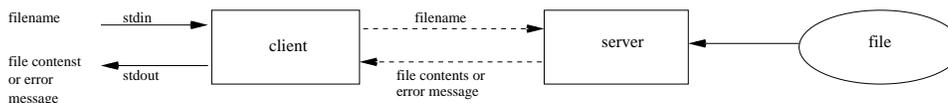
Für die Prozesse kann es völlig transparent sein, wo der jeweilige Partnerprozeß abläuft.

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()`

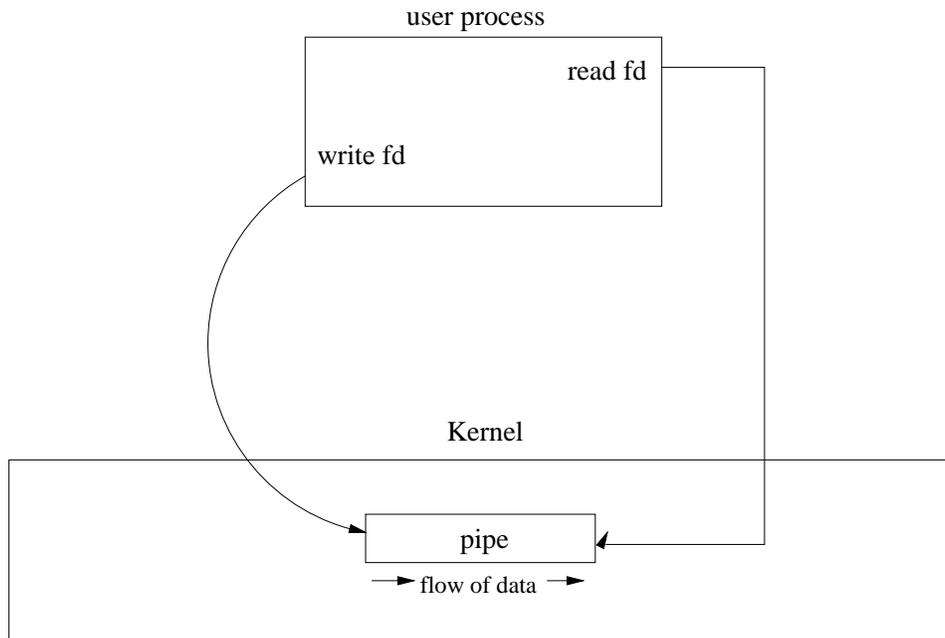
```
int pipe( int pipefd[2] ) /* create a pipe */

/* 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.



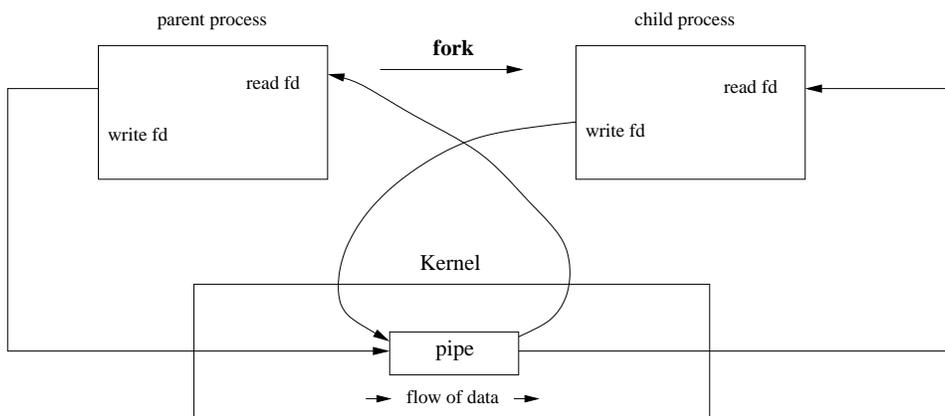
- Deadlock

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

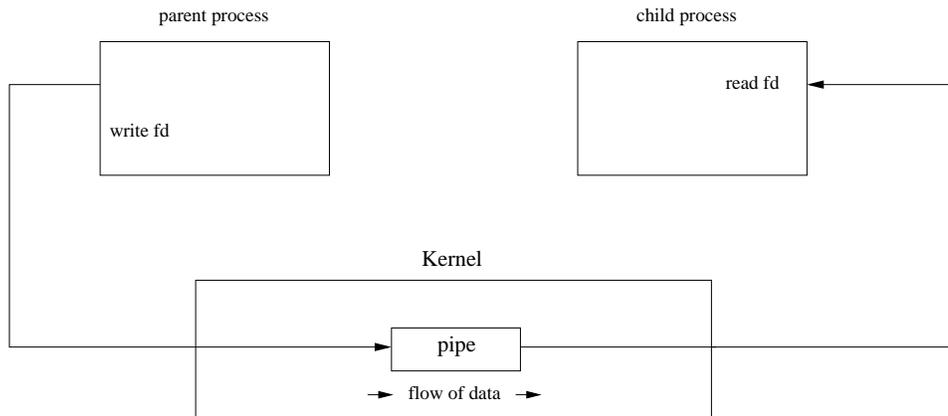
VORSICHT! Der Kernel synchronisiert Prozesse, die in Pipes schreiben oder aus Pipes lesen – *Producer / Consumer* Modell. Sollte hier in dem Ein-Prozess-Beispiel ein `read()` oder `write` System Call blockieren, entsteht ein *Deadlock*, denn der blockierte Prozess kann den befreienden, komplementären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

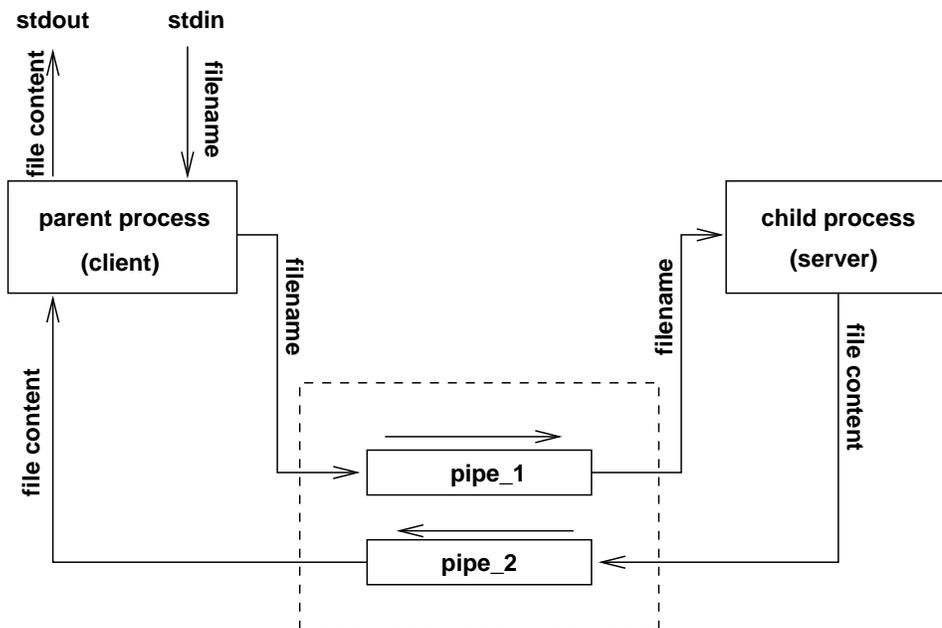
        do /* wait for child */ {
            if ( (gone = wait( (int *) 0 )) < 0 ) {
                /*no interest for exit_status*/
                perror("wait");
                exit(2);
            }
        } while ( gone != childpid );
    }
}
exit(0);
}

```

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 Maxi-Shell: Eine fast wirkliche Shell

Von der Mini-Shell über die Midi-Shell zur einer etwas umfassenderen Implementierung:

2.7.1 Merkmale

- **Zeilenende:** \n oder & (Hintergrundaufführung)
- **einfache Kommandos:**
 Kommandoname gefolgt optional von Argumenten, getrennt durch Leer-/Tab-Zeichen;
 ein Argument ist ein einzelnes Wort oder eine in **Doppelapostroph** eingeschlossene Zeichenfolge, in der Sonderzeichen wie | ; & > < oder Leer-/Tab-Zeichen oder *newline* enthalten sein dürfen.
 Ein Doppelapostroph darin muß ebenso wie \ durch \ (also \" resp. \\) geschützt sein.
 Anzahl der Argumente auf 20, Länge eines Arguments auf 200 beschränkt.
- **einfache Kommandos mit I/O-Umlenkung**
 < > >> → i.P. wie bei der "richtigen" Shell

- **Kommando-Pipeline()**
- **Einschränkung:**
`cmd1 args | cmd2 args`
 In `cmd1` darf keine Ausgabeumlenkung, in `cmd2` keine Eingabeumlenkung erfolgen!
- Während der Ausführung von Hintergrundkommandos werden *Interrupt*-Signale ignoriert.
- **Kommando-Sequenz:** durch Semikolon getrennte Folge einfacher Kommandos / einfacher Kommandos mit I/O-Umlenkung / Kommando-Pipeline, die nacheinander ausgeführt werden
- **eingebaute Kommando's (builtins):** `cd`, `set`, `export`

NICHT: Dateinamen-Substitution, Kommando-Substitution, History, ...

2.7.2 Analyse der Kommandozeile

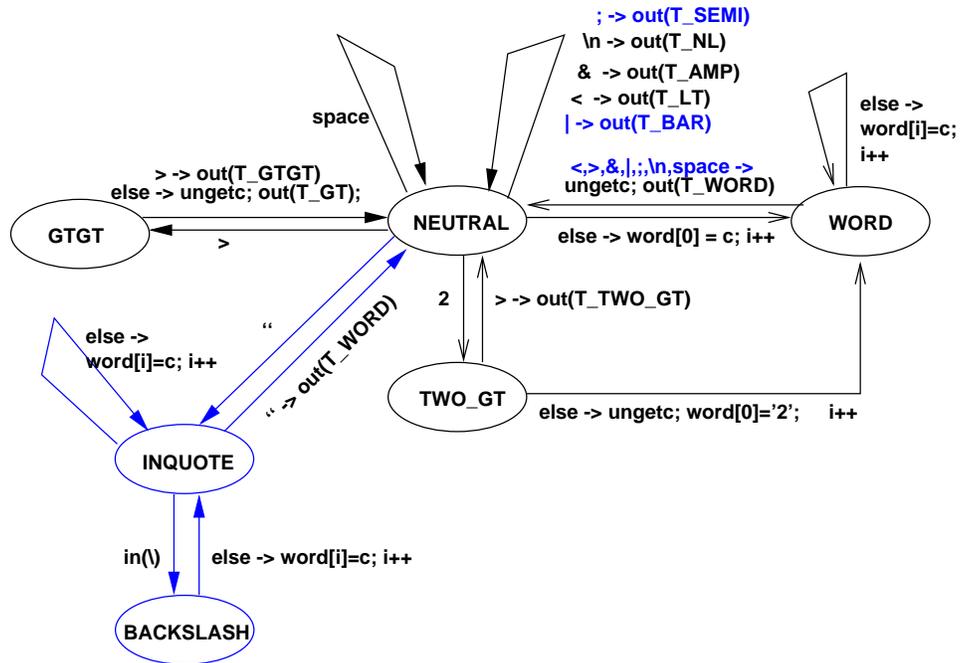
Hierzu sind nur geringfügige Ergänzungen zur Midi-Shell nötig, werden hier aber der Vollständigkeit wegen alle noch einmal aufgeführt:

Zerlegung in Symbole (*Token*), die syntaktische Einheit bilden

Symb. Konst.	Erläuterung
T_WORD	Argument oder Dateiname; falls quotiert, werden Apostrophen nach Erkennen des Symbols entfernt
T_BAR	
T_AMP	&
T_SEMI	;
T_GT	>
T_GTGT	>>
T_TWO_GT	2 >
T_LT	<
T_NL	<i>newline</i>
T_EOF	Dateiende oder bei stdin EOT (<i>end of transmission, ctrl-d</i>)

Anm.: Zur lexikalischen Analyse gibt es in UNIX das "Standardwerkzeug" `lex`, das hier wegen der geringen Komplexität der Aufgabe aber nicht verwendet wird.

Endlicher Automat zur TOKEN-Bestimmung



Die Datei *defs.h* enthält für das Folgende einige nützliche Vereinbarungen:

```
/* ----- defs.h ----- */
#include <stdio.h>
#include <unistd.h>

typedef enum {FALSE, TRUE} BOOLEAN;

typedef enum {T_WORD, T_BAR, T_AMP,
T_SEMI, T_GT, T_GTGT, T_TWO_GT,
T_LT, T_NL, T_EOF} TOKEN;

#define BADFD -2
#define MAXARG 20
#define MAXWORD 200
#define MAXFNAME 14
```

Anm.: BADFD=-2: -2 ist kein Returnwert von open()!

Realisierung des Automaten durch die Funktion `gettoken()`:

```
/* ----- gettoken.c ----- */
#define GET_H
#include "gettoken.h"

/* lexikalische Analyse der Kommandozeile */

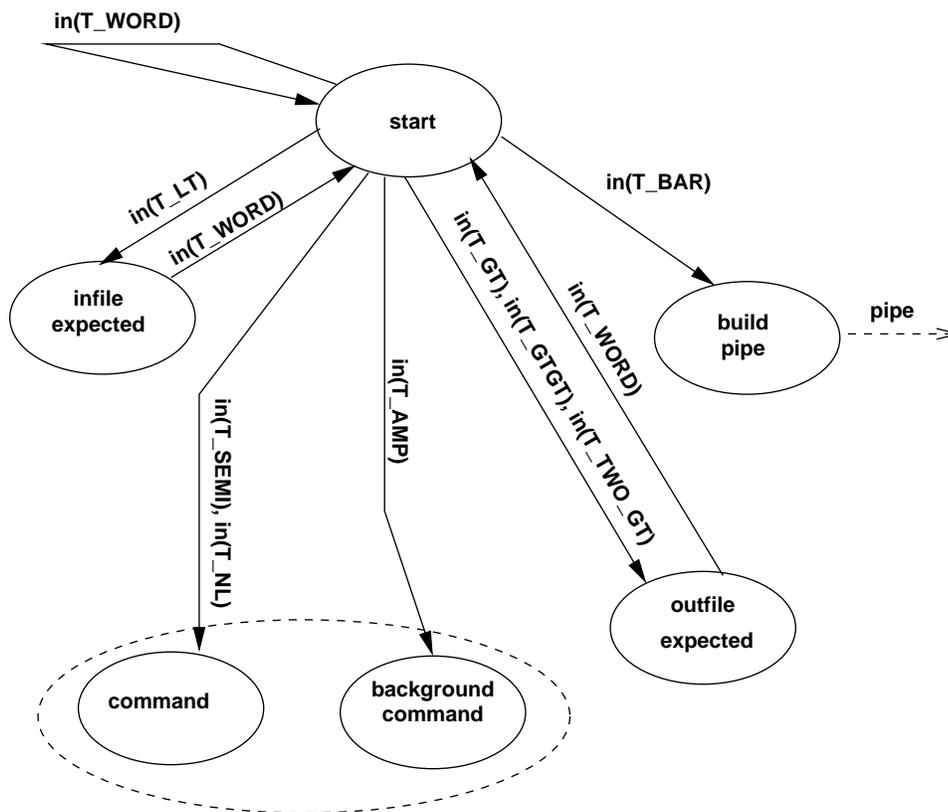
TOKEN gettoken( char * word) {
    int c;
    char * w;
```



```
        *w = '\\0';
        return (T_WORD);
    default:
        * w++ = c;
        continue;
    }
case INWORD:
    switch (c) {
        case ';' :
        case '&' :
        case '|' :
        case '<' :
        case '>' :
        case '\\n' :
        case ' ' :
        case '\\t' :
            ungetc(c,stdin);
            * w = '\\0';
            return (T_WORD);
        default:
            * w++ =c;
            continue;
    }
}
}
return (T_EOF);
}
```

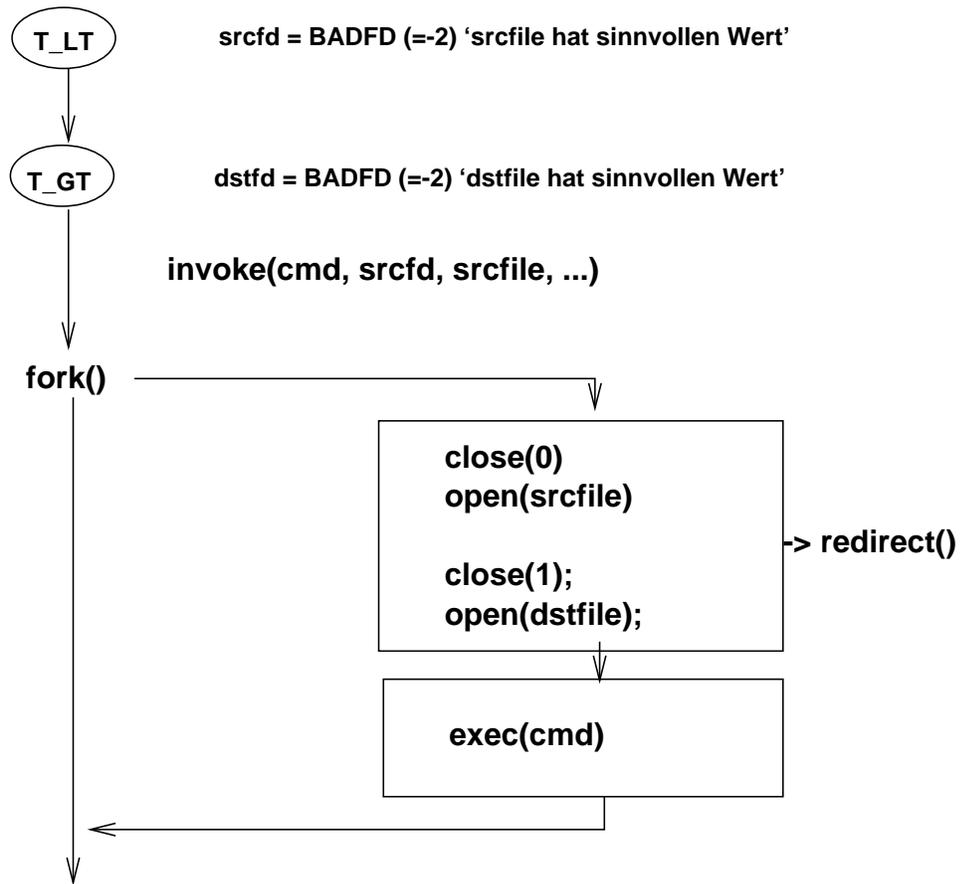
2.7.3 Bearbeiten und Ausführen von Kommandos

- einfaches Kommando aufbauen



- einfaches Kommando mit I/O-Umlenkung:

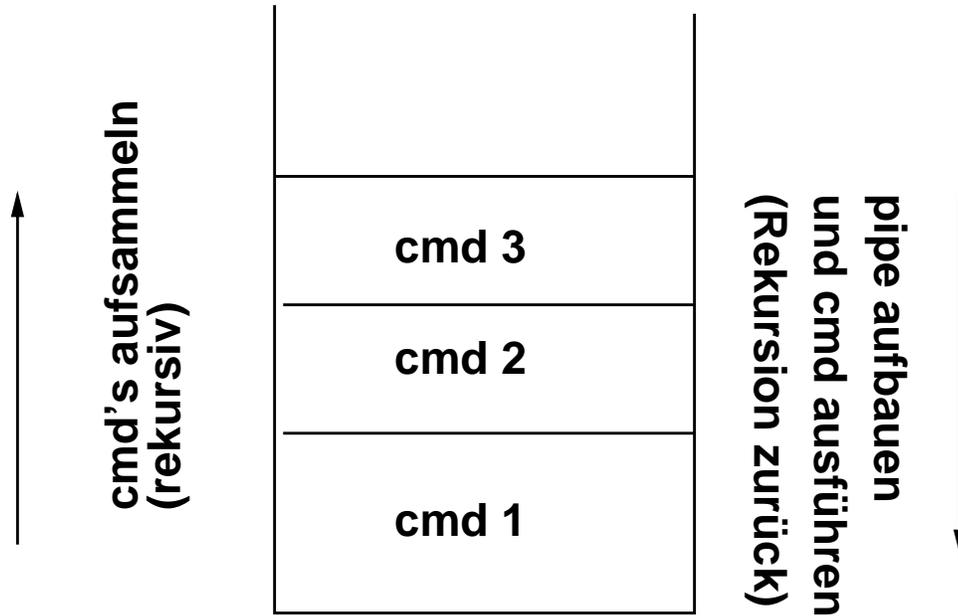
```
$ <infile cmd args > outfile
```



- Pipeline aufbauen (von vorne oder von hinten her ??)
hier: von 'hinten' her! → Funktion `command`

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

Stack via Rekursion



term=command(&pid,makepipe=FALSE,pipefdp=NULL)

reads cmd1, terminated by T_BAR

term=command(waitpid,makepipe=TRUE,&dstfd)

reads cmd2, terminated by T_BAR

term=command(waitpid,makepipe=TRUE,&dstfd)

reads cmd3, terminated by T_NL;

term = T_NL;

pipe()

returns write_end via dstfd to calling function

fork()

parent:

returns pid of child to
command argument

child:

close(0)
dup(reading end of pipe)
exec(cmd3)

back to calling function

got fd for write end of created pipe

pipe()

returns write_end via dstfd to calling function

fork()

the same as above

close(0); close(1)

dup(read_end of last created pipe)

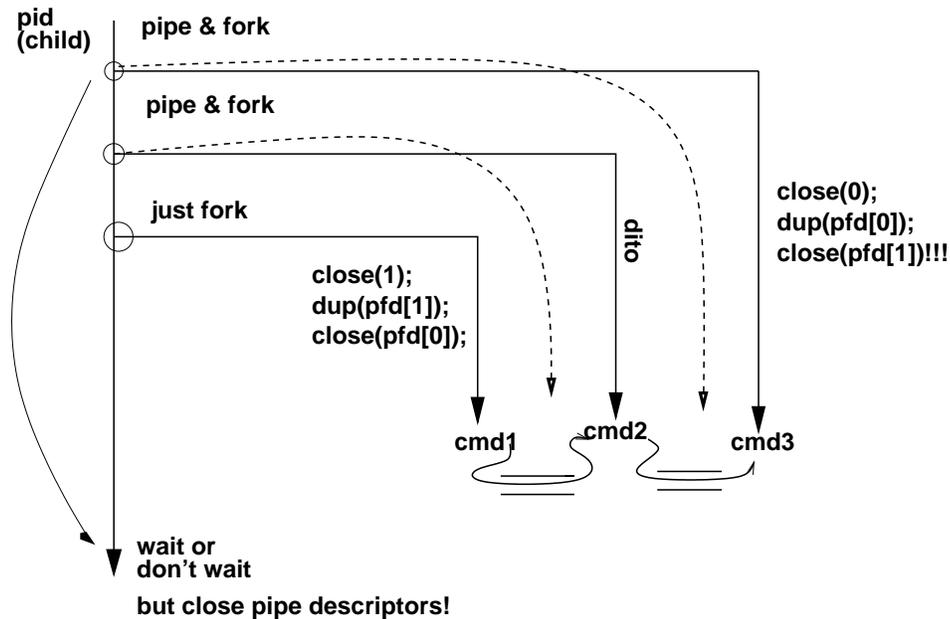
dup(fd - see above)

exec(cmd2)

got fd for write end of created pipe → fork() →

close(1); dup(fd); exec(cmd1)

vereinfacht:



- **command()**

bearbeitet ein einfaches Kommando, das durch | oder & oder ; oder *newline* abgeschlossen ist. Argumente werden im Vektor `argv` zur späteren Benutzung für `execvp()` hinterlegt.

Bestimmung der E/A-Kanäle:

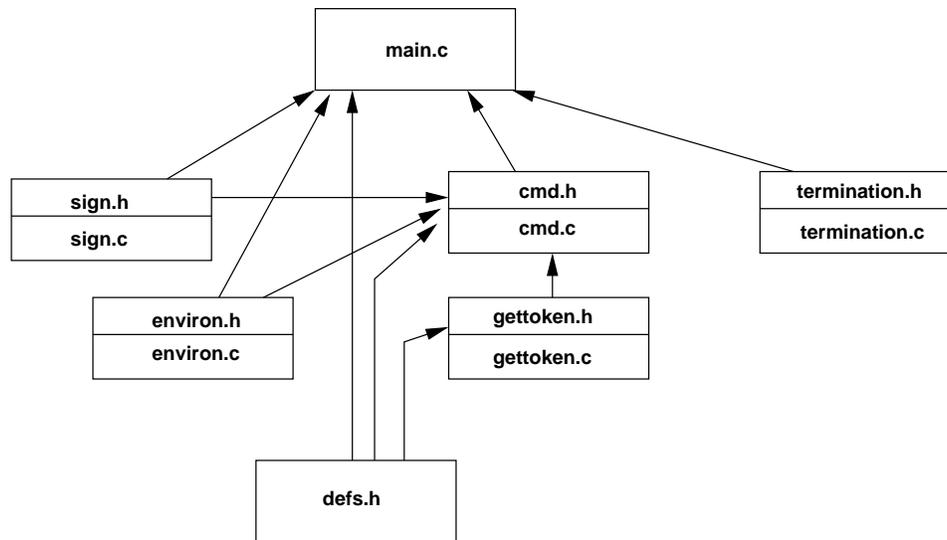
- Eingabe
 - * Voreinstellung, Filedeskriptor 0
 - * Datei, falls "<" erkannt (`gettoken()`)
 - * lesbares Ende einer Pipe (vor dem Kommando stand |).
- Ausgabe
 - * Voreinstellung (Filedeskriptor 1)
 - * zu erzeugende / verkürzende Datei (Symbol >).
 - * zu erzeugende / verlängernde Datei (Symbol >>).
 - * schreibbares Ende einer Pipe (nach dem Kommando steht |).

Mittels folgender Variablen wird die jeweilige Situation festgehalten:

- **srcfd** (source file descriptor): falls Quelle eine Datei, so wird in **srcfile** der Dateiname festgehalten
- **dstfd** (destination file descriptor): falls Ziel eine Datei, so wird deren Name in **dstfile** festgehalten.
- **APPEND** anfügen?
- **makepipe**
Aufforderung für **command()**, eine Pipe zu erzeugen, deren lesbares Ende als Standardeingabe dient.

- **invoke()**
Ausführen des Kommandos ggf. im Hintergrund ("fork/exec")
- **builtin()**
Kommando als "builtin" ausführen
- **redirect()**
Ein-/Ausgabe umlenken, Quelle (Ziel) mit Filedeskriptor 0 (1)

2.7.4 Gesamtstruktur



```

# eine kleine Shell
mysh:    main.o cmd.o sign.o termination.o environ.o gettoken.o
gcc -Wall -o mysh main.o cmd.o sign.o termination.o envi-
ron.o gettoken.o
# executable: mysh
gcc -Wall -o sleepwell sleepwell.c
# long runner: sleepwell
gcc -Wall -o read_something read_something.c
# sleep and read from stdin: read_something
gcc -Wall -o write_something write_something.c
# sleep and from write to stdout: write_something

main.o:   main.c defs.h sign.h cmd.h environ.h termination.h
gcc -Wall -c main.c

cmd.o:    cmd.c cmd.h gettoken.h environ.h defs.h sign.h
gcc -Wall -c cmd.c

gettoken.o:  gettoken.c gettoken.h defs.h
gcc -Wall -c gettoken.c

sign.o:    sign.c sign.h

```

```

gcc -Wall -c sign.c

termination.o:    termination.c termination.h
gcc -Wall -c termination.c

environ.o:    environ.c environ.h defs.h
gcc -Wall -c environ.c

clean:
    rm -f *.o core
realclean:
    rm -f mysh *.o core sleepwell read_something write_something

```

2.7.5 Kommando-Berarbeitung

cmd.h:

```

/* ----- cmd.h -----*/

#ifdef CMD_H
#include <stdio.h>
#include <strings.h>
#include <fcntl.h>
#include <errno.h>
#include "defs.h"
#include "sign.h"
#include "environ.h"
#include "gettoken.h"
#else
extern BOOLEAN builtin(int argc, char *argv[],
                      int srcfd, int dstfd);
/* execute built-in-commands */

extern void redirect(int srcfd, char *srcfile,
                   int dstfd, char *dstfile,
                   int errfd, char *errfile,
                   BOOLEAN append, BOOLEAN bckgrnd);
/* Redirection of I/O */

extern int invoke(int argc, char *argv[],
                int srcfd, char *srcfile, int dstfd,
                char *dstfile, int errfd, char * errfile,
                BOOLEAN append, BOOLEAN bckgrnd);
/* invoke() - execute simple command */

extern TOKEN command(int *waitpid,
                   BOOLEAN makepipe,
                   int *pipefdp);
/* collect a simple command */
#endif

```

cmd.c:

```

/* ----- cmd.c -----*/

#define CMD_H
# include "cmd.h"
# include <stdlib.h>
# include <string.h>

static BOOLEAN builtin(int argc, char * argv[],
                      int srcfd, int dstfd) {
/*recognize and execute built-ins*/

    char * path;

    if (strchr(argv[0], '=') != NULL)
        /* strchr: locate character in string
         * strchr returns pointer to "=" or NULL
         *
         * assignment ?
         * example: x=3
         * argc must be 1, argv[0] contains "x=3", e.g.
         */

        asg(argc,argv);

    else if (strcmp(argv[0],"export") == 0)
        /* export of a variable?
         * example: export x or export x y x or
         * just export (list environment)
         */
        export(argc,argv);

    else if (strcmp(argv[0], "set") == 0)
        /* print environment */
        set(argc, argv);

    else if (strcmp(argv[0], "cd") == 0) {
        /* change directory */
        if (argc > 1)
            path = argv[1];
        else
            if ( (path = EVget("HOME")) == NULL)
                path = "."; /*HOME not defined*/
            /*else: path has got value from HOME*/
            if (chdir(path) == -1)
                fprintf(stderr, "%s: bad dir\n", path);
    }
    else
        return (FALSE);
    if (srcfd != 0 || dstfd != 1)
        /*we allow only from stdin and to stdout*/
        fprintf(stderr, "Illegal redir. or pipeline\n");
}

```

```

    return (TRUE);
}

/* Redirection of I/O: */
/* after redirect the caller has file descriptor 0 for input,
 * fd 1 for output, fd 2 for errors!
 */
static BOOLEAN redirect(int srcfd, char * srcfile,
                       int dstfd, char * dstfile,
                       int errfd, char * errfile,
                       BOOLEAN append, BOOLEAN bckgrnd) {
    int flags, fd;

    /* we expect for srcfd:
     * 0: nothing to do
     * -2: redirect 0 to file
     * >0: redirect 0 to pipe
     *
     * we expect for dstfd:
     * 1: nothing to do
     * -2: redirect 1 to file
     *      (with respect to parameter 'append'
     * >1: redirect 1 to pipe
     *
     * we expect for errfd:
     * 2: nothing to do
     * -2: redirect to file
     */

    if (srcfd == 0 && bckgrnd) {
        strcpy(srcfile, "/dev/null");
        /* /dev/null ->
         * there is nothing to read, only EOF
         * a background command couldn't get any
         * input from stdin;
         */
        srcfd = BADFD;
        /* so redirect 0 to srcfile
         * set to /dev/null above
         */
    }

    if (srcfd != 0) {
        /* 0 should point to file or
         * pipe for input
         */
        if (close(0) == -1)
            perror("close");
        if (srcfd > 0) {
            /* it's a pipe */
            if (dup(srcfd) != 0) {
                perror("dup");
                return FALSE;
            }
        }
    }
}

```

```

    }
  }
  else if (open(srcfile, O_RDONLY, 0) == -1) {
    perror("open");
    return FALSE;
  }
}

/*now file or pipe is referenced
 *by file descriptor 0
 */

/* stderr: */
if(errfd != 2) {
  if( close(2) == -1)
    perror("close");
  if(open(errfile, O_WRONLY | O_CREAT | O_TRUNC, 0664) == -
1) {
    perror("open");
    return FALSE;
  }
}

/* now the same for output */
if (dstfd != 1) {
  /* output to pipe (dstfd>1) or
   * file (>,>> (dstfd==2))
   */
  if (close(1) == -1)
    perror("close");
  if (dstfd > 1) {
    /* to pipe */
    if (dup(dstfd) != 1) {
      perror("dup");
      return FALSE;
    }
  }
}
else {
  /* to file */
  flags = O_WRONLY | O_CREAT;
  if (!append) /* > file */
    flags |= O_TRUNC;
  else
    flags |= O_APPEND;

  if (open(dstfile, flags, 0666) == -1) {
    /* open returns the smallest
     * free file descriptor
     */
    fprintf(stderr, "can't create %s\n", dstfile);
    return FALSE;
  }
}
}
}

```



```

        exit(3);

        execvp(argv[0], argv);
        /* this shouldn't be reached */
        fprintf(stderr, "can't execute %s\n", argv[0]);
        exit(4);
default:
    /* PARENT */
    if (srcfd > 0 && close(srcfd) == -1)
        perror("close src");
    if (dstfd > 1 && close(dstfd) == -1)
        perror("close dst");
    if (bckgrnd)
        printf("%d\n", pid);
    return (pid);
}
}

/* recursive collection of the command line and
 * recursive execution of all of the simple commands
 * 'from back' in the pipeline
 */

TOKEN command(int * waitpid, BOOLEAN makepipe,
              int * pipefdp) {
    /* int * waitpid: perhaps we have to wait
     * for the last command in the pipeline
     *
     * makepipe==TRUE: we are in a pipeline and
     * have to return the write-pipefdp to the caller
     *
     * return values: T_NL or T_SEMI or T_BAR (!)
     *
     * uses: gettoken(), invoke()
     */

    TOKEN token, term;
    int argc, srcfd, dstfd, errfd, pid, pfd[2];
    char * argv[MAXARG+1];
    char srcfile[MAXFNAME+1];
    char dstfile[MAXFNAME+1];
    char errfile[MAXFNAME+1];
    char word[MAXWORD], *malloc();
    BOOLEAN append;

    argc = 0; srcfd = 0; dstfd = 1; errfd = 2;
    /* defaults */

    while (1) {
        switch (token = gettoken(word)) {
            case T_WORD:
                if (argc == MAXARG) {
                    fprintf(stderr, "Too many args\n");
                    break;
                }

```

```

    }
    if ((argv[argc]=malloc(strlen(word)+1))==NULL) {
        fprintf(stderr,"Out of arg memory\n");
        break;
    }
    strcpy(argv[argc],word);
    argc++;
    continue;
case T_LT:
    /* more than 1 '>' is here allowed (???)
     * we just take the last one
     */

    /* after a T_BAR an in/out redirection
     * isn't allowed:
     */
    if (makepipe) {
        fprintf(stderr, "Extra <\n");
        break;
    }

    /* after T_LT we expect an input-file-name */
    if (gettoken(srcfile) != T_WORD) {
        fprintf(stderr, "Illegal <\n");
        break;
    }
    /* we have to redirect 0 to a file */
    srcfd = BADFD;
    continue;
case T_GT:
case T_GTGT:
    if (dstfd != 1) {
        fprintf(stderr, "EXTRA > or >>\n");
        break;
    }
    if (gettoken(dstfile) != T_WORD) {
        fprintf(stderr, "Illegal > or >>\n");
        break;
    }
    dstfd = BADFD;
    append = (token == T_GTGT);
    continue;
case T_TWO_GT:
    if (errfd !=2) {
        fprintf(stderr, "EXTRA 2>\n");
        break;
    };
    if (gettoken(errfile) != T_WORD) {
        fprintf(stderr, "Illegal 2>\n");
        break;
    };
    errfd = BADFD;
    continue;
case T_BAR:

```

```

case T_AMP:
case T_SEMI:
case T_NL:
    /* one simple command is read */
    argv[argc] = NULL;
    if (token == T_BAR) {
        if (dstfd != 1) {
            fprintf(stderr, "> or >> conflicts with |\n");
            break;
        }

        /* RECURSION */
        term = command(waitpid, TRUE, &dstfd);
    }
    else
        /* AFTER RECURSION or just one simple command */

        term = token;
    /* T_BAR for all cmd's except the last
    * in the pipeline
    */

    if (makepipe) {
        if (pipe(pfd) == -1)
            perror("pipe");
        * pipefdp = pfd[1];
        srcfd = pfd[0];
    }
    pid = invoke(argc,argv,srcfd, srcfile,dstfd,
                dstfile, errfd, errfile, append,
                term == T_AMP);
    if (token != T_BAR)
        /* last command in pipeline */
        * waitpid = pid;
    if (argc == 0 && (token != T_NL || srcfd > 1))
        fprintf(stderr, "Missing command\n");
    while (--argc >= 0)
        free(argv[argc]);
    return (term);
case T_EOF:
    exit(0);
}
}
}

```

2.7.6 Behandlung des Einvironments

environ.h

```
/* ----- environ.h ----- */
```

```

#ifndef ENVIRON_H
#define EXTERN extern
EXTERN BOOLEAN EVset(char *name, char * val);
    /* add to environment*/

EXTERN BOOLEAN EVexport(char *name);
    /* set variable to be exported*/

EXTERN char * EVget(char *name);
    /* get value */

EXTERN BOOLEAN EVinit();
    /* init. symtable from environment*/

EXTERN BOOLEAN EVupdate();
    /*build environment from symtable*/

EXTERN void EVprint();
    /* print environment*/

EXTERN void asg(int argc, char *argv[]);
    /*assign*/

EXTERN void export(int argc, char *argv[]);
    /*export command*/

EXTERN void set(int argc, char *argv[]);
    /*set command*/

#else
#include "defs.h"
#endif

```

environ.c

```

/* ----- environ.c ----- */
#define ENVIRON_H

#include <string.h>
#include <stdio.h>
#include <stdlib.h>
#include "environ.h"

/* table of symbols: how much and how long
 * is one (var=value) ???
 */
#define MAXVAR 100
#define VARLENGTH 200

static struct varslot {
    char * name;    /*name of variable */
    char * val;    /*value*/
    BOOLEAN exported; /* exported ?*/

```

```

} sym[MAXVAR];

char * malloc(), * realloc();

extern char ** environ;

/*initialize*/
static BOOLEAN assign(char **p, char *s) {
    int size;
    size = strlen(s) + 1;
    if (*p == NULL) {
        if ((*p = malloc(size)) == NULL)
            return(FALSE);
    }
    else if ((*p = realloc(*p,size)) == NULL)
        return(FALSE);
    strcpy(*p, s);
    return(TRUE);
}

/*find slot: */
static struct varslot * find(char * name) {
    int i;
    struct varslot * v;

    v = NULL;
    for ( i=0; i < MAXVAR; i++)
        if (sym[i].name == NULL) {
            if (v == NULL)
                v = &sym[i];
        }
        else if (strcmp(sym[i].name, name) == 0) {
            v = &sym[i];
            break;
        }
    return(v);
}

/* add to environment:*/
BOOLEAN EVset(char * name, char * val) {
    struct varslot * v;

    if ( (v = find(name)) == NULL)
        return FALSE;
    return(assign(&v->name,name) && assign(&v->val,val));
}

/* set variable to be exported*/
BOOLEAN EVexport(char * name) {
    struct varslot * v;

    if ((v = find(name)) == NULL)
        return(FALSE);
    if (v->name == NULL)

```

```

        if (!assign(&v->name,name) || !assign(&v->val,""))
            return(FALSE);
        v->exported = TRUE;
        return(TRUE);
    }

/* get value */
char * EVget(char * name)    {
    struct varslot * v;

    if ((v=find(name))==NULL || v->name==NULL)
        return(NULL);
    return(v->val);
}

/* init. symtable from environment*/
BOOLEAN EVinit()    {
    int i, namelen;
    char name[VARLENGTH];

    for(i=0; (environ[i] != NULL) && (i < MAXVAR) ; i++) {
        namelen = strcspn(environ[i],"=");
        strncpy(name, environ[i],namelen);
        name[namelen] = '\0';
        if (!EVset(name,&environ[i][namelen+1]) ||
            !EVexport(name))
            return(FALSE);
    }
    return(TRUE);
}

/*build environment from symtable*/
BOOLEAN EVupdate()    {
    int i, envi, nvlen;
    struct varslot *v;
    static BOOLEAN updated = FALSE;

    if (!updated)
        if ((environ=(char **) malloc((MAXVAR+1)*
            sizeof(char *)))==NULL)
            return(FALSE);
    envi = 0;
    for (i=0; i < MAXVAR; i++) {
        v = &sym[i];
        if (v->name == NULL || !v->exported)
            continue;
        nvlen = strlen(v->name) + strlen(v->val) + 2;
        if (!updated) {
            if ((environ[envi] = malloc(nvlen)) == NULL)
                return(FALSE);
        }
        else
            if ((environ[envi]=realloc(environ[envi],nvlen))==NULL)
                return(FALSE);
    }
}

```

```

        sprintf(environ[envi], "%s=%s", v->name, v->val);
        envi++;
    }
    environ[envi] = NULL;
    updated = TRUE;
    return(TRUE);
}

/* print environment*/
void EVprint()    {
    int i;

    for (i=0; i<MAXVAR; i++)
        if (sym[i].name != NULL)
            printf("%3s %s=%s\n", sym[i].exported ? "[E]" : "",
                sym[i].name, sym[i].val);
}

/*assign*/
void asg(int argc, char *argv[])    {
    char * name, * val, *strtok();

    if (argc != 1)
        fprintf(stderr, "Extra args\n");
    else {
        name = strtok(argv[0], "=");
        val = strtok(NULL, "\\1");
        /*take the rest - \\1 doesn't occur*/
        if (!EVset(name, val))
            fprintf(stderr, "can't assign\n");
    }
}

/*set command*/
void set(int argc, char *argv[])    {
    if (argc != 1)
        fprintf(stderr, "Extra args\n");
    else
        EVprint();
}

/*export command*/
void export(int argc, char *argv[])    {
    int i;

    if (argc == 1) {
        set(argc, argv);
        return;
    }

    for (i=1; i < argc; i++)
        if (!EVexport(argv[i])) {
            fprintf(stderr, "can't export %s\n", argv[i]);
        }
}

```

```

        return;
    }
}

```

2.7.7 Signalbehandlung

sign.h

```

/* ----- 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

```

/* ----- 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;
    newact.sa_handler = shell_handler;
    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;
    }
}

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;
}

```

2.7.8 main()

```

/*main.c : eine fast wirkliche Shell*/
#include <stdio.h>
#include <sys/types.h>
#include <sys/wait.h>
#include "defs.h"
#include "sign.h"
#include "cmd.h"
#include "environ.h"
#include "termination.h"

int main() {
    char * prompt;
    int pid, fd, status;
    TOKEN term;

    if ( (ignoresig(SIGINT) == SIG_ERR) ||
         (ignoresig(SIGCHLD) == SIG_ERR)
        ) {
        perror("signal");
        exit(1);
    };

    if (!EVinit()) {
        fprintf(stderr, "can't init environment\n");
        exit(2);
    }
    if ( (prompt = EVget("PS2")) == NULL )
        prompt = "mysh> ";

    printf("%s", prompt);

    while (1) {
        term = command(&pid, FALSE, NULL);

        if (term != T_AMP && pid != 0) {
            waitpid(pid, &status, 0);
        }
        if (term == T_NL)
            printf("%s", prompt);
        for (fd = 3; fd < 20; fd++)
            (void) close(fd);
    }
    exit(0);
}

```

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*

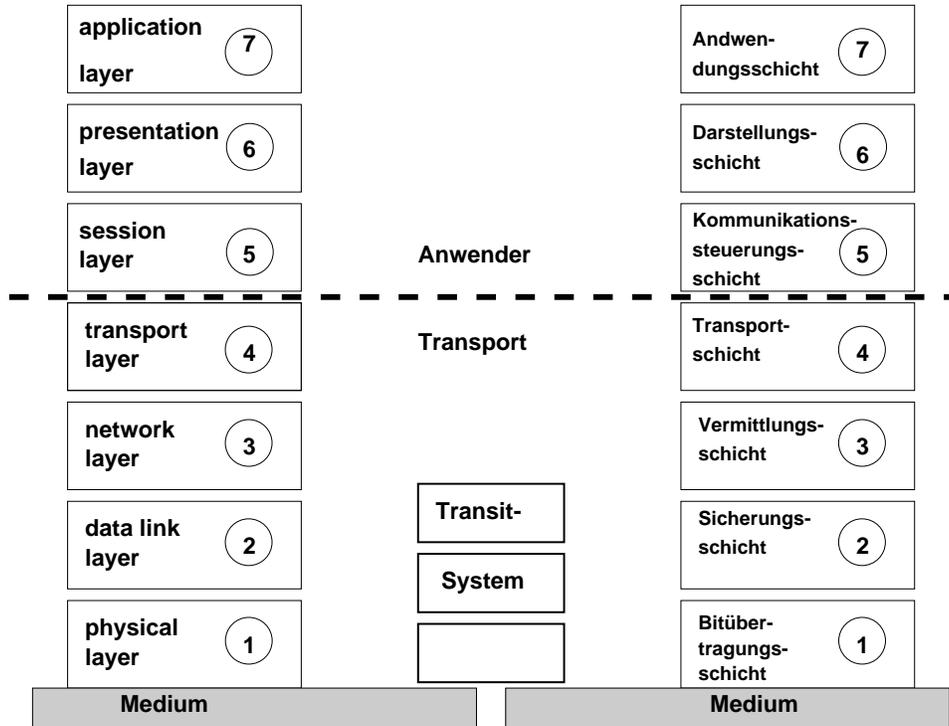


Abbildung 3.1: ISO-OSI-Referenzmodell

- 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**
Bereitstellung von Sprachmitteln zur Steuerung der Kommunikations-
beziehung (*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

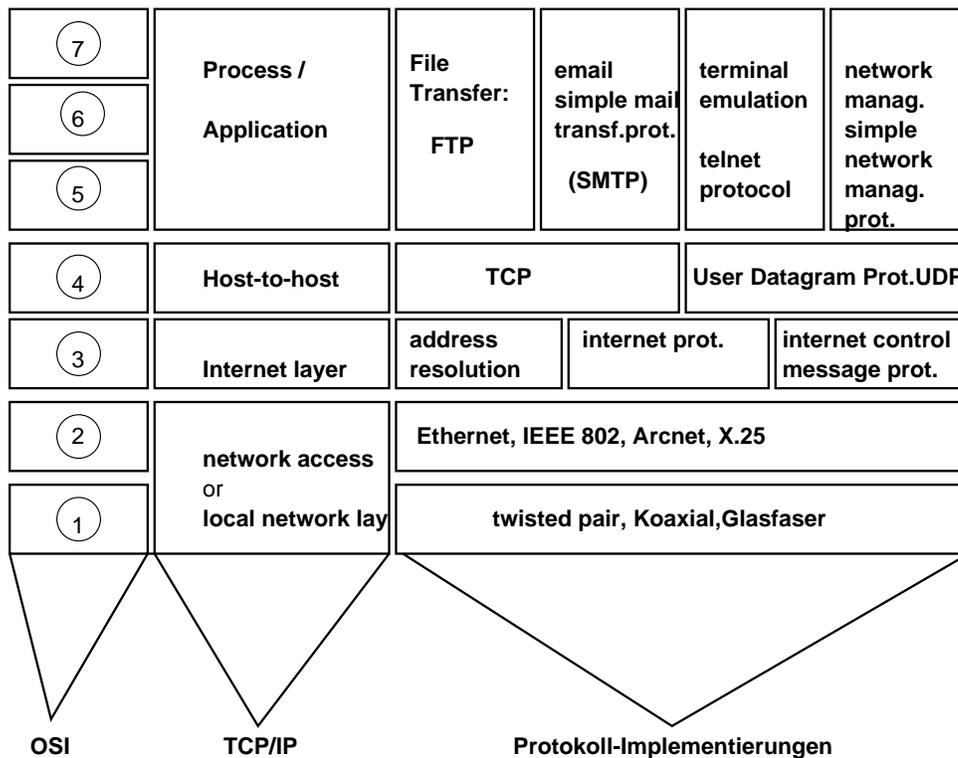


Abbildung 3.2: TCP/IP

Kopplung von Netzen

- Kopplungseinheiten zur Verbindung von Netzen (*internetworking units*)
Adreßumwandlung, Wegewahl (*routing*), Flußkontrolle, Fragmentierung und Wiedierzusammenfügung 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 10Mbits/s

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 (1-4Mbits/s)
- Koaxialkabel (4-40Mbits/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 Mbits/s), meist als *Backbone* eingesetzt

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:

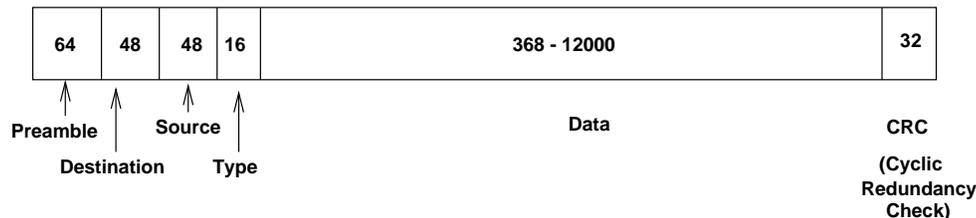


Abbildung 3.3: Ethernet Frame Format

- 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

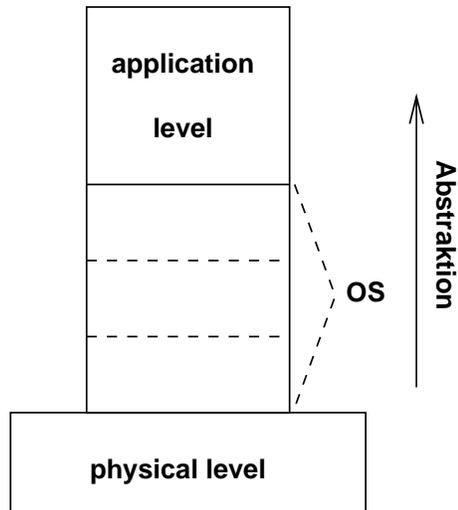


Abbildung 3.4: Layered System Architecture

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

□ *The TCP/IP internet protocols treat all networks equally (independent 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

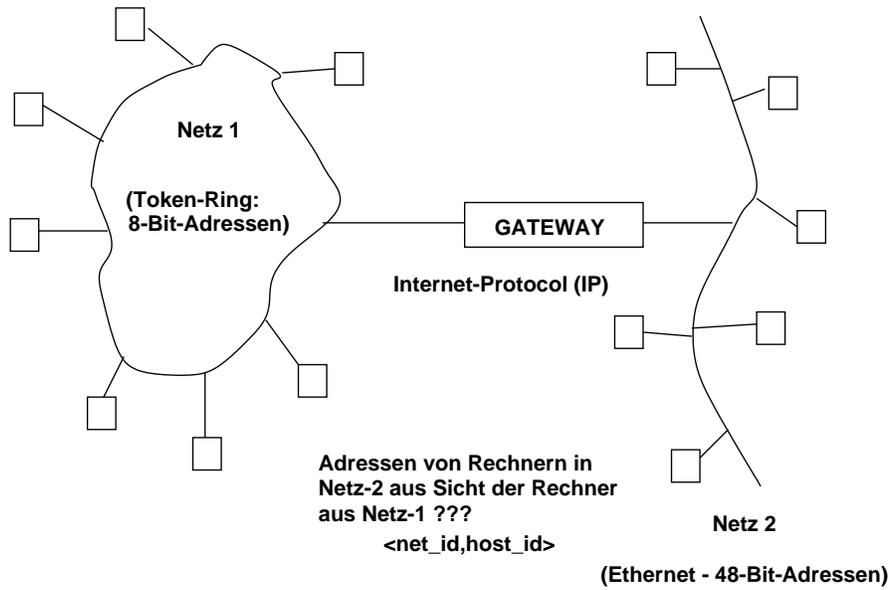


Abbildung 3.5: Internet Architecture

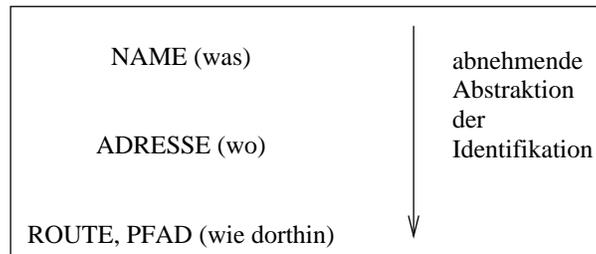


Abbildung 3.6: Abstraktion der Identifikation

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

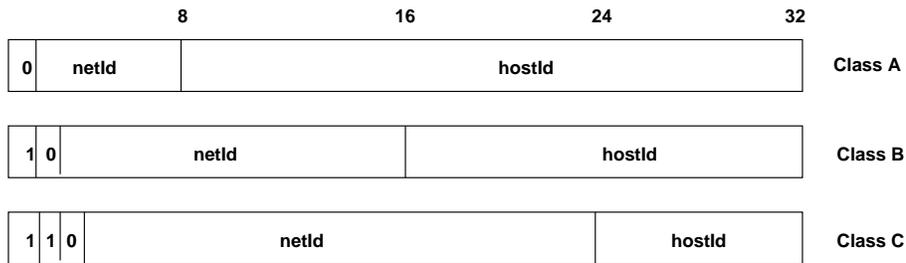


Abbildung 3.7: 32-Bit IP-Adresse

- 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)

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)

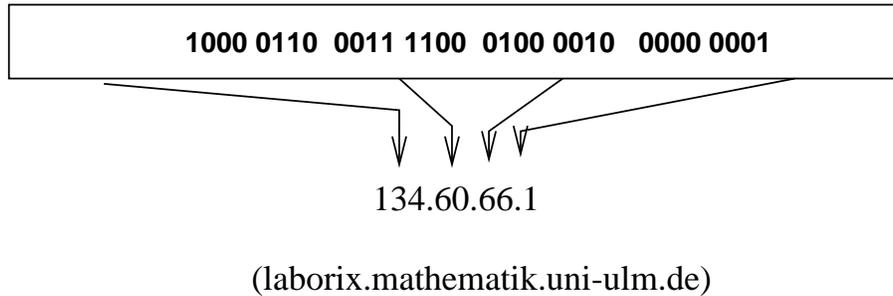


Abbildung 3.8: IP-Adresse: dotted decimal form

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

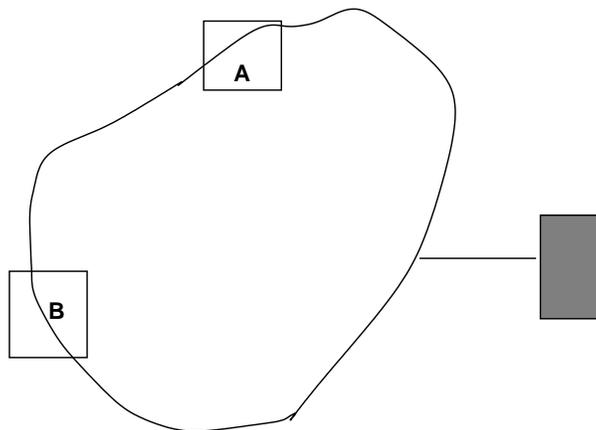


Abbildung 3.9: Adress-Auflösung

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

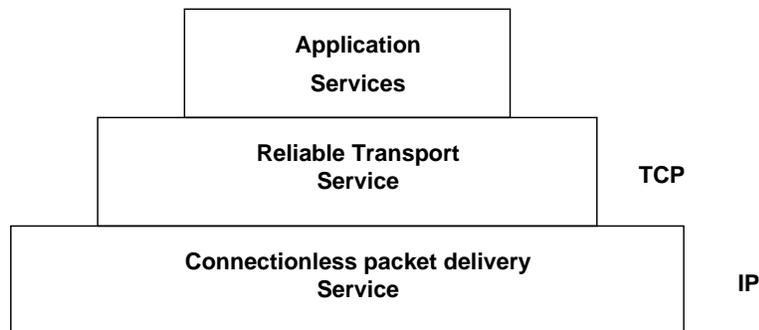


Abbildung 3.10: TCP/IP-Schichtenmodell

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:

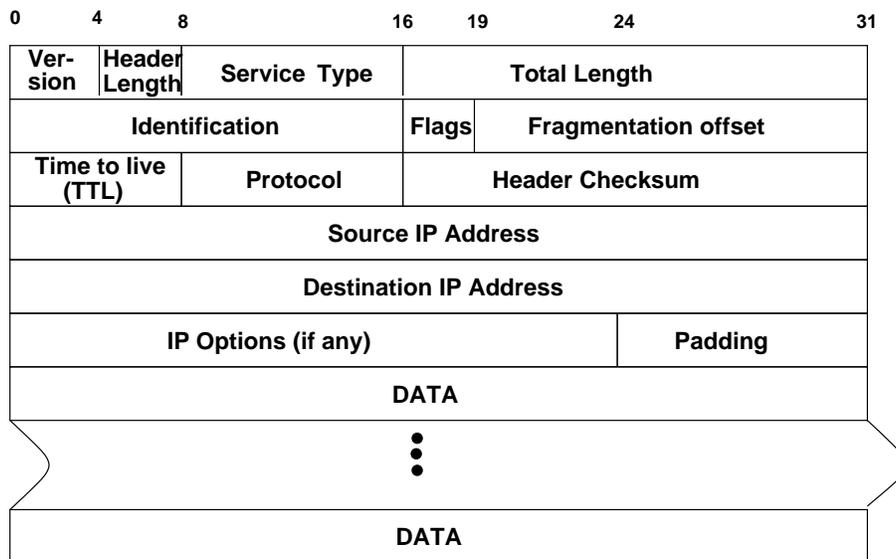


Abbildung 3.12: IP Datagram - im Detail

Erläuterungen siehe unten 'Felder im Datagram'

Einschub: Datagram Encapsulation

- Physikalische Network Frames → erkannt von Hardware
- Datagram → erkannt von Software

Ein Datagram wird in ein (Idealfall) oder mehrere physikalische Frames gepackt, abhängig vom physikalischen Netz

- Aufbrechen eines Datagrams in Fragmente → fragmentation

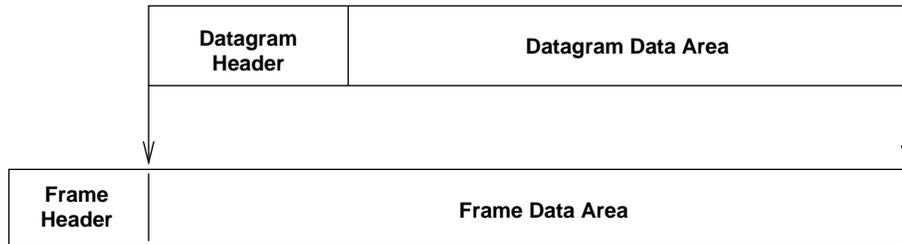


Abbildung 3.13: Datagram Encapsulation

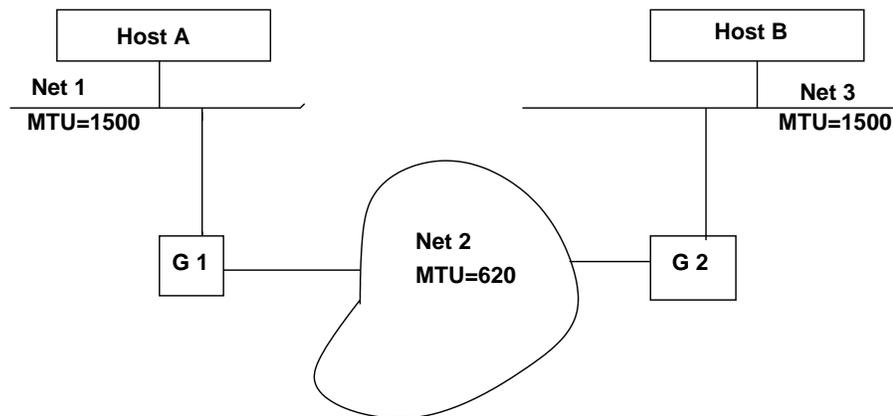


Abbildung 3.14: Maximum Transfer Unit

- MTU (*maximum transfer unit*) in Bytes im physikalischen Netz
- ein Fragment == 1 phys. Frame
- Fragmentierung ↔ Reassembly (Wiederzusammensetzen beim Empfänger!)

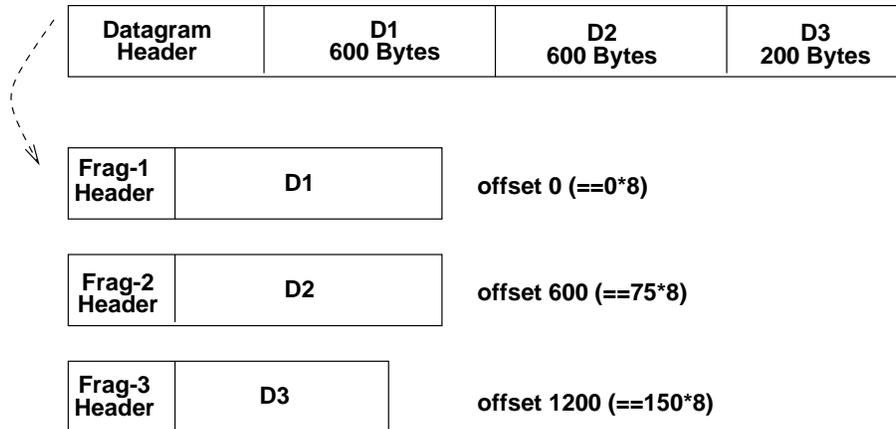


Abbildung 3.15: Datagramm Fragmentation / Reassembly

Felder im Datagramm:

- Version:
IP Software Version muß mit der angegebenen Version übereinstimmen
 - Header length:
in 32-Bit-Einheiten
 - Total length:
in Bytes (bei 16 Bits: max. Datagramm-Länge 2^{16} Bytes)
 - Service Type:
"Wie soll das Paket behandelt werden?"
Substruktur:
- Fragmentation Control:
- Identification
Eindeutige Identifikation des Datagramm für Wiederausammenbau, in alle Fragmente kopiert
 - Fragment Offset
in Einheiten von 8 Bytes
 - Flags
z.B.
"Do not fragment" → nur Datagramm als Ganzes ist für Empfänger brauchbar (z.B. Bootstrap-Code)
"more fragments" → bei Fragmentierung wird *total length* auf Länge des Fragments gesetzt, d.h. Empfänger weiß nichts über Gesamtlänge des Datagramms; wenn dieses Bit auf 0 gesetzt ist, erkennt Empfänger, dass dies das letzte Fragment ist - aus *fragment offset* und *total length* aller angekommenen Fragmente kann der Empfänger die Vollständigkeit der "Lieferung" ermitteln.

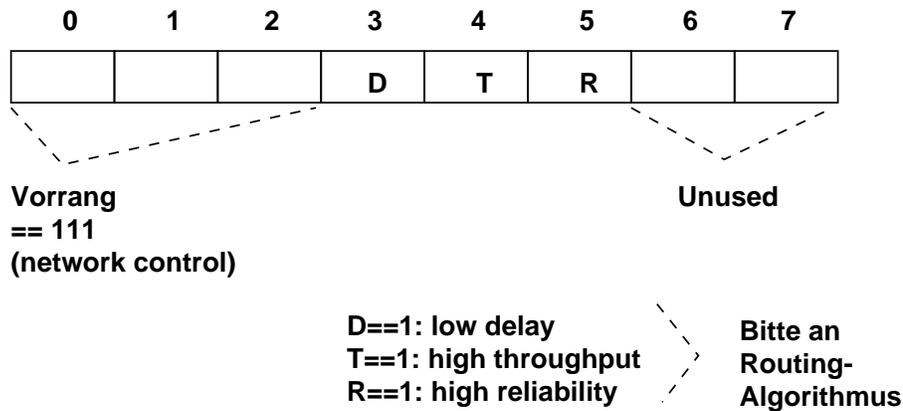


Abbildung 3.16: Service Type (IP)

- Time To Live (TTL):
 - maximale Zeit, die das Datagram im Internet bleiben soll
 - wird auf dem Weg entsprechend dekrementiert
 - TTL auf 0 → entfernen (und Fehlermeldung an Sender)
- Protocol
spezifiziert das Protokoll für den Datenbereich
- IP-Options
i.w. für Debugging u.dgl., z.B. Aufzeichnen der Route; u.U. mehrere Worte lang, wieviele: PADDING

ICMP (Internet Control Messages Protocol)

- In einem verbindungslosen System operieren Gateways autonom (*routing and delivering datagrams*) ohne Koordination mit dem ursprünglichen Sender
- Fehler können auftreten
 - in Verbindungsstrecken
 - bei Prozessoren (Gateways, Rechner)
 - TTL (*time2live*) ist abgelaufen
 - Überlast bei Gateways
- auch Fehlermeldungen gehen in *datagrams* im Datenteil übers Netz

The ICMP allows gateways to send error or control messages to other gateways or hosts; ICMP provides communication between the IP software on one machine and that of another machine

ICMP only reports error conditions to the original source; the source must relate errors to individual application programs and take action to correct the problem

- ICMP ist (notwendiger) Bestandteil der IP-Software (kein eigenständiges *higher level protocol*)
 - Fehler-/Kontrollmeldungen: Ein Datagram, bei dem der Datenteil einen speziellen Aufbau hat (Bestandteil der IP-Protokolle)

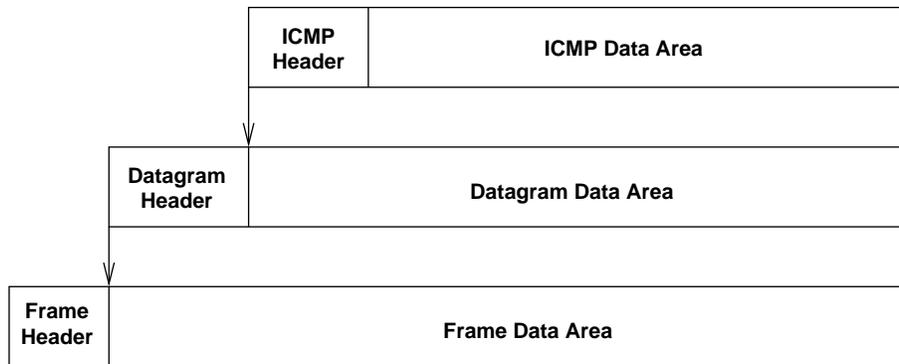


Abbildung 3.17: ICMP Encapsulation

- jede ICMP-Meldungen hat eigenes Format, alle beginnen jedoch mit den gleichen drei Feldern:
 - das erste Feld (8-Bit-Integer) - *TYPE-field* - identifiziert die Meldung (Bedeutung und Aufbau)
 - das zweite Feld (8-Bit-Integer) - *CODE-field* - enthält weitere Informationen über den Meldungstyp
 - das dritte Feld (16-Bit-Integer) - *CHECKSUM-field*
- ICMP-Meldungen, die Fehler berichten, enthalten neben dem *header* die ersten 64 Bits des Datenteils, die die Fehlerursache beschreiben - dadurch kann der Empfänger der Fehlermeldung bestimmen, welches Protokoll und welches Anwendungsprogramm verantwortlich sind.

IP-Protokoll-Schichten

Probleme bei der Rechnerkommunikation im Netz:

- Hardware-Fehler / -Ausfall
- Netzwerk-Überlast
- Verzögerung / Verlust von Paketen
- Datenverfälschung
- Datenverdoppelung oder Fehler in der Abfolge von Datenpaketen

TCP/IP Internet 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.

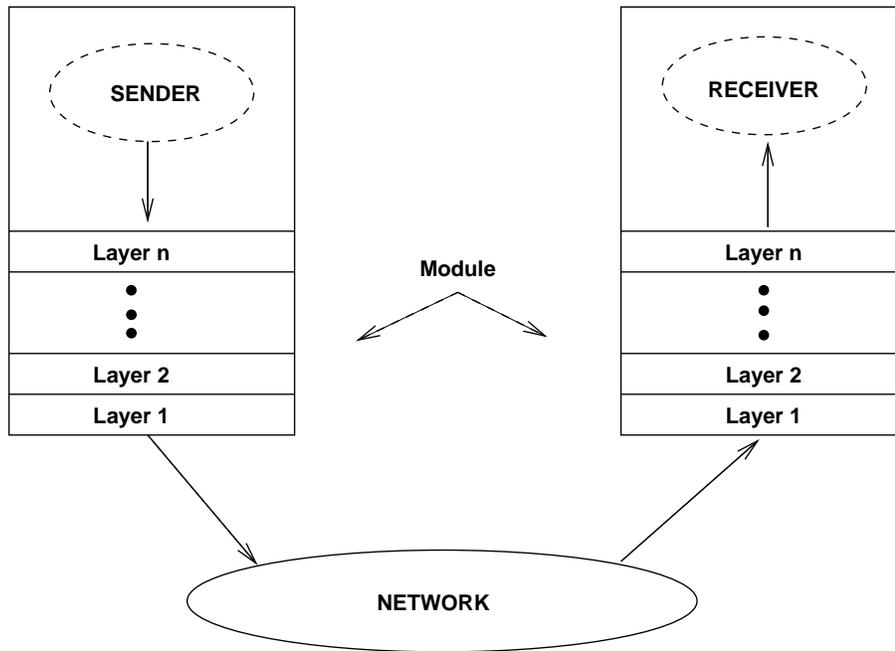


Abbildung 3.18: Module Layering

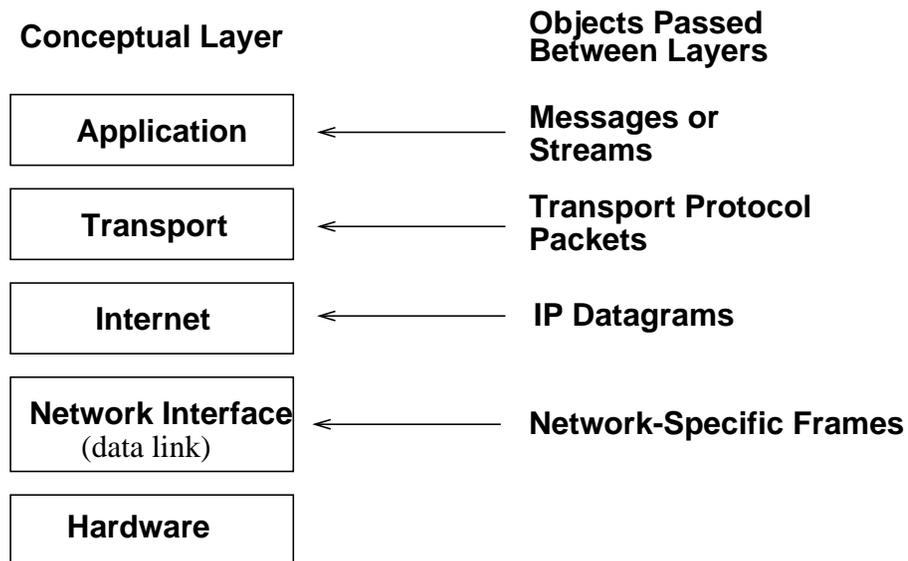


Abbildung 3.19: TCP/IP Layering Model

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

Protocol Layering Principle

☞ *Layered protocols are designed so that layer n at the destination receives exactly the same object sent by layer n at the source.*

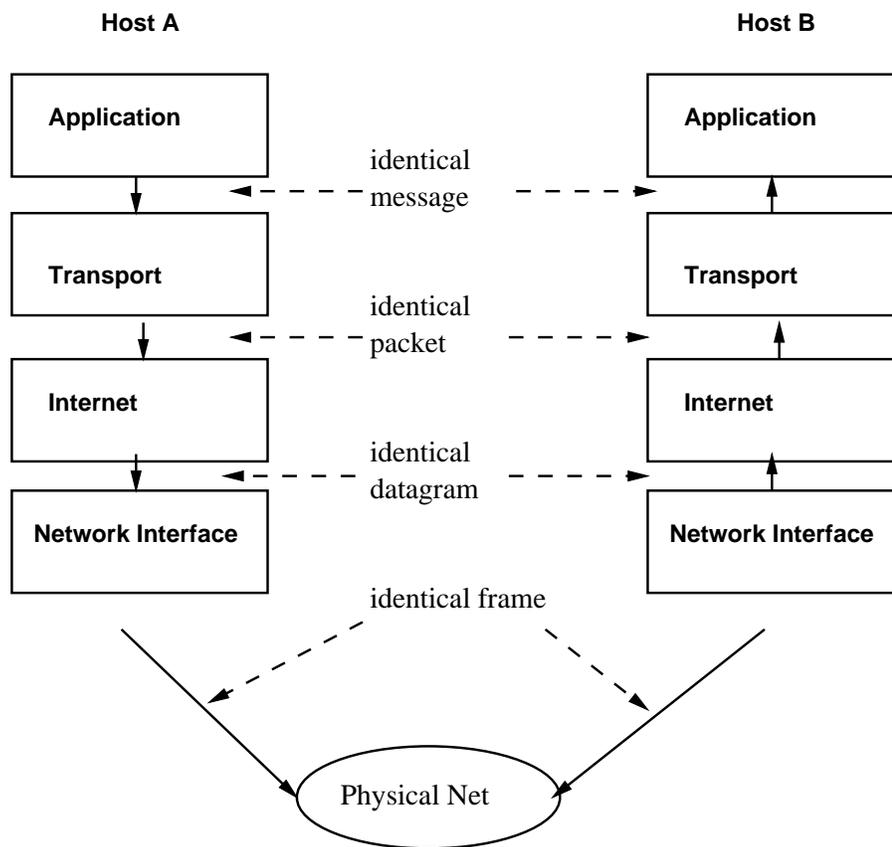


Abbildung 3.20: Protocol Layering

Mit Gateway(s):

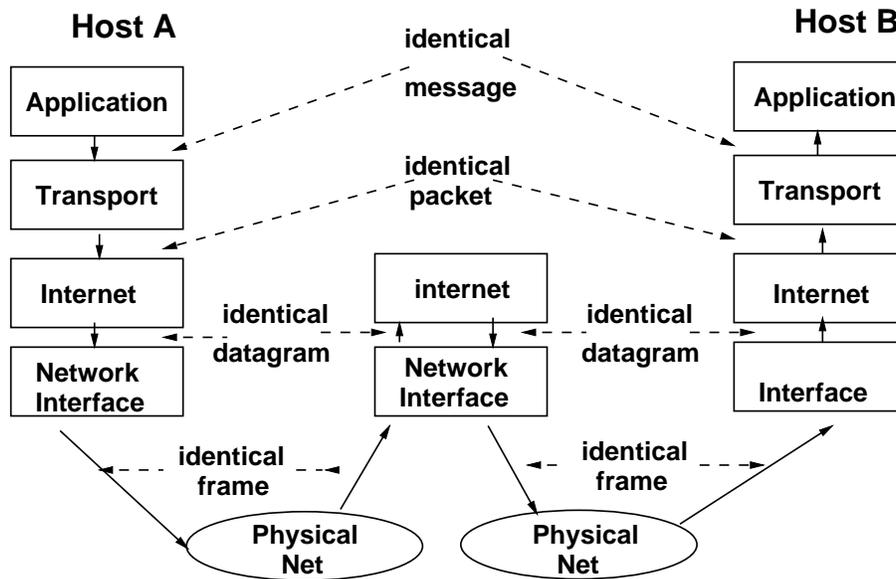


Abbildung 3.21: Protocol Layering with Gateways

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:

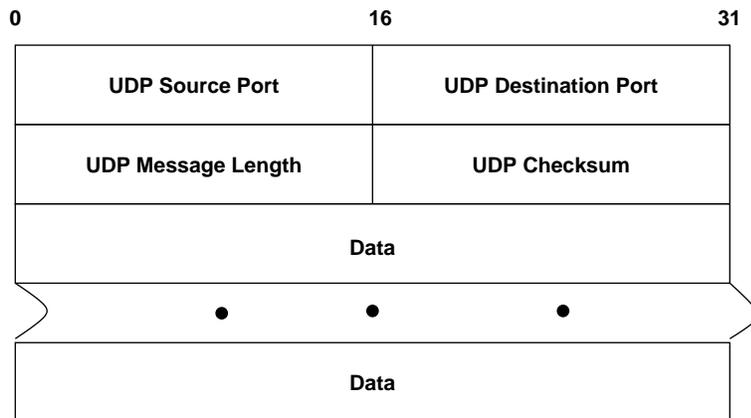


Abbildung 3.22: UDP - Format

□ *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 keiner 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!

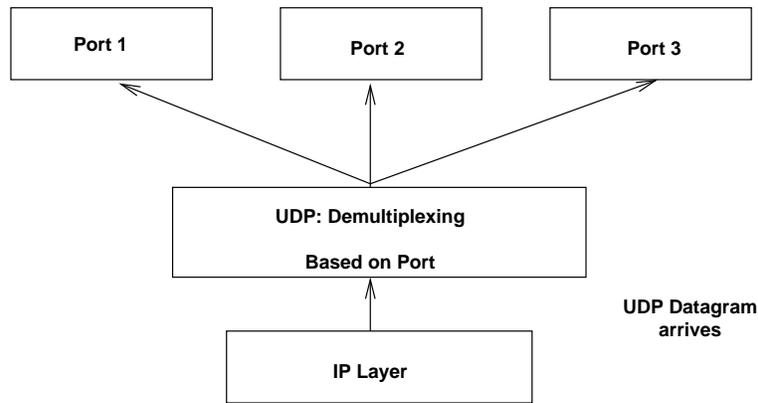


Abbildung 3.23: UDP-Demultiplexing

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	sysstat	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
.			
.			
.			

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.

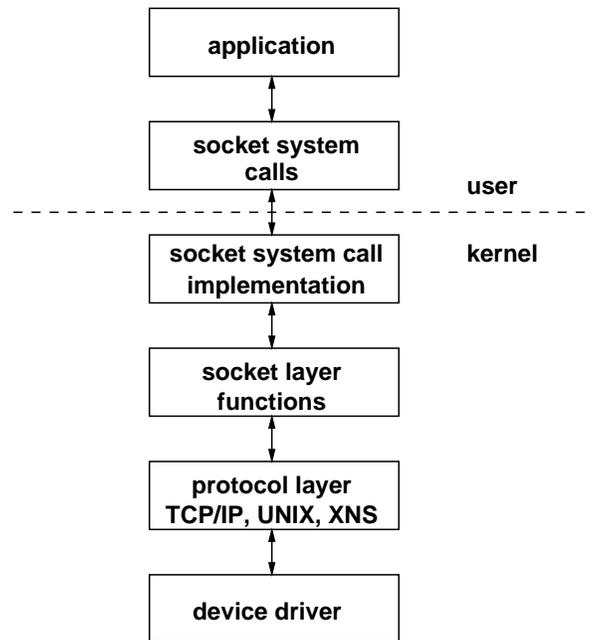


Abbildung 4.1: Vereinf. Modell der Implementierung von Sockets unter BSD

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

4.2 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.

4.3 Die Socket-Programmierschnittstelle

4.3.1 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.

4.3.2 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.

4.3.3 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 sockaddr * 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 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` 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 *)&sun_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));
```

4.3.4 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.

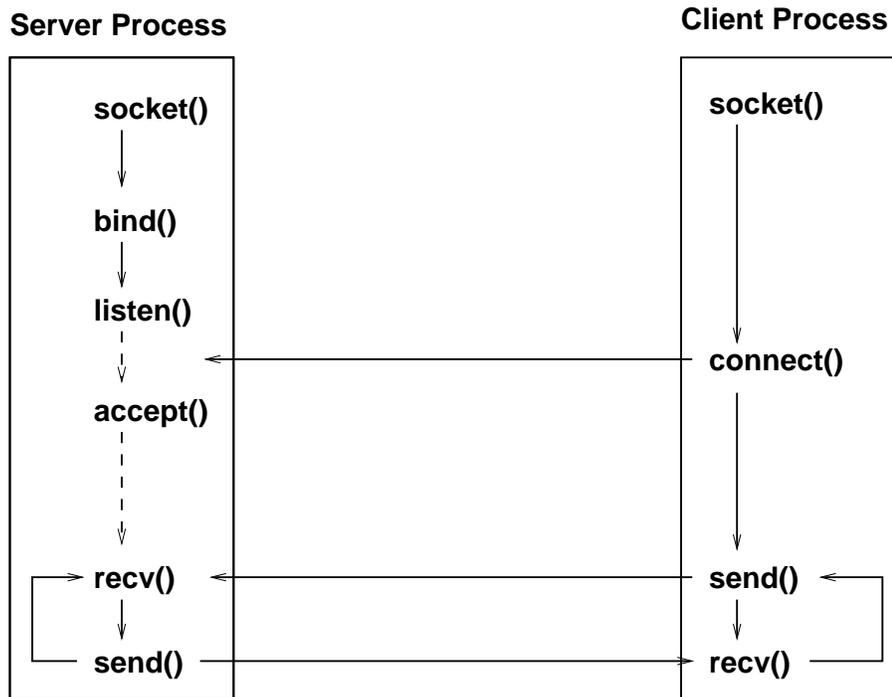


Abbildung 4.2: Verbindungsorientierte Client-Server-Kommunikation

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.

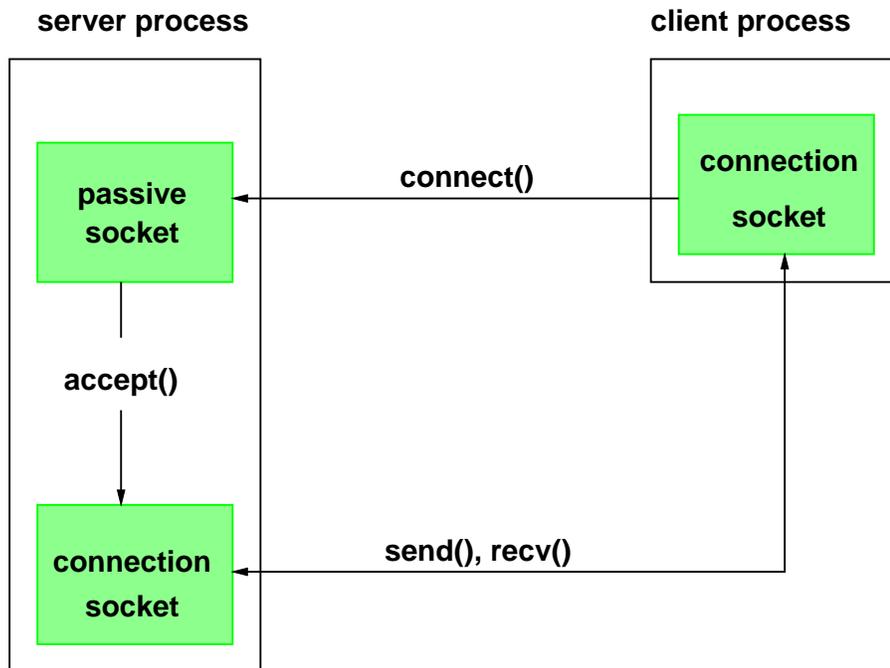


Abbildung 4.3: Aufbau einer verbindungsorientierten Client/Server-Kommunikation

- 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);

```

4.3.5 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()**, **writv()** (*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!

4.3.6 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

4.3.7 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:

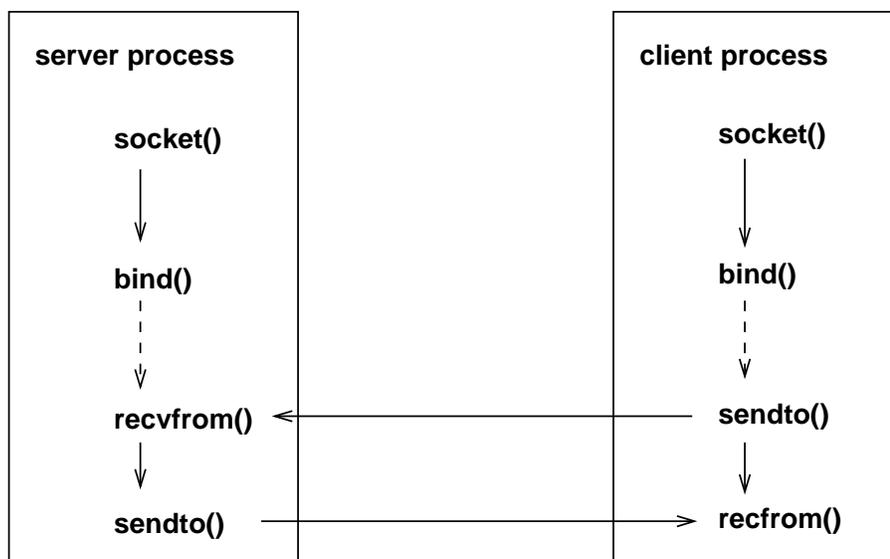


Abbildung 4.4: Aufbau einer verb.-losen Client/Server-Kommunikation

- 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.

4.3.8 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.

4.3.9 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
 - **writenv():**
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

4.4 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.

4.4.1 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 Längenangabe 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).

4.4.3 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**:

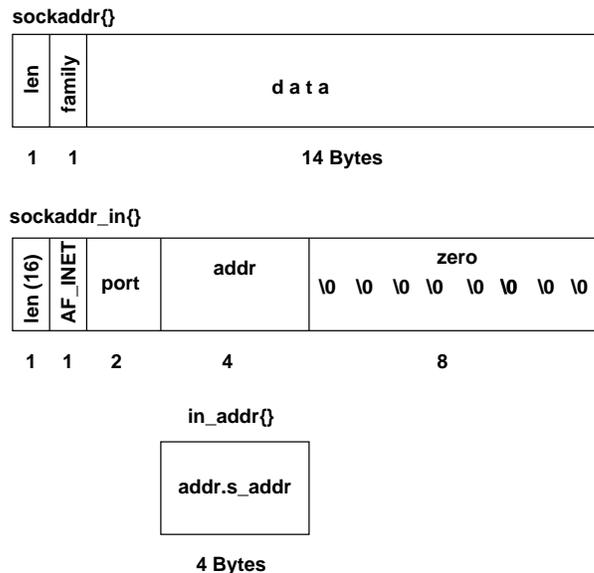


Abbildung 4.5: Organisation der Adress-Struktur `sockaddr_in`

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.

4.4.4 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*).

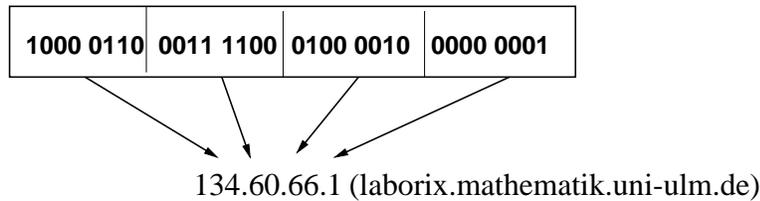


Abbildung 4.6: dotted-decimal notation

4.4.5 Spezifikation von Internet-Adressen

Die 32-Bit-IP-Adressen werden meist in der sog. **punktiertes Dezimalformat (dotted decimal notation)** angegeben; diese entsteht dadurch, dass byteweise 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));
```

4.4.6 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:

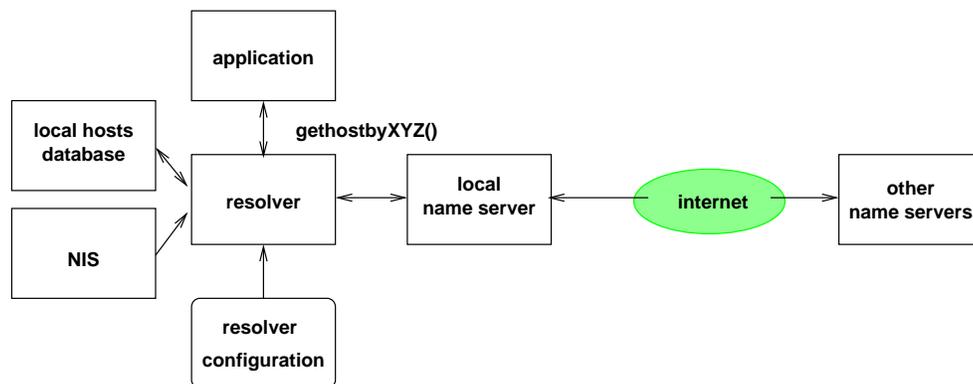


Abbildung 4.7: Methoden der Adress-Resolution

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
-----
```

4.4.7 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$
```

4.4.8 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 Superuser-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$

```

4.4.9 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>

4.4.10 Zusammenfassung der Netzwerkinformationen

Die Erzeugung und Benennung eines Kommunikationsendpunktes erfordert

- die Selektion eines Kommunikationsprotokolls und
- die Konstruktion einer protokollspezifischen Adresse, die in der Internetdomäne 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!**

4.4.11 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 Datenflusskontrolle von Echtzeitanwendungen (z.B. Übermittlung von multimedialen 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

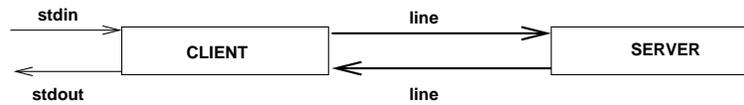


Abbildung 5.1: echo: Client/Server

- 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( (ignoresig(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.

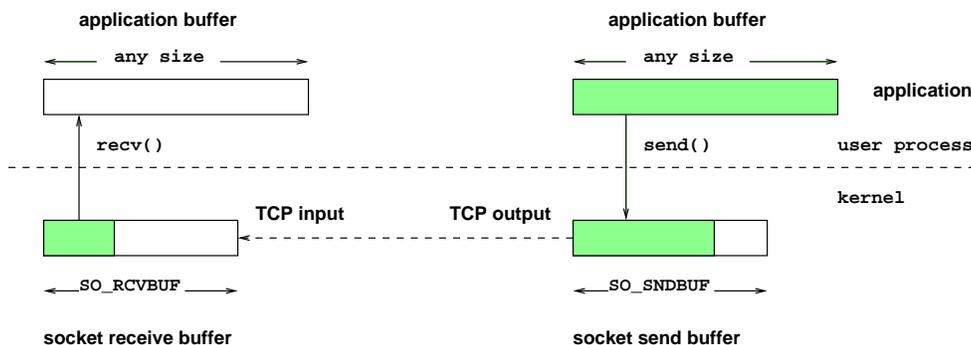


Abbildung 5.2: Socket-Pufferung (vereinfacht)

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 Netzwerkanwendungen 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 `recvn()` 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 `recvn()` 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_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_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	96
3.2	TCP/IP	97
3.3	Ethernet Frame Format	100
3.4	Layered System Architecture	101
3.5	Internet Architecture	102
3.6	Abstraktion der Identifikation	102
3.7	32-Bit IP-Adresse	103
3.8	IP-Adresse: dotted decimal form	104
3.9	Adress-Auflösung	104
3.10	TCP/IP-Schichtenmodell	105
3.11	IP Datagramm - grober Aufbau	106
3.12	IP Datagram - im Detail	106
3.13	Datagram Encapsulation	107
3.14	Maximum Transfer Unit	107
3.15	Datagram Fragmentation / Reassembly	108
3.16	Service Type (IP)	109
3.17	ICMP Encapsulation	110
3.18	Module Layering	111
3.19	TCP/IP Layering Model	111
3.20	Protocol Layering	113
3.21	Protocol Layering with Gateways	114
3.22	UDP - Format	115
3.23	UDP-Demultiplexing	116
4.1	Vereinf. Modell der Implementierung von Sockets unter BSD	118
4.2	Verbindungsorientierte Client-Server-Kommunikation	123
4.3	Aufbau einer verbind.-orient. Client/Server-Kommunikation	125
4.4	Aufbau einer verb.-losen Client/Server-Kommunikation	129
4.5	Organisation der Adress-Struktur sockaddr_in	135
4.6	dotted-decimal notation	137
4.7	Methoden der Adress-Resolution	140
5.1	echo: Client/Server	150
5.2	Socket-Pufferung (vereinfacht)	164

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