



ulm university universität
uulm

**Vorlesungsbegleiter zu
Systemnahe Software II
SS 2015**

Andreas F. Borchert

Fakultät für Mathematik und Wirtschaftswissenschaften
Institut für Numerische Mathematik

Hinweise:

- Auf eine detaillierte Unterscheidung zwischen *BSD-Unix*, *System-V-Unix* oder *Linux* wird hier verzichtet. Stattdessen dient der IEEE Standard 1003.1 (POSIX) weitgehend als Grundlage.
- Die enthaltenen Beispiel-Programme wurden zum großen Teil unter Linux entwickelt und sind weitestgehend unter Solaris getestet (für konstruktive Hinweise sind die Autoren dankbar).
- Die Beispiele sollen jeweils gewisse Aspekte verdeutlichen und erheben nicht den Anspruch von Robustheit und Zuverlässigkeit. Man kann alles anders und besser machen.
- Details zu den behandelten bzw. verwendeten Systemaufrufen sollten jeweils im **Manual** bzw. den entsprechenden Header-Files nachgelesen werden.
- Die Sprache C dient in erster Linie als *Werkzeug* zur Darstellung systemnaher Konzepte!

Inhaltsverzeichnis

1	Prozesse unter UNIX	1
1.1	Virtueller Adressraum	1
1.2	Ausführungskontext	2
1.3	Die Prozess-ID	3
1.4	Terminierung eines Prozesses	3
1.5	Das Erzeugen neuer Prozesse	3
1.6	Synchronisierung bei der Prozessterminierung	5
1.7	Sonderfälle: Zombies und der <i>init</i> -Prozess	9
1.8	Der Wechsel zu einem anderen Programm	11
1.9	Das Zusammenspiel von <i>fork</i> , <i>exec</i> , <i>exit</i> und <i>wait</i>	12
2	Signale	21
2.1	Einführung	21
2.2	Signalbehandler	22
2.3	Wecksignale mit <i>alarm</i>	24
2.4	Das Versenden von Signalen	26
2.5	Die Zustellung von Signalen	29
2.6	Signale als Indikatoren für terminierte Prozesse	32
2.7	Signalbehandlung in einer Shell	35
2.8	Überblick der Signale aus dem POSIX-Standard	41
3	Pipelines	43
4	Einführung in Netzwerkdienste und TCP/IP	59
4.1	Netzwerkdienste	59
4.2	IP-Adressen	60
4.3	Berkeley Sockets	62
4.4	Aufbau einer Netzwerk-Verbindung	63
4.5	Gepufferte Ein- und Ausgabe für Netzwerkverbindungen	67
4.6	Ein- und Ausgabe von Paketen für Netzwerkverbindungen	72
4.7	Parallele Sitzungen	77

Kapitel 1

Prozesse unter UNIX

1.1 Virtueller Adressraum

Bei einfachen Prozessoren gibt es keinen Unterschied zwischen physischen und virtuellen Speicheradressen und damit auch keine scharfe Trennlinie zwischen Betriebssystem und Anwendung. Dies macht nicht nur die Speicheraufteilung unübersichtlich, es führt auch dazu, dass ein Absturz einer Anwendung, beispielsweise hervorgerufen durch einen ungültigen Zeigerzugriff, das gesamte System beeinträchtigen kann.

Besser ausgestattete Prozessoren besitzen eine Speicherverwaltung (*memory management unit*, abgekürzt MMU), die die Einrichtung virtueller Speicherumgebungen ermöglicht. Anwendungen verwenden dann virtuelle Speicheradressen, die mittels einer vom Betriebssystem konfigurierten Funktion in physische Speicheradressen abgebildet werden.

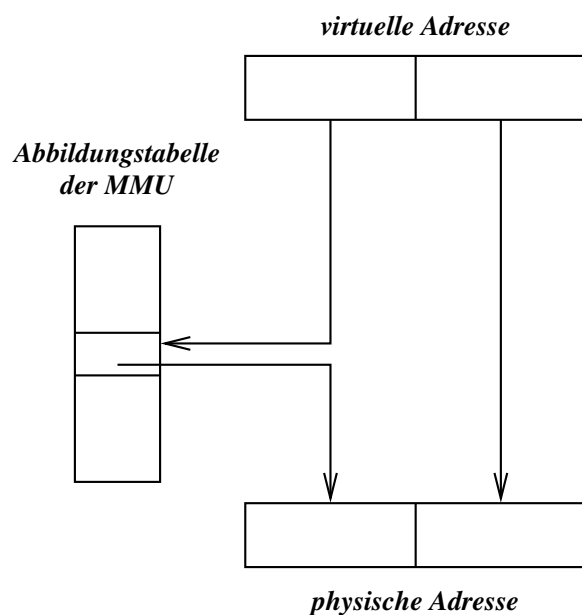


Abbildung 1.1: Virtuelle und physische Adressen

Abbildung 1.1 zeigt schematisch, wie typischerweise virtuelle Adressen in physische Adressen abgebildet werden. Um die Abbildung effizient (in der Hardware) durchführen

zu können, wird der gesamte Speicher in Form von Seiten organisiert. Typische Seitengrößen liegen zwischen 4 und 64 Kilobyte. Entsprechend lassen sich virtuelle Adressen in zwei Teile zerlegen: Der niedrigwertige Teil spezifiziert nur die Position innerhalb einer Seite, der höherwertige gibt die Seitennummer an. In einer Tabelle der MMU können nun die Adressen physischer Seiten den Seitennummern virtueller Adressen zugeordnet werden.

Mehrere Abbildungstabellen können parallel nebeneinander existieren: Eine für das Betriebssystem selbst und eine für jede Anwendung. Ein Wechsel lässt sich effizient realisieren, indem ein spezielles Hardware-Register verwendet wird, das auf die derzeitig gültige Abbildungstabelle verweist.

Die Abbildungstabellen einer MMU bieten typischerweise noch Platz für Zugriffsrechte. So kann geregelt werden, welche Seiten aus der Sicht einer gegebenen Abbildungstabelle vorhanden, lesbar, schreibbar und ausführbar sind. Die Gesamtheit aus Abbildungstabelle und Zugriffsrechten wird als *virtueller Adressraum* bezeichnet. Kommt es bei einem Zugriff zu einem Fehler durch die MMU, kann das Betriebssystem entscheiden, ob es sich dabei um einen Fehler der Anwendung handelt oder ob es den Zugriff nach einigen Manipulationen doch noch ermöglicht.

Viele Tricks sind auf diese Weise möglich:

- Zwei oder mehr Adressräume können gemeinsam auf die gleiche Bibliothek zugreifen. Hierbei ist es sinnvoll, für den gemeinsamen Bereich nur Lese- und Ausführungsrechte einzuräumen.
- Teile eines belegten Adressraumes können auf die Platte ausgelagert werden. Bei einem Zugriff kommt es zuerst zu einem Fehler der MMU, der dann vom Betriebssystem abgefangen wird, um den Bereich wieder von Platte einzulesen und zur Verfügung zu stellen.
- Kopien von Speicherbereichen werden verzögert angefertigt, indem der zunächst noch gemeinsam gehaltene Bereich mit einem Schreibschutz versehen wird. Sobald eine Schreiboperation erfolgt, wird die betroffene Seite dupliziert und bei beiden Kopien werden dann die Schreibrechte wieder zurückgegeben.

1.2 Ausführungskontext

Ein Programm wird in einem bestimmten Kontext ausgeführt. Zu diesem Kontext gehören

- der Adressraum, in dem unter anderem der Programmtext (als Maschinencode) und die Daten untergebracht sind,
- ein Satz Maschinenregister einschließlich der Stackverwaltung (Stack-Zeiger, Frame-Zeiger) und dem PC (*program counter*, verweist auf die nächste auszuführende Instruktion) und
- weitere Statusinformationen, die vom Betriebssystem verwaltet werden wie beispielsweise Informationen über geöffnete Dateien.

Die Einrichtung eines Kontexts und der Wechsel zwischen verschiedenen Kontexten gehört zu den Aufgaben des Betriebssystems. Interessant ist die Frage, ob Ausführungskontexte nur isoliert existieren oder ob es mehrere Kontexte geben kann, die Teile gemeinsam haben.

Die Kombination aus Adressraum und Statusinformationen des Betriebssystems wird unter UNIX als *Prozess* bezeichnet. Zu einem Prozess gehört mindestens ein Satz Maschinenregister einschließlich einem PC. Es können aber auch mehrere sein. In diesem Falle

wird von *Threads* gesprochen, d.h. es existieren mehrere Ausführungsfäden, die parallel abgearbeitet werden. Allerdings verwaltet UNIX auch Statusinformationen für einzelne *Threads*, so dass bei einem Prozess unter UNIX auch gelegentlich von einer Rechtegemeinschaft gesprochen wird.

Teile der Statusinformationen wie beispielsweise die Dateiverbindungen einschließlich der aktuellen Zugriffsposition können auf dem Wege der Vererbung zwischen UNIX-Prozessen geteilt oder auch explizit übergeben werden.

1.3 Die Prozess-ID

Jeder Prozess hat unter UNIX eine gleichbleibende identifizierende positive ganze Zahl, die mit *getpid()* abgefragt werden kann:

Programm 1.1: Ausgabe der eigenen Prozess-ID (*printpid.c*)

```
1 #include <stdio.h>
2 #include <unistd.h>
3
4 int main() {
5     printf("%d\n", (int) getpid());
6 }
```

Bei der Mehrheit der UNIX-Systeme liegt die Prozess-ID im Bereich von 1 bis 32767. Die Eindeutigkeit ist jedoch nur zu Lebzeiten garantiert. Sobald ein Prozess beendet wird, kann die gleiche Prozess-ID später einem neuen Prozess zugeordnet werden. Alle gängigen UNIX-Systeme vergeben Prozess-IDs reihum, wobei bereits vergebene Prozess-IDs übersprungen werden.

1.4 Terminierung eines Prozesses

Ein Prozess kann sich jederzeit mit *exit()* beenden und dabei einen Statuswert im Bereich von 0 bis 255 angeben. Die *exit*-Funktion kann in C-Programmen auch implizit aufgerufen werden: Ein **return** in der *main*-Funktion führt zu einem entsprechenden *exit* und wenn das Ende der *main*-Funktion erreicht wird, entspricht dies einem *exit(0)*. Ein Exit-Wert von 0 deutet dabei eine erfolgreiche Terminierung an, andere Werte, insbesondere *EXIT_FAILURE*, werden als Mißerfolg gewertet. Diese Konventionen orientieren sich zwar an UNIX, sind aber auch Bestandteil des ISO-Standards 9899-1999.

1.5 Das Erzeugen neuer Prozesse

Neue Prozesse können nur in Form eines Klon-Vorganges mit Hilfe des Systemaufrufs *fork()* erzeugt werden:

- Der Adressraum, die Maschinenregister und fast der gesamte Status des Betriebssystems für den erzeugenden Prozess werden dupliziert. Das bedeutet, dass beide Prozesse (der *fork()* aufrufende Prozess und der neu erzeugte Prozess) einen zu Beginn gleich aussehenden Adressraum vorfinden. Änderungen werden jedoch nur bei jeweils einem der beiden Prozesse wirksam, so dass hier die Verzögerungstechnik beim Kopieren von Speicherbereichen zum Zuge kommt (*copy on write*).
- Einige Statusinformationen beim Betriebssystem betreffen beide Prozesse. So werden offene Dateiverbindungen vererbt und können gemeinsam genutzt werden. Dies

bezieht sich aber nur auf Dateiverbindungen, die zum Zeitpunkt des *fork*-Aufrufs eröffnet waren und nicht auf Dateien, die später von einem der beiden Prozesse neu eröffnet werden.

- Einige Statusinformationen des Betriebssystems werden *nicht* weitergegeben. Dazu gehören beispielsweise Locks und anhängige Signale.

Das überraschende ist hier, dass ein neuer Prozess nicht irgendwo mit einem neuen Programm bei *main()* beginnt, sondern wir nach dem *fork()* zwei weitgehend übereinstimmende Kopien eines Prozesses vorfinden, die alle das Programm hinter dem Aufruf von *fork()* fortsetzen. Der folgende Programmtext demonstriert bei der Ausführung, dass das zweite *printf* doppelt ausgeführt wird:

Programm 1.2: Aus einem Prozess werden zwei (*clones.c*)

```

1 #include <stdio.h>
2 #include <unistd.h>
3 int main() {
4     printf("I am feeling lonely!\n");
5     fork();
6     printf("Hey, I am cloned!\n");
7 }
```

```

doolin$ clones
I am feeling lonely!
Hey, I am cloned!
Hey, I am cloned!
doolin$
```

Dieses einfache Beispiel ist gut geeignet, eine Falle von *fork()* zu demonstrieren:

```

doolin$ clones | cat
I am feeling lonely!
Hey, I am cloned!
I am feeling lonely!
Hey, I am cloned!
doolin$
```

Warum erhalten wir jetzt die Ausgabe "I am feeling lonely!" nun doppelt? Die Antwort ist in der Pufferung der *stdio*-Bibliothek zu suchen. Erfolgt die Ausgabe direkt auf ein Terminal, wird zeilenweise gepuffert. In diesem Falle erfolgt die Ausgabe des ersten *printf()* noch vor dem Aufruf von *fork()*. Falls jedoch voll gepuffert wird — dies ist bei der Ausgabe in eine Datei oder in eine Pipeline der Fall — dann erfolgt vor dem *fork()* noch keine Ausgabe. Stattdessen wird der Puffer von *stdout* durch *fork()* dupliziert, womit die doppelte Ausgabe der ersten Zeile provoziert wird. Dieser Effekt lässt sich durch die rechtzeitige Leerung des Puffers mit Hilfe von *fflush()* vermeiden, wie folgender Programmtext zeigt:

Programm 1.3: Prozessduplizierung mit vorheriger Pufferleerung (*clones2.c*)

```

1 #include <stdio.h>
2 #include <unistd.h>
3 int main() {
4     printf("I am feeling lonely!\n"); fflush(stdout);
5     fork();
6     printf("Hey, I am cloned!\n");
7 }
```

Es bleibt die Frage, ob die beiden Prozesse in Abhängigkeit davon, ob sie der alte oder der neue Prozess sind, unterschiedliche Dinge tun können. Folgendes Beispiel zeigt, wie dies mit Hilfe von `getpid()` geschehen kann:

Programm 1.4: Unterscheidung von Vorfahre und Nachfahre (*clones3.c*)

```

1 #include <stdio.h>
2 #include <unistd.h>
3 int main() {
4     pid_t parent;
5
6     printf("I_am_feeling_lonely!\n"); fflush(stdout);
7     parent = getpid();
8     fork();
9     if (getpid() == parent) {
10        printf("I_am_the_parent_process!\n");
11    } else {
12        printf("I_am_the_child_process!\n");
13    }
14 }
```

Die Unterscheidung kann aber auch mit Hilfe des Rückgabewertes von `fork()` selbst erfolgen. `fork()` liefert -1 im Falle von Fehlern, 0 für den neu erzeugten Prozess und die Prozess-ID des neu erzeugten Prozesses beim alten Prozess. Folgende Variante demonstriert nicht nur die Nutzung des Rückgabewertes, sondern auch die Vermeidung eines gemeinsamen Pfades nach dem `fork()` durch ein explizites `exit()`.

Programm 1.5: Der Rückgabewert von `fork()` (*fork.c*)

```

1 #include <stdio.h>
2 #include <unistd.h>
3 int main() {
4     pid_t pid;
5
6     pid = fork();
7     if (pid == -1) {
8         perror("unable_to_fork"); exit(1);
9     }
10    if (pid == 0) {
11        /* child process */
12        printf("I_am_the_child_process:%d.\n", (int) getpid());
13        exit(0);
14    }
15    /* parent process */
16    printf("The_pid_of_my_child_process_is_%d.\n", (int) pid);
17 }
```

1.6 Synchronisierung bei der Prozessterminierung

Es mag Fälle geben, bei denen neue Prozesse erzeugt und dann „vergessen“ werden. Im Normalfall jedoch stößt das weitere Schicksal des neuen Prozesses auf Interesse und insbesondere ist es nicht unüblich, dass der erzeugende Prozess auf das Ende der von ihm erzeugten Prozesse warten möchte.

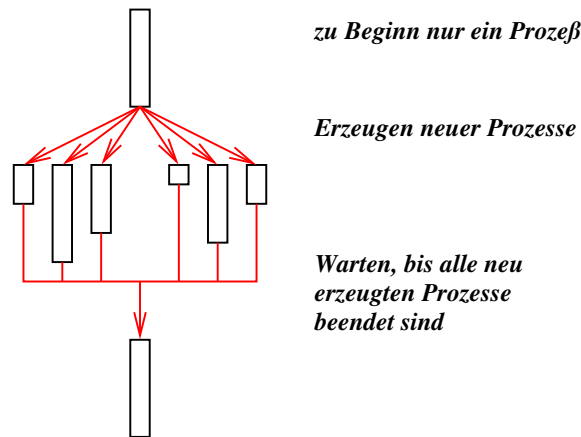


Abbildung 1.2: „Fork and Join“

Dies macht insbesondere dann Sinn, wenn mehrere Prozesse erzeugt werden, die parallel Teilprobleme des Gesamtproblems lösen. Dann wartet der erzeugende Prozess nach Erzeugung aller Unterprozesse, bis sie alle ihre Teilaufgaben erledigt haben. Dieses Muster wird „Fork and Join“ genannt (siehe Abbildung 1.2).

Der Systemaufruf, der es erlaubt, auf das Ende eines Prozesses zu warten, nennt sich `wait()`. Dabei ist zu beachten, dass nur das Warten auf Unterprozesse des aufrufenden Prozesses möglich ist. In der einfachsten Variante wartet `wait()`, bis der nächste Unterprozess beendet ist und liefert die zugehörige Prozess-ID zurück. Zusätzlich wird auch noch der bei `exit()` angegebene Wert geliefert.

Der folgende Programmtext demonstriert `wait()` an einem einfachen Beispiel:

Programm 1.6: Warten auf die Terminierung (`forkandwait.c`)

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <sys/wait.h>
5
6  int main() {
7      pid_t child, pid;
8      int stat;
9
10     child = fork();
11     if (child == -1) {
12         perror("unable_to_fork"); exit(1);
13     }
14     if (child == 0) {
15         /* child process */
16         srand(getpid());
17         char randval = rand();
18         exit(randval);
19     }
20
21     /* parent process */
22     pid = wait(&stat);
23     if (pid == child) {

```

```

24     if (WIFEXITED(stat)) {
25         printf("exit_code_of_child_=%d\n", WEXITSTATUS(stat));
26     } else {
27         printf("child_terminated_abnormally\n");
28     }
29 } else {
30     perror("wait");
31 }
32 }

```

In den Zeilen 16 bis 18 gibt der neu erzeugte Prozess einen etwas abwechslungsreichen Rückgabewert bei `exit()` an. In Zeile 22 wartet der erzeugende Prozess auf die Terminierung des Unterprozesses. `wait()` liefert die Prozess-ID des terminierten Prozesses oder -1, falls es keine Unterprozesse mehr gibt, auf die gewartet werden könnte. (Prinzipiell kann ein Prozess auch gestoppt werden. Gestoppte Prozesse sind ebenfalls bei `wait` zu sehen, aber das beschränkt sich typischerweise auf Szenarien im Kontext eines Debuggers.) Der in `stat` abgelegte Status des Unterprozesses besteht aus mehreren Komponenten, die angeben,

- wie ein Prozess sein Leben beendete (durch `exit()` oder durch ein Signal (bei einem Crash oder Verwendung von `kill()`) oder ob der Prozess nur gestoppt wurde,
- welcher Wert bei `exit()` angegeben wurde, falls `exit()` benutzt wurde und
- welches Signal das Leben des Prozesses terminierte bzw. stoppte, falls der Prozess nicht mit `exit()` endete.

Um das Zerpfücken des Status-Wertes zu vereinfachen, wurden im Rahmen des POSIX-Standards einige Makros definiert, die sich in `<sys/wait.h>` befinden und die Abfragen erleichtern. So ermöglicht `WIFEXITED` die Fallunterscheidung, wie der Prozess terminierte und `WEXITSTATUS` liefert den Exit-Wert zurück.

Ein etwas trickreicheres Beispiel ist im folgenden Programmtext zu finden:

Programm 1.7: Rekursion mit Unterprozessen (*forkingqueens.c*)

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <sys/wait.h>
4  #include <unistd.h>
5  #define INSET(member,set) ((1<<(member))&set)
6  #define INCL(set,member) ((set) | = 1<<(member))
7  int main() {
8      const int size = 6; /* square size of the board */
9      struct { int row, col; } pos[size]; /* queen positions */
10     /* a queen on (row, col) threatens a row, a column,
11        and 2 diagonals;
12        rows and columns are characterized by their number (0..n-1),
13        the diagonals by row-col+n-1 and row+col,
14        (n is a shorthand for the square size of the board)
15     */
16     int rows = 0, cols = 0; /* bitmaps of [0..n-1] */
17     int diags1 = 0; /* bitmap of [0..2*(n-1)] used for row-col+n-1 */
18     int diags2 = 0; /* bitmap of [0..2*(n-1)] used for row+col */
19     int row = 0; int col;
20

```

```

21  setnextqueen: for (col = 0; col < size; ++col) {
22      if (INSET(row, rows)) continue;
23      if (INSET(col, cols)) continue;
24      if (INSET(row - col + size - 1, diags1)) continue;
25      if (INSET(row + col, diags2)) continue;
26      int child = fork();
27      if (child == -1) {
28          perror("fork"); exit(255);
29      }
30      if (child == 0) {
31          INCL(rows, row); INCL(cols, col);
32          INCL(diags1, row - col + size - 1);
33          INCL(diags2, row + col);
34          pos[row].row = row; pos[row].col = col;
35          ++row; /* set next queen in next row */
36          if (row == size) exit(1); /* solution found */
37          goto setnextqueen;
38      }
39  }
40
41  /* count the results of all children */
42  int stat; pid_t child; int nofsolutions = 0;
43  while ((child = wait(&stat)) > 0) {
44      if (!WIFEXITED(stat)) continue;
45      int exitval = WEXITSTATUS(stat);
46      if (exitval == 255) exit(255);
47      nofsolutions += exitval;
48  }
49  exit(nofsolutions);
50 }

```

Beim n -Damen-Problem geht es darum, n Damen auf einem $n \times n$ Schachbrett so unterzubringen, dass sie sich gegenseitig nicht bedrohen. Typischerweise wird das Problem im Backtracking-Verfahren gelöst, bei dem rekursiv alle Varianten durchprobiert werden. Klappt es mit einem Weg nicht, werden die bereits gesetzten Damen wieder sukzessive abgebaut, bis sich andere, bislang noch nicht getestete Varianten eröffnen.

Diese sequentielle Vorgehensweise lässt sich auch parallelisieren. Auf der ersten Reihe auf dem Schachbrett gibt es n mögliche Positionen für die erste Dame. Entsprechend können n Prozesse erzeugt werden, die sich jeweils um einen Teilbaum kümmern. Das Verfahren kann auch weiter mit Parallelisierung fortgesetzt werden. Bei größeren Brettgrößen stößt dies jedoch rasch an die Grenze der Prozesstabelle, so dass die vorgestellte Lösung nicht wirklich skalierbar ist.

Um das Problem zu vereinfachen, soll das Programm nur die Zahl der gefundenen Lösungen zählen. Eine Ausgabe der gefundenen Lösungen ist nicht trivial, da hier eine Synchronisierung stattfinden müsste. Sonst würden viele Prozesse konkurrierend versuchen, Ausgabe zu produzieren mit dem Resultat einer wilden Textmischung. Über den Exit-Wert kann für nicht zu große Werte von n bequem die Zahl der gefundenen Lösungen eines Teilbaumes zurückgeliefert werden, so dass der übergeordnete Prozess die Gelegenheit hat, die Einzelresultate zusammenzuzählen.

Zur Lösung selbst: Der gesamte Stand auf dem Schachbrett wird in einer Reihe lokaler Variablen von `main()` verwaltet. Abgesehen von `col` wird keine dieser Variablen im ersten Prozess modifiziert. Stattdessen erfolgen sämtliche Änderungen nur bei Unterprozessen, die dank der Magie des `fork()`-Aufrufes auf Kopien der Variablen arbeiten.

Die entscheidende Schleife beginnt in Zeile 21. Auf den Reihen 0 bis *nofqueens* - 1 sind die Damen bereits gesetzt. Als nächstes ist nun eine Dame auf Reihe *nofqueens* zu plazieren. Dazu geht die Schleife sämtliche Spaltenpositionen durch und überprüft in den Zeilen 22 bis 25, ob die einzelnen Varianten im Konflikt zu den bereits plazierten Damen stehen.

Falls eine weitere Dame gefahrlos gesetzt werden kann, dann wird diese Variante einem Prozess überlassen, der in Zeile 26 neu erzeugt wird. Dieser setzt die Dame innerhalb der Datenstruktur in den Zeilen 31 bis 35. Wenn damit *size* Damen gesetzt sind, ist eine Lösung gefunden worden und der Unterprozess signalisiert dies mit einem *exit(1)* in Zeile 36. Ansonsten gibt es einen Sprung zur Schleife auf Zeile 21, die dann die nächste Zeile zu besetzen versucht.

Die Anweisungen ab Zeile 41 werden von allen Prozessen ausgeführt mit Ausnahme derjenigen, die entweder bei *fork()* scheitern oder höchstselbst eine Lösung gefunden haben. Die Aufgabe besteht hier darin, all die Zahlen der gefundenen Lösungen der einzelnen Unterprozesse zusammenzuzählen. Dies erledigt die **while**-Schleife auf Zeile 43, die *wait()* solange aufruft, bis -1 zurückgeliefert wird, d.h. alle Unterprozesse berücksichtigt worden sind. Die Exit-Werte werden in Zeile 47 einfach aufaddiert. Es findet nur eine zusätzliche Überprüfung statt, ob der Unterprozess normal terminierte und selbst keine Probleme mit der Erzeugung weiterer Unterprozesse hatte. So ein Problem wird mit einem Exit-Wert von 255 in Zeile 28 signalisiert.

Dieses Beispiel ist nicht zur Nachahmung geeignet. Schon ab einer Brettgröße von 7 scheiterte es auf meiner Workstation am Mangel zur Verfügung stehender Einträge in der Prozesstabelle. Dennoch ist der Ansatz brauchbar, wenn es darum geht, auf einer Mehrprozessor-Maschine die vorhandenen Ressourcen auszunutzen. Allerdings macht es dann nicht Sinn, mehr Prozesse zu erzeugen als tatsächlich Prozessoren vorhanden sind. Ein pragmatischer Ansatz könnte darin bestehen, die erste Stufe der Rekursion (also hier die erste Reihe auf dem Schachbrett) zu parallelisieren und ansonsten sequentiell weiterzuarbeiten. Elegant wird der Programmtext jedoch durch so einen Mischansatz nicht.

Ein weiterer Problempunkt ist die Rückgabe gefundener Lösungen. Dies geht entweder nur unter Verwendung von Dateien (leicht zu programmieren, jedoch nicht sehr effizient) oder der direkten Kommunikation zwischen den Prozessen (effizienter, jedoch leider nicht einfach zu programmieren).

Weitere Abzugspunkte ergeben sich aus der Verwendung einer **goto**-Anweisung und gemeinsamer Programmpfade nach dem *fork()*.

1.7 Sonderfälle: Zombies und der *init*-Prozess

Was geschieht mit dem Rückgabewert bei *exit()* und dem sonstigen Endstatus eines Prozesses, wenn der übergeordnete Prozess nicht zeitig *wait()* aufruft? Das UNIX-System lässt solche toten Prozesse noch in seiner Verwaltung weiterleben, so dass der Endstatus noch bewahrt wird, aber die nicht mehr benötigten Ressourcen freigegeben werden. Prozesse, die sich in diesem Stadium befinden, werden als Zombies bezeichnet.

Das folgende Beispiel demonstriert die Erzeugung eines Zombies:

Programm 1.8: Erzeugung eines Zombies (*genzombie.c*)

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <unistd.h>
4
5 int main() {
6     pid_t child = fork();
7     if (child == -1) {
```

```

8     perror("fork"); exit(1);
9     }
10    if (child == 0) exit(0);
11    printf("%d\n", child);
12    sleep(60);
13 }

```

Der neu erzeugte Prozess verabschiedet sich sofort mit `exit()`, während der übergeordnete Prozess mit Hilfe eines `sleep()`-Aufrufes sich für 60 Sekunden zur Ruhe legt. Während dieser Zeit verbleibt der Unterprozeß im Zombie-Status, wie das `ps`-Kommando belegt:

```

doolin$ genzombie&
[1] 24489
doolin$ 24490

doolin$ ps -y lp 24489,24490
 S  UID  PID  PPID  C  PRI  NI   RSS   SZ   WCHAN TTY        TIME CMD
 S  120 24489 23591  0   64  28   616   936          ? pts/31    0:00 genzombi
 Z  120 24490 24489  0    0                0:00 <defunct>
doolin$

```

In der ersten Spalte gibt `ps` bei dieser Aufrufvariante den Status eines Prozesses an. „Z“ steht dabei für Zombie, „S“ für schlafend. Weitere Varianten sind „O“ für gerade arbeitend, „R“ für arbeitsbereit und „T“ für gestoppt.

Hält der Zombie-Status ewig an? Was geschieht, wenn der übergeordnete Prozess — wie in diesem Beispiel — sich verabschiedet, ohne jemals mit `wait()` sich den Status abgeholt zu haben? Solange der übergeordnete Prozess lebt, bleibt der Zombie-Status tatsächlich beliebig lange bestehen. Wenn jedoch ein Unterprozeß verwaist, weil sein übergeordneter Prozess sich verabschiedet, dann wird ihm der Prozess mit der Prozess-ID 1 als neuer übergeordneter Prozess zugewiesen. Dies geschieht unabhängig davon, ob der untergeordnete Prozess noch aktiv ist oder bereits ein Zombie geworden ist. Folgender Programmtext demonstriert die Erzeugung eines Waisenkindes, indem es `getppid()` vor und nach dem Ende des übergeordneten Prozesses ausgibt. `getppid()` steht dabei für „get parent process id“. So sieht eine beispielhafte Ausführung aus:

Programm 1.9: Ein Prozess wird zum Waisenkind (*orphan.c*)

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4
5  int main() {
6      pid_t child;
7      child = fork();
8      if (child == -1) {
9          perror("fork"); exit(1);
10     }
11     if (child == 0) {
12         printf("Hi,_my_parent_is_%d\n", (int) getppid());
13         sleep(5);
14         printf("My_parent_is_now_%d\n", (int) getppid());
15     }
16     sleep(3);
17     exit(0);
18 }

```

```
doolin$ orphan
Hi, my parent is 24583
doolin$ My parent is now 1

doolin$
```

Der Prozess mit der Prozess-ID 1 spielt eine besondere Rolle unter UNIX. Es ist der erste Prozess, der vom Betriebssystem selbst erzeugt wird und den Programmtext, der unter */etc/init* zu finden ist, ausführt. Dieser Prozess startet weitere Prozesse anhand einer Konfigurationsdatei (bei uns unter */etc/inittab*) und ruft ansonsten *wait()* auf, um den Status der von ihm selbst erzeugten Prozesse oder den von Waisenkindern entgegenzunehmen. Auf diese Weise wird dann auch der Zombie-Status eines Prozesses beendet, wenn es zum Waisenkind wird.

1.8 Der Wechsel zu einem anderen Programm

Mit *fork()* ist es möglich, neue Prozesse zu erzeugen. Allerdings teilen die neuen Prozesse sich den Programmtext mit ihrem Erzeuger. Wie ist nun der Wechsel zu einem anderen Programmtext möglich? Die Lösung dafür ist der Systemaufruf *exec()*, der

- den gesamten virtuellen Adressraum des aufrufenden Prozesses auflöst,
- an seiner Stelle einen neuen einrichtet mit einem angegebenen Programmtext,
- sämtliche Maschinenregister für den Prozess neu initialisiert und
- Statusinformationen des Betriebssystems weitgehend unverändert belässt

Der folgende Programmtext dient als erstes Beispiel, bei dem der laufende Prozess seinen Programmtext durch den von *date* austauscht:

Programm 1.10: Wechsel zum Programm *date* (*datum.c*)

```
1 #include <unistd.h>
2
3 int main() {
4     execl(
5         "/usr/bin/date", /* path of the program */
6         "/usr/bin/date", /* name of the program, i.e. argv[0] */
7         "+%d.%m.%Y", /* first argument, i.e. argv[1] */
8         0 /* terminate list of arguments */
9     );
10    /* not reached except if execl failed */
11    perror("/usr/bin/date");
12 }
```

Dabei ergibt sich die Möglichkeit, Kommandozeilenargumente zu übergeben, die dann vom neuen Programm in *argv[]* vorgefunden werden. Hierbei ist zu beachten, dass der Pfadname des auszuführenden Programms (erster Parameter von *execl()*) getrennt angegeben wird von dem Parameter, der später unter *argv[0]* zu finden ist. Normalerweise sind beide gleich — der Name kann aber auch vom Pfad beliebig abweichen. Bei *execl()* dürfen beliebig viele Argumente angegeben werden. Die Liste wird mit einem Nullzeiger beendet.

Normalerweise geht es im Programmtext nach einem Aufruf von *execl()* nicht weiter, weil im Erfolgsfalle das Programm ausgetauscht wurde. Nur bei einem Fehler (weil z.B.

das `date`-Kommando nicht gefunden wurde) wird das Programm hinter dem Aufruf von `execl()` fortgesetzt.

1.9 Das Zusammenspiel von `fork`, `exec`, `exit` und `wait`

Auf den ersten Blick erscheinen diese vier Systemaufrufe seltsam. Warum ist eine Kombination aus `fork()` und `exec()` notwendig, um einen neuen Prozess mit einem neuen Programmtext in Gang zu setzen? Wäre es nicht besser und einfacher, nur einen einzigen Systemaufruf dafür zu haben? Die Frage verschärft sich, wenn berücksichtigt wird, dass in der Zeit der frühen UNIX-Implementierungen die Technik des „copy on write“ noch nicht zur Verfügung stand. Stattdessen war es bei `fork()` notwendig, den gesamten Speicher zu kopieren. Bei BSD wurde deswegen zeitweise `fork1()` eingeführt, das diesen Kopiervorgang unterdrückte, um die typische Kombination von `fork()` und `exec()` nicht zu teuer werden zu lassen.

UNIX ist keinesfalls das erste Betriebssystem, das Prozesse unterstützte. Die älteren Systeme boten in der Tat die Kombination aus `fork()` und `exec()` in einem Systemaufruf an. Überraschenderweise zeigt sich jedoch, dass dies in der Mehrheit der Fälle viel komplizierter ist. Der Haken liegt darin, dass Prozesse häufig eine Umgebung erwarten, die mehr umfaßt als eine Kommandozeile. Wichtiger Bestandteil der Umgebung sind bereits im Vorfeld eingerichtete Ein- und Ausgabeverbindungen und die Zuteilung von Ressourcen.

So sieht die traditionelle Erzeugung eines Prozesses aus:

- Erzeuge einen neuen Prozess mit einem gegebenen Programmtext mit einem Systemaufruf der `fork()` und `exec()` kombiniert.
- Einrichtung der Umgebung für den neuen Prozess.
- Start des neuen Prozesses.

Entsprechend ist es notwendig, alle wichtigen Systemaufrufe für die Einrichtung einer Umgebung einschließlich dem Öffnen von Ein- und Ausgabeverbindungen in zwei Varianten zu unterstützen: Die eine Variante bezieht sich auf den eigenen Prozess, die andere für einen untergeordneten Prozess, der noch nicht gestartet wurde.

Die Trennung in `fork()` und `exec()` eröffnet die Möglichkeit, dass der Programmtext des übergeordneten Prozesses direkt im neu erzeugten Prozess die Umgebung vorbereitet, die dann bei `exec()` von dem gleichen Prozess mit dem neuen Programmtext vorgefunden wird. Wie Abbildung 1.3 zeigt, wird genau dies von den UNIX-Shells genutzt, um die Umgebung für Anwendungen einzurichten. Der übergeordnete Prozess der Shell wartet dann normalerweise mit `wait()` auf das Ende des untergeordneten Prozesses.

Das folgende Beispiel zeigt die Grundstruktur einer einfachen Shell. Jede Eingabezeile ist entweder leer oder enthält genau ein Kommando, das innerhalb der **while**-Bedingung mit `readline()` eingelesen wird. Die eingelesene Zeile wird dann mit der Hilfe des `tokenizer()` in einzelne Wörter zerlegt, die dann als Kommandoname und als weitere Argumente des Kommandos interpretiert werden.

Programm 1.11: Hauptprogramm einer einfachen Shell (`tinys.c`)

```

1 /*
2
3 =head1 NAME
4
5 tinys -- a tiny shell with a minimal set of features
6
```

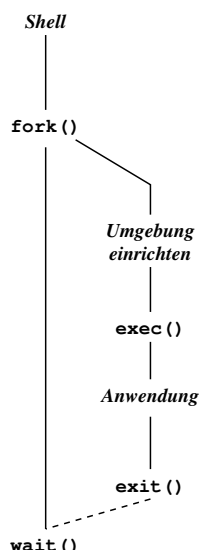



Abbildung 1.3: Start einer Anwendung von der Shell

7 =head1 SYNOPSIS
8
9 B<tinysh>
10
11 =head1 DESCRIPTION
12
13 B<tinysh> reads command lines from its standard input
14 and executes them. % " is given as prompt.
15
16 Each command line consists of space-separated tokens. Special tokens
17 begin with <"', indicating the input file, >"', indicating the output
18 file, and >>"', specifying an output file that is to be extended. The
19 first non-special token specifies the command name which must be either
20 an absolute path of an executable file or locatable within the list of
21 directories provided by the environment variable B<PATH>.
22
23 =head1 DIAGNOSTICS
24
25 B<tinysh> prints the exit code of a terminated program if it is
26 non-zero. It tells that a program terminated abnormally" if it
27 did not exit. An exit code of 255 is used by subprocesses if they
28 are unable to start the command for some reason.
29
30 =head1 BUGS
31
32 No signal handling, no support of backgrounded commands, no pipelines,
33 no builtins, no shell variables.
34
35 =head1 AUTHOR
36
37 Andreas Borchert
38

```

39 =cut
40
41 */
42
43 #include <fcntl.h>
44 #include <stdio.h>
45 #include <stdlib.h>
46 #include <unistd.h>
47 #include <sys/wait.h>
48 #include "sareadline.h"
49 #include "strlist.h"
50 #include "tokenizer.h"
51
52 /*
53  * assign an opened file with the given flags and mode to fd
54  */
55 void fassign(int fd, char* path, int oflags, mode_t mode) {
56     int newfd = open(path, oflags, mode);
57     if (newfd < 0) {
58         perror(path); exit(255);
59     }
60     if (dup2(newfd, fd) < 0) {
61         perror("dup2"); exit(255);
62     }
63     close(newfd);
64 } // fassign
65
66 int main() {
67     stralloc line = {0};
68     while (printf("%%_\n", readline(stdin, &line)) {
69         strlist tokens = {0};
70         stralloc_0(&line); /* required by tokenizer() */
71         if (!tokenizer(&line, &tokens)) break;
72         if (tokens.len == 0) continue;
73         pid_t child = fork();
74         if (child == -1) {
75             perror("fork"); continue;
76         }
77         if (child == 0) {
78             strlist argv = {0}; /* list of arguments */
79             char* cmdname = 0; /* first argument */
80             char* path; /* of output files */
81             int oflags;
82
83             for (int i = 0; i < tokens.len; ++i) {
84                 switch (tokens.list[i][0]) {
85                     case '<':
86                         fassign(0, &tokens.list[i][1], O_RDONLY, 0);
87                         break;
88                     case '>':
89                         path = &tokens.list[i][1];
90                         oflags = O_WRONLY | O_CREAT;
91                         if (*path == '>') {

```

```

92         ++path; oflags |= O_APPEND;
93     } else {
94         oflags |= O_TRUNC;
95     }
96     fassign(1, path, oflags, 0666);
97     break;
98     default:
99         strlist_push(&argv, tokens.list[i]);
100        if (cmdname == 0) cmdname = tokens.list[i];
101    }
102 }
103 if (cmdname == 0) exit(0);
104 strlist_push0(&argv);
105 execvp(cmdname, argv.list);
106 perror(cmdname);
107 exit(255);
108 }
109
110 /* wait for termination of child */
111 int stat;
112 pid_t pid = wait(&stat);
113 if (pid == child) {
114     if (WIFEXITED(stat)) {
115         int code = WEXITSTATUS(stat);
116         if (code && code != 255) {
117             printf("terminated_with_exit_code_%d\n", code);
118         }
119     } else {
120         printf("terminated_abnormally\n");
121     }
122 } else {
123     perror("wait");
124 }
125 }
126 } // main

```

Die folgende Implementierung der Funktion *readline()* basiert auf der *stralloc*-Bibliothek:

Programm 1.12: Zeilenweises Einlesen mit der *stralloc*-Bibliothek (*sareadline.c*)

```

1  /*
2  * Read a string of arbitrary length from a
3  * given file pointer. LF is accepted as terminator.
4  * 1 is returned in case of success, 0 in case of errors.
5  * afb 4/2003
6  */
7
8  #include <stralloc.h>
9  #include <stdio.h>
10
11 int readline(FILE* fp, stralloc* sa) {
12     if (!stralloc_copys(sa, "")) return 0;
13     for(;;) {
14         if (!stralloc_readyplus(sa, 1)) return 0;

```

```

15     if (fread(sa->s + sa->len, sizeof(char), 1, fp) <= 0) return 0;
16     if (sa->s[sa->len] == '\n') break;
17     ++sa->len;
18 }
19 return 1;
20 } // readline

```

Wenn die Zahl der Kommandozeilenargumente variabel ist, empfiehlt sich die Verwendung von `execvp()` oder `execv()` anstelle von `execl()`. Sowohl `execvp()` als auch `execv()` akzeptieren eine mit einem 0-Zeiger terminierte Argumentliste in der Form `char**` wie sie auch an `main()` übergeben wird. `execvp()` durchsucht (anders als die einfachere Variante `execv()`) die Umgebungsvariable `PATH`, um den Programmtext zu finden.

Um eine obere Schranke in der Zahl der Kommandozeilenargumente zu vermeiden, ist es angemessen, analog zur `stralloc`-Bibliothek eine ähnliche Bibliothek für Listen von Zeichenketten einzurichten. Der folgende Programmtext zeigt, wie dies erfolgen kann:

Programm 1.13: Datenstruktur und Schnittstelle für Listen von Zeichenketten (`strlist.h`)

```

1  /*
2  * Data structure for dynamic string lists that works
3  * similar to the stralloc library.
4  * Return values: 1 if successful, 0 in case of failures.
5  * afb 4/2003
6  */
7
8  #ifndef STRLIST_H
9  #define STRLIST_H
10
11 typedef struct strlist {
12     char** list;
13     unsigned int len; /* # of strings in list */
14     unsigned int allocated; /* allocated length for list */
15 } strlist;
16
17 /* assure that there is at least room for len list entries */
18 int strlist_ready(strlist* list, unsigned int len);
19
20 /* assure that there is room for len additional list entries */
21 int strlist_readyplus(strlist* list, unsigned int len);
22
23 /* truncate the list to zero length */
24 int strlist_clear(strlist* list);
25
26 /* append the string pointer to the list */
27 int strlist_push(strlist* list, const char* string);
28 #define strlist_push0(list) strlist_push((list), 0)
29
30 /* free the strlist data structure but not the strings */
31 int strlist_free(strlist* list);
32
33 #endif

```

Programm 1.14: Verwaltung von Listen für Zeichenketten (`strlist.c`)

```

1  /*
2  * Data structure for dynamic string lists that works
3  * similar to the stralloc library.
4  * Return values: 1 if successful, 0 in case of failures.
5  * afb 4/2003
6  */
7  #include <stdlib.h>
8  #include "strlist.h"
9
10 /* assure that there is at least room for len list entries */
11 int strlist_ready(strlist* list, unsigned int len) {
12     if (list->allocated < len) {
13         unsigned int wanted = len + (len>>3) + 8;
14         char** newlist = (char**) realloc(list->list,
15             sizeof(char*) * wanted);
16         if (newlist == 0) return 0;
17         list->list = newlist;
18         list->allocated = wanted;
19     }
20     return 1;
21 }
22
23 /* assure that there is room for len additional list entries */
24 int strlist_readyplus(strlist* list, unsigned int len) {
25     return strlist_ready(list, list->len + len);
26 }
27
28 /* truncate the list to zero length */
29 int strlist_clear(strlist* list) {
30     list->len = 0;
31     return 1;
32 }
33
34 /* append the string pointer to the list */
35 int strlist_push(strlist* list, const char* string) {
36     if (!strlist_ready(list, list->len + 1)) return 0;
37     list->list[list->len++] = string;
38     return 1;
39 }
40
41 /* free the strlist data structure but not the strings */
42 int strlist_free(strlist* list) {
43     free(list->list); list->list = 0;
44     list->allocated = 0;
45     list->len = 0;
46     return 1;
47 }

```

Der Wortzerleger *tokenizer()* arbeitet auf einer als *stralloc* repräsentierten Zeichenkette und generiert eine Liste vom Typ *strlist*, die auf die einzelnen Wörter verweist. Dabei wird ein Umkopieren der Wörter vermieden. Stattdessen verweisen die Zeiger in die originale Zeichenkette und die Leerzeichen werden in Nullbytes verwandelt, um als Begrenzer dienen zu können. Abbildung 1.4 zeigt die resultierende Datenstruktur an einem Beispiel.

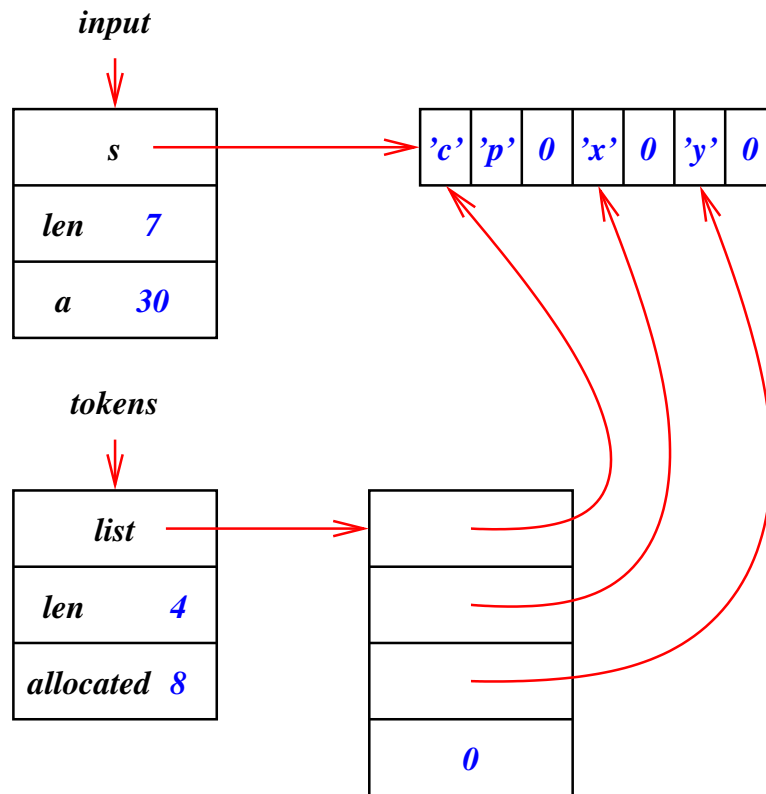


Abbildung 1.4: Resultierende Datenstruktur des Wortzerlegers am Beispiel „cp x y“

Programm 1.15: Zerlegung einer Zeile in Wörter (*tokenizer.c*)

```

1  /*
2  * Simple tokenizer: Take a 0-terminated stralloc object and return a
3  * list of pointers in tokens that point to the individual tokens.
4  * Whitespace is taken as token-separator and all whitespaces within
5  * the input are replaced by null bytes.
6  * afb 4/2003
7  */
8
9  #include <ctype.h>
10 #include <stdlib.h>
11 #include <stralloc.h>
12 #include "strlist.h"
13 #include "tokenizer.h"
14
15 int tokenizer(stralloc* input, strlist* tokens) {
16     char* cp;
17     int white = 1;
18
19     strlist_clear(tokens);
20     for (cp = input->s; *cp && cp < input->s + input->len; ++cp) {
21         if (isspace(*cp)) {
22             *cp = '\0'; white = 1; continue;

```

```
23     }
24     if (!white) continue;
25     white = 0;
26     if (!strlist_push(tokens, cp)) return 0;
27 }
28 return 1;
29 }
```

Damit die Einrichtung einer Umgebung für den neuen Prozess demonstriert werden kann, soll *tinys* einfache Ein- und Ausgabeumleitungen unterstützen. Vorgesehen sind dabei „<“, „>“ und „>>“, wobei diese Symbole unmittelbar dem jeweiligen Dateinamen vorangehen müssen, so dass sie mit ihm zusammen als ein Wort betrachtet werden. Der Programmtext innerhalb der *if*-Anweisung *if (child == 0)* wird vom neuen Prozess ausgeführt. Die *for*-Schleife geht dabei durch die Liste aller Wörter aus der Eingabezeile und fügt dabei die normalen Kommandozeilenargumente zur Liste *argv* hinzu, während die Wörter, die Umleitungen spezifizieren, sofort interpretiert werden.

Die Feinarbeit übernimmt hier die Funktion *fassign*, die die angegebene Datei mit den gewünschten Modi eröffnet und dafür den als ersten vorgegebenen Dateideskriptor verwendet. Hierbei kommt die Systemfunktion *dup2()* zum Einsatz, die *fd* zum Alias von *newfd* macht. Danach kann *newfd* geschlossen werden. Falls *fd* zuvor bereits auf eine geöffnete Datei verwies, wurde diese alte Verbindung zuerst gekappt.

Kapitel 2

Signale

2.1 Einführung

Signale werden für vielfältige Zwecke eingesetzt. Sie können verwendet werden,

- um den normalen Ablauf eines Prozesses für einen wichtigen Hinweis zu unterbrechen,
- um die Terminierung eines Prozesses zu erbitten oder zu erzwingen und
- um schwerwiegende Fehler bei der Ausführung zu behandeln wie z.B. den Verweis durch einen invaliden Zeiger.

Signale ersetzen keine Interprozeßkommunikation, da sie fast keine Informationen mit sich führen. In Abhängigkeit von der jeweiligen Systemumgebung gibt es mehr oder weniger fest definierte Signale, die über natürliche Zahlen identifiziert werden.

Der ISO-Standard 9899-1999 für die Programmiersprache C definiert eine einfache und damit recht portable Schnittstelle für die Behandlung von Signalen. Hier gibt es neben der Signalnummer selbst keine weiteren Informationen. Der IEEE Standard 1003.1 (POSIX) bietet eine Obermenge der Schnittstelle des ISO-Standards an, bei der wenige zusätzliche Informationen (wie z.B. die Angabe des invaliden Zeigers) dabei sein können und der insbesondere eine sehr viel feinere Kontrolle der Signalbehandlung erlaubt.

Signale können von verschiedenen Parteien ausgelöst werden: Von anderen Prozessen, die die dafür notwendige Berechtigung haben (entweder der gleiche Benutzer oder der Super-User), durch den Prozess selbst entweder indirekt (durch einen schwerwiegenden Fehler) oder explizit oder auch durch das Betriebssystem.

Typisch für letzteres ist die Terminalschnittstelle unter UNIX. Diese wurde ursprünglich für ASCII-Terminals mit serieller Schnittstelle entwickelt, die nur folgende Eingabemöglichkeiten anboten:

- Einzelne ASCII-Zeichen, jeweils ein Byte (zusammen mit etwas Extra-Kodierung wie Prüf- und Stop-Bits).
- Ein BREAK, das als spezielles Signal repräsentiert wird, das länger als die Kodierung für ein ASCII-Zeichen währt.
- Ein HANGUP, bei dem ein Signal wegfällt, das zuvor die Existenz der Leitung bestätigt hat. Dies benötigt einen weiteren Draht in der seriellen Leitung.

Diese Eingaben werden auf der Seite des Betriebssystems vom Terminal-Treiber bearbeitet, der in Abhängigkeit von den getroffenen Einstellungen

- die eingegebenen Zeichen puffert und das Editieren der Eingabe ermöglicht (beispielsweise mittels BACKSPACE, CTRL-u und CTRL-w) und
- bei besonderen Eingaben Signale an alle Prozesse schickt, die mit diesem Terminal verbunden sind.

Ziel war es, dass im Normalfall ein BREAK zu dem Abbruch oder zumindest der Unterbrechung der gerade laufenden Anwendung führt. Und ein HANGUP sollte zu dem Abbruch der gesamten Sitzung führen, da bei einem Wegfall der Leitung keine Möglichkeit eines regulären Abmeldens besteht.

Heute sind serielle Terminals rar geworden, aber das Konzept wurde dennoch beibehalten. Zwischen einem virtuellen Terminal (beispielsweise einem xterm) und den Prozessen, die zur zugehörigen Sitzung gehören, ist ein sogenanntes Pseudo-Terminal im Betriebssystem geschaltet, das der Sitzung die Verwendung eines klassischen Terminals vorspielt. Da es BREAK in diesem Umfeld nicht mehr gibt, wird es durch ein beliebiges Zeichen ersetzt wie beispielsweise CTRL-c. Und wenn das virtuelle Terminal wegfällt (z.B. durch eine gewaltsame Beendigung der xterm-Anwendung), dann gibt es weiterhin ein HANGUP für die Sitzung.

Auf fast alle Signale können Prozesse, die sie erhalten, auf dreierlei Weise reagieren:

- Voreinstellung: Terminierung des Prozesses.
- Ignorieren.
- Bearbeitung durch einen Signalbehandler.

Es mag harsch erscheinen, dass die Voreinstellung zur Terminierung eines Prozesses führt. Aber genau dies führt bei normalen Anwendungen genau zu den gewünschten Effekten wie Abbruch des laufenden Programms bei BREAK (die Shell ignoriert das Signal) und Abbau der Sitzung bei HANGUP. Wenn ein Prozess diese Signale ignoriert, sollte es genau wissen, was es tut, da der Nutzer auf diese Weise eine wichtige Kontrollmöglichkeit seiner Sitzung verliert. Sinnvoll ist es natürlich, eine Anwendung mit einem Signalbehandler zu ergänzen, wenn dadurch Datenverluste bei einer Terminierung vermieden werden können.

2.2 Signalbehandler

Der folgende Programmtext demonstriert die Behandlung des Signals *SIGINT*. Als Signalbehandler operiert die Funktion *signal_handler()*. Signalbehandler erhalten als Argument eine ganze Zahl, worüber sie die Nummer des Signals erhalten, das sie gerade bearbeiten. Einen Rückgabewert gibt es nicht.

Programm 2.1: Behandlung eines *SIGINT*-Signals (*sigint.c*)

```

1 #include <signal.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4
5 volatile sig_atomic_t signal_caught = 0;
6
7 void signal_handler(int signal) {
8     signal_caught = signal;
9 }
10
11 int main() {
```

```

12  if (signal(SIGINT, signal_handler) == SIG_ERR) {
13      perror("unable_to_setup_signal_handler_for_SIGINT");
14      exit(1);
15  }
16  printf("Try_to_send_a_SIGINT_signal!\n");
17  int counter = 0;
18  while (!signal_caught) {
19      for (int i = 0; i < counter; ++i);
20      ++counter;
21  }
22  printf("Got_signal_%d_after_%d_steps!\n", signal_caught, counter);
23  }

```

Der Signalbehandler in diesem Beispiel setzt nur eine globale Variable auf den Wert des erhaltenen Signals (Zeile 8). Die Verwendung des Typs `volatile sig_atomic_t` wird später diskutiert. `main()` richtet auf Zeile 12 die Funktion `signal_handler()` als Signalbehandler ein. Dies geschieht mit der Funktion `signal()`, die eine Signalnummer (hier: `SIGINT`) und eine Reaktion (entweder `SIG_DFL` für die Prozeßterminierung, `SIG_IGN` für das Ignorieren oder einen Funktionszeiger) als Parameter erwartet. Im Erfolgsfalle liefert `signal()` die frühere Einstellung zurück, ansonsten `SIG_ERR`.

Weiter geht es auf Zeile 18 mit einer Schleife, die erst dann abbricht, wenn sich der Wert von `signal_caught` ändert. Dies kann in diesem Beispiel nur durch den Signalbehandler passieren. Wenn dies geschieht, wird auf Zeile 22 die Nummer des eingetroffenen Signals ausgegeben und das Programm beendet. So könnte ein Aufruf dieses Programms aussehen, wenn das Versenden des Signals `SIGINT` durch den Terminaltreiber durch die Eingabe von CTRL-c recht schnell geschieht:

```

doolin$ sigint
Try to send a SIGINT signal!
^CGot signal 2 after 2074 steps!
doolin$

```

Die 2 ist dabei die Nummer des Signals `SIGINT`:

```

doolin$ grep SIGINT /usr/include/sys/iso/signal_iso.h
#define SIGINT 2 /* interrupt (rubout) */
doolin$

```

Leider sind mit der Bearbeitung von Signalen große Probleme verbunden. Wenn ein optimierender Übersetzer den obigen Programmtext analysiert, könnten folgende Punkte auffallen:

- Die Schleife in den Zeilen 18 bis 21 ruft keine externen Funktionen auf.
- Innerhalb der Schleife wird `signal_caught` nirgends verändert.

Daraus könnte vom Übersetzer der Schluß gezogen werden, dass die Schleifenbedingung nur zu Beginn einmal überprüft werden muss. Findet der Eintritt in die Schleife statt, könnte der weitere Test der Bedingung ersatzlos wegfallen. Analysen wie diese sind für heutige optimierende Übersetzer Pflicht, um guten Maschinen-Code erzeugen zu können. Es wäre also fatal, wenn darauf nur wegen der Existenz von asynchron aufgerufenen Signalbehandlern verzichtet werden würde.

Um beides zu haben, die fortgeschrittenen Optimierungstechniken und die Möglichkeit, Variablen innerhalb von Signalbehandlern setzen zu können, wurde in C die Speicherklasse `volatile` eingeführt. Damit lassen sich Variablen kennzeichnen, deren Wert sich jederzeit ändern kann — selbst dann, wenn dies aus dem vorliegenden Programmtext nicht

ersichtlich ist. Entsprechend gilt dann auch in C, dass alle anderen Variablen, die nicht als **volatile** klassifiziert sind, sich nicht durch „magische“ Effekte verändern dürfen.

Daraus folgt, dass korrekte Signalbehandler in ihren Möglichkeiten stark eingeschränkt sind. So ist es nur zulässig,

- lokale Variablen zu verwenden,
- mit **volatile** deklarierte Variablen zu benutzen und
- Funktionen aufzurufen, die sich an die gleichen Spielregeln halten.

Letzteres schließt insbesondere die Verwendung von Ein- und Ausgabe innerhalb eines Signalbehandlers aus. Der ISO-Standard 9899-1999 nennt nur *abort()*, *_Exit()*¹ und *signal()* als zulässige Bibliotheksfunktionen. Beim POSIX-Standard werden noch zahlreiche weitere Systemaufrufe genannt. Auf den Manualseiten von Solaris wird dies dokumentiert durch die Angabe „Async-Signal-Safe“ bei „MT-Level“.²

Wozu der Datentyp *sig_atomic_t*? Dieser Datentyp ist ein ganzzahliger Typ, bei dem Lese- und Schreiboperationen garantiert atomar ablaufen. Es handelt sich also um einen Datentyp, den der verwendete Prozessor mit einer einzigen ununterbrechbaren Instruktion laden oder speichern kann. Bei allen anderen Datentypen ist es nicht ausgeschlossen, dass mitten in einer Zuweisung ein asynchrones Signal eintreffen kann und der zugehörige Signalbehandler dann eine partiell modifizierte Variable vorfindet.

2.3 Wecksignale mit *alarm*

Zu den Signalen, den der POSIX-Standard definiert, gehört auch *SIGALRM*, das sich als Wecksignal verwenden lässt. Der Wecker wird mit *alarm()* unter Angabe einer relativen Weckzeit in Sekunden gestellt. Am Ende dieser Frist kommt es zum Eintreffen des Signals *SIGALRM*.

Programm 2.2: *read*-Operation mit Zeitlimit (*tread.c*)

```

1 /*
2  * Timed read operation. timed_read() works like read() but
3  * returns 0 if no input was received within the given number
4  * of seconds.
5  */
6
7 #include <signal.h>
8 #include <unistd.h>
9 #include "tread.h"
10
11 static volatile sig_atomic_t time_exceeded = 0;
12
13 static void alarm_handler(int signal) {
14     time_exceeded = 1;
15 }
16
17 int timed_read(int fd, void* buf, size_t nbytes, unsigned seconds) {
18     if (seconds == 0) return 0;
19     /*

```

¹*_Exit()* unterlässt im Vergleich zu *exit()* sämtliche Aufräumarbeiten.

²„MT“ steht hier für Multi-Threading, da dabei ähnliche Probleme auftreten, wenngleich in einem noch größeren Umfang.

```

20  * setup signal handler and alarm clock but
21  * remember the previous settings
22  */
23  void (*previous_handler)(int) = signal(SIGALRM, alarm_handler);
24  if (previous_handler == SIG_ERR) return -1;
25  time_exceeded = 0;
26  int remaining_seconds = alarm(seconds);
27  if (remaining_seconds > 0) {
28      if (remaining_seconds <= seconds) {
29          remaining_seconds = 1;
30      } else {
31          remaining_seconds -= seconds;
32      }
33  }
34
35  int bytes_read = read(fd, buf, nbytes);
36
37  /* restore previous settings */
38  if (!time_exceeded) alarm(0);
39  signal(SIGALRM, previous_handler);
40  if (remaining_seconds) alarm(remaining_seconds);
41
42  if (time_exceeded) return 0;
43  return bytes_read;
44  }

```

Der vorstehende Programmtext zeigt, wie *alarm()* verwendet werden kann, um eine Operation zeitlich zu befristen. Die Funktion *timed_read()* arbeitet genauso wie *read()*, wobei jedoch die Wartezeit auf die gegebene Zahl von Sekunden begrenzt wird. Wenn die Zeit verstreicht, ohne dass eine Eingabe erfolgte, wird 0 zurückgeliefert.

Was geschieht, wenn mehrere Wecksignale nebeneinander eingerichtet werden? Da *alarm()* nur die Verwaltung einer einzigen Weckzeit unterstützt, ist eine kooperative Vorgehensweise notwendig. Dies wird dadurch erleichtert, indem *alarm()* bei der Einrichtung einer neuen Weckzeit entweder 0 zurückgibt (vorher war keine Weckzeit aktiv) oder eine positive Zahl von Sekunden als Restzeit zum vorher eingestellten Wecksignal. Entsprechend wird auf Zeile 26 in der Variablen *remaining_seconds* die alternative Weckzeit notiert und später auf Zeile 40 findet eine Wiedereinsetzung statt, nachdem zuvor auf Zeile 31 die verbleibende Restzeit neu berechnet worden ist. Analog wird in Zeile 23 der frühere Signalbehandler für *SIGALRM* in der Variablen *previous_handler* gesichert, so dass er in Zeile 39 wieder restauriert werden kann.

Bei der Restaurierung ist es wichtig, dass kein Fenster offenbleibt, das zum Verlust einer Signalbehandlung führt oder den alten Signalbehandler auf das noch nicht eingetretene Signal für das von *tread()* eingerichtete Zeitlimit reagieren lässt. Wenn der alte Signalbehandler *SIG_DFL* war, also die Voreinstellung, dann würde das sogar unerwartet zur Terminierung unseres Prozesses führen. Deswegen wird in Zeile 38 der Wecker zuerst deaktiviert, wenn er bislang sein Signal noch nicht von sich gegeben hat. Danach kann in Zeile 39 der alte Signalbehandler eingesetzt werden, ohne dass wir Gefahr laufen, dass er sofort verwendet wird. In Zeile 40 wird dann der Wecker neu aufgesetzt, falls er vor dem Aufruf von *tread()* eingeschaltet war.

Es bleibt hier noch anzumerken, dass es noch bessere Wege gibt, Leseoperationen mit Zeitbeschränkungen zu versehen. Es empfiehlt sich hier entweder die Verwendung der Systemaufrufe *poll()* oder *select()* oder der Einsatz asynchroner Lesetechniken. Dennoch ist die Verwendung von *alarm()* sinnvoll, da es genügend Operationen gibt, bei denen es

keine Variante mit Zeitbeschränkung gibt.

2.4 Das Versenden von Signalen

Der ISO-Standard für C sieht nur eine Funktion *raise()* vor, die es erlaubt, ein Signal an den eigenen Prozess zu versenden. Im POSIX-Standard kommt die Funktion *kill()* hinzu, die es erlaubt, ein Signal an einen anderen Prozess zu verschicken, sofern die dafür notwendigen Privilegien vorliegen. Das folgende Programm zeigt ein Beispiel, bei dem ein neu erzeugter Prozess ein Signal an den Erzeuger sendet.

Programm 2.3: Versenden eines Signals an den übergeordneten Prozess (*killparent.c*)

```

1 #include <signal.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <unistd.h>
5 #include <sys/wait.h>
6
7 void sigterm_handler(int signo) {
8     const char msg[] = "Goodbye,_cruel_world!\n";
9     write(1, msg, sizeof msg - 1);
10    _Exit(1);
11 }
12
13 int main() {
14     if (signal(SIGTERM, sigterm_handler) == SIG_ERR) {
15         perror("signal"); exit(1);
16     }
17
18     pid_t child = fork();
19     if (child == 0) {
20         kill(getppid(), SIGTERM);
21         exit(0);
22     }
23     int wstat;
24     wait(&wstat);
25     exit(0);
26 }

```

Trotz ihres geringen Informationsgehalts dienen Signale gelegentlich zur Kommunikation. So gibt es eine weitverbreitete Konvention, dass bei langfristig laufenden Diensten *SIGHUP* das erneute Einlesen der Konfiguration veranlasst und *SIGTERM* eine geordnete Terminierung einleiten soll. Gelegentlich sind für Dienste auch Reaktionen für *SIGUSR1* und *SIGUSR2* definiert. So wird der Apache-Webserver beispielsweise bei *SIGUSR1* veranlasst, bei nächster Gelegenheit auf sanfte Weise neu zu starten. Anders als bei *SIGHUP* kommt es dann nicht zur Unterbrechung aktueller Netzwerkverbindungen.

Der Systemaufruf *kill()* erfüllt aber auch noch einen weiteren Zweck. Bei einer Signalnummer von 0 wird nur die Zulässigkeit des Signalversendens überprüft. Das folgende Beispiel demonstriert, wie dies dazu verwendet werden kann, um die Existenz eines Prozesses zu überprüfen:

Programm 2.4: Verwendung von *kill* zur Überprüfung der Existenz eines Prozesses (*waitfor.c*)

```

1 #include <errno.h>
2 #include <signal.h>
3 #include <stdio.h>
4 #include <stdlib.h>
5 #include <unistd.h>
6
7 int main(int argc, char** argv) {
8     char* cmdname = *argv++; --argc;
9     if (argc != 1) {
10         fprintf(stderr, "Usage:_%s_pid\n", cmdname);
11         exit(1);
12     }
13
14     /* convert first argument to pid */
15     char* endptr = argv[0];
16     pid_t pid = strtol(argv[0], &endptr, 10);
17     if (endptr == argv[0]) {
18         fprintf(stderr, "%s:_integer_expected_as_argument\n",
19             cmdname);
20         exit(1);
21     }
22
23     while (kill(pid, 0) == 0) sleep(1);
24
25     if (errno == ESRCH) exit(0);
26     perror(cmdname); exit(1);
27 }

```

Gelegentlich kommt es vor, dass Prozesse nur auf das Eintreffen eines Signals warten möchten und sonst nichts zu tun haben. Theoretisch könnte ein Prozess dann in eine Dauerschleife mit leerem Inhalt treten (auch *busy loop* bezeichnet), aber dies wäre nicht sehr fair auf einem System mit mehreren Prozessen, da dadurch Rechenzeit vergeudet würde. Abhilfe schafft hier der Systemaufruf *pause()*, der einen Prozess schlafen legt, bis ein Signal eintrifft.

Der folgende Programmtext demonstriert das Warten auf ein Signal anhand eines virtuellen Ballspiels zweier Prozesse. Das Programm beginnt in Zeile 37 mit der Erzeugung eines weiteren Prozesses, um einen Spielpartner zu haben. Beide Prozesse rufen anschließend *playwith()* auf, um das Spiel durchzuführen. Der neu erzeugte Prozess hat zuerst den virtuellen Ball und darf mit dem Spiel beginnen.

Programm 2.5: Virtuelles Ballspiel zweier Prozesse (*pingpong.c*)

```

1 #include <signal.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <unistd.h>
5
6 static volatile sig_atomic_t sigcount = 0;
7
8 void sighandler(int sig) {
9     ++sigcount;
10    if (signal(sig, sighandler) == SIG_ERR) _Exit(1);
11 }
12

```

```

13 static void playwith(pid_t partner, int start) {
14     if (signal(SIGUSR1, sighandler) == SIG_ERR) {
15         perror("signal_SIGUSR1"); exit(1);
16     }
17     /* give our partner some time for preparation */
18     if (start) sleep(1);
19     /* start the ping pong game */
20     if (start) sigcount = 1;
21     for(int i = 0; i < 10; ++i) {
22         if (!sigcount) pause();
23         printf("[%d]_send_signal_to_%d\n",
24             (int) getpid(), (int) partner);
25         if (kill(partner, SIGUSR1) < 0) {
26             printf("[%d]_%d_is_no_longer_alive\n",
27                 (int) getpid(), (int) partner);
28             return;
29         }
30         --sigcount;
31     }
32     printf("[%d]_finishes_playing\n", (int) getpid());
33 }
34
35 int main() {
36     pid_t parent = getpid();
37     pid_t child = fork();
38
39     if (child < 0) {
40         perror("fork"); exit(1);
41     }
42     if (child == 0) {
43         playwith(parent, 1);
44     } else {
45         playwith(child, 0);
46     }
47 }

```

Der Besitzer des virtuellen Balles sendet in Zeile 25 den Ball mit Hilfe des Signals *SIGUSR1* an den Partner und legt sich anschließend in Zeile 22 mittels *pause()* schlafen, um auf das Wiedereintreffen des Balls zu warten, das wiederum durch *SIGUSR1* signalisiert wird.

Die Variable *sigcount* in Zeile 6 repräsentiert die Zahl der Bälle, die sich im Augenblick im Besitz des Prozesses sind. Diese Zahl wird von dem Signalbehandler *sighandler()* in Zeile 9 hochgezählt, wenn ein *SIGUSR1*-Signal eintrifft und in Zeile 30 wieder heruntergezählt, nachdem das *SIGUSR1*-Signal an den Partner verschickt wurde.

Zu Beginn setzen beide Spieler in Zeile 14 *sighandler()* als Signalbehandler ein. Derjenige, der mit dem Spiel beginnt, wartet dann in Zeile 18 noch eine Sekunde, um zu vermeiden, dass ein Signal geschickt wird, bevor der Spielpartner die Gelegenheit hatte, seinen Signalbehandler aufzusetzen.

Zu beachten ist dabei, dass der Signalbehandler in Zeile 10 sich selbst wieder einsetzt, da der POSIX-Standard nicht festlegt, ob diese Einstellung nach Eintreffen eines Signals erhalten bleibt. *signal()* gehört mit zu den vom POSIX-Standard genannten Funktionen, die auch innerhalb eines Signalbehandlers aufgerufen werden dürfen.

2.5 Die Zustellung von Signalen

Die vorangegangenen Beispiele werfen die Frage auf, wie UNIX bei der Zustellung von Signalen vorgeht, wenn

- der Prozess zur Zeit nicht aktiv ist,
- gerade ein Systemaufruf für den Prozess abgearbeitet wird oder
- gerade ein Signalbehandler bereits aktiv ist.

Vom ISO-Standard 9899-1999 für C wird in dieser Beziehung nichts festgelegt. Der POSIX-Standard geht jedoch genauer darauf ein.³

- Wenn ein Prozess ein Signal erhält, wird dieses Signal zunächst in den zugehörigen Verwaltungsstrukturen des Betriebssystems vermerkt. Signale, die für einen Prozess vermerkt sind, jedoch noch nicht zugestellt worden sind, werden als *anhängige* Signale bezeichnet.
- Wenn mehrere Signale mit der gleichen Nummer anhängig sind, ist nicht festgelegt, ob eine Mehrfachzustellung erfolgt. Es können also Signale wegfallen.
- Nur aktiv laufende Prozesse können Signale empfangen. Prozesse werden normalerweise durch die Existenz eines anhängigen Signals aktiv — aber dieses kann auch längere Zeit in Anspruch nehmen, wenn dem zwischenzeitlich mangelnde Ressourcen entgegenstehen.
- Für jeden Prozess gibt es eine Menge blockierter Signale, die im Augenblick nicht zugestellt werden sollen. Dies hat nichts mit dem Ignorieren von Signalen zu tun, da blockierte Signale anhängig bleiben, bis die Blockierung aufgehoben wird.
- Der POSIX-Standard legt nicht fest, was mit der Signalbehandlung geschieht, wenn ein Signalbehandler aufgerufen wird. Möglich ist das Zurückfallen auf *SIG_DFL* (Voreinstellung mit Prozeßterminierung) oder die temporäre automatische Blockierung des Signals bis zur Beendigung des Signalbehandlers. Alle modernen UNIX-Systeme wählen die zweite Variante. Dies lässt sich aber gemäß dem POSIX-Standard auch erzwingen, indem die umfangreichere Schnittstelle *sigaction()* anstelle von *signal()* verwendet wird. Allerdings ist *sigaction()* nicht mehr Bestandteil des ISO-Standards für C.
- UNIX unterscheidet zwischen unterbrechbaren und unterbrechungsfreien Systemaufrufen. Zur ersteren Kategorie gehören weitgehend alle Systemaufrufe, die zu einer längeren Blockierung eines Prozesses führen können. Ist ein nicht blockiertes Signal anhängig, kann ein unterbrechbarer Systemaufruf aufgrund des Signals mit einer Fehlerindikation beendet werden. *errno* wird dann auf *EINTR* gesetzt. Dabei ist zu beachten, dass der unterbrochene Systemaufruf nach Beendigung der Signalbehandlung *nicht* fortgesetzt wird, sondern manuell erneut gestartet werden muss. Dies kann leider zu unerwarteten Überraschungseffekten führen, weil insbesondere auch die *stdio*-Bibliothek keinerlei Vorkehrungen trifft, Systemaufrufe automatisch erneut aufzusetzen, falls es zu einer Unterbrechung kam. Dies ist eine wesentliche Schwäche sowohl des POSIX-Standards als auch der *stdio*-Bibliothek und ein Grund mehr dafür, auf die Verwendung der *stdio* in kritischen Anwendungen völlig zu verzichten.

³Siehe „Signal Concepts“, im Web unter http://www.opengroup.org/onlinepubs/009695399/functions/xsh_chap02_04.html

Datentyp	Feldname	Beschreibung
void(*) (int)	<i>sa_handler</i>	Funktionszeiger (wie bisher)
void(*) (int, <i>siginfo_t*</i> , void*)	<i>sa_sigaction</i>	alternativer Zeiger auf einen Signalbehandler, der mehr Informationen zum Signal erhält
<i>sigset_t</i>	<i>sa_mask</i>	Menge von Signalen, die während der Signalbehandlung dieses Signals zu blockieren sind
int	<i>sa_flags</i>	Menge von Boolean-wertigen Optionen

Tabelle 2.1: Felder der **struct** *sigaction*

Für die genauere Regulierung der Signalbehandlung bietet POSIX (jedoch nicht ISO-C) den Systemaufruf *sigaction()* an. Während bei *signal()* zur Spezifikation der Signalbehandlung nur ein Funktionszeiger genügte, kommen bei der **struct** *sigaction*, die *sigaction()* verwendet, die in Tabelle 2.1 genannten Felder zum Einsatz.

Ein wesentlicher Unterschied zwischen *sigaction()* und *signal()* besteht bereits darin, dass per Voreinstellung das Signal, das eine Signalbehandlung auslöst, während der Bearbeitung automatisch blockiert wird. Ferner findet (solange nichts Gegenteiliges in den Optionen angegeben wurde) keine implizite Veränderung des Signalhandlers auf *SIG_DFL* nach der Signalbehandlung statt. Bei *signal()* ist dies ebenfalls möglich, jedoch nicht garantiert.

Das folgende Beispiel demonstriert den möglichen Verlust von Signalen trotz umfangreicher Vorsichtsmaßnahmen. Das Experiment besteht hier im Versenden von 1000 *SIGUSR1*-Signalen, die beim Empfänger nachgezählt werden. In den Zeilen 26 bis 30 setzt der neu erzeugte Prozess die Funktion *count_signals()* als Signalbehandler für *SIGUSR1* auf. Dabei wird hier *sigaction()* anstelle von *signal()* verwendet. Dies garantiert uns, dass während der Bearbeitung des Signals *SIGUSR1* weitere eintreffende *SIGUSR1*-Signale aufgeschoben und nicht sofort in verschachtelter Form bearbeitet werden.

Programm 2.6: Verlust von Signalen (*sigfire.c*)

```

1 #include <signal.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <unistd.h>
5
6 static const int NOF_SIGNALS = 1000;
7 static volatile sig_atomic_t received_signals = 0;
8 static volatile sig_atomic_t terminated = 0;
9
10 static void count_signals(int sig) {
11     ++received_signals;
12 }
13
14 void termination_handler(int sig) {
15     terminated = 1;
16 }
17
18 int main() {
19     sighold(SIGUSR1); sighold(SIGTERM);
20

```

```

21  pid_t child = fork();
22  if (child < 0) {
23      perror("fork"); exit(1);
24  }
25  if (child == 0) {
26      struct sigaction action = {0};
27      action.sa_handler = count_signals;
28      if (sigaction(SIGUSR1, &action, 0) != 0) {
29          perror("sigaction"); exit(1);
30      }
31      action.sa_handler = termination_handler;
32      if (sigaction(SIGTERM, &action, 0) != 0) {
33          perror("sigaction"); exit(1);
34      }
35      sigrelse(SIGUSR1); sigrelse(SIGTERM);
36      while (!terminated) pause();
37      printf("[%d]_received_%d_signals\n",
38             (int) getpid(), received_signals);
39      exit(0);
40  }
41
42  sigrelse(SIGUSR1); sigrelse(SIGTERM);
43  for (int i = 0; i < NOF_SIGNALS; ++i) {
44      kill(child, SIGUSR1);
45  }
46  printf("[%d]_sent_%d_signals\n",
47         (int) getpid(), NOF_SIGNALS);
48  kill(child, SIGTERM); wait(0);
49  }

```

Auf diese Weise ist die Atomizität des Hochzählens der Variable *received_signals* garantiert. Zu bedenken ist dabei, dass der Datentyp *sig_atomic_t* selbst nur die Atomizität einer einzelnen Lese- oder Schreiboperation gewährleistet. Bei einem Inkrement liegt jedoch eine Lese- und eine Schreib-Operation vor. Zwischen diesen Operationen wäre eine Unterbrechung wegen einer Signalbehandlung denkbar.

In diesem Beispiel versendet der übergeordnete Prozess 1000 Signale und der neu erzeugte Prozess zählt die eingetroffenen Signale. Wann darf der übergeordnete Prozess davon ausgehen, dass der neu erzeugte Prozess seinen Signalbehandler fertig aufgesetzt hat, so dass er mit dem Zählen beginnen kann? Wenn der übergeordnete Prozess zu früh Signale versendet, während noch die voreingestellte Reaktion für *SIGUSR1* eingerichtet ist, würde dies nur zur vorzeitigen Terminierung des erzeugten Prozesses führen. Diese Problematik lässt sich vermeiden, indem die Signale, für die der erzeugte Prozess Behandler aufsetzt, vor der Prozeßerzeugung blockiert werden. Das geht am einfachsten mit der Funktion *sighold()* (auf Zeile 19), die zum POSIX-Standard gehört. Zu jedem Prozess unterhält das Betriebssystem eine Menge blockierter Signale. Mit *sighold()* wird das angegebene Signal zu der Menge hinzugefügt. Entscheidend ist hier, dass die Menge der blockierten Signale an den neu erzeugten Prozess vererbt wird. So können nach dem *fork()* in aller Ruhe auf den Zeilen 26 bis 34 die Signalbehandler für *SIGUSR1* und *SIGTERM* aufgesetzt werden, bevor auf Zeile 35 diese Signale wieder mit Hilfe von *sigrelse()* (eine unglückliche Abkürzung von *signal release*) wieder aus der Menge der blockierten Signale entfernt werden. Auch der übergeordnete Prozess nimmt die beiden Signale wieder heraus auf der Zeile 42.

Wie wird das Experiment beendet? Da, wie das Experiment zeigen soll, Signale verloren gehen können, sollte der erzeugte Prozess nicht darauf warten, bis *NOF_SIGNALS* eingetroffen sind. Stattdessen wird *SIGTERM* vom übergeordneten Prozess an den erzeugten Prozess verwendet, um das Ende des Experiments zu signalisieren.

Hier sind einige Läufe des Experiments, die demonstrieren, wie sehr die Werte voneinander abweichen können:

```
doolin$ sigfire
[22073] sent 1000 signals
[22074] received 264 signals
doolin$ sigfire
[22075] sent 1000 signals
[22076] received 227 signals
doolin$ sigfire
[22077] sent 1000 signals
[22078] received 481 signals
doolin$ sigfire
[22079] sent 1000 signals
[22080] received 136 signals
doolin$
```

Wenn anstelle von nur *SIGUSR1* zwei Signale *SIGUSR1* und *SIGUSR2* abwechselnd verwendet werden, können höhere Werte erzielt werden, wobei der Erfolg keinesfalls garantiert ist:

```
doolin$ sigfire2
[22142] sent 1000 signals
[22143] received 495 signals
doolin$ sigfire2
[22144] sent 1000 signals
[22145] received 462 signals
doolin$ sigfire2
[22146] sent 1000 signals
[22147] received 468 signals
doolin$ sigfire2
[22151] sent 1000 signals
[22152] received 688 signals
doolin$
```

Dieses Experiment untermauert die Regel, dass einzelne Signale zuverlässig zugestellt werden, während es bei dem Mehrfachen Eintreffen des gleichen Signals zu Verlusten kommen kann. In der Praxis zeigen sich jedoch Signalverluste nur bei härteren Rahmenbedingungen, sei es durch ein explizites Dauerfeuer (wie in diesem Experiment) oder durch eine hohe Belastung der Maschine.

2.6 Signale als Indikatoren für terminierte Prozesse

Mit Hilfe der Funktionen *wait()* oder *waitpid()* wird die Terminierung erzeugter Prozesse *synchron* abgewickelt. Gelegentlich ist es auch sinnvoll, sich die Terminierung über Signale *asynchron* mitteilen zu lassen. Dies geht mit dem Signal *SIGCHLD*, das an den Erzeuger versendet wird, sobald eine der von ihm erzeugten Prozesse terminiert. Per Voreinstellung wird dieses Signal ignoriert.

Programm 2.7: Demonstration von *SIGCHLD* (*sigchld.c*)

```
1 #include <signal.h>
```

```

2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <unistd.h>
5 #include <sys/wait.h>
6 #include "processlist.h"
7
8 static processlist alive, dead;
9
10 void child_term_handler(int sig) {
11     pid_t pid; int wstat;
12     while ((pid = waitpid((pid_t)-1, &wstat, WNOHANG)) > 0) {
13         if (pl_move(&alive, &dead, pid)) {
14             pl_modify(&dead, pid, wstat);
15         }
16     }
17 }
18
19 int main() {
20     struct sigaction action = {0};
21     action.sa_handler = child_term_handler;
22     if (sigaction(SIGCHLD, &action, 0) != 0) {
23         perror("sigaction");
24     }
25     pl_alloc(&alive, 4); pl_alloc(&dead, 4);
26     sighold(SIGCHLD);
27     for (int i = 0; i < 10; ++i) {
28         fflush(0);
29         pid_t child = fork();
30         if (child < 0) {
31             perror("fork"); exit(1);
32         }
33         if (child == 0) {
34             srand(getpid()); sleep(rand() % 5);
35             exit((char) rand());
36         }
37         pl_add(&alive, child, 0);
38     }
39     sigrelse(SIGCHLD);
40     while (pl_length(&alive) > 0) {
41         if (pl_length(&dead) == 0) pause();
42         while (pl_length(&dead) > 0) {
43             int wstat;
44             sighold(SIGCHLD);
45             pid_t pid = pl_pick(&dead, &wstat);
46             sigrelse(SIGCHLD);
47             printf("[%d]_ %d\n", (int) pid, WEXITSTATUS(wstat));
48         }
49     }
50 }

```

Der vorstehende Programmtext demonstriert die Verwendung von `SIGCHLD`. In den Zeilen 27 bis 38 werden zehn Prozesse erzeugt, die nach einem zufällig bestimmten Zeitintervall mit einem zufälligen `exit`-Wert enden. Die Prozesse werden in der Zeile 37 in die

Liste der noch lebenden Prozesse *alive* eingetragen. Der folgende Programmtext zeigt die Schnittstelle der Prozeßverwaltung:

Programm 2.8: Schnittstelle der Bibliothek für Listen von Prozessen (*processlist.h*)

```

1  #ifndef PROCESSLIST_H
2  #define PROCESSLIST_H
3
4  typedef struct process {
5      pid_t pid; int wstat;
6      struct process* next;
7  } process;
8  typedef struct processlist {
9      unsigned int size, length;
10     process** bucket; /* hash table */
11     unsigned int it_index;
12     process* it_entry;
13 } processlist;
14
15 // All functions with the exception of pl_length, pl_next,
16 // and pl_pick return 1 on success, 0 in case of failures.
17
18 /* allocate a hash table for processes with the given bucket size */
19 int pl_alloc(processlist* pl, unsigned int size);
20
21 /* add tuple (pid,wstat) to the process list, pid must be unique */
22 int pl_add(processlist* pl, pid_t pid, int wstat);
23
24 /* modify wstat for a given pid */
25 int pl_modify(processlist* pl, pid_t pid, int wstat);
26
27 /* delete tuple by pid */
28 int pl_remove(processlist* pl, pid_t pid);
29
30 /* move entry for pid to another list */
31 int pl_move(processlist* from, processlist* to, pid_t pid);
32
33 /* return number of elements */
34 unsigned int pl_length(processlist* pl);
35
36 /* lookup wstat by pid */
37 int pl_lookup(processlist* pl, pid_t pid, int* wstat);
38
39 /* start iterator */
40 int pl_start(processlist *pl);
41
42 /* fetch next pid from iterator; returns 0 on end */
43 pid_t pl_next(processlist *pl);
44
45 /* pick and remove one element out of the list */
46 pid_t pl_pick(processlist *pl, int* wstat);
47
48 /* free allocated memory */
49 int pl_free(processlist* pl);

```

50 #endif

Der Signalbehandler *child_term_handler* auf den Zeilen 10 bis 17 verzeichnet alle Eingänge des *SIGCHLD*-Signals, indem mit Hilfe von *waitpid()* der Status abgeholt wird und der terminierte Prozess von der Liste *alive* nach *dead* verschoben wird. Das Problem ist, dass bei einer zeitgleichen Terminierung mehrerer Prozesse es wiederum zu Verlusten des *SIGCHLD*-Signals kommen kann. Somit können wir uns nicht darauf verlassen, dass für jeden terminierten Prozess der Signalbehandler genau einmal aufgerufen wird. Deswegen empfiehlt es sich, den Status aller bereits terminierter Prozesse abzurufen. Dies wird mit *waitpid()* unter der Verwendung der Option *WNOHANG* erreicht. Diese verhindert ein Blockieren des Systemaufrufs *waitpid()* und somit kann *waitpid()* gefahrlos solange aufgerufen werden, bis *waitpid()* 0 oder einen negativen Wert zurückliefert.

Um die Prozesse und deren Status zu erfassen, sind umfangreichere Datenstrukturen erforderlich, die nicht mehr ohne weiteres mit dem Attribut **volatile** versehen werden können. Außerdem sollte es nicht dazu kommen, dass das Entfernen eines Eintrages aus der Liste der toten Prozesse in Zeile 45 von einer Signalbehandlung unterbrochen wird, die in Zeile 13 in genau die gleiche Liste, einen Eintrag hinzuzufügen versucht. Nur ein zuverlässiger gegenseitiger Ausschluß bewahrt uns hier vor Überraschungen.

Dies wird erreicht, indem gleich zu Beginn in Zeile 26 das Signal *SIGCHLD* blockiert wird. Das ist schon aus dem Grund wichtig, dass die Liste der noch lebenden Prozesse gefüllt werden kann, bevor diese Prozesse die Gelegenheit haben, sich allzu frühzeitig zu verabschieden. Nachdem alle zehn Prozesse erzeugt worden sind, wird in Zeile 39 die Blockierung wieder aufgehoben. Später, wenn die Liste *dead* modifiziert werden soll, wird die kritische Operation mit Hilfe von *sighold()* und *sigrelse()* wiederum geschützt.

Wenn auf gemeinsame Datenstrukturen von mehreren Seiten in asynchroner Form zugegriffen werden kann, dann sind die zugreifenden Programmbereiche sogenannte *kritische Regionen*. Nur durch den gegenseitigen Ausschluß wird verhindert, dass die betroffene Datenstruktur durch einen ungünstigen Unterbrechungszeitpunkt inkonsistent wird. Da *sigaction()* anstelle von *signal()* verwendet wird, ist bereits sichergestellt, dass *child_term_handler()* nicht mehrfach in verschachtelter Form aufgerufen wird. Somit müssen nur alle verbliebenen Programmbereiche, die auf die gleiche Datenstruktur zugreifen, in *sighold()* und *sigrelse()* geklammert werden.

2.7 Signalbehandlung in einer Shell

Die im vorherigen Kapitel vorgestellte einfache Shell kümmerte sich nicht um die Signalbehandlung. Dies kann zu überraschenden Effekten führen, wenn der Versuch unternommen wird, aufgerufene Prozesse beispielsweise mit *SIGINT* zu unterbrechen:

```
doolin$ tinysh
% cat >OUT
Some input...
^Cdoolin$
```

Hier wurde zunächst *cat* aufgerufen und nach der ersten eingegebenen Zeile CTRL-c eingegeben, welches bei den aktuellen Einstellungen zu einem *SIGINT* an alle Prozesse der aktuellen Sitzung führte. Zur aktuellen Sitzung gehört jedoch nicht nur das gerade laufende *cat*-Kommando, sondern natürlich auch *tinysh*. Da *tinysh* keinerlei Vorkehrungen traf, wurde es genauso wie *cat* einfach terminiert, weil dies die voreingestellte Reaktion ist. Wäre *tinysh* die Login-Shell gewesen, wäre damit die gesamte Sitzung beendet. Hier in diesem Beispiel wurde *SIGINT* offensichtlich von der normalen Shell ignoriert.

Wie muss also die Signalbehandlung einer Shell aussehen?

- Wenn ein Kommando *im Vordergrund* läuft, muss die Shell die Signale *SIGINT* und *SIGQUIT* ignorieren.
- Wenn ein Kommando **im Hintergrund** läuft, müssen für diesen Prozess *SIGINT* und *SIGQUIT* ignoriert werden.
- Wenn die Shell ein Kommando einliest, sollten *SIGINT* und *SIGQUIT* die Neu-Eingabe des Kommandos ermöglichen.
- Bezüglich *SIGHUP* muss nichts unternommen werden.

Der folgende Programmtext zeigt das Hauptprogramm einer einfachen Shell mit Signalbehandlung. Das zeilenweise Einlesen wurde in die Funktion *getline()* ausgelagert, die Unterbrechungen berücksichtigt. Bei jeder sich bietenden Gelegenheit werden nicht-blockierend mit *waitpid()* die Statusinformationen der inzwischen beendeten Prozesse abgeholt. Das Parsieren einer Kommandozeile wurde ausgelagert, wobei als weiteres unterstütztes Token noch "&" hinzugekommen ist, das für die Ausführung im Hintergrund steht.

Programm 2.9: Einfache Shell mit Signalbehandlung (*tinys2.c*)

```

1 /*
2
3 =head1 NAME
4
5 tinys2 -- a tiny shell with a minimal set of features
6
7 =head1 SYNOPSIS
8
9 B<tinys2>
10
11 =head1 DESCRIPTION
12
13 B<tinys2> reads command lines from its standard input
14 and executes them. tinys2> " is given as prompt.
15
16 Each command line consists of space-separated tokens. Special tokens
17 begin with <" , indicating the input file, >" , indicating the output
18 file, >>" , specifying an output file that is to be extended, and
19 &" , indicating a background execution. The first non-special token
20 specifies the command name which must be either an absolute path of an
21 executable file or locatable within the list of directories provided by
22 the environment variable B<PATH>.
23
24 =head1 DIAGNOSTICS
25
26 B<tinys2> prints the exit code of a terminated program if it is non-zero.
27 An exit code of 255 is used by subprocesses if they are unable to start
28 the command for some reason.
29
30 =head1 BUGS
31
32 No pipelines, no builtins, no shell variables.
33
34 =head1 AUTHOR

```



```
35
36 Andreas Borchert
37
38 =cut
39
40 */
41
42 #include <errno.h>
43 #include <signal.h>
44 #include <stdio.h>
45 #include <stdlib.h>
46 #include <unistd.h>
47 #include <sys/wait.h>
48 #include "command.h"
49 #include "sareadline.h"
50 #include "strlist.h"
51 #include "tokenizer.h"
52
53 static sig_atomic_t interrupted = 0;
54
55 void interrupt_handler(int sig) {
56     interrupted = 1;
57 }
58
59 void print_child_status(pid_t pid, int wstat) {
60     printf("[%d]_", (int) pid);
61     if (WIFEXITED(wstat)) {
62         printf("exit_%d", WEXITSTATUS(wstat));
63     } else if (WIFSIGNALED(wstat)) {
64         printf("terminated_with_signal_%d", WTERMSIG(wstat));
65         if (WCOREDUMP(wstat)) printf("_ (core_dump)");
66     } else {
67         printf("???");
68     }
69     printf("\n");
70 }
71
72 void status_report(void) {
73     pid_t pid; int wstat;
74     while ((pid = waitpid((pid_t)-1, &wstat, WNOHANG)) > 0) {
75         print_child_status(pid, wstat);
76     }
77 }
78
79 int getline(stralloc* line) {
80     int first = 1;
81     interrupted = 0;
82     for(;;) {
83         if (interrupted) {
84             interrupted = 0;
85             printf("\n");
86             first = 1;
87         }
```

```

88     if (first) {
89         status_report();
90         printf("%%_%");
91         first = 0;
92     }
93     errno = 0;
94     if (readline(stdin, line)) return 1;
95     if (errno != EINTR) return 0;
96 }
97 }
98
99 int main() {
100     struct sigaction action = {0};
101     action.sa_handler = interrupt_handler;
102     if (sigaction(SIGINT, &action, 0) != 0 ||
103         sigaction(SIGQUIT, &action, 0) != 0) {
104         perror("sigaction");
105     }
106
107     stralloc line = {0};
108     while (getline(&line) {
109         strlist tokens = {0};
110         stralloc_0(&line); /* required by tokenizer() */
111         if (!tokenizer(&line, &tokens)) break;
112         if (tokens.len == 0) continue;
113         command cmd = {0};
114         if (!scan_command(&tokens, &cmd)) continue;
115
116         sighold(SIGINT); sighold(SIGQUIT);
117         pid_t child = fork();
118         if (child == -1) {
119             perror("fork"); continue;
120         }
121         if (child == 0) {
122             sigrelse(SIGINT); sigrelse(SIGQUIT);
123             if (cmd.background) {
124                 sigignore(SIGINT); sigignore(SIGQUIT);
125             }
126             exec_command(&cmd);
127             perror(cmd.cmdname);
128             exit(255);
129         }
130
131         if (cmd.background) {
132             printf("%d\n", (int)child);
133         } else {
134             int wstat;
135             pid_t pid = waitpid(child, &wstat, 0);
136             if (!WIFEXITED(wstat) || WEXITSTATUS(wstat)) {
137                 print_child_status(pid, wstat);
138             }
139         }
140         sigrelse(SIGINT); sigrelse(SIGQUIT);

```

```

141     }
142 }

```

Programm 2.10: Schnittstelle für die Behandlung von Kommandos (*command.h*)

```

1  #ifndef COMMAND_H
2  #define COMMAND_H
3
4  #include <fcntl.h>
5  #include "strlist.h"
6
7  typedef struct fd_assignment {
8      char* path;
9      int oflags;
10     mode_t mode;
11 } fd_assignment;
12
13 typedef struct command {
14     char* cmdname;
15     strlist argv;
16     int background;
17     /* for file descriptors 0 and 1 */
18     fd_assignment assignments[2];
19 } command;
20
21 /* convert list of tokens into a command record */
22 int scan_command(strlist* tokens, command* cmd);
23
24 /*
25  * open input and output files, if required, and
26  * exec to the given command
27  */
28 void exec_command(command* cmd);
29
30 #endif

```

Programm 2.11: Funktionen für die Behandlung von Kommandos (*command.c*)

```

1  #include "command.h"
2  #include <stdio.h>
3  #include <stdlib.h>
4  #include <unistd.h>
5
6  int scan_command(strlist* tokens, command* cmd) {
7      char* path; int oflags;
8      strlist_clear(&cmd->argv);
9      for (int fd = 0; fd < 2; ++fd) {
10         cmd->assignments[fd].path = 0;
11     }
12     cmd->background = 0;
13     for (int i = 0; i < tokens->len; ++i) {
14         switch (tokens->list[i][0]) {
15             case '<':

```

```

16     cmd->assignments[0].path = &tokens->list[i][1];
17     cmd->assignments[0].oflags = O_RDONLY;
18     cmd->assignments[0].mode = 0;
19     break;
20     case '>':
21         path = &tokens->list[i][1];
22         oflags = O_WRONLY | O_CREAT;
23         if (*path == '>') {
24             ++path; oflags |= O_APPEND;
25         } else {
26             oflags |= O_TRUNC;
27         }
28         cmd->assignments[1].path = path;
29         cmd->assignments[1].oflags = oflags;
30         cmd->assignments[1].mode = 0666;
31         break;
32     case '&':
33         cmd->background = 1;
34         break;
35     default:
36         strlist_push(&cmd->argv, tokens->list[i]);
37         if (cmd->cmdname == 0) {
38             cmd->cmdname = tokens->list[i];
39         }
40     }
41 }
42 strlist_push0(&cmd->argv);
43 return cmd->cmdname != 0;
44 }
45
46 /* assign an opened file with the given flags and mode to fd */
47 void fassign(int fd, char* path, int oflags, mode_t mode) {
48     int newfd = open(path, oflags, mode);
49     if (newfd < 0) {
50         perror(path); exit(255);
51     }
52     if (dup2(newfd, fd) < 0) {
53         perror("dup2"); exit(255);
54     }
55     close(newfd);
56 }
57
58 void exec_command(command* cmd) {
59     for (int fd = 0; fd < 2; ++fd) {
60         if (cmd->assignments[fd].path != 0) {
61             fassign(fd, cmd->assignments[fd].path,
62                 cmd->assignments[fd].oflags,
63                 cmd->assignments[fd].mode);
64         }
65     }
66     execvp(cmd->cmdname, cmd->argv.list);
67 }

```

2.8 Überblick der Signale aus dem POSIX-Standard

Signal	Voreinstellung	Beschreibung
<i>SIGABRT</i>	A	Process abort signal.
<i>SIGALRM</i>	T	Alarm clock.
<i>SIGBUS</i>	A	Access to an undefined portion of a memory object.
<i>SIGCHLD</i>	I	Child process terminated, stopped, or continued.
<i>SIGCONT</i>	C	Continue executing, if stopped.
<i>SIGFPE</i>	A	Erroneous arithmetic operation.
<i>SIGHUP</i>	T	Hangup.
<i>SIGILL</i>	A	Illegal instruction.
<i>SIGINT</i>	T	Terminal interrupt signal.
<i>SIGKILL</i>	T	Kill (cannot be caught or ignored).
<i>SIGPIPE</i>	T	Write on a pipe with no one to read it.
<i>SIGQUIT</i>	A	Terminal quit signal.
<i>SIGSEGV</i>	A	Invalid memory reference.
<i>SIGSTOP</i>	S	Stop executing (cannot be caught or ignored).
<i>SIGTERM</i>	T	Termination signal.
<i>SIGTSTP</i>	S	Terminal stop signal.
<i>SIGTTIN</i>	S	Background process attempting read.
<i>SIGTTOU</i>	S	Background process attempting write.
<i>SIGUSR1</i>	T	User-defined signal 1.
<i>SIGUSR2</i>	T	User-defined signal 2.
<i>SIGPOLL</i>	T	Pollable event.
<i>SIGPROF</i>	T	Profiling timer expired.
<i>SIGSYS</i>	A	Bad system call.
<i>SIGTRAP</i>	A	Trace/breakpoint trap.
<i>SIGURG</i>	I	High bandwidth data is available at a socket.
<i>SIGVTALRM</i>	T	Virtual timer expired.
<i>SIGXCPU</i>	A	CPU time limit exceeded.
<i>SIGXFSZ</i>	A	File size limit exceeded.

Voreinstellung	Beschreibung
T	Abbruch des Prozesses. Bei dem bei <i>wait()</i> zurückgelieferten Status ist <i>WIFSIGNALED</i> wahr und über <i>WTERMSIG</i> lässt sich das Signal ermitteln.
A	Analog zu T. Hinzu kommt möglicherweise noch die Erzeugung eines Speicherauszugs (in der Datei <i>core</i>). Letzteres lässt sich mit <i>WCOREDUMP</i> untersuchen.
I	Das Signal wird ignoriert.
S	Der Prozess wird gestoppt.
C	Der Prozess wird fortgesetzt.

Tabelle 2.2: Im POSIX-Standard genannte Signale (Quelle: www.opengroup.org)

Die Tabelle 2.2 liefert einen Überblick aller vom POSIX-Standard genannten Signale. Einzelne Implementierungen können noch weitere Signale unterstützen. Die Signale lassen sich dabei in mehrere Gruppen aufteilen:

- Programmierfehler: *SIGBUS*, *SIGFPE*, *SIGILL*, *SIGSEGV* und *SIGSYS*.
- Ressourcenverbrauch: *SIGVTALRM*, *SIGXCPU* und *SIGXFSZ*.
- Prozeßkontrolle: *SIGCONT*, *SIGKILL*, *SIGSTOP*, *SIGTERM* und *SIGTRAP*.

- Sitzungskontrolle: *SIGHUP*, *SIGINT*, *SIGQUIT* und *SIGTSTP*.
- Ereignis-Indikatoren: *SIGALRM*, *SIGCHLD*, *SIGPIPE*, *SIGPOLL*, *SIGPROF*, *SIGTTIN*, *SIGTTOU*, *SIGURG*, *SIGUSR1*, *SIGUSR2* und *SIGVTALRM*.

Kapitel 3

Pipelines

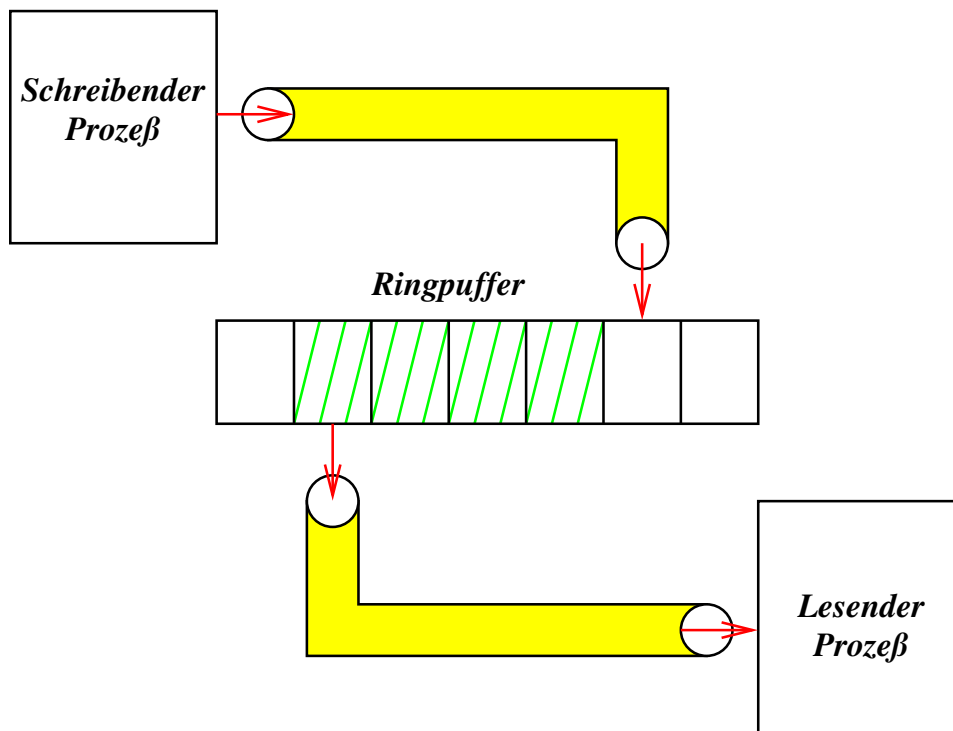


Abbildung 3.1: Aufbau einer Pipeline

Eine Pipeline ist ein unidirektionaler Kommunikationskanal (siehe Abbildung 3.1). Die schreibende und die lesende Seite einer Pipeline werden über verschiedene Dateiverbindungen angesprochen. Zur Pipeline gehört ein vom Betriebssystem verwalteter Ringpuffer, dessen Größenordnung typischerweise bei 4 Kilobyte liegt. Wenn der Puffer vollständig gefüllt ist, wird ein Prozess, der ihn weiter zu füllen versucht, blockiert, bis wieder genügend Platz zur Verfügung steht. Wenn der Puffer leer ist, wird ein lesender Prozess blockiert, bis der Puffer sich zumindest partiell füllt. Dies ist vergleichbar mit der Datenstruktur einer FIFO-Queue (*first in, first out*) mit explizit begrenzter Kapazität.

Der POSIX-Standard unterstützt sowohl benannte Pipelines als auch solche, die mit Hilfe des Systemaufrufs *pipe()* erzeugt werden. Die benannten Pipelines sind aber kaum noch in Gebrauch, da die bidirektionalen UNIX-Domain-Sockets (mehr dazu später) nor-

malerweise bevorzugt werden.

Der folgende Programmtext zeigt den Aufbau und die Verwendung einer Pipeline an einem einfachen Beispiel. In Zeile 10 wird die Pipeline mit Hilfe des Systemaufrufs *pipe()* angelegt. Dabei werden zwei Dateideskriptoren zurückgegeben — einen Dateideskriptor für die lesende und einen für die schreibende Seite.

Programm 3.1: Aufbau und Verwendung einer Pipeline (*pipehello.c*)

```

1 #include <unistd.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4
5 const int PIPE_READ = 0;
6 const int PIPE_WRITE = 1;
7
8 int main() {
9     int pipefds[2];
10    if (pipe(pipefds) < 0) {
11        perror("pipe"); exit(1);
12    }
13    pid_t child = fork();
14    if (child < 0) {
15        perror("fork"); exit(1);
16    }
17    if (child == 0) {
18        close(pipefds[PIPE_WRITE]);
19        char buf[32];
20        ssize_t nbytes;
21        while ((nbytes = read(pipefds[PIPE_READ],
22            buf, sizeof buf)) > 0) {
23            if (write(1, buf, nbytes) < nbytes) exit(1);
24        }
25        exit(0);
26    }
27    close(pipefds[PIPE_READ]);
28    const char message[] = "Hello!\n";
29    write(pipefds[PIPE_WRITE], message, sizeof message - 1);
30    close(pipefds[PIPE_WRITE]);
31    wait(0);
32 }

```

Nach dem *fork()* auf Zeile 13 sind beide Dateideskriptoren für die beiden Enden der Pipeline von beiden Prozessen aus zugänglich. Es empfiehlt sich, unmittelbar nach dem *fork()* das jeweils unbenützte Ende zu schließen. In diesem Beispiel soll das schreibende Ende beim erzeugenden Prozess verbleiben und das lesende Ende beim Kind-Prozess. Entsprechend schließt auf Zeile 18 der Kind-Prozess das schreibende Ende und auf Zeile 27 der Erzeuger das lesende Ende. Danach hat jeder der beiden Prozesse nur die eigene Seite zur Pipeline geöffnet. Dies ist sehr wichtig, da ein Eingabeende bei einer Pipeline nur dann erkannt werden kann, wenn alle Dateiverbindungen geschlossen sind, über die etwas in die Pipeline gefüllt werden kann. Sobald auf Zeile 30 der Erzeuger seine Schreib-Verbindung zur Pipeline schließt, kann die Schleife des Kind-Prozesses in Zeile 21 terminieren, die die komplette Eingabe aus der Pipeline zum Dateideskriptor 1 (Standard-Ausgabe) kopiert. Wenn die *close()*-Operation auf Zeile 18 fehlen würde, dann käme es zu einem endlosen Hänger bei der *read()*-Operation in Zeile 21.

Was geschieht im umgekehrten Falle, wenn versucht wird, in eine Pipeline zu schreiben, bei der alle lesenden Verbindungen geschlossen sind? Glücklicherweise wird der schreibende Prozess nicht endlos blockiert, sondern erhält das Signal *SIGPIPE*. Zudem wird bei einer *write()*-Operation, die wegen einer geschlossenen Pipeline fehlschlägt, *EPIPE* als Fehlerindikation zurückgegeben. Aber auch hierfür ist es notwendig, dass alle überflüssigen Verbindungen zur Pipeline geschlossen worden sind. Der folgende Programmtext demonstriert die Berücksichtigung des *SIGPIPE*-Signals. Zu beachten ist auch hier die voreingestellte Reaktion: Wenn es keine Signalbehandlung für *SIGPIPE* gibt, wird der schreibende Prozess einfach terminiert.

Programm 3.2: Demonstration von *SIGPIPE* (*sigpipe.c*)

```

1  #include <signal.h>
2  #include <stdio.h>
3  #include <stdlib.h>
4  #include <unistd.h>
5
6  const int PIPE_READ = 0;
7  const int PIPE_WRITE = 1;
8
9  volatile sig_atomic_t sigpipe_received = 0;
10
11 void sigpipe_handler(int sig) {
12     sigpipe_received = 1;
13 }
14
15 int main() {
16     int pipefds[2];
17     if (pipe(pipefds) < 0) {
18         perror("pipe"); exit(1);
19     }
20     pid_t child = fork();
21     if (child < 0) {
22         perror("fork"); exit(1);
23     }
24     if (child == 0) {
25         close(pipefds[PIPE_WRITE]);
26         char buf[32];
27         ssize_t nbytes = read(pipefds[PIPE_READ],
28             buf, sizeof buf);
29         if (nbytes > 0) {
30             if (write(1, buf, nbytes) < nbytes) exit(1);
31         }
32         exit(0);
33     }
34     close(pipefds[PIPE_READ]);
35     struct sigaction action = {0};
36     action.sa_handler = sigpipe_handler;
37     if (sigaction(SIGPIPE, &action, 0) < 0) {
38         perror("sigaction"); exit(1);
39     }
40     while (!sigpipe_received) {
41         const char message[] = "Hello!\n";
42         write(pipefds[PIPE_WRITE], message, sizeof message - 1);

```

```

43     }
44     close(pipefds[PIPE_WRITE]);
45     wait(0);
46 }

```

Wozu können Pipelines sinnvoll eingesetzt werden? Praktisch kann die Verwendung in einem Fork-and-Join-Szenario sein (siehe Abbildung 1.2), wenn es darum geht, Ergebnisse von den Unterprozessen abzuholen. Das hätte beispielsweise bei dem Programmtext 1.7 genutzt werden können, um gefundene Lösungen an den übergeordneten Prozess wieder zurückzugeben.

Sehr viel häufiger ist allerdings der Einsatz von Pipelines, um die Standard-Eingabe oder Standard-Ausgabe eines Prozesses mit einem anderen Prozess zu verbinden. Hierfür bietet POSIX die Funktionen `popen()` und `pclose()`, die allerdings unvermeidbar die Shell mit all ihren Sicherheitsrisiken und die Verwendung der `stdio` mit sich bringen. Das folgende Beispiel zeigt eine alternative Schnittstelle, die analog zu `execvp()` mit Argumentlisten arbeitet:

Programm 3.3: Schnittstelle für Pipelines zu Kommandos (`pconnect.h`)

```

1  /*
2   * Create and manage pipelines to given commands.
3   * In comparison to popen(), neither stdio nor the shell are used.
4   * afb 6/2003
5   */
6  #ifndef PCONNECT_H
7  #define PCONNECT_H
8
9  #include <unistd.h>
10
11 #define PIPE_READ 0
12 #define PIPE_WRITE 1
13
14 typedef struct pipe_end {
15     int fd;
16     pid_t pid;
17     int wstat;
18 } pipe_end;
19
20 /*
21  * create a pipeline to the given command;
22  * mode should be either PIPE_READ or PIPE_WRITE;
23  * return a filled pipe_end structure and 1 on success
24  * and 0 in case of failures
25  */
26 int pconnect(const char* path, char* const* argv,
27             int mode, pipe_end* pipe_con);
28
29 /*
30  * close pipeline and wait for the forked-off process to exit;
31  * the wait status is returned in pipe->wstat;
32  * 1 is returned if successful, 0 otherwise
33  */
34 int phangup(pipe_end* pipe_end);
35

```

36 #endif

Nicht zu vergessen ist dabei eine Verwaltungsstruktur, die sich die Prozess-ID merkt und die es ermöglicht, auf das Ende des Prozesses an der anderen Seite der Pipeline zu warten. `pconnect()` initialisiert im Erfolgsfalle eine entsprechende Struktur vom Typ `pipe_end`, und `phangup()` erlaubt es, die Pipeline zu schließen und auf das Ende der anderen Seite zu warten. So könnte eine zugehörige Implementierung aussehen:

Programm 3.4: Pipelines zu Kommandos (`pconnect.c`)

```

1 #include <fcntl.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <sys/wait.h>
5 #include <unistd.h>
6 #include "pconnect.h"
7
8 /*
9  * create a pipeline to the given command;
10  * mode should be either PIPE_READ or PIPE_WRITE;
11  * return a filled pipe_end structure and 1 on success
12  * and 0 in case of failures
13  */
14 int pconnect(const char* path, char* const* argv,
15             int mode, pipe_end* pipe_con) {
16     int pipefds[2];
17     if (pipe(pipefds) < 0) return 0;
18     int myside = mode; int otherside = 1 - mode;
19     fflush(0);
20     pid_t child = fork();
21     if (child < 0) {
22         close(pipefds[0]); close(pipefds[1]);
23         return 0;
24     }
25     if (child == 0) {
26         close(pipefds[myside]);
27         dup2(pipefds[otherside], otherside);
28         close(pipefds[otherside]);
29         execvp(path, argv); exit(255);
30     }
31     close(pipefds[otherside]);
32     int flags = fcntl(pipefds[myside], F_GETFD);
33     flags |= FD_CLOEXEC;
34     fcntl(pipefds[myside], F_SETFD, flags);
35     pipe_con->pid = child;
36     pipe_con->fd = pipefds[myside];
37     pipe_con->wstat = 0;
38     return 1;
39 }
40
41 /*
42  * close pipeline and wait for the forked-off process to exit;
43  * the wait status is returned in wstat (if non-null);
44  * 1 is returned if successful, 0 otherwise

```

```

45  */
46  int phangup(pipe_end* pipe) {
47      if (close(pipe->fd) < 0) return 0;
48      if (waitpid(pipe->pid, &pipe->wstat, 0) < 0) return 0;
49      return 1;
50  }

```

Nach der Erzeugung der Pipeline in Zeile 17 wird festgelegt, welche der beiden Dateideskriptoren welcher Seite gehört. In Zeile 19 werden alle Puffer der *stdio*-Bibliothek geleert. Der Kind-Prozess, der die gegenüberliegende Seite vertritt, schließt in Zeile 26 das unbenutzte Ende der Pipeline und sorgt in Zeile 27 dafür, dass der entsprechende Standard-Deskriptor (0 für Standardeingabe und 1 für Standardausgabe) der Pipeline zugeordnet wird. Der zuvor verwendete Dateideskriptor für das Ende der Pipeline wird dann in Zeile 28 geschlossen.

Der übergeordnete Prozess schließt nicht nur auf Zeile 31 das andere Ende, sondern setzt für das eigene Ende die Option *FD_CLOEXEC*, die bewirkt, dass der betreffende Dateideskriptor nicht über einen *exec()*-Systemaufruf hinweg vererbt, sondern geschlossen wird. Dies ist zwingend notwendig, wenn weitere Prozesse erzeugt werden, um zu verhindern, dass versehentlich Zugänge zu einer Pipeline offenbleiben, wenn die eigentlichen Beteiligten sie nicht mehr benötigen. Besser als ein *close-on-exec* wäre natürlich ein implizites *close()* bei *fork()*. Leider sieht dies der POSIX-Standard nicht vor. Entsprechend ist die Verwendung von Pipelines nur bei *fork()/exec()*-Kombinationen sicher, jedoch nicht bei der Erzeugung länger laufender Unterprozesse, die nicht mit *exec()* zu einem anderen Programmtext wechseln.

Programm 3.5: Verarbeitung der Ausgabe eines Kommandos (*rwhousers.c*)

```

1  #include <stralloc.h>
2  #include <string.h>
3  #include <unistd.h>
4  #include "pconnect.h"
5  #include "rwhousers.h"
6  #include "strlist.h"
7
8  const char rwho_path[] = "/usr/bin/rwho";
9
10 /*
11  * invoke rwho and get list of users that are currently logged in;
12  * return 1 in case of success, otherwise 0
13  */
14  int get_rwho_users(strlist* users) {
15      strlist argv = {0};
16      strlist_push(&argv, rwho_path);
17      strlist_push0(&argv);
18      pipe_end pipe;
19      int ok = pconnect(rwho_path, argv.list, PIPE_READ, &pipe);
20      strlist_free(&argv);
21      if (!ok) return 0;
22
23      stralloc rwho_output = {0};
24      ssize_t nbytes;
25      char buf[32];
26      while ((nbytes = read(pipe.fd, buf, sizeof buf)) > 0) {
27          stralloc_catb(&rwho_output, buf, nbytes);

```

```

28     }
29     phangup(&pipe);
30
31     strlist_clear(users);
32     char* user = rwho_output.s;
33     for (int i = 0; i < rwho_output.len; ++i) {
34         switch (rwho_output.s[i]) {
35             case '_':
36                 if (user != 0) {
37                     rwho_output.s[i] = 0;
38                     strlist_push(users, strdup(user));
39                     user = 0;
40                 }
41                 break;
42             case '\n':
43                 user = rwho_output.s + i + 1;
44                 break;
45         }
46     }
47     stralloc_free(&rwho_output);
48     return 1;
49 }

```

Der vorstehende Programmtext demonstriert, wie `pconnect()` verwendet werden kann, um die Ausgabe eines Kommandos einzulesen. In diesem Beispiel geht es um die Ausgabe des `rwho`-Kommandos, das angibt, welche Nutzer sich im gleichen Netzwerk angemeldet haben:¹

```

turing# rwho
borchert theseus:pts/10  May  9 09:39 :19
borchert theseus:pts/13  Mar  9 13:39 :20
borchert theseus:pts/36  May  9 09:37 :06
borchert theseus:pts/38  Apr 26 13:36 :19
borchert turing:pts/31   Apr 26 13:37 :14
borchert xylopi:pts/4    Apr  8 10:38 :19
jens_k   theseus:pts/97   Jun 14 13:51 :02
rudib    centaurus:pts/2   Jun 14 14:57 :02
rudib    centaurus:pts/4   Jun 14 14:57 :24
francis  kleemann:pts/2       Jun 14 11:02 :17
francis  turing:pts/56       Jun 14 11:02 :18
hschmidt turing:pts/24        Jun 14 14:34
rmack    theseus:pts/49      Jun 14 15:15
fkraft   turing:pts/5         Jun 14 13:53 :52
prince   theseus:pts/84       Jun 14 15:01 :03
drtiger  turing:pts/75       Jun 14 11:51 :01
turing#

```

Die Funktion `get_rwho_users()` aus dem Programmtext `rwhousers.c` soll das Kommando `rwho` aufrufen und aus dessen Ausgabe nur die Namen der Benutzer extrahieren. In dem angezeigten Beispiel würden wir also 6 Mal `borchert` erhalten, gefolgt von `jens_k`, zwei Mal `rudib` und so fort. Zu Beginn eröffnet `get_rwho_users()` in den Zeilen 15 bis 21 die Pipeline. Dann wird in den Zeilen 23 bis 28 die gesamte Ausgabe des `rwho`-Kommandos in einen Puffer vom Typ `stralloc` eingelesen. In Zeile 29 wird die Pipeline geschlossen. Danach wird

¹Die Benutzernamen sind hier mit der Ausnahme meines eigenen frei erfunden.

in den Zeilen 31 bis 46 das jeweils erste Wort jeder Zeile aus dem Puffer in die Liste *users* eingefügt.

Der folgende Programmtext zeigt den umgekehrten Fall, bei dem die Ausgabe für ein Kommando generiert wird:

Programm 3.6: Generierung der Eingabe für ein Kommando (*sendmail.c*)

```

1  #include <stralloc.h>
2  #include <unistd.h>
3  #include "pconnect.h"
4  #include "sendmail.h"
5  #include "strlist.h"
6
7  const char sendmail_path[] = "/usr/lib/sendmail";
8  /*
9   * return a pipeline opened to /usr/lib/sendmail on the
10  * local system; return the opened pipeline and 1 in
11  * case of success; 0 in case of failures
12  */
13  int sendmail(char* recipient, char* subject, pipe_end* pipe_con) {
14      strlist argv = {0};
15      strlist_push(&argv, sendmail_path);
16      strlist_push(&argv, "-t");
17      strlist_push0(&argv);
18      int ok = pconnect(sendmail_path, argv.list, PIPE_WRITE, pipe_con);
19      strlist_free(&argv);
20      if (!ok) return 0;
21      stralloc header = {0};
22      stralloc_cats(&header, "To:");
23      stralloc_cats(&header, recipient);
24      stralloc_cats(&header, "\n");
25      stralloc_cats(&header, "Subject:");
26      stralloc_cats(&header, subject);
27      stralloc_cats(&header, "\n\n");
28      ssize_t written = 0; ssize_t left = header.len;
29      while (left > 0) {
30          ssize_t nbytes = write(pipe_con->fd, header.s + written, left);
31          if (nbytes < 0) {
32              stralloc_free(&header);
33              phangup(pipe_con);
34              return 0;
35          }
36          written += nbytes; left -= nbytes;
37      }
38      stralloc_free(&header);
39      return 1;
40  }

```

Bei diesem Beispiel geht es um die Erstellung einer E-Mail, wobei *sendmail()* den Aufbau einer Pipeline zu */usr/lib/sendmail* übernimmt, womit lokal E-Mails versendet werden können. Obwohl */usr/lib/sendmail* nicht zum POSIX-Standard gehört, ist es auf allen gängigen UNIX-Systemen zu finden, unabhängig von der eingesetzten Mail-Software. Die Option *-t* sorgt dafür, dass die Empfängeradresse der entsprechenden Kopfzeile der E-Mail entnommen wird und nicht auf der Kommandozeile angegeben werden muß. Die

Funktion `sendmail()` übernimmt hier nur das Schreiben der Kopfzeilen. Diese werden zunächst auf den Zeilen 21 bis 27 in einen Puffer geschrieben, um sie danach in den Zeilen 29 bis 37 an die Pipeline weiterzugeben. Hier ist zu beachten, dass das `write()` keinesfalls verpflichtet ist, den Puffer in der gewünschten Gesamtlänge in einem Schwung zu übernehmen. Stattdessen steht es `write()` frei, nur einen Teil zu schreiben und die entsprechende Zahl der geschriebenen Bytes zurückzuliefern. Dies ist insbesondere bei Pipelines möglich, wenn die Zahl der zu schreibenden Bytes die zur Verfügung stehende Kapazität im Ringpuffer übersteigt.

Programm 3.7: Anwendung mit zwei Pipeline-Verbindungen (*bigbrother.c*)

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include "pconnect.h"
4  #include "sendmail.h"
5  #include "strhash.h"
6  #include "strlist.h"
7  #include "rwhousers.h"
8
9  int main(int argc, char** argv) {
10     if (argc <= 2) {
11         fprintf(stderr, "Usage: _%s_email_login...\n", argv[0]);
12         exit(1);
13     }
14     char* email = ++argv; --argc;
15     strhash friends = {0};
16     strhash_alloc(&friends, 4);
17     while (--argc > 0) {
18         if (!strhash_add(&friends, ++argv, 0)) exit(1);
19     }
20
21     strlist users = {0};
22     if (!get_rwho_users(&users)) exit(1);
23     strhash found = {0};
24     strhash_alloc(&found, 4);
25     for (int i = 0; i < users.len; ++i) {
26         if (strhash_exists(&found, users.list[i])) continue;
27         if (!strhash_exists(&friends, users.list[i])) continue;
28         if (!strhash_add(&found, users.list[i], 0)) exit(1);
29     }
30     if (strhash_length(&found) == 0) exit(0);
31
32     pipe_end pipe_con;
33     if (!sendmail(email, "Your_Friends_Are_Online!", &pipe_con))
34         exit(1);
35     if (dup2(pipe_con.fd, 1) < 0) exit(1);
36     printf("Hi, ");
37     if (strhash_length(&found) == 1) {
38         printf("one_of_your_friends_is");
39     } else {
40         printf("some_of_your_friends_are");
41     }
42     printf("_online:\n");
43     strhash_start(&found);

```

```

44     char* key;
45     while (strhash_next(&found, &key)) {
46         printf("%s\n", key);
47     }
48     fclose(stdout);
49     if (!phangup(&pipe_con)) exit(1);
50 }

```

Der vorstehende Programmtext verwendet `get_rwho_users()` und `sendmail()`, um per E-Mail einen Hinweis zu verschicken, wenn auf der Kommandozeile benannte Freunde angemeldet sind. Dabei kommt eine Hash-Organisation für Zeichenketten zum Einsatz, um Duplikate zu eliminieren und das Abgleichen der Namenslisten zu erleichtern:

Programm 3.8: Schnittstelle einer Hash-Organisation für Zeichenketten (*strhash.h*)

```

1  /*
2   * Hashes for strings, i.e. keys and values are of type char*
3   * afb 6/2003
4   */
5
6  #ifndef STRHASH_H
7  #define STRHASH_H
8
9  typedef struct strhash_entry {
10     char* key;
11     char* value;
12     struct strhash_entry* next;
13 } strhash_entry;
14
15 typedef struct strhash {
16     unsigned int size, length;
17     strhash_entry** bucket; /* hash table */
18     unsigned int it_index;
19     strhash_entry* it_entry;
20 } strhash;
21
22 /* allocate a hash table with the given bucket size */
23 int strhash_alloc(strhash* hash, unsigned int size);
24
25 /* add tuple (key,value) to the hash, key must be unique */
26 int strhash_add(strhash* hash, char* key, char* value);
27
28 /* remove tuple with the given key from the hash */
29 int strhash_remove(strhash* hash, char* key);
30
31 /* return number of elements */
32 unsigned int strhash_length(strhash* hash);
33
34 /* check existance of a key */
35 int strhash_exists(strhash* hash, char* key);
36
37 /* lookup value by key */
38 int strhash_lookup(strhash* hash, char* key, char** value);
39

```

```

40 /* start iterator */
41 int strhash_start(strhash *hash);
42
43 /* fetch next key from iterator; returns 0 on end */
44 int strhash_next(strhash *hash, char** key);
45
46 /* free allocated memory */
47 int strhash_free(strhash* hash);
48
49 #endif

```

Die Zeile 34 bei *bigbrother.c* zeigt noch einen kleinen Trick. Wenn die Absicht besteht, die *stdio* zu verwenden, obwohl zunächst nur ein Dateideskriptor vorliegt, dann besteht die Möglichkeit, *fdopen()* zu verwenden oder — wie in diesem Beispiel — mit einem *dup2()* die bisherige Standardausgabe zu schließen und sie danach auf eine andere geöffnete Verbindung verweisen zu lassen. Danach kann *printf()* verwendet werden, um den Inhalt der zu versendenden E-Mail zu schreiben.

Interessant ist der Aufbau größerer Pipeline-Ketten, wie sie z.B. von der Shell unterstützt werden. Abbildung 3.2 zeigt die Struktur bei einer Kette von drei Kommandos, demonstriert am Beispiel „*who | grep borchert | wc -l*“. Zu beachten ist dabei, dass die Kette mit dem letzten Kommando beginnt. Der unmittelbare Unterprozeß der Shell ruft also das Kommando „*wc -l*“ auf. Bevor es dazu kommt, wird aber eine Pipeline erzeugt, die an einen neu erzeugten Prozess als Standardausgabe weitervererbt wird, der dann die Ausführung von „*grep borchert*“ übernimmt. Doch bevor es dazu kommt, erzeugt der Prozess analog wie zuvor eine neue Pipeline und einen weiteren Prozess, der für die Ausführung von „*who*“ zuständig ist.

Der folgende Programmtext zeigt eine einfache Datenstruktur für Pipeline-Ketten. Die Felder *cmdname* und *argv* sind direkt an *execvp()* übergebene Parameter. Das Feld *input* ist entweder 0 (Standard-Eingabe wird einfach übernommen) oder ein Verweis auf die Pipeline-Kette, die die Eingabe für dieses Kommando generiert.

Programm 3.9: Datenstruktur für Kommandos einer Pipeline (*pipecmd.h*)

```

1 #ifndef PIPECMD_H
2 #define PIPECMD_H
3
4 #include "strlist.h"
5
6 typedef struct pipecmd {
7     char* cmdname;
8     strlist* argv;
9     struct pipecmd* input;
10 } pipecmd;
11
12 /* construct pipeline chain out of tokens */
13 pipecmd* pipe_constructor(strlist* tokens);
14
15 /* free space allocated by pipe_constructor() */
16 void pipe_free(pipecmd* cmd);
17
18 #endif

```

Abbildung 3.3 zeigt die zugehörige Datenstruktur für das genannte Beispiel. Die Datenstruktur entspricht dabei genau der benötigten Prozessstruktur, bei der die Reihenfol-

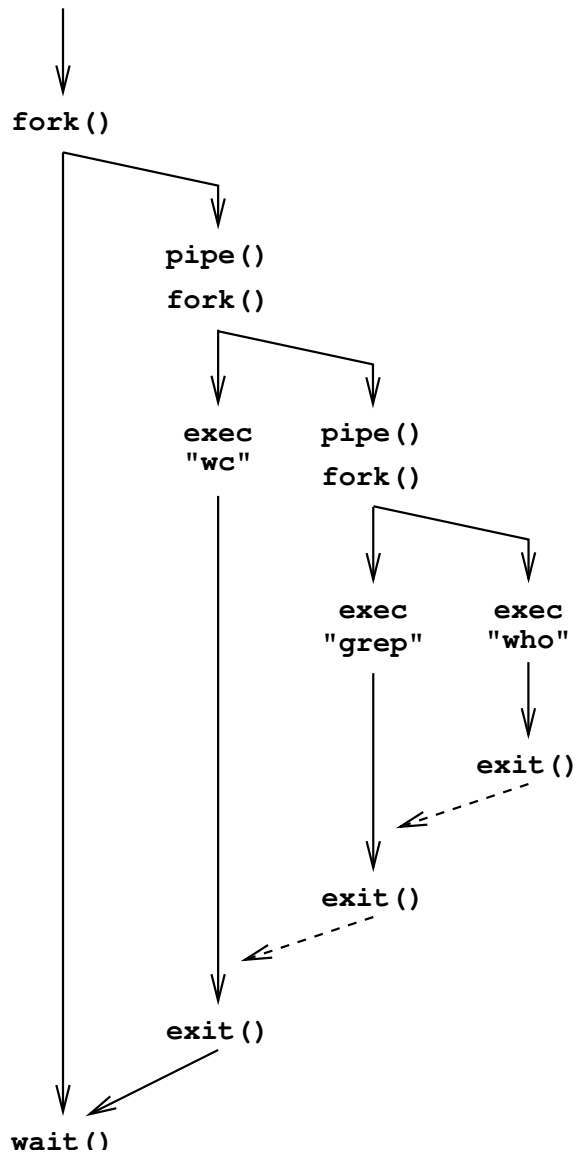


Abbildung 3.2: Prozessstruktur bei einer Kette von Pipelines

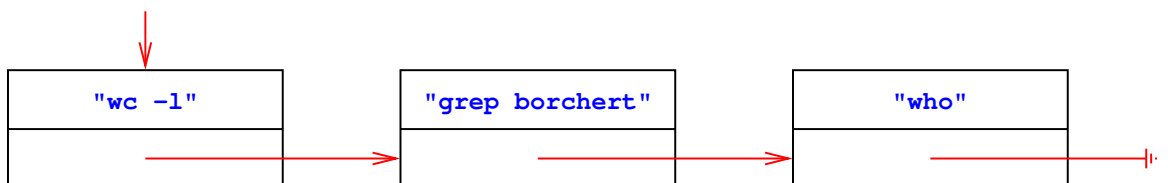


Abbildung 3.3: Datenstruktur für Kommandos einer Pipeline

ge im Vergleich zur Kommandozeile reversiert ist. Die gestrichelten Linien, die den *exit()*-Aufrufen folgen, deuten die Synchronisierung durch das Erkennen des Eingabe-Endes an. Nur beim Kommando, das am Ende der Zeile stand, erfolgt eine Synchronisierung durch *wait()*.

Programm 3.10: Aufbau einer Pipeline-Liste ausgehend von der Wortliste (*pipecmd.c*)

```

1  #include <assert.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include "pipecmd.h"
5  #include "strlist.h"
6
7  /* generate new pipeline command node */
8  static pipecmd* new_command(char* cmdname,
9      strlist* argv, pipecmd* input) {
10     assert(cmdname != 0);
11     pipecmd* cmd = (pipecmd*) malloc(sizeof(pipecmd));
12     if (cmd == 0) return 0;
13     cmd->cmdname = cmdname;
14     strlist_push0(argv);
15     cmd->argv = argv;
16     cmd->input = input;
17     return cmd;
18 }
19
20 /* construct pipeline chain out of tokens */
21 pipecmd* pipe_constructor(strlist* tokens) {
22     pipecmd* input = 0;
23     char* cmdname;
24     strlist* argv = 0;
25     for (int index = 0; index < tokens->len; ++index) {
26         if (strcmp(tokens->list[index], "|") == 0) {
27             if (argv == 0) return 0; /* syntax error */
28             input = new_command(cmdname, argv, input);
29             cmdname = 0; argv = 0;
30         } else {
31             if (argv == 0) {
32                 argv = (strlist*) malloc(sizeof(strlist));
33                 if (argv == 0) return 0;
34                 argv->list = 0; argv->len = 0; argv->allocated = 0;
35             }
36             if (cmdname == 0) cmdname = tokens->list[index];
37             if (!strlist_push(argv, tokens->list[index])) return 0;
38         }
39     }
40     if (cmdname != 0) input = new_command(cmdname, argv, input);
41     return input;
42 }
43
44 /* free space allocated by pipe_constructor() */
45 void pipe_free(pipecmd* cmd) {
46     if (cmd->input) pipe_free(cmd->input);
47     strlist_free(cmd->argv); free(cmd->argv);

```

```

48     free(cmd);
49 }

```

Vorstehender Programmtext demonstriert, wie eine Sequenz von Worten, die aus Kommandoargumenten und möglichen „|“-Zeichen besteht, in eine Pipeline-Liste konvertiert werden kann. Diese Aufgabe wird von der Funktion *pipe_constructor()* erledigt, die in der **for**-Schleife in den Zeilen 25 bis 39 durch die Liste der Tokens geht. Wird das Ende eines Kommandos gefunden, sei es durch ein Pipeline-Symbol in Zeile 26 oder ganz am Ende der Token-Liste in Zeile 40, so wird mit der Funktion *new_command()* ein neues Element der Kette angelegt. Da das neueste Element der Kette jeweils auf alle früheren Elemente verweist, erhalten wir genau die gewünschte Reversion.

Der folgende Programmtext zeigt, wie sich die Datenstruktur unmittelbar in die zugehörige Prozessstruktur umsetzen lässt.

Programm 3.11: Ausführung einer Pipeline-Kette (*execpipe.c*)

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include "execpipe.h"
5  #include "pipecmd.h"
6
7  #define PIPE_READ 0
8  #define PIPE_WRITE 1
9
10 /*
11  * exec() to the given pipeline command;
12  * does not return in case of success
13  */
14 void exec_pipecmd(pipecmd* command) {
15     if (command->input != 0) {
16         int pipefds[2];
17         if (pipe(pipefds) < 0) {
18             perror("pipe"); return;
19         }
20         pid_t child = fork();
21         if (child < 0) {
22             perror("fork"); return;
23         }
24         if (child == 0) {
25             close(pipefds[PIPE_READ]);
26             if (dup2(pipefds[PIPE_WRITE], 1) < 0) {
27                 perror("dup2"); exit(255);
28             }
29             close(pipefds[PIPE_WRITE]);
30             exec_pipecmd(command->input);
31             exit(255);
32         }
33         close(pipefds[PIPE_WRITE]);
34         if (dup2(pipefds[PIPE_READ], 0) < 0) {
35             perror("dup2"); return;
36         }
37         close(pipefds[PIPE_READ]);
38     }

```



```
37         if (code && code != 255) {
38             printf("terminated_with_exit_code_%d\n", code);
39         }
40     } else {
41         printf("terminated_abnormally\n");
42     }
43 } else {
44     perror("wait");
45 }
46 }
47 }
```

Kapitel 4

Einführung in Netzwerkdienste und TCP/IP

4.1 Netzwerkdienste

Ein Netzwerkdienst ist ein Prozess, der unter einer Netzwerkadresse einen Dienst anbietet. Ein Klient, der die Netzwerkadresse kennt, kann einen bidirektionalen Kommunikationskanal zu dem Netzwerkdienst eröffnen und über diesen mit dem Dienst kommunizieren, beispielsweise indem Kommandos auf dem Hinweg übermittelt werden und der Dienst auf dem Rückweg des Kommunikationskanals die Antwort überträgt.

Im Unterschied zu Pipelines müssen die beiden Kommunikationspartner nicht miteinander verwandt sein. Sie müssen nicht einmal auf dem gleichen Rechner laufen. Da der Kommunikationskanal bidirektional ist, wird ein echter Dialog zwischen den beiden Prozessen möglich. Der Aufbau einer Verbindung ist jedoch schwieriger, da die Netzwerkadresse des gewünschten Partners ermittelt werden muss.

Wenn Dienste über das Netzwerk angeboten und in Anspruch genommen werden, ergeben sich viele Vorteile:

- Der Dienst kann allen offenstehen, und ein direkter Zugang zu dem Rechner, auf dem der Dienst angeboten wird, ist nicht notwendig.
- Viele Parteien können in kooperativer Weise einen Dienst gleichzeitig nutzen.
- Der Dienste-Anbieter hat weniger Last, da die Benutzerschnittstelle auf anderen Rechnern laufen kann.

Andererseits bergen Netzwerkdienste eine Reihe von Risiken und Problematiken:

- Der Kreis derjenigen, die auf einen Netzwerkdienst zugreifen können, ist möglicherweise ziemlich umfangreich (normalerweise das gesamte Internet).
- Somit muss jeder Netzwerkdienst Zugriffsberechtigungen einführen und überprüfen und kann sich dabei nicht wie traditionelle Applikationen auf die des Betriebssystems verlassen.
- Dienste, die gleichzeitig von vielen genutzt werden können, haben vielerlei zusätzliche Konsistenz- und Synchronisierungsprobleme, für die nicht jede Art von Datenhaltung geeignet ist.
- Netzwerke bringen neue Arten von Ausfällen mit sich, wenn eine Netzwerkverbindung zusammenbricht oder es zu längeren „Hängern“ kommt.

Es gibt eine unüberschaubare Vielfalt an Netzwerk-Hardware, Transport-Protokollen und Schnittstellen zu deren Benutzung. Einen exzellenten Überblick hierfür bietet das Standardwerk *Computer Networks* von Andrew S. Tanenbaum (Prentice Hall, Third Edition, 1996, ISBN 0-13-349945-6). Im weiteren konzentrieren wir uns jedoch auf die Transport-Protokolle des Internets (TCP/IP).

4.2 IP-Adressen

Hier ist ein triviales Beispiel für den Aufbau einer Netzwerkverbindung. Als Klient wird hier *telnet* verwendet, das eine interaktive Nutzung von Netzwerkdiensten ermöglicht. Als trivialer Netzwerkdienst kommt hier *daytime* zum Einsatz, das die Eingabe auf dem Kommunikationskanal ignoriert und nur die aktuelle Zeit analog zum *date*-Kommando zurückgibt:

```
theseus# telnet 134.60.166.2 13
Trying 134.60.166.2...
Connected to 134.60.166.2.
Escape character is '^]'.
Sat Jun 21 14:48:34 2003
Connection to 134.60.166.2 closed by foreign host.
theseus#
```

In diesem Beispiel besteht die Netzwerkadresse aus *134.60.54.12* (einer sogenannten IP-Adresse) und 13 (der Port-Nummer). IP-Adressen dienen dazu, einen Rechner zu adressieren, während Port-Nummern einem Dienst, der auf einem Rechner läuft, zugeordnet werden. In diesem Beispiel adressiert *134.60.54.12* unsere Theseus und die Port-Nummer 13 üblicherweise den *daytime*-Dienst.

Beim Internet gibt es für die Adressierung von Rechnern zwei Adreßräume. IPv4 arbeitet mit 32-Bit-Adressen, ist seit dem 1. Januar 1983 in Benutzung und dominiert bis heute. Daneben gibt es auch IPv6, das mit 128-Bit-Adressen arbeitet, um einen Ausweg zu bieten für die drohende Knappheit mit Adressen bei IPv4. Im weiteren beschäftigen wir uns hier allerdings nur mit IPv4-Adressen, die kurz als IP-Adressen bezeichnet werden.

Die sogenannte *dotted-decimal*-Notation für IP-Adressen spezifiziert eine 32-Bit-Adresse als Folge von vier dezimalen Werten, die jeweils ein Byte angeben. *134.60.54.12* ist somit eine etwas übersichtlichere Schreibweise als 2252092940. Port-Nummern liegen im Bereich von 1 bis 65535, wobei der Bereich von 1 bis 1023 typischerweise für wohldefinierte Dienste reserviert ist. Diese Zuordnungen werden von der IANA (*Internet Assigned Numbers Authority*) vergeben und sind unter <http://www.iana.org/assignments/port-numbers> zu finden.

Beide Adreßräume werden in hierarchisierter Form verwaltet. Ganz oben steht IANA, das große Teile des Adreßraumes an regionale Institutionen weitergibt (ARIN für Amerika, RIPE für Europa, den Mittleren Osten und Zentralasien, APNIC für Asien, Australien und Ozeanien, AfriNIC für Afrika und LACNIC für Lateinamerika einschließlich Teile der Karibik). Die Universität Ulm hat seit 1989 den Adreßbereich *134.60.0.0/16*¹ (damals noch von ARIN) zugeteilt bekommen.

Die Abbildung von Rechnernamen wie *theseus.mathematik.uni-ulm.de* in IP-Adressen wie *134.60.54.12* erfolgt dabei durch sogenannte Domain-Server, die ebenfalls hierarchisch organisiert sind. Wenn jemand von außerhalb den Namen *theseus.mathematik.uni-ulm.de* angibt, wird zunächst untersucht, wer für die Domain „de“ zuständig ist. Dies geschieht

¹Diese Schreibweise wird als CIDR-Notation bezeichnet, wobei CIDR für *Classless Inter-Domain Routing* steht. Die Zahl hinter dem Schrägstrich gibt dabei die Zahl der führenden Bits an, die den Netzwerkteil ausmachen. Bei 16 sind das die ersten beiden Bytes, also *134.60*. Alle Adressen in der Form *134.60.*.** gehören somit zur Universität Ulm.

durch eine Anfrage an einen der sogenannten Root-Server, beispielsweise *198.41.0.4*, der selbst den Namen *a.root-servers.net* trägt. Dieser liefert unter anderem den Name-Server *194.0.0.53* mit dem Namen *a.nic.de*. Wenn *a.nic.de* nach *uni-ulm.de* gefragt wird, wird u.a. auf *134.60.1.111* mit dem Namen *dns1.uni-ulm.de* verwiesen. Letzterer Name-Server kann dann endlich auch *theseus.mathematik.uni-ulm.de* in *134.60.54.12* abbilden.

IP-Adressen wie *134.60.54.12* werden nur auf einer abstrakten Ebene zur Verfügung gestellt. IP-Adressen werden auf der darunterliegenden physischen Ebene und denen damit verbundenen Protokollen nicht verstanden. So wird beispielsweise beim Ethernet, das bei uns weitgehend an der Universität zum Einsatz kommt, mit 6-Byte-Adressen gearbeitet. Die Theseus hat beispielsweise die Ethernet-Adresse *0:14:4f:3e:a1:f0* (Bytes werden hier in Form von Hexzahlen angegeben). Diese Adressen sind jedoch nur lokal auf einem Ethernet-Segment von Bedeutung.

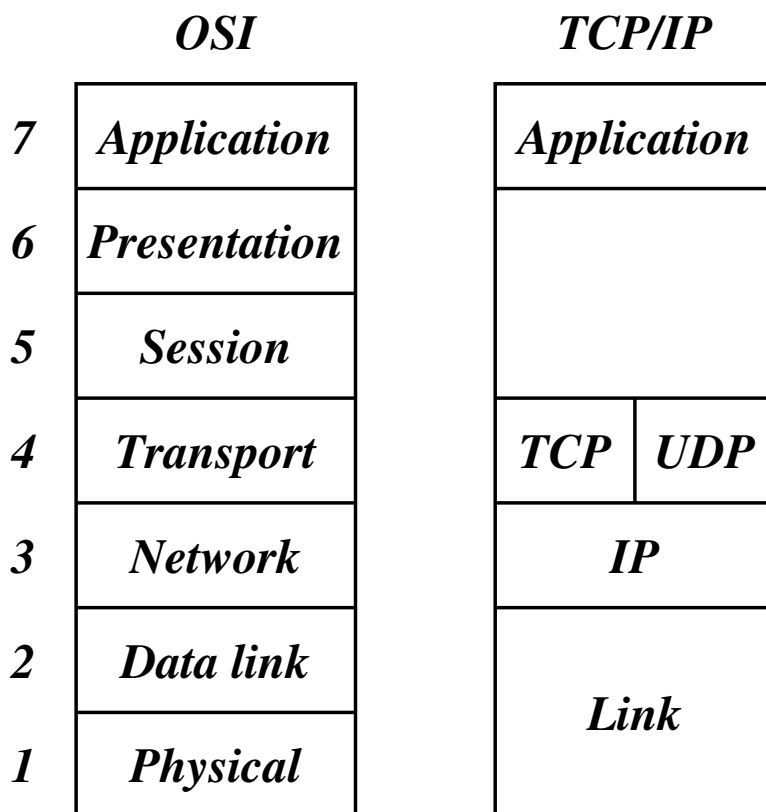


Abbildung 4.1: Schichtenmodell von OSI und TCP/IP

Wir haben also eine Schicht mit IP-Adressen und darunter eine Schicht, die von der verwendeten Netzwerk-Hardware abhängt. Aufbauend auf der Schicht mit IP-Adressen (IP-Protokoll) gibt es alternative Transport-Schichten, über die Pakete versendet werden können: Mittels UDP (*User Datagram Protocol*) können einzelne Pakete sehr effizient, aber unzuverlässig versendet werden, während TCP (*Transmission Control Protocol*) eine sichere Verbindung gewährleistet, die jedoch weniger effizient ist. Abbildung 4.1 zeigt auf der rechten Seite die vier Schichten, die beim Internet eine Rolle spielen: Anwendungen (wie etwa *telnet* und *daytime* im obigen Beispiel), Transport (UDP oder TCP), IP und die physische Ebene. Parallel dazu entstand 1983 das OSI-Referenz-Modell² der ISO (*International Standards Organization*), das eine etwas feinere Schichtung vorsieht. Die Präsentations-

²OSI steht für *Open Systems Interconnection*.

oder Sitzungs-Ebene fand jedoch nie ihren Weg in die Protokollhierarchie von TCP/IP.

4.3 Berkeley Sockets

Für TCP/IP gibt es zwei Schnittstellen, die beide zum POSIX-Standard gehören: Die Berkeley Sockets, die mit BSD 4.2 im Jahr 1983 eingeführt worden sind, und TLI (*Transport Layer Interface*), das auf Streams basiert und zuerst durch UNIX System V Release 3.0 im Jahr 1987 veröffentlicht wurde. Die Berkeley-Socket-Schnittstelle hat sich weitgehend durchgesetzt, da sie bei einigen UNIX-Systemen die einzige zur Verfügung stehende Schnittstelle ist.

Das Standardwerk über BSD 4.3 (*The Design and Implementation of the 4.3 BSD UNIX Operating System*, Samuel J. Leffler et al, Addison Wesley, 1989, ISBN 0-201-06196-1) nennt folgende Ziele:

- **Transparenz:** Die Kommunikation zwischen zwei Prozessen soll nicht davon abhängen, ob sie auf dem gleichen Rechner laufen oder nicht.
- **Effizienz:** Zu Zeiten von BSD 4.2 (also 1983) war dies ein außerordentlich wichtiges Kriterium wegen der damals noch sehr geringen Rechenleistung. Aus diesem Grund werden insbesondere keine weiteren System-Prozesse zur Kommunikation eingesetzt, obwohl dies zu mehr Flexibilität und Modularität hätte führen können.
- **Kompatibilität:** Viele bestehende Applikationen und Bibliotheken wissen nichts von Netzwerken und sollen dennoch in einem verteilten Umfeld eingesetzt werden können. Dies wurde dadurch erreicht, daß nach einem erfolgten Verbindungsaufbau (der z.B. von einem anderen Prozess durchgeführt werden kann) Ein- und Ausgabe in gewohnter Weise (wie bei Dateien, Pipelines oder Terminal-Verbindungen) erfolgen können.

Um den Anforderungen gerecht zu werden, wurden folgende Abstraktionen entwickelt:

- Die Abstraktion eines Kommunikationsbereiches macht es möglich, nicht nur TCP/IP zu unterstützen, sondern auch viele weitere Netzwerke (z.B. Appletalk, DECnet, IPX von Novell). Zu jedem Kommunikationsbereich gibt es unterschiedliche Namen (bzw. Adressen) für Kommunikationsendpunkte. Bei TCP/IP sind das die bekannten 32 Bit langen IP-Adressen (z.B. 134.60.166.2 für die Theseus) kombiniert mit der Port-Nummer des einzelnen Dienstes (z.B. 13 für *daytime*).
- Die Abstraktion eines Kommunikationsendpunktes (daher der Name „socket“), der mit der eines Dateideskriptors verbunden wird und über den eine bidirektionale Kommunikation möglich ist.
- Die Semantik einer Kommunikation.

Die Semantik einer Kommunikation umschließt bei jeder Verbindung eine Teilmenge der folgenden Punkte:

1. Daten werden in der Reihenfolge empfangen, in der sie abgeschickt worden sind.
2. Daten kommen nicht doppelt an.
3. Daten werden zuverlässig übermittelt.
4. Einzelne Pakete kommen in der originalen Form an (d.h. sie werden weder zerstückelt noch mit anderen Paketen kombiniert).

5. Nachrichten außerhalb des normalen Kommunikationsstromes (*out-of-band messages*) werden unterstützt.
6. Die Kommunikation erfolgt verbindungs-orientiert, womit die Notwendigkeit entfällt, sich bei jedem Paket identifizieren zu müssen.

Die folgende Tabelle zeigt die Varianten, die von der Berkeley-Socket-Schnittstelle unterstützt werden:

Name	1	2	3	4	5	6
<i>SOCK_STREAM</i>	*	*	*		*	*
<i>SOCK_DGRAM</i>				*		
<i>SOCK_SEQPACKET</i>	*	*	*	*	*	*
<i>SOCK_RDM</i>	*	*	*	*		

Von diesen Varianten kommt *SOCK_STREAM* den Pipelines am nächsten, wenn davon abgesehen wird, daß die Verbindungen bei Pipelines nur unidirektional sind. Die Variante *SOCK_STREAM* lässt sich ziemlich direkt auf TCP abbilden, während UDP ziemlich genau durch *SOCK_DGRAM* widerspiegelt wird. Die Varianten *SOCK_SEQPACKET* (TCP-basiert) und *SOCK_RDM* (UDP-basiert) fügen hier noch weitere Funktionalitäten hinzu. Allerdings fand *SOCK_RDM* nicht den Weg in den POSIX-Standard und wird auch von einigen Implementierungen nicht angeboten. Im weiteren Verlauf werden wir uns nur mit *SOCK_STREAM*-Sockets beschäftigen.

4.4 Aufbau einer Netzwerk-Verbindung

Bis zu einem gewissen Grad mag hier zur Einführung die Analogie mit unserem Telefonsystem hilfreich sein. Bevor Sie Telefonanrufe entgegennehmen oder selbst anrufen können, benötigen Sie einen Telefonanschluss. Dieser Anschluss wird mit dem Systemaufruf `socket()` erzeugt:

```
int sfd = socket(domain, type, protocol);
```

Bei *domain* wird hier normalerweise *PF_INET* angegeben, um das IPv4-Protokoll auszuwählen. *PF* steht dabei für *protocol family*. Bei *type* kann eine der unterstützten Semantiken ausgewählt werden, also beispielsweise *SOCK_STREAM*. Der dritte Parameter *protocol* erlaubt in einigen Fällen eine weitere Selektion. Normalerweise wird hier schlicht 0 angegeben.

Nachdem der Anschluss existiert, fehlt noch eine zugeordnete Telefonnummer. Um bei der Analogie zu bleiben, haben wir eine Vorwahl (IP-Adresse) und eine Durchwahl (Port-Nummer). Auf einem Rechner können mehrere IP-Adressen zur Verfügung stehen. Es ist dabei möglich, nur eine dieser IP-Adressen zu verwenden oder alle, die zur Verfügung stehen. Bei den Port-Nummern ist eine automatische Zuteilung durch das Betriebssystem möglich. Alternativ ist es auch möglich, sich selbst eine Port-Nummer auszuwählen. Diese darf aber noch nicht vergeben sein und muss bei nicht-privilegierten Prozessen eine Nummer jenseits des Bereiches der wohldefinierten Port-Nummern sein, also typischerweise mindestens 1024 betragen. Die Verknüpfung eines Anschlusses mit einer vollständigen Adresse erfolgt mit dem Systemaufruf `bind()`:

```
struct sockaddr_in address = {0};
address.sin_family = AF_INET;
address.sin_addr.s_addr = htonl(INADDR_ANY);
address.sin_port = htons(port);
bind(sfd, (struct sockaddr *) &address, sizeof address);
```

Die Datenstruktur `struct sockaddr_in` repräsentiert Adressen für IPv4, die aus einer IP-Adresse und einer Port-Nummer bestehen. Das Feld `sin_family` legt den Adreßraum fest. Hier gibt es passend zur Protokollfamilie `PF_INET` nur `AF_INET` (`AF` steht hier für *address family*). Bei dem Feld `sin_addr.s_addr` lässt sich die IP-Adresse angeben. Mit `INADDR_ANY` übernehmen wir alle IP-Adressen, die zum eigenen Rechner gehören. Das Feld `sin_port` spezifiziert die Port-Nummer. Da Netzwerkadressen grundsätzlich nicht von der Byte-Anordnung eines Rechners abhängen dürfen, wird mit `htonl` (*host to network long*) der 32-Bit-Wert der IP-Adresse in die standardisierte Form konvertiert. Analog konvertiert `htons`() (*host to network short*) den 16-Bit-Wert `port` in die standardisierte Byte-Reihenfolge. Wenn die Port-Nummer vom Betriebssystem zugeteilt werden soll, kann bei `sin_port` auch einfach 0 angegeben werden. Der Datentyp `struct sockaddr_in` ist eine spezielle Variante des Datentyps `struct sockaddr`. Letzterer sieht nur ein Feld `sin_family` vor und ein generelles Datenfeld `sa_data`, das umfangreich genug ist, um alle unterstützten Adressen unterzubringen. Bei `bind`() wird der von `socket`() erhaltene Deskriptor angegeben (hier `sfd`), ein Zeiger, der auf eine Adresse vom Typ `struct sockaddr` verweist, und die tatsächliche Länge der Adresse, die normalerweise kürzer ist als die des Typs `struct sockaddr`. Schön sind diese Konstruktionen nicht, aber C bietet eben keine objekt-orientierten Konzepte, wenngleich die Berkeley-Socket-Schnittstelle sehr wohl polymorph und damit objekt-orientiert ist.

Grundsätzlich kann auf `bind`() auch verzichtet werden. Wenn dies geschieht, entspricht dies der Angabe von `INADDR_ANY` bei dem Feld `sin_addr.s_addr` und einer 0 bei `sin_port`. Das bedeutet dann, daß die Adresse vollständig von dem Betriebssystem vergeben wird. Die auf diese Weise erhaltene Adresse lässt sich mit der Funktion `getsockname`() ermitteln, sobald eine Verbindung eröffnet wird oder eingehende Verbindungen angenommen werden können.

Damit eingehende Verbindungen (oder Anrufe in unserer Telefon-Analogie) entgegenommen werden können, muss `listen`() aufgerufen werden:

```
listen(sfd, SOMAXCONN);
```

Wohlgermerkt, nach `listen`() kann der Anschluss „klingeln“, aber noch sind keine Vorbereitungen getroffen, das Klingeln zu hören oder den Hörer abzunehmen. Der zweite Parameter bei `listen`() gibt an, wieviele Kommunikationspartner es gleichzeitig klingeln lassen dürfen. `SOMAXCONN` ist hier das Maximum, das die jeweilige Implementierung erlaubt. Der folgende Programmtext zeigt, wie mit `socket`() und `listen`() ein Anschluss angelegt werden kann:

Programm 4.1: Einrichtung eines Anschlusses (*newsocket.c*)

```

1 #include <sys/socket.h>
2 #include <netinet/in.h>
3 #include <stdio.h>
4 #include <stdlib.h>
5
6 void print_ip_addr(in_addr_t ipaddr) {
7     if (ipaddr == INADDR_ANY) {
8         printf("INADDR_ANY");
9     } else {
10        uint32_t addr = ntohl(ipaddr);
11        printf("%d.%d.%d.%d",
12            addr>>24, (addr>>16)&0xff,
13            (addr>>8)&0xff, addr&0xff);
14    }
15 }
16
```

```

17 int main() {
18     int sfd = socket(PF_INET, SOCK_STREAM, 0);
19     if (sfd < 0) exit(1);
20     if (listen(sfd, SOMAXCONN) < 0) exit(2);
21     struct sockaddr address;
22     socklen_t len = sizeof address;
23     if (getsockname(sfd, &address, &len) < 0) exit(3);
24     struct sockaddr_in *inaddr = (struct sockaddr_in *) &address;
25     printf("This is the address of my new socket: \n");
26     printf("IP Address: \n"); print_ip_addr(inaddr->sin_addr.s_addr);
27     printf("\n");
28     printf("Port Number: \n", (int) ntohs(inaddr->sin_port));
29 }

```

Die zugeteilte Port-Nummer wird dann mit `getsockname()` ermittelt. Interessanterweise gibt es keine von POSIX unterstützte Methode, alle verfügbaren IP-Adressen zu ermitteln. Dies ist allerdings kein gravierendes Problem, da typischerweise `INADDR_ANY` oder eine explizit (z.B. über die Kommandozeile) angegebene IP-Adresse verwendet wird. Die Entgegennahme eines Anrufes erfolgt mit `accept()`:

```

struct sockaddr client_addr;
socklen_t client_addr_len = sizeof client_addr;
int fd = accept(sfd, &client_addr, &client_addr_len);

```

Liegt noch kein Anruf vor, blockiert `accept()` bis zum nächsten Anruf. Wenn mit `accept()` ein Anruf eingeht, wird ein Dateideskriptor auf den bidirektionalen Verbindungskanal zurückgeliefert. Normalerweise speichert `accept()` die Adresse des Klienten beim angegebenen Zeiger ab. Wenn als Zeiger 0 angegeben wird, entfällt dies.

Das folgende Beispiel zeigt einen einfachen Zeitdienst analog zu `daytime`, der die Port-Nummer 11011 verwendet:

Programm 4.2: Ein einfacher Zeitdienst (`timeserver.c`)

```

1 #include <netinet/in.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <string.h>
5 #include <sys/socket.h>
6 #include <sys/time.h>
7 #include <time.h>
8 #include <unistd.h>
9
10 #define PORT 11011
11
12 int main () {
13     struct sockaddr_in address = {0};
14     address.sin_family = AF_INET;
15     address.sin_addr.s_addr = htonl(INADDR_ANY);
16     address.sin_port = htons(PORT);
17
18     int sfd = socket(PF_INET, SOCK_STREAM, 0);
19     int optval = 1;
20     if (sfd < 0 ||
21         setsockopt(sfd, SOL_SOCKET, SO_REUSEADDR,
22                 &optval, sizeof optval) < 0 ||

```

```

23     bind(sfd, (struct sockaddr *) &address,
24           sizeof address) < 0 ||
25     listen(sfd, SOMAXCONN) < 0) {
26     perror("socket"); exit(1);
27 }
28
29 int fd;
30 while ((fd = accept(sfd, 0, 0)) >= 0) {
31     char timebuf[32]; time_t clock; time(&clock);
32     ctime_r(&clock, timebuf);
33     write(fd, timebuf, strlen(timebuf)); close(fd);
34 }
35 }

```

Zusätzlich kommt hier noch `setsockopt()` in den Zeilen 21 und 22 zum Einsatz, um die Option `SO_REUSEADDR` einzuschalten. Dies empfiehlt sich immer, wenn eine feste Port-Nummer verwendet wird. Fehlt diese Option, kann es passieren, daß bei einem Neustart des Dienstes die Port-Nummer nicht sofort wieder zur Verfügung steht, da noch alte Verbindungen nicht vollständig abgewickelt worden sind. In den Zeilen 31 bis 33 wird jeweils die aktuelle Zeitangabe in eine Zeichenkette konvertiert und mit Hilfe von `write()` an die zuvor geöffnete Verbindung verschickt. Mit `close()` wird die Verbindung dann auf der eigenen Seite aufgehängt.

Die Verbindungsaufnahme erfolgt mit `connect()`, wobei der bei `socket()` zurückgelieferte Deskriptor anzugeben ist und die Adresse des Kommunikationspartners:

```
connect(fd, (struct sockaddr *) &addr, sizeof addr);
```

Der folgende Programmtext zeigt einen Klienten, der den vorgestellten Zeitdienst kontaktiert, die Zeitangabe entgegennimmt und sie ausgibt:

Programm 4.3: Ein Klient des Zeitdienstes (*timeclient.c*)

```

1 #include <netdb.h>
2 #include <netinet/in.h>
3 #include <stdio.h>
4 #include <stdlib.h>
5 #include <string.h>
6 #include <sys/socket.h>
7 #include <unistd.h>
8
9 #define PORT 11011
10
11 int main (int argc, char** argv) {
12     char* cmdname = *argv++; --argc;
13     if (argc != 1) {
14         fprintf(stderr, "Usage:_%s_host\n", cmdname);
15         exit(1);
16     }
17     char* hostname = *argv;
18     struct hostent* hp;
19     if ((hp = gethostbyname(hostname)) == 0) {
20         fprintf(stderr, "unknown_host:_%s\n", hostname);
21         exit(1);
22     }
23     char* hostaddr = hp->h_addr_list[0];

```

```

24  struct sockaddr_in addr = {0};
25  addr.sin_family = AF_INET;
26  memmove((void *) &addr.sin_addr, (void *) hostaddr, hp->h_length);
27  addr.sin_port = htons(PORT);
28
29  int fd;
30  if ((fd = socket(PF_INET, SOCK_STREAM, 0)) < 0) {
31      perror("socket"); exit(1);
32  }
33  if (connect(fd, (struct sockaddr *) &addr, sizeof addr) < 0) {
34      perror("connect"); exit(1);
35  }
36
37  char buffer[BUFSIZ]; ssize_t nbytes;
38  while((nbytes = read(fd, buffer, sizeof buffer)) > 0 &&
39      write(1, buffer, nbytes) == nbytes);
40  }

```

timeclient

Der Klient erhält über die Kommandozeile den Namen des Rechners, auf dem der Zeitdienst zur Verfügung steht. Für die Abbildung eines Rechnernamens in eine IP-Adresse wird die Funktion `gethostbyname()` benötigt, die im Erfolgsfall eine oder mehrere IP-Adressen liefert, unter denen sich der Rechner erreichen lässt. In Zeile 23 wird die erste IP-Adresse ausgewählt, die bereits in der richtigen Byte-Anordnung vorliegt. Mit Hilfe von `memmove()` wird sie zu `addr.sin_addr` in Zeile 26 kopiert. Danach kann in Zeile 33 mit der zuvor erzeugten Socket die gewünschte Verbindung mit Hilfe von `connect()` aufgebaut werden. Die `while`-Schleife in Zeile 38 gibt dann alles aus, was sich über die Netzwerkverbindung einlesen ließ.

4.5 Gepufferte Ein- und Ausgabe für Netzwerkverbindungen

Die Ein- und Ausgabe über Netzwerkverbindungen bringt in Vergleich zur Behandlung von Dateien und interaktiven Benutzern einige Veränderungen mit sich. Wenn eine Verbindung des Typs `SOCK_STREAM` zum Einsatz gelangt, so kommen die Daten zwar in der korrekten Reihenfolge an, jedoch nicht in der ursprünglichen Paketisierung. Als ursprüngliche Pakete werden hier die Daten betrachtet, die mit Hilfe eines einzigen Aufrufs von `write()` geschrieben werden:

```

const char greeting[] = "Hi,_how_are_you?\r\n";
ssize_t nbytes = write(sfd, greeting, sizeof greeting);

```

Wenn beispielsweise bei einer Netzwerkverbindung immer vollständige Zeilen mit `write()` geschrieben werden, so ist es möglich, daß die korrespondierende `read()`-Operation nur einen Teil einer Zeile zurückliefert oder auch ein Fragment, das sich über mehr als eine Zeile erstreckt. Diese Problematik legt es nahe, nur zeichenweise einzulesen, wenn genau eine einzelne Zeile eingelesen werden soll:

```

char ch;
stralloc line = {0};
while (read(fd, &ch, sizeof ch) == 1 && ch != '\n') {
    stralloc_append(&line, &ch);
}

```

Diese Vorgehensweise ist jedoch außerordentlich ineffizient, weil Systemaufrufe wie `read()` zu einem Kontextwechsel zwischen dem aufrufenden Prozess und dem Betriebssystem führen. Wenn ein Kontextwechsel für jedes einzulesende Byte initiiert wird, dann ist der betroffene Rechner mehr mit Kontextwechseln als mit sinnvollen Tätigkeiten beschäftigt. Wenn jedoch mit

```
char buf[512];
ssize_t nbytes = read(fd, buf, sizeof buf)
```

eingelassen wird, ist möglicherweise mehr als nur die gewünschte Zeile in `buf` zu finden. Oder auch nur ein Teil der Zeile.

Entsprechend ist eine gepufferte Eingabe notwendig, bei der die Eingabe-Operationen aus einem Puffer versorgt werden, der, wenn er leer wird, mit Hilfe einer `read()`-Operation aufzufüllen ist. Die Datenstruktur für einen Eingabe-Puffer benötigt entsprechend einen Dateideskriptor, einen Puffer und einen Positionszeiger innerhalb des Puffers:

```
typedef struct inbuf {
    int fd;
    stralloc buf;
    unsigned int pos;
} inbuf;
```

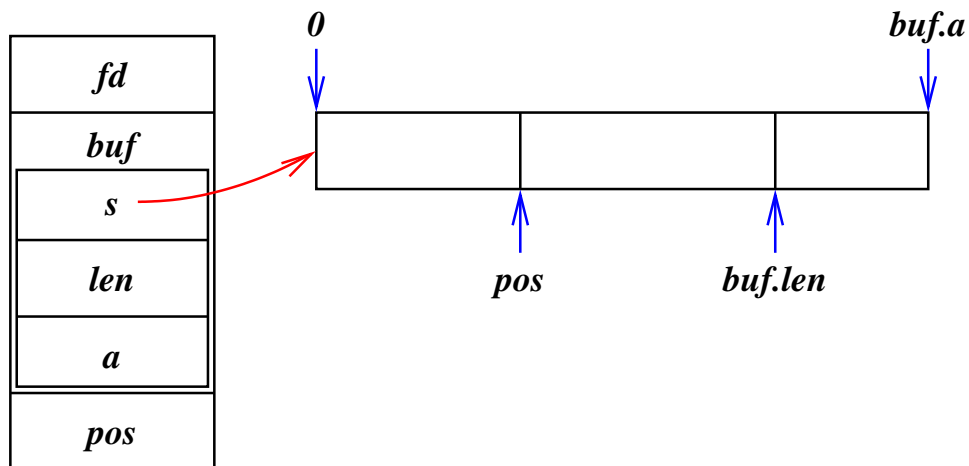


Abbildung 4.2: Struktur eines Eingabe-Puffers

Abbildung 4.2 dient der Illustrierung des Eingabe-Puffers. Als Puffer-Bereich dient `buf.s` mit einem zulässigen Indextbereich von `0` bis `buf.a-1`. Entsprechend dem letzten erfolgreichen Aufruf von `read()` ergibt sich ein Füllgrad des Puffers, der von `buf.len` repräsentiert wird. Der Positionszeiger `pos` begann unmittelbar nach der `read()`-Operation auf Position `0` und wandert bei jeder Einlese-Operation aus dem Puffer der Grenze von `buf.len` entgegen. Wird die Grenze erreicht, so ist die nächste `read()`-Operation fällig.

Der folgende Programmtext zeigt die entsprechende Schnittstelle:

Programm 4.4: Schnittstelle für den Eingabe-Puffer (`inbuf.h`)

```
1 #ifndef INBUF_H
2 #define INBUF_H
3
4 #include <stralloc.h>
5 #include <unistd.h>
```



```

6
7 typedef struct inbuf {
8     int fd;
9     stralloc buf;
10    unsigned int pos;
11 } inbuf;
12
13 /* set size of input buffer */
14 int inbuf_alloc(inbuf* ibuf, unsigned int size);
15
16 /* works like read(2) but from ibuf */
17 ssize_t inbuf_read(inbuf* ibuf, void* buf, size_t size);
18
19 /* works like fgetc but from ibuf */
20 int inbuf_getchar(inbuf* ibuf);
21
22 /* move backward one position */
23 int inbuf_back(inbuf* ibuf);
24
25 /* release storage associated with ibuf */
26 void inbuf_free(inbuf* ibuf);
27
28 #endif

```

Die Funktion `inbuf_alloc()` dient dazu, die Größe des Puffers einzurichten, wobei eine sinnvolle Voreinstellung automatisch gewählt wird, wenn der Aufruf dieser Funktion unterbleibt. Als Einlese-Operationen vom Puffer dienen `inbuf_read()` und `inbuf_getchar()`, die sich in ihrer Aufrufsemantik an `read()` bzw. `fgetc()` orientieren. Ein zuviel gelesenes Zeichen kann mit `inbuf_back()` wieder zum erneuten Einlesen zur Verfügung gestellt werden. Mit `inbuf_free()` wird der Puffer deallokiert.

Programm 4.5: Implementierung des Eingabe-Puffers (*inbuf.c*)

```

1 #include <errno.h>
2 #include <string.h>
3 #include <unistd.h>
4 #include "inbuf.h"
5
6 /* set size of input buffer */
7 int inbuf_alloc(inbuf* ibuf, unsigned int size) {
8     return stralloc_ready(&ibuf->buf, size);
9 }
10
11 /* works like read(2) but from ibuf */
12 ssize_t inbuf_read(inbuf* ibuf, void* buf, size_t size) {
13     if (size == 0) return 0;
14     if (ibuf->pos >= ibuf->buf.len) {
15         if (ibuf->buf.a == 0 && !inbuf_alloc(ibuf, 512)) return -1;
16         /* fill input buffer */
17         ssize_t nbytes;
18         do {
19             errno = 0;
20             nbytes = read(ibuf->fd, ibuf->buf.s, ibuf->buf.a);
21         } while (nbytes < 0 && errno == EINTR);

```

```

22     if (nbytes <= 0) return nbytes;
23     ibuf->buf.len = nbytes;
24     ibuf->pos = 0;
25 }
26 ssize_t nbytes = ibuf->buf.len - ibuf->pos;
27 if (size < nbytes) nbytes = size;
28 memcpy(buf, ibuf->buf.s + ibuf->pos, nbytes);
29 ibuf->pos += nbytes;
30 return nbytes;
31 }
32
33 /* works like fgetc but from ibuf */
34 int inbuf_getchar(inbuf* ibuf) {
35     char ch;
36     ssize_t nbytes = inbuf_read(ibuf, &ch, sizeof ch);
37     if (nbytes <= 0) return -1;
38     return ch;
39 }
40
41 /* move backward one position */
42 int inbuf_back(inbuf* ibuf) {
43     if (ibuf->pos == 0) return 0;
44     ibuf->pos--;
45     return 1;
46 }
47
48 /* release storage associated with ibuf */
49 void inbuf_free(inbuf* ibuf) {
50     stralloc_free(&ibuf->buf);
51 }

```

Der vorstehende Programmtext zeigt die zugehörige Implementierung, wobei hier insbesondere *inbuf_read()* interessant ist. In Zeile 14 wird untersucht, ob der Puffer bereits geleert ist, d.h. ob *pos* bereits *buf.len* erreicht hat. Falls ja, wird in den Zeilen 15 bis 24 der Puffer neu gefüllt. In Zeile 15 wird zunächst untersucht, ob der Puffer bereits allokiert worden ist. Falls nicht, wird dies mit der Standardgröße von 512 Bytes versucht. In Zeile 20 erfolgt die *read()*-Operation, bei der grundsätzlich versucht wird, den gesamten Puffer zu füllen. Allerdings ist damit zu rechnen, daß die Zahl der tatsächlich gelesenen Bytes *nbytes* niedriger als *buf.a* ist. Dies liegt daran, daß das Betriebssystem uns bereits vorhandene Daten sofort zur Verfügung stellt, selbst wenn es sich um eine geringere Quantität als angefordert handelt. Dies stellt sicher, daß effizientes Einlesen mit großen Puffergrößen ohne unnötiges Blockieren möglich ist. Die *read()*-Operation selbst ist in eine Schleife in den Zeilen 18 bis 21 eingebettet, die sicherstellt, daß es zu einem erneuten Versuch kommt, falls *read()* wegen einer Signalunterbrechung nicht erfolgreich sein konnte.

Sobald sichergestellt ist, daß mindestens ein Byte in dem Puffer verfügbar ist, wird *nbytes* in Zeile 26 auf die maximal mögliche Rückgabequantität gesetzt. Wurden weniger verlangt, so wird *nbytes* in Zeile 27 entsprechend zurückgesetzt. In Zeile 28 wird die zurückzuliefernde Quantität an Bytes aus dem Puffer in *buf* mit Hilfe von *memcpy()* kopiert. Der erste Parameter (*buf*) zeigt dabei auf das Ziel, der zweite Parameter (*buf.s* + *pos*) auf die Quelle und der dritte Parameter gibt die Zahl der zu kopierenden Bytes an (*nbytes*). Nach der Kopieraktion wird *pos* entsprechend aktualisiert.

Die Ausgabe sollte ebenfalls gepuffert erfolgen, um die Zahl der Systemaufrufe zu minimieren. Ein Positionszeiger ist nicht erforderlich, wenn Puffer grundsätzlich vollständig

an `write()` übergeben werden. Hier ist das einzige Problem, daß die `write()`-Operation unter Umständen nicht den gesamten gewünschten Umfang akzeptiert und nur einen Teil der zu schreibenden Bytes akzeptiert und entsprechend eine geringere Quantität als Wert zurückgibt.

Der folgende Programmtext zeigt die Schnittstelle für den Ausgabe-Puffer:

Programm 4.6: Schnittstelle für den Ausgabe-Puffer (*outbuf.h*)

```

1 #ifndef OUTBUF_H
2 #define OUTBUF_H
3
4 #include <stralloc.h>
5 #include <unistd.h>
6
7 typedef struct outbuf {
8     int fd;
9     stralloc buf;
10 } outbuf;
11
12 /* works like write(2) but to obuf */
13 ssize_t outbuf_write(outbuf* obuf, void* buf, size_t size);
14
15 /* works like fputc but to obuf */
16 int outbuf_putchar(outbuf* obuf, char ch);
17
18 /* write contents of obuf to the associated fd */
19 int outbuf_flush(outbuf* obuf);
20
21 /* release storage associated with obuf */
22 void outbuf_free(outbuf* obuf);
23
24 #endif

```

Die Funktion `outbuf_write()` schreibt in den gegebenen Puffer und entspricht ansonsten dem Systemaufruf `write()`. Mit Hilfe von `outbuf_putchar()` können bequem einzelne Zeichen in den Puffer ausgegeben werden. Beide Schreiboperationen führen nur zur Verlängerung des Pufferinhalts, ohne daß dieser mit Hilfe einer `write()`-Operation geleert wird. Letzteres ist nur durch den Aufruf von `outbuf_flush()` möglich. Wenn der Puffer nicht mehr benötigt wird, kann er durch `outbuf_free()` freigegeben werden. Es folgt die zugehörige Implementierung:

Programm 4.7: Implementierung des Ausgabe-Puffers (*outbuf.c*)

```

1 #include <errno.h>
2 #include <stralloc.h>
3 #include <string.h>
4 #include "outbuf.h"
5
6 /* works like write(2) but to obuf */
7 ssize_t outbuf_write(outbuf* obuf, void* buf, size_t size) {
8     if (size == 0) return 0;
9     if (!stralloc_readyplus(&obuf->buf, size)) return -1;
10    memcpy(obuf->buf.s + obuf->buf.len, buf, size);
11    obuf->buf.len += size;
12    return size;

```

```

13 }
14
15 /* works like fputc but to obuf */
16 int outbuf_putchar(outbuf* obuf, char ch) {
17     if (outbuf_write(obuf, &ch, sizeof ch) <= 0) return -1;
18     return ch;
19 }
20
21 /* write contents of obuf to the associated fd */
22 int outbuf_flush(outbuf* obuf) {
23     ssize_t left = obuf->buf.len; ssize_t written = 0;
24     while (left > 0) {
25         ssize_t nbytes;
26         do {
27             errno = 0;
28             nbytes = write(obuf->fd, obuf->buf.s + written, left);
29         } while (nbytes < 0 && errno == EINTR);
30         if (nbytes <= 0) return 0;
31         left -= nbytes; written += nbytes;
32     }
33     obuf->buf.len = 0;
34     return 1;
35 }
36
37 /* release storage associated with obuf */
38 void outbuf_free(outbuf* obuf) {
39     stralloc_free(&obuf->buf);
40 }

```

In `outbuf_write()` wird in Zeile 9 darauf geachtet, daß der Puffer genügend Platz für den aufzunehmenden Inhalt aufweist, wonach in Zeile 10 der Kopiervorgang mit Hilfe von `memcpy()` durchgeführt werden kann. Danach muss nur noch `buf.len` in Zeile 11 angepaßt werden. In der Funktion `outbuf_flush()` gibt es zwei Schleifen. Die äußere Schleife in den Zeilen 24 bis 32 sorgt dafür, daß der gesamte Puffer-Inhalt geschrieben wird, da einzelne `write()`-Operationen die Freiheit haben, nur einen Teil umzusetzen. Mit Hilfe der Variablen `left` und `written` wird vermerkt, wieviel noch zu schreiben ist bzw. wieviel bereits geschrieben wurde. Die innere Schleife in den Zeilen 26 bis 29 wiederholt die `write()`-Operation im Falle von Unterbrechungen.

4.6 Ein- und Ausgabe von Paketen für Netzwerkverbindungen

Zwischen Dienste-Anbietern und ihren Klienten auf dem Netzwerk besteht häufig ein ähnliches Verhältnis wie zwischen einer Shell und dem zugehörigen Benutzer: Der Klient gibt ein Kommando, das typischerweise mit dem Zeilentrenner CR LF³ beendet wird, und der Dienste-Anbieter sendet darauf eine Antwort zurück,

- die zum Ausdruck bringt, ob das Kommando erfolgreich verlief oder fehlschlug, und
- einen Antworttext über eine oder mehrere Zeilen bringt.

³carriage return, gefolgt von line feed.

Es gibt keine zwingende Notwendigkeit, bei einem Protokoll Zeilentrenner zu verwenden. Es gibt auch Alternativen wie die von Dan Bernstein vorgeschlagenen Net-Strings⁴, jedoch erlaubt die Konvention mit dem Zeilentrenner CR LF die interaktive Benutzung eines Dienstes mit dem *telnet*-Kommando.

Bei jedem Protokoll lohnt es sich, eine gewisse Grundstruktur über die generelle Syntax von Kommandos und deren Antworten festzulegen. Hier ist ein Beispiel für eine (nicht sehr konstruktive) Sitzung mit einem Dienst, über den E-Mails zugestellt werden können. Zum Einsatz kommt hier das SMTP-Protokoll⁵:

```
clonard$ telnet mail.rz.uni-ulm.de smtp
Trying 134.60.1.11...
Connected to mail.rz.uni-ulm.de.
Escape character is '^]'.
220 mail.uni-ulm.de ESMTP Sendmail 8.14.2/8.14.2; Mon, 2 Jun 2008 10:18:51 +0200 (MEST)
help
214-2.0.0 This is sendmail version 8.14.2
214-2.0.0 Topics:
214-2.0.0      HELO      EHLO      MAIL      RCPT      DATA
214-2.0.0      RSET      NOOP      QUIT      HELP      VRFY
214-2.0.0      EXPN      VERB      ETRN      DSN       AUTH
214-2.0.0      STARTTLS
214-2.0.0 For more info use "HELP <topic>".
214-2.0.0 To report bugs in the implementation see
214-2.0.0      http://www.sendmail.org/email-addresses.html
214-2.0.0 For local information send email to Postmaster at your site.
214 2.0.0 End of HELP info
huhu
500 5.5.1 Command unrecognized: "huhu"
helo clonard.mathematik.uni-ulm.de
250 mail.uni-ulm.de Hello borchert@clonard.mathematik.uni-ulm.de [134.60.66.13], pleased to meet y
quit
221 2.0.0 mail.uni-ulm.de closing connection
Connection to mail.rz.uni-ulm.de closed by foreign host.
clonard$
```

Beim SMTP-Protokoll erfolgt zunächst eine Begrüßung des Dienste-Anbieters. Die Begrüßung oder auch eine andere Antwort des Anbieters besteht aus einer dreistelligen Nummer, einem Leerzeichen oder einem Minus und beliebigem Text, der durch CR LF abgeschlossen wird. Die erste Ziffer der dreistelligen Nummer legt hier fest, ob ein Erfolg oder ein Problem vorliegt. Die beiden weiteren Ziffern werden zur feineren Unterscheidung der Rückmeldung verwendet. Eine führende 2 bedeutet Erfolg, eine 4 signalisiert ein temporäres Problem und eine 5 signalisiert einen permanenten Fehler.

In der Beispielsitzung ist das erste Kommando ein „help“, gefolgt von CR LF. Da die Antwort sich über mehrere Zeilen erstreckt, werden alle Zeilen, hinter der noch mindestens eine folgt, mit einem Minuszeichen hinter der dreistelligen Zahl gekennzeichnet. Danach wurde mit dem unbekanntem Kommando „huhu“ eine Fehlermeldung provoziert, die durch den Code 500 signalisiert wurde. Das SMTP-Protokoll erlaubt auch eine Fortsetzung des Dialogs nach Fehlern, so daß dann noch ein „helo“-Kommando akzeptiert wurde. Die Verbindung wurde mit dem „quit“-Befehl beendet.

Ungeachtet seines Namens ist das SMTP-Protokoll nicht mehr sehr einfach. Um ein überschaubareres Beispiel zu erhalten, wird im weiteren Text von folgenden Konventionen ausgegangen:

- Kommandos bestehen aus einem Kommandonamen (nur aus Kleinbuchstaben), ei-

⁴Siehe <http://cr.yp.to/proto/netstrings.txt>

⁵SMTP steht für *simple mail transfer protocol*, siehe RFC 2821

nem Leerzeichen und genau einem Parameter. Abgeschlossen wird ein Kommando durch die Sequenz CR LF.

- Antworten beginnen mit „S“, „F“ oder „C“, gefolgt von einem beliebigen Text, der durch CR LF abgeschlossen wird. „S“ deutet einen Erfolg an (*success*), „F“ einen Fehlschlag (*failure*) und „C“ gibt an, daß noch weitere Zeilen folgen.

Sobald dieser syntaktische Rahmen festliegt, ist es möglich, entsprechende Datenstrukturen für Anfragen und deren Antworten zu deklarieren und zugehörige Ein- und Ausgabefunktionen anzubieten wie die folgenden Programmtexte zeigen.⁶ Es ist dabei sinnvoll, jeweils die Ein- und Ausgabefunktionen zusammen zu belassen, da beide Funktionen jeweils genau zusammenpassen müssen.

Programm 4.8: Schnittstelle für MXP-Anfrage-Pakete (*mxprequest.h*)

```

1 #ifndef MXP_REQUEST_H
2 #define MXP_REQUEST_H
3
4 #include <stralloc.h>
5 #include "inbuf.h"
6 #include "outbuf.h"
7
8 typedef struct mxp_request {
9     stralloc keyword;
10    stralloc parameter;
11 } mxp_request;
12
13 /* read one request from the given input buffer */
14 int read_mxp_request(inbuf* ibuf, mxp_request* request);
15
16 /* write one request to the given outbuf buffer */
17 int write_mxp_request(outbuf* obuf, mxp_request* request);
18
19 /* release resources associated with request */
20 void free_mxp_request(mxp_request* request);
21
22 #endif

```

Programm 4.9: Schnittstelle für MXP-Antwort-Pakete (*mxpresponse.h*)

```

1 #ifndef MXP_RESPONSE_H
2 #define MXP_RESPONSE_H
3
4 #include "inbuf.h"
5 #include "outbuf.h"
6
7 typedef enum mxp_status {
8     MXP_SUCCESS = 'S',
9     MXP_FAILURE = 'F',
10    MXP_CONTINUATION = 'C',
11 } mxp_status;
12
13 typedef struct mxp_response {

```

⁶ MXP steht hierbei für das *mutual exclusion protocol*, das im übernächsten Abschnitt vorgestellt wird.

```

14     mxp_status status;
15     stralloc message;
16 } mxp_response;
17
18 /* write one (possibly partial) response to the given output buffer */
19 int write_mxp_response(outbuf* obuf, mxp_response* response);
20
21 /* read one (possibly partial) response from the given input buffer */
22 int read_mxp_response(inbuf* ibuf, mxp_response* response);
23
24 void free_mxp_response(mxp_response* response);
25
26 #endif

```

Programm 4.10: Ein- und Ausgabe von MXP-Anfrage-Paketen (*mxprequest.c*)

```

1 #include "inbuf.h"
2 #include "outbuf.h"
3 #include "mxprequest.h"
4
5 static int read_keyword(inbuf* ibuf, stralloc* keyword) {
6     int ch;
7     if (!stralloc_copys(keyword, "")) return 0;
8     while ((ch = inbuf_getchar(ibuf)) >= 0 &&
9           ch >= 'a' && ch <= 'z') {
10        if (!stralloc_readyplus(keyword, 1)) return 0;
11        keyword->s[keyword->len++] = ch;
12    }
13    if (ch >= 0) inbuf_back(ibuf);
14    return keyword->len > 0;
15 }
16
17 static int read_parameter(inbuf* ibuf, stralloc* parameter) {
18     int ch;
19     if (!stralloc_copys(parameter, "")) return 0;
20     while ((ch = inbuf_getchar(ibuf)) >= 0 &&
21           ch != '\r' && ch != '\n' && ch != '/' && ch != '\0') {
22        if (!stralloc_readyplus(parameter, 1)) return 0;
23        parameter->s[parameter->len++] = ch;
24    }
25    if (ch >= 0) inbuf_back(ibuf);
26    return parameter->len > 0;
27 }
28
29 static int expect_delimiter(inbuf* ibuf, char delimiter) {
30     int ch = inbuf_getchar(ibuf);
31     if (ch < 0) return 0;
32     return ch == delimiter;
33 }
34
35 /* read one request from the given input buffer */
36 int read_mxp_request(inbuf* ibuf, mxp_request* request) {
37     return

```

```

38     read_keyword(ibuf, &request->keyword) &&
39     expect_delimiter(ibuf, '_') &&
40     read_parameter(ibuf, &request->parameter) &&
41     expect_delimiter(ibuf, '\r') &&
42     expect_delimiter(ibuf, '\n');
43 }
44
45 /* write one request to the given outbuf buffer */
46 int write_mxp_request(outbuf* obuf, mxp_request* request) {
47     return
48         outbuf_write(obuf, request->keyword.s,
49                     request->keyword.len) == request->keyword.len &&
50         outbuf_putchar(obuf, '_') == '_' &&
51         outbuf_write(obuf, request->parameter.s,
52                     request->parameter.len) == request->parameter.len &&
53         outbuf_putchar(obuf, '\r') == '\r' &&
54         outbuf_putchar(obuf, '\n') == '\n';
55 }
56
57 /* release resources associated with request */
58 void free_mxp_request(mxp_request* request) {
59     stralloc_free(&request->keyword);
60     stralloc_free(&request->parameter);
61 }

```

Programm 4.11: Ein- und Ausgabe von MXP-Antwort-Paketen (*mxpresponse.c*)

```

1 #include "inbuf.h"
2 #include "mxpresponse.h"
3 #include "outbuf.h"
4
5 static int read_message(inbuf* ibuf, stralloc* message) {
6     int ch;
7     if (!stralloc_copys(message, "")) return 0;
8     while ((ch = inbuf_getchar(ibuf)) >= 0 &&
9           ch != '\r' && ch != '\n') {
10        if (!stralloc_readyplus(message, 1)) return 0;
11        message->s[message->len++] = ch;
12    }
13    if (ch >= 0) inbuf_back(ibuf);
14    return 1;
15 }
16
17 static int expect_delimiter(inbuf* ibuf, char delimiter) {
18     int ch = inbuf_getchar(ibuf);
19     if (ch < 0) return 0;
20     return ch == delimiter;
21 }
22
23 /* read one (possibly partial) response from the given input buffer */
24 int read_mxp_response(inbuf* ibuf, mxp_response* response) {
25     int ch = inbuf_getchar(ibuf);
26     switch (ch) {

```



```

27     case MXP_SUCCESS:
28     case MXP_FAILURE:
29     case MXP_CONTINUATION:
30         response->status = ch;
31         break;
32     default:
33         return 0;
34     }
35     if (!read_message(ibuf, &response->message)) return 0;
36     if (!expect_delimiter(ibuf, '\r')) return 0;
37     if (!expect_delimiter(ibuf, '\n')) return 0;
38     return 1;
39 }
40
41 /* write one (possibly partial) response to the given output buffer */
42 int write_mxp_response(outbuf* obuf, mxp_response* response) {
43     if (obuf_putchar(obuf, response->status) < 0) return 0;
44     if (response->message.len > 0) {
45         ssize_t nbytes = obuf_write(obuf, response->message.s,
46             response->message.len);
47         if (nbytes != response->message.len) return 0;
48     }
49     return
50         obuf_putchar(obuf, '\r') >= 0 &&
51         obuf_putchar(obuf, '\n') >= 0;
52 }
53
54 void free_mxp_response(mxp_response* response) {
55     stralloc_free(&response->message);
56 }

```

Da eine Anfrage aus einem Kommandonamen, einem Leerzeichen, dem Parameter und den Zeilenterminator besteht, wurde das Einlesen in *mxrequest.c* entsprechend aufgeteilt. Die Funktionen *read_keyword()* und *read_parameter()* lesen jeweils ihren Teil in ein *stralloc*-Objekt ein und achten mit Hilfe von *inbuf_back()* darauf, daß das jeweilige terminierende Zeichen anschließend zum Einlesen wieder zur Verfügung steht. Die Funktion *expect_delimiter()* liest ein Zeichen ein und vergleicht es mit dem erwarteten Trennzeichen. Auf diese Weise lassen sich Leerzeichen und der Zeilentrenner überlesen. Übertriebene Toleranz bezüglich mehreren Leerzeichen oder verschiedenen Varianten bei Zeilentrennern ist hier fehl am Platze, da dies nur den Programmtext verkompliziert, ohne einen wahren Vorteil zu bringen, da primär nur Programme über dieses Protokoll kommunizieren. Bei der zugehörigen Schreibfunktion ist zu beachten, daß *obuf_flush()* nicht von *write_mxp_request()* aufgerufen wird. Die Kontrolle darüber, wann wirklich der Ausgabepuffer geleert wird, verbleibt so beim Kern des Programmtexts, der mit dem Protokoll umgeht.

4.7 Parallele Sitzungen

Der zuvor mit dem Programmtext *timeserver.c* (auf Seite 65) vorgestellte Zeitdienst wartete mit *accept()* auf den nächsten Aufruf und kümmerte sich dann ausschließlich um diesen. Da eine Sitzung des Zeitdienstes ohne Interaktionen auskommt und nur aus der Zeitansage besteht, führte dies nicht zu Problemen bei mehreren parallelen Anfragen. Entsprechend der bei *listen()* angegebenen Zahl (mit *SOMAXCONN* wird normalerweise das

Maximum gewählt) ist es möglich, weitere Anrufe in eine Warteschlange einzureihen, bis die derzeitige laufende Sitzung beendet ist und `accept()` erneut aufgerufen wird. Dies ist jedoch bei längeren Sitzungen nicht mehr akzeptabel.

Es gibt drei Ansätze, um parallele Sitzungen zu ermöglichen:

- Für jede neue Sitzung wird mit Hilfe von `fork()` ein neuer Prozess erzeugt, der sich um die Verbindung zu genau einem Klienten kümmert.
- Für jede neue Sitzung wird ein neuer Thread gestartet.
- Sämtliche Ein- und Ausgabe-Operationen werden asynchron abgewickelt mit Hilfe der `O_NONBLOCK`-Option und der Verwendung von `poll()` oder `select()`. Hierbei blockieren Ein- und Ausgabe-Operationen nicht mehr, sondern geben eine Fehlerindikation zurück, falls noch keine Eingabe vorliegt bzw. eine Ausgabe momentan noch nicht akzeptiert werden kann. Mit Hilfe von `poll()` oder `select()` ist es dann möglich, bei einer Vielzahl von Dateideskriptoren darauf zu warten, daß eine Lese- oder Schreib-Operation möglich ist.

Die erste Methode ist am einfachsten umzusetzen. Normalerweise wird nur dann eine der beiden anderen Varianten gewählt, wenn die parallel laufenden Sitzungen auf gemeinsame Datenstrukturen zugreifen möchten.

Die erste Variante ist so gängig, daß es sich lohnt, sie in verallgemeinerter Form zur Verfügung zu stellen. Von Dan J. Bernstein gibt es hierfür ein Werkzeug namens `tcpserver`⁷, das die gesamte Sequenz von `socket()` bis `accept()` so verpackt, daß der eigentliche Dienst die gesamte Kommunikation über die Standard-Ein- und Ausgabe abwickeln kann. Auf diese Weise lassen sich sogar relativ einfach kleine Dienste auf Basis bestehender Kommandos einrichten. Hier ist ein Beispiel auf Basis des `factor`-Kommandos, das ganze Zahlen in ihre Primteiler zerlegt:

```
cordelia$ tcpserver 0 23456 /bin/sh -c \
> 'while read x; do echo $x | factor; done'
^Ccordelia$
```

Das Kommando `tcpserver` interpretiert die beiden ersten Argumente als Adresse und Portnummer, bei der der Dienst zu etablieren ist. Die 0 ist hier als Adreßangabe zulässig und entspricht dem `INADDR_ANY`. Für jede eingehende Verbindung wird ein neuer Prozess erzeugt, bei dem Standard-Ein- und Ausgabe mit der Netzwerkverbindung verknüpft werden. Anschließend erfolgt ein `execvp()` zum angegebenen Kommando. Als Kommando kommt hier die Shell zum Einsatz, bei der die Option „-c“ es erlaubt, ein Shell-Kommando einschließlich der Verwendung von Kontrollstrukturen anzugeben. In diesem Beispiel handelt es sich dabei um eine **while**-Schleife, die jede eingelesene Zeile jeweils an das Kommando `factor` weiterleitet.

Prinzipiell wäre es natürlich auch möglich, `factor` direkt von `tcpserver` aufrufen zu lassen. Wegen der gepufferten Ausgabe von `factor` würden wir die Faktorisierung jedoch zu spät oder überhaupt nicht zu sehen bekommen. Wenn hingegen für jede Eingabezeile ein Kommandoaufruf erfolgt, ist damit auch ein implizites Leeren des Ausgabe-Puffers verbunden, sobald das Kommando beendet ist. Das Problem tritt bei einer normalen interaktiven Verwendung nicht auf, da die `stdio` implizit bei der Verwendung eines Terminals auf eine zeilenweise Pufferung umschaltet. Dies geschieht jedoch nicht bei Netzwerkverbindungen.

Hier ist eine Beispielsitzung des Faktorisierungsdienstes:

⁷Siehe <http://cr.yp.to/ucspi-tcp/tcpserver.html>

```
cordelia$ telnet cordelia 23456
Trying 134.60.66.8...
Connected to cordelia.
Escape character is '^]'.
8714
    2
    4357

9726
    2
    3
    1621

^]
telnet> quit
Connection to cordelia closed.
cordelia$
```

Die Kontrollstruktur, die *tcpserver* auf der Ebene der Kommandozeile liefert, lässt sich ebenso in genereller Form als Bibliotheksfunktion repräsentieren:

Programm 4.12: Schnittstelle für die Kontrollstruktur von Netzwerkdiensten (*service.h*)

```
1 #ifndef SERVICE_H
2 #define SERVICE_H
3
4 #include <netinet/in.h>
5
6 typedef void (*session_handler)(int fd, int argc, char** argv);
7
8 /*
9  * listen on the given port and invoke the handler for each
10  * incoming connection
11  */
12 void run_service(in_port_t port, session_handler handler,
13                 int argc, char** argv);
14
15 #endif
```

Programm 4.13: Implementierung für die Kontrollstruktur von Netzwerkdiensten (*service.c*)

```
1 #include <netinet/in.h>
2 #include <signal.h>
3 #include <stdio.h>
4 #include <stdlib.h>
5 #include <string.h>
6 #include <sys/socket.h>
7 #include <sys/time.h>
8 #include <time.h>
9 #include <unistd.h>
10 #include "service.h"
11
12 /*
13  * listen on the given port and invoke the handler for each
```

```

14  * incoming connection
15  */
16  void run_service(in_port_t port, session_handler handler,
17                 int argc, char** argv) {
18      struct sockaddr_in address = {0};
19      address.sin_family = AF_INET;
20      address.sin_addr.s_addr = htonl(INADDR_ANY);
21      address.sin_port = htons(port);
22
23      int sfd = socket(PF_INET, SOCK_STREAM, 0);
24      int optval = 1;
25      if (sfd < 0 ||
26          setsockopt(sfd, SOL_SOCKET, SO_REUSEADDR,
27                    &optval, sizeof optval) < 0 ||
28          bind(sfd, (struct sockaddr *) &address,
29              sizeof address) < 0 ||
30          listen(sfd, SOMAXCONN) < 0) {
31          return;
32      }
33
34      /* our childs shall not become zombies */
35      struct sigaction action = {0};
36      action.sa_handler = SIG_IGN;
37      action.sa_flags |= SA_NOCLDWAIT;
38      if (sigaction(SIGCHLD, &action, 0) < 0) return;
39
40      int fd;
41      while ((fd = accept(sfd, 0, 0)) >= 0) {
42          pid_t child = fork();
43          if (child == 0) {
44              handler(fd, argc, argv);
45              exit(0);
46          }
47          close(fd);
48      }
49  }

```

Genauso wie *tcpserver* benötigt die Funktion *run_service()* die lokal zu verwendende Adresse. Der Einfachheit halber wird hier allerdings nur die Portnummer als Parameter übergeben und implizit *INADDR_ANY* verwendet. Für jeden eingehenden Anruf soll *run_service()* einen neuen Prozess mit *fork()* erzeugen und in diesem die als Parameter übergebene Funktion *handler* aufrufen. Diese Funktion erhält als Parameter den Dateideskriptor der offenen Netzwerkverbindung und — analog zu *tcpserver* — die verbliebenen Kommandozeilenparameter. Genau wie *tcpserver* endet *run_service()* nur im Fehlerfalle und läuft ansonsten endlos.

In der zugehörigen Implementierung gleichen die Zeilen 18 bis 32 dem Aufbau des Anschlusses beim Zeitdienst (siehe Programmtext *timeserver.c* auf Seite 65). Der entscheidende Kern der Kontrollstruktur ist in der Schleife in den Zeilen 41 bis 48 zu finden, bei der jeder erfolgreiche Aufruf von *accept()* zu einem neuen Prozess führt, der *handler* aufruft. Der übergeordnete Prozess schließt nur seinen Dateideskriptor zur Netzwerkverbindung und startet den nächsten Aufruf von *accept()*.

Im Rahmen dieses Verfahrens entstehen zahlreiche Kind-Prozesse, für die *wait()* nicht aufgerufen wird. Um nicht endlos Zombies zu akkumulieren, wird in den Zeilen 35 bis

38 das Signal `SIG_CHLD` ignoriert, wobei die Option `SA_NOCLDWAIT` gesetzt wird. Diese Option bewirkt, daß der Endstatus beendeter Kind-Prozesse sofort entsorgt wird, ohne daß diese zu Zombies werden oder später die Möglichkeit besteht, den Status noch mit `wait()` abzurufen. Alternativ wäre es natürlich auch möglich, `wait()` innerhalb eines Signalhandlers aufzurufen oder regelmäßig in der **while**-Schleife `waitpid()` mit der `WNOHANG`-Option aufzurufen (siehe Programmtext `sigchld.c` auf Seite 32).