

- Ein- und Ausgabe-Operationen blockieren normalerweise, bis sie durchgeführt werden können.
- Dies erschwert die Parallelisierung solcher Operationen bzw. die Möglichkeit, auf unterschiedliche Ein- und Ausgabe-Ereignisse zu reagieren.
- Mit den Systemaufrufen *poll* und *select* gibt es die Möglichkeit, zu warten, bis wir mindestens eine von beliebig vielen geplanten Ein- und Ausgabe-Operationen durchführen können, ohne blockiert zu werden.
- Der Vorteil dieser Schnittstelle liegt darin, dass wir die synchrone Arbeitsweise nicht aufgeben müssen.
- Wir betrachten hier im weiten *poll*, da dieser Systemaufruf eine etwas elegantere Schnittstelle als *select* bietet.

multiplexor.c

```
if (poll(mpx.pollfds, count, -1) <= 0) return;
```

- *poll* erhält drei Parameter:
  - ▶ Einen Zeiger auf ein Array mit Einträgen des Datentyps **struct pollfd**,
  - ▶ einer natürlichen Zahl, die die Länge des Arrays angibt, und
  - ▶ einer zeitlichen Beschränkung in Millisekunden. (Hier wird -1 angegeben, wenn keine Befristung gewünscht wird.)
- Der Datentyp **struct pollfd** umfasst folgende Felder:
  - fd*        Dateideskriptor
  - events*   Menge der Ereignisse, auf die gewartet wird
  - revents*   Menge der Ereignisse, die eingetreten sind
- Im Erfolgsfall liefert *poll* die Zahl der eingetretenen Ereignisse zurück. Falls die zeitliche Beschränkung erreicht wurde, ohne dass eines der Ereignisse eintrat, wird 0 zurückgeliefert. Im Falle von Fehlern wird -1 zurückgegeben.

- Relevant sind nur *POLLIN* und *POLLOUT*. Prinzipiell kann *poll* noch Unterscheidungen treffen, ob priorisierte Pakete über die Netzwerkverbindung ankamen, aber das wird normalerweise nicht verwendet.
- Das Ereignis *POLLIN* bedeutet, dass ein *read*-Systemaufruf für den Dateideskriptor abgesetzt werden kann, ohne dass der Prozess blockiert wird.
- Analog bedeutet *POLLOUT*, dass ein *write*-Systemaufruf abgesetzt werden kann, ohne Gefahr zu laufen, blockiert zu werden.
- Bei mit *listen* vorbereiteten Sockets kann ebenfalls *POLLIN* verwendet werden. Das Ereignis tritt dann ein, sobald sich eine neue Netzwerkverbindung anbahnt und *accept* blockierungsfrei aufgerufen werden kann.

- Die Umsetzung des Prefork-Modells lässt sich mit Hilfe von *poll* verbessern, da wir dann keinen Wächterprozess pro Prozess benötigen, der bereit ist, eine Verbindung mit *accept* entgegenzunehmen.
- Bei  $n$  Prozessen, die bereit sein sollen, eine Sitzung entgegenzunehmen, werden jetzt nur noch insgesamt  $n + 1$  Prozesse benötigt, d.h. es kommt nur noch der Hauptprozess hinzu.
- Der Hauptprozess erzeugt selbst alle weiteren Prozesse und beobachtet dann mit Hilfe von *poll* die Pipeline-Verbindungen zu den einzelnen Prozessen.
- Sobald die letzte offene Schreibverbindung einer Pipeline geschlossen wird, tritt auf der lesenden Seite das *POLLIN*-Ereignis ein, damit das Eingabe-Ende erkannt werden kann. (Ein *read* würde dann blockierungsfrei eine 0 zurückliefern.)

preforked\_service.c

```
static pid_t spawn_preforked_process(int sfd, int pipefds[2],
    session_handler handler, int argc, char** argv) {
    if (pipe(pipefds) < 0) return -1;
    pid_t child = fork();
    if (child) {
        close(pipefds[1]);
        return child;
    }
    close(pipefds[0]);

    int fd = accept(sfd, 0, 0); close(sfd);
    if (fd < 0) exit(1);
    /* now close the writing side of the pipe to indicate that
       we are busy with running a session */
    close(pipefds[1]);
    /* run the session and exit */
    handler(fd, argc, argv);
    exit(0);
}
```

- Die Funktion *spawn\_preforked\_process* vereinfacht sich, da nur noch ein Prozess erzeugt wird.

preforked\_service.c

```
/* create preforked processes */
pid_t child_pid[number_of_processes];
struct pollfd pollfds[number_of_processes];
for (int i = 0; i < number_of_processes; ++i) {
    /* a pipe is used to signal that one of the
       preforked processes accepted a connection */
    int pipefds[2];
    pid_t pid = spawn_preforked_process(sfd, pipefds, handler,
                                       argc, argv);
    pollfds[i] = (struct pollfd) { .fd = pipefds[0], .events = POLLIN};
    if (pid < 0) return;
    child_pid[i] = pid;
}
```

- Der Hauptprozess erzeugt hier zu Beginn die gewünschte Zahl von Prozessen.
- Dabei wird gleichzeitig die *pollfds*-Datenstruktur aufgebaut, um all die Pipelines gleichzeitig beobachten zu können.

preforked\_service.c

```
while (!terminate) {
    if (poll(pollfds, number_of_processes, -1) <= 0) break;
    for (int i = 0; i < number_of_processes; ++i) {
        if (pollfds[i].revents == 0) continue;
        close(pollfds[i].fd);
        int pipefds[2];
        pid_t pid = spawn_preforked_process(sfd, pipefds, handler,
            argc, argv);
        if (pid < 0) return;
        pollfds[i] = (struct pollfd) {
            .fd = pipefds[0], .events = POLLIN};
        child_pid[i] = pid;
    }
}
```

- Mit *poll* warten wir darauf, dass die schreibende Seite eine der Pipes geschlossen wird.
- Dies ist das Signal, dass ein neuer Prozess zu starten ist, dessen Pipeline dann in *pollfds* ersatzweise eingetragen wird.

- In manchen Fällen ist es vorteilhaft, wenn alle Sitzungen einen gemeinsamen Adressraum verwenden, damit sitzungsübergreifende Datenstrukturen leichter verwaltet werden können.
- Prinzipiell lässt sich das mit Hilfe des Systemaufrufs *poll* erreichen, mit dem auf das Eintreten eines Ein- oder Ausgabe-Ereignisses gewartet werden kann.
- Dies führt zu einem grundlegend anderen Programmierstil, bei dem Ein- und Ausgaben ereignisgesteuert abgewickelt werden.
- Da bei jedem Ereignis entsprechende Behandler neu aufgerufen werden, kann der Sitzungskontext nicht in lokalen Variablen verwaltet werden. Stattdessen sind dafür dynamische Datenstrukturen zu verwenden, die bei jedem Aufruf erst lokalisiert werden müssen.



multiplexor.h

```
typedef void (*multiplexor_handler)(connection* link);
void run_multiplexor(int socket, multiplexor_handler open_handler,
    multiplexor_handler input_handler, multiplexor_handler close_handler,
    void* mpx_handle);
bool write_to_link(connection* link, char* buf, unsigned int len);
ssize_t read_from_link(connection* link, char* buf, unsigned int len);
void close_link(connection* link);
```

- Es ist sinnvoll, die Verwendung von *poll* in eine geeignete Bibliothek zu verpacken.
- Die Funktion *run\_multiplexor* läuft dann permanent und übernimmt somit die vollständige Kontrolle des Programms. Es werden nur noch Behandler aufgerufen, wenn
  - ▶ neue Netzwerkverbindungen eröffnet werden,
  - ▶ neue Eingaben vorliegen oder
  - ▶ eine Verbindung beendet wird.
- Eine Rückkehr von *run\_multiplexor* gibt es nur im Fehlerfalle.

multiplexor.h

```
typedef void (*multiplexor_handler)(connection* link);
void run_multiplexor(int socket, multiplexor_handler open_handler,
    multiplexor_handler input_handler, multiplexor_handler close_handler,
    void* mpx_handle);
bool write_to_link(connection* link, char* buf, unsigned int len);
ssize_t read_from_link(connection* link, char* buf, unsigned int len);
void close_link(connection* link);
```

- Konkret ruft *run\_multiplexor* den Behandler *open\_handler* für neue Verbindungen, *input\_handler* für neue Eingaben und *close\_handler* für beendete Verbindungen auf.
- Die Behandler dürfen selbst nichts direkt auf eine Netzwerkverbindung ausgeben, da dies zu längeren Blockaden führen könnte. Stattdessen muss dies durch *write\_to\_link* erfolgen, das dafür Warteschlangen unterhält.
- Der Parameter *mpx\_handle* dient als Zeiger auf eine eigene Datenstruktur, die den Behandlern unter *connection->mpx\_handle* zur Verfügung gestellt wird.

multiplexor.h

```
typedef struct connection {
    int fd;
    void* handle; /* may be freely used by the application */
    void* mpx_handle; /* corresponding parameter from run_multiplexor */
    bool eof;
    struct output_queue_member* oqhead;
    struct output_queue_member* oqtail;
    struct connection* next;
    struct connection* prev;
} connection;
```

- Für jede Netzwerkverbindung gibt es eine zugehörige Datenstruktur.
- Neben der Netzwerkverbindung *fd* und den beiden benutzerdefinierten Zeigern *handle* und *mpx\_handle*, kommen noch folgende Felder hinzu:
  - eof* wird auf *true* gesetzt, sobald ein Eingabeende erkannt wurde
  - oqhead* und *oqtail* Zeiger auf das erste und letzte Element der Warteschlange mit den auszugehenden Puffern
  - next* und *prev* doppelt verkettete Liste aller Netzwerkverbindungen

multiplexor.c

```
typedef struct output_queue_member {
    char* buf;
    unsigned int len;
    unsigned int pos;
    struct output_queue_member* next;
} output_queue_member;
// ...
int write_to_link(connection* link, char* buf, unsigned int len);
```

- Jedes Element der Warteschlange weist auf einen Puffer.
- Zu Beginn ist die Position *pos* gleich 0 und *len* entspricht der Länge, die an *write\_to\_link* übergeben worden ist.
- Wenn jedoch der entsprechende Aufruf von *write* nicht vollständig umgesetzt werden kann, dann wird *pos* um die übertragene Quantität erhöht und *len* entsprechend gesenkt.
- Sobald die Schreiboperation abgeschlossen ist, wird nicht nur das Warteschlangen-Element, sondern auch der Puffer freigegeben.

```
typedef struct multiplexor {
    /* parameters passed to run_multiplexor */
    int socket;
    multiplexor_handler ohandler, ihandler, chandler;
    void* mpx_handle;
    /* additional administrative fields */
    bool socketok; /* becomes false when accept() fails */
    connection* head; /* double-linked linear list of connections */
    connection* tail; /* its last element */
    int count; /* number of connections */
    struct pollfd* pollfds; /* parameter for poll() */
    unsigned int pollfdslen; /* allocated len of pollfds */
    connection** pollcs; /* of the same len as pollfds */
} multiplexor;
```

- Es gibt nur ein Objekt dieser Datenstruktur, das von *run\_multiplexor* zu Beginn angelegt wird.
- Neben den Parametern von *run\_multiplexor* werden in der doppelt verketteten Liste mit *head* und *tail* alle offenen Verbindungen verwaltet. In *count* findet sich deren Zahl.
- *pollfds* zeigt auf ein dynamisch belegtes Feld mit *pollfdslen* Elementen. Dies dient der Verwaltung der *poll* zu übergebenden Datenstruktur.

multiplexor.c

```
/* prepare fields pollfds and pollfdslen in mpx in
dependence of the current set of connections */
static int setup_polls(multiplexor* mpx) {
    int len = mpx->count;
    if (mpx->socketok) ++len;
    if (len == 0) return 0;
    /* weed out links which have been closed
and where our output queue is empty */
    connection* link = mpx->head;
    while (link) {
        connection* next = link->next;
        if (link->eof && link->oqhead == 0) remove_link(mpx, link);
        link = next;
    }
    /* allocate or enlarge pollfds, if necessary */
    if (mpx->pollfdslen < len) {
        mpx->pollfds = realloc(mpx->pollfds, sizeof(struct pollfd) * len);
        if (mpx->pollfds == 0) return 0;
        mpx->pollcs = realloc(mpx->pollcs, sizeof(connection*) * len);
        if (mpx->pollcs == 0) return 0;
        mpx->pollfdslen = len;
    }

    /* ... */
}
```

multiplexor.c

```
/* prepare fields pollfds and pollfdslen in mpx in
   dependence of the current set of connections */
static int setup_polls(multiplexor* mpx) {
    /* ... */

    int index = 0;
    /* look for new network connections as long accept()
       returned no errors so far */
    if (mpx->socketok) {
        mpx->pollcs[index] = 0;
        mpx->pollfds[index++] = (struct pollfd) {mpx->socket, POLLIN};
    }
    /* look for incoming network connections and
       check whether we can write any pending output packets
       without blocking */
    link = mpx->head;
    while (link) {
        short events = 0;
        if (!link->eof) events |= POLLIN;
        if (link->oqhead) events |= POLLOUT;
        mpx->pollcs[index] = link;
        mpx->pollfds[index++] = (struct pollfd) {link->fd, events};
        link = link->next;
    }
    return index;
}
```

```
static bool add_connection(multiplexor* mpx) {
    int newfd;
    if ((newfd = accept(mpx->socket, 0, 0)) < 0) {
        mpx->socketok = false; return true;
    }
    connection* link = malloc(sizeof(connection));
    if (link == 0) return false;
    *link = (connection) {
        .fd = newfd, .handle = 0, .mpx = mpx,
        .mpx_handle = mpx->mpx_handle,
        .eof = false, .oqhead = 0, .oqtail = 0,
        .next = 0, .prev = mpx->tail,
    };
    if (mpx->tail) {
        mpx->tail->next = link;
    } else {
        mpx->head = link;
    }
    mpx->tail = link; ++mpx->count;
    if (mpx->ohandler) (*mpx->ohandler)(link);
    return true;
}
```



multiplexor.c

```
/* remove a connection from the double-linked linear
   list of connections
*/
static void remove_link(multiplexor* mpx, connection* link) {
    close(link->fd);
    if (link->prev) {
        link->prev->next = link->next;
    } else {
        mpx->head = link->next;
    }
    if (link->next) {
        link->next->prev = link->prev;
    } else {
        mpx->tail = link->prev;
    }
    if (mpx->chandler) (*mpx->chandler)(link);
    free(link);
    --mpx->count;
}
```

multiplexor.c

```
/* read one input packet from the given network connection */
ssize_t read_from_link(connection* link, char* buf, unsigned int len) {
    if (link->eof) return 0;
    ssize_t nbytes = read(link->fd, buf, len);
    if (nbytes <= 0) {
        link->eof = true;
        if (link->oqhead == 0) remove_link((multiplexor*)link->mpx, link);
    }
    return nbytes;
}
```

- Wenn *poll* signalisiert hat, dass wir von einer Verbindung einlesen dürfen, dann wird der entsprechende Behandler aufgerufen, der wiederum *read\_from\_link* aufruft, um die Eingabe in den eigenen Puffer einzulesen.

```
/* write one pending output packet to the given network connection */
static void write_to_socket(multiplexor* mpx, connection* link) {
    ssize_t nbytes = write(link->fd,
        link->oqhead->buf + link->oqhead->pos,
        link->oqhead->len - link->oqhead->pos);
    if (nbytes <= 0) {
        remove_link(mpx, link);
    } else {
        link->oqhead->pos += nbytes;
        if (link->oqhead->pos == link->oqhead->len) {
            output_queue_member* old = link->oqhead;
            link->oqhead = old->next;
            if (link->oqhead == 0) {
                link->oqtail = 0;
            }
            free(old->buf); free(old);
            if (link->oqhead == 0 && link->eof) {
                remove_link(mpx, link);
            }
        }
    }
}
```

```
bool write_to_link(connection* link, char* buf, unsigned int len) {
    assert(len >= 0);
    if (len == 0) {
        free(buf); return true;
    }
    output_queue_member* member = malloc(sizeof(output_queue_member));
    if (!member) return false;
    member->buf = buf; member->len = len; member->pos = 0;
    member->next = 0;
    if (link->oqtail) {
        link->oqtail->next = member;
    } else {
        link->oqhead = member;
    }
    link->oqtail = member;
    return true;
}
```

- Diese Funktion ist von den Behandlern aufzurufen, wenn etwas auf eine der Netzwerkverbindungen auszugeben ist.
- Der Ausgabepuffer wird dann in die entsprechende Warteschlange eingereiht.

multiplexor.c

```
void close_link(connection* link) {
    link->eof = 1;
    shutdown(link->fd, SHUT_RD);
}
```

- Bei bidirektionalen Netzwerkverbindungen ist es möglich, nur eine Seite zu schließen.
- Dies geht nicht mit *close*, das sofort beide Seiten schließen würde, sondern mit *shutdown*, mit dem eine spezifizierte Seite geschlossen werden kann.
- Hier wird aus der Sicht des Aufrufers die lesende Seite geschlossen, also die Verbindung vom Klienten zum Dienst. Danach können keine weiteren Anfragen mehr eintreffen, aber die Warteschlange der abzuarbeitenden Ausgabe-Puffer kann noch abgearbeitet werden.
- Erst wenn die Warteschlange ganz leer ist, dann wird (von *remove\_link*) die Verbindung vollständig geschlossen.

multiplexor.c

```
void run_multiplexor(int socket, multiplexor_handler open_handler,
    multiplexor_handler input_handler, multiplexor_handler close_handler,
    void* mpx_handle) {
    multiplexor mpx = {
        .socket = socket, .ohandler = open_handler,
        .ihandler = input_handler, .chandler = close_handler,
        .mpx_handle = mpx_handle, .socketok = true,
    };
    int count;
    while ((count = setup_polls(&mpx)) > 0) {
        if (poll(mpx.pollfds, count, -1) <= 0) return;
        for (int index = 0; index < count; ++index) {
            if (mpx.pollfds[index].revents == 0) continue;
            int fd = mpx.pollfds[index].fd;
            if (fd == mpx.socket) {
                if (!add_connection(&mpx)) return;
            } else {
                connection* link = mpx.pollcs[index]; assert(link);
                if (mpx.pollfds[index].revents & POLLIN) {
                    (*mpx.ihandler)(link);
                }
                if (mpx.pollfds[index].revents & POLLOUT) {
                    write_to_socket(&mpx, link);
                }
            }
        }
    }
}
```

- Der *input\_handler* wird für jedes eingehende Paket aufgerufen.
- Da Pakete fragmentiert sein können, sind dies möglicherweise Bruchstücke einer Anfrage oder auch Teile mehrerer Anfragen.
- Entsprechend muss die Eingabe wieder gepuffert und zerlegt werden, da normalerweise eine Reaktion erst bei einer vollständig übermittelten Anfrage erfolgen sollte.
- Eine ereignisgesteuerte Behandlung wäre daher aus Anwendungssicht leichter zu programmieren, wenn sie auf vollständigen Anfragen beruhen würde.
- Die Erkennung einer vollständigen Anfrage ist im allgemeinen Fall nicht ganz trivial zu spezifizieren. Im folgenden wird eine Lösung auf Basis regulärer Ausdrücke vorgestellt, die für textbasierte Protokolle gut geeignet ist.

mpx\_session.h

```
typedef void (*mpx_handler)(session* s);

int mpx_session_scan(session* s, ...);
int mpx_session_printf(session* s, const char* restrict format, ...);
void close_session(session* s);

void run_mpx_service(hostport* hp, const char* regexp,
    mpx_handler ohandler, mpx_handler rhandler, mpx_handler hhandler,
    void* global_handle);
```

- *run\_mpx\_service* erhält einen regulären Ausdruck, der eine Anfrage spezifiziert.
- Dieser reguläre Ausdruck darf mit Hilfe runder Klammern beliebig viele Elemente der Anfrage herausgreifen – analog zu *inbuf\_scan*.
- Der *rhandler* (*request handler*) wird dann für jede vollständig vorliegende Anfrage aufgerufen und kann dann mit *mpx\_session\_scan* die herausgegriffenen Elemente in *stralloc*-Objekte hineinkopieren lassen.



`mxprequest.h`

```
#define MXP_REQUEST_RE "([a-z]+) (.*)\r?\n"
```

`mutexd.c`

```
run_mpx_service(&hp, MXP_REQUEST_RE,  
               mpx_session_open, mpx_session_read, mpx_session_hangup,  
               locks);
```

- Beim Aufruf von `run_mpx_service` wird der reguläre Ausdruck zum Erkennen einer Anfrage mitgegeben.

mxpession.c

```
void mxp_session_read(session* s) {
    struct mxp_session* ms = s->handle; assert(ms);
    if (!read_mxp_request(s, &ms->request)) {
        close_session(s); return;
    }
    /* ... process request and generate response ... */
    if (!write_mxp_response(s, &ms->response)) {
        close_session(s);
    }
}
```

- Der Behandler `mxp_session_read` wird jetzt nur aufgerufen, wenn eine vollständige Anfrage vorliegt. Entsprechend sollte `read_mxp_request` eine Anfrage einlesen können.

mxprequest.c

```
bool read_mxp_request(session* s, mxp_request* request) {
    return
        mpx_session_scan(s, &request->keyword, &request->parameter) == 2;
}
```

mxresponse.c

```
/* write one (possibly partial) response to */
bool write_mxp_response(session* s, mxp_response* response) {
    return mpx_session_printf(s, "%c%.*s\r\n", response->status,
        response->message.len, response->message.s) > 0;
}
```

- Die Einlese-Operation für Anfragen und die Ausgabe-Operation für Antworten verwenden hier die entsprechenden Funktionen aus *mpx\_session.h*