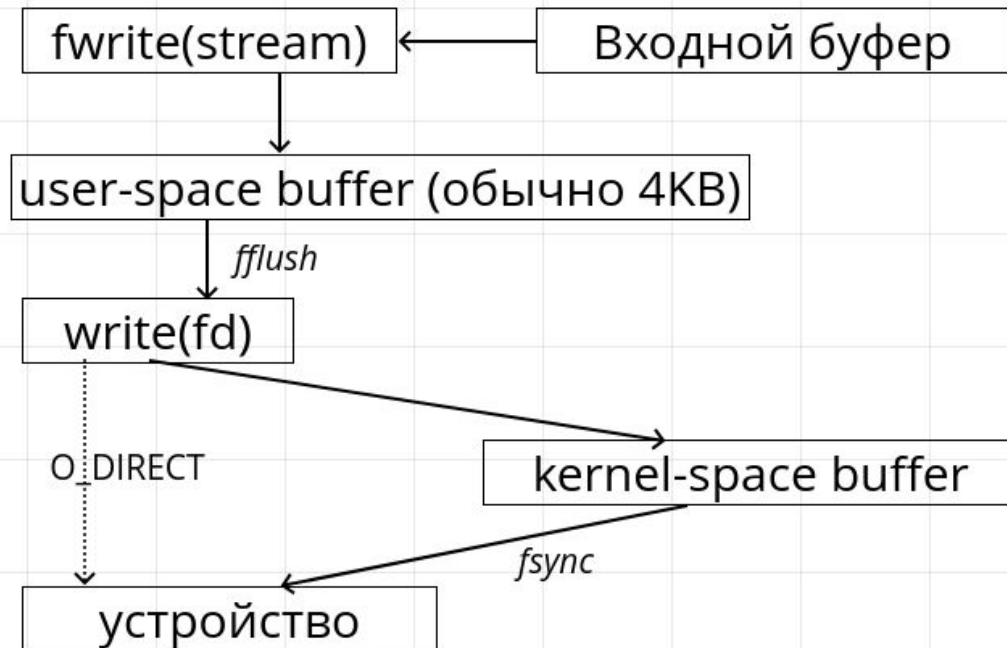


# Программирование в Linux

Файлы. IO.

```
1 #include <unistd.h>
2
3 ssize_t read(int fd, void* buf, size_t count);
4 ssize_t write(int fd, const void* buf, size_t count);
```



## Конкурентный доступ к файлу

- Системные read/write атомарны
- Порядок определяет планировщик
- В каждом открытом файле есть мьютекс (см. flock)



# Блокирующий ввод/вывод

Как прочитать N байт?

read(fd, buf, N)?

```
1 ssize_t read_all(int fd, char* buff, size_t len) {
2     ssize_t actual_read = 0;
3     while (len != 0) {
4         auto read_cnt = read(fd, buff, len);
5         if (read_cnt == 0) {
6             // файл кончился
7             break;
8         }
9         if (read_cnt < 0) {
10             if (errno == EINTR) {
11                 // нас прервали, попробуй еще раз
12                 continue;
13             }
14             // серьезная ошибка
15             return -1;
16         }
17         actual_read += read_cnt;
18         buff += read_cnt;
19         len -= read_cnt;
20     }
21     return actual_read;
22 }
```

## Векторные чтение/запись

```
1 ssize_t readv(int fd, const struct iovec *iov, int iovcnt);
2
3 ssize_t writev(int fd, const struct iovec *iov, int iovcnt);
4
5 struct iovec {
6     void *iov_base;    /* Starting address */
7     size_t iov_len;    /* Number of bytes to transfer */
8 };
9
```

```
char *str0 = "hello ";
char *str1 = "world\n";
struct iovec iov[2];
ssize_t nwritten;

iov[0].iov_base = str0;
iov[0].iov_len = strlen(str0);
iov[1].iov_base = str1;
iov[1].iov_len = strlen(str1);

nwritten = writev(STDOUT_FILENO, iov, 2);
```

# Позиционирование

- Внутренний указатель связан с fd
- read/write двигают указатель
- pread/pwrite — не двигают

```
1 ssize_t pread(int fd, void *buf, size_t count, off_t offset);  
2  
3 ssize_t pwrite(int fd, const void *buf, size_t count, off_t offset);
```

- lseek — подвигать вручную
- lseek за конец + write = дыры в файле

# Неблокирующий ввод/вывод

open с флагом O\_NONBLOCK

```
1 start:
2 read_cnt = read(fd, buff, len);
3 if (read_cnt < 0) {
4     if (errno == EINTR) {
5         goto start;
6     }
7     if (errno == EAGAIN) {
8         // читать рано -- надо ждать,
9         // займитесь чем-то другим
10    } else {
11        // что-то более серьезное
12        return -1;
13    }
14 }
```

# Ожидание событий IO

- while (!ready) { ready = try\_read(...); } // активное ожидание ест CPU
- while (!ready) { sleep(1ms); ready = try\_read(...); } // частота опроса может быть недостаточной
- лучше попросить ОС об оповещении

```
1 #include <sys/select.h>
2
3 int select(int nfds, fd_set *readfds, fd_set *writefds,
4           fd_set *exceptfds, struct timeval *timeout);
5
6 void FD_CLR(int fd, fd_set *set);
7 int  FD_ISSET(int fd, fd_set *set);
8 void FD_SET(int fd, fd_set *set);
9 void FD_ZERO(fd_set *set);
10
11
12 int pselect(int nfds, fd_set *readfds, fd_set *writefds,
13            fd_set *exceptfds, const struct timespec *timeout,
14            const sigset_t *sigmask);
```



```

1 #include <sys/select.h>
2 #include <sys/unistd.h>
3 #include <sys/types.h>
4
5 #include <cstdlib>
6 #include <stdio>
7
8 int main() {
9     fd_set readfds; FD_ZERO(&readfds);
10    FD_SET(STDIN_FILENO, &readfds);
11    timeval timeout; timeout.tv_sec = 5;
12                    timeout.tv_usec = 0;
13
14    int ready_fd_cnt = select(STDIN_FILENO + 1, &readfds, nullptr, nullptr,
15                             &timeout);
16    if (ready_fd_cnt < 0) {
17        perror("select");
18        return EXIT_FAILURE;
19    }
20    if (ready_fd_cnt == 0) {
21        printf("Timeout");
22        return EXIT_SUCCESS;
23    }
24    if (FD_ISSET(STDIN_FILENO, &readfds)) {
25        const int kBufSize = 1024;
26        char buf[kBufSize + 1] = {0};
27        auto read_cnt = read(STDIN_FILENO, buf, kBufSize);
28        if (read_cnt < 0) {
29            perror("read");
30            return EXIT_FAILURE;
31        }
32        if (read_cnt) {
33            printf("read: %s\n", buf);
34        }
35    }
36 }

```

# Альтернатива select — poll

```
1 #include <poll.h>
2 #include <signal.h>
3
4
5 int poll(struct pollfd *fds, nfds_t nfds, int timeout_ms);
6
7 int ppoll(struct pollfd *fds, nfds_t nfds,
8           const struct timespec *tmo_p, const sigset_t *sigmask);
9
10
11 struct pollfd {
12     int    fd;           /* file descriptor */
13     short  events;       /* requested events */
14     short  revents;      /* returned events */
15 };
```

```

1 #include <poll.h>
2 #include <unistd.h>
3
4 #include <cstdlib>
5 #include <stdio.h>
6
7 int main() {
8
9     pollfd request;
10    request.fd = STDIN_FILENO;
11    request.events = POLLIN;
12
13    int ready_cnt = poll(&request, 1, 5000);
14
15    if (ready_cnt < 0) {
16        perror("poll");
17        return EXIT_FAILURE;
18    }
19    if (ready_cnt == 0) {
20        printf("Timeout");
21        return EXIT_SUCCESS;
22    }
23    if (request.revents & POLLIN) {
24        const int kBufSize = 1024;
25        char buf[kBufSize + 1] = {0};
26        auto read_cnt = read(STDIN_FILENO, buf, kBufSize);
27        if (read_cnt < 0) {
28            perror("read");
29            return EXIT_FAILURE;
30        }
31        if (read_cnt) {
32            printf("read: %s\n", buf);
33        }
34    }
35 }

```

# epoll — linux specific

```
int listen_sock; // каким-то образом полученный
// дескриптор серверного сокета,
// для которого ждем событий

// создаем ожидальку
const int kEpollFlags = 0;
int epollfd = epoll_create1(kEpollFlags);
if (epollfd == -1) {
    perror("epoll_create1");
    std::exit(EXIT_FAILURE);
}

{
    epoll_event ev;
    ev.events = EPOLLIN; // будем входящих данных
    ev.data.fd = listen_sock;
    // регистрируем событие для ожидания
    if (epoll_ctl(epollfd, EPOLL_CTL_ADD, listen_sock, &ev) == -1) {
        perror("epoll_ctl: listen_sock");
        exit(EXIT_FAILURE);
    }
}
```

```
#define MAX_EVENTS 10
epoll_event events_buffer[MAX_EVENTS];
for (;;) {
    int events_cnt = epoll_wait(epollfd,
                                events_buffer, MAX_EVENTS,
                                -1); // Timeout
    if (events_cnt == -1) {
        perror("epoll_wait");
        exit(EXIT_FAILURE);
    }
    // дождался событий
    for (int idx = 0; idx < events_cnt; ++idx) {
        if (events_buffer[idx].data.fd == listen_sock) {
            // для примера сервер подключает новых клиентов →
            // новые файловые дескрипторы
            int conn_sock = accept(listen_sock, &addr, &addrlen);
            if (conn_sock == -1) {
                perror("accept");
                exit(EXIT_FAILURE);
            }
            // включаем неблокирующий режим
            fcntl(conn_sock,
                  F_SETFL,
                  fcntl(conn_sock, F_GETFL, 0) | O_NONBLOCK);
            // и будем ждать на них событий — добавляем к ожидальке
            epoll_event ev;
            ev.events = EPOLLIN | EPOLLET;
            ev.data.fd = conn_sock;
            if (epoll_ctl(epollfd,
                          EPOLL_CTL_ADD,
                          conn_sock, &ev) == -1) {
                perror("epoll_ctl: conn_sock");
                exit(EXIT_FAILURE);
            }
        } else {
            // дескриптор клиента готов к неблокирующим операциям
            // events_buffer[idx].data.fd;
        }
    }
}

close(epollfd);
```

флаг EPOLLET — детект события по фронту сигнала, а не по уровню



# Posix aio

```
1 #include <aio.h>
2 #include <aio.h>
3
4 struct aiocb {
5     /* The order of these fields is implementation-dependent */
6
7     int         aio_fildes;    /* File descriptor */
8     off_t       aio_offset;    /* File offset */
9     volatile void *aio_buf;    /* Location of buffer */
10    size_t      aio_nbytes;    /* Length of transfer */
11    int         aio_reqprio;    /* Request priority */
12    struct sigevent aio_sigevent; /* Notification method */
13    int         aio_lio_opcode; /* Operation to be performed;
14                                lio_listio() only */
15
16    /* Various implementation-internal fields not shown */
17 };
18
19 /* Operation codes for 'aio_lio_opcode': */
20
21 enum { LIO_READ, LIO_WRITE, LIO_NOP };
```

и пачка методов *aio\_OPNAME*: *aio\_read*, *aio\_write*, *aio\_cancel*...

# Ссылки

1. Epoll <https://habr.com/ru/post/416669/>
2. <https://idea.popcount.org/2017-01-06-select-is-fundamentally-broken/>
3. <https://habr.com/ru/company/badoo/blog/439972/>