

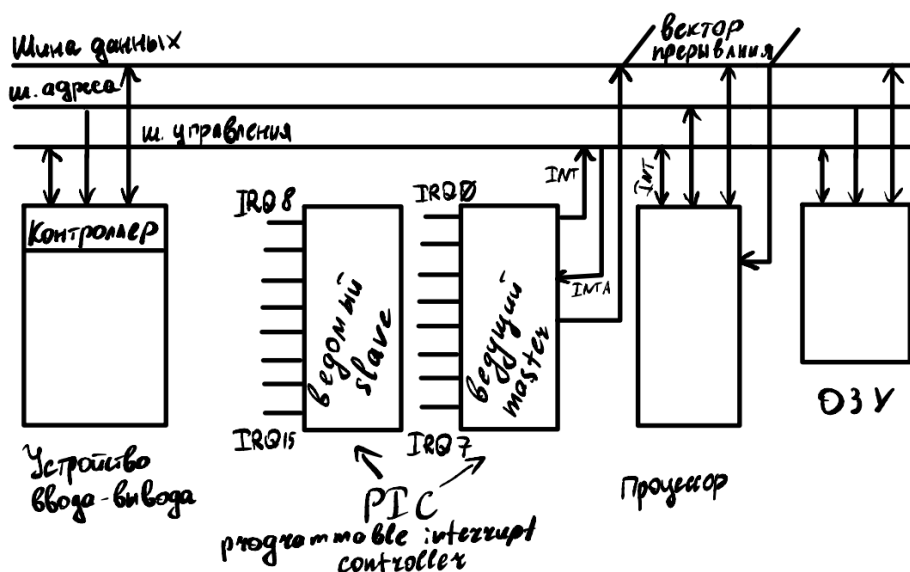
# Оглавление

|           |  |           |
|-----------|--|-----------|
| <b>1</b>  | <b>Ликбез по 1 семестру</b>                            | <b>2</b>  |
| 1.1       | Концептуальная трехшинная архитектура . . . . .        | 2         |
| 1.2       | Канальная архитектура . . . . .                        | 2         |
| 1.3       | Прерывания в последовательности ввода-вывода . . . . . | 3         |
| 1.4       | Классификация прерываний . . . . .                     | 3         |
| <b>2</b>  | <b>Возможно стоит знать</b>                            | <b>4</b>  |
| 2.1       | Архитектуры . . . . .                                  | 4         |
| 2.2       | Мультиплексирование . . . . .                          | 4         |
| 2.3       | Монитор Хоара . . . . .                                | 4         |
| 2.4       | Семафоры . . . . .                                     | 5         |
| 2.5       | Системные вызовы . . . . .                             | 5         |
| 2.6       | Способы взаимoisключения . . . . .                     | 5         |
| <b>3</b>  | <b>struct task_struct</b>                              | <b>6</b>  |
| <b>4</b>  | <b>Сокеты</b>  | <b>7</b>  |
|           | Сетевой стек . . . . .                                 | 8         |
| <b>5</b>  | <b>Модели ввода-вывода</b>                             | <b>10</b> |
| 5.1       | Модель блокирующего ввода-вывода . . . . .             | 10        |
| 5.2       | Модель неблокирующего ввода-вывода . . . . .           | 10        |
| 5.3       | Мультиплексирование ввода-вывода . . . . .             | 11        |
| 5.4       | Управляемый сигналом ввод-вывод . . . . .              | 11        |
| 5.5       | Асинхронный ввод-вывод . . . . .                       | 12        |
| <b>6</b>  | <b>Файлы</b>   | <b>13</b> |
|           | VFS . . . . .  | 13        |
|           | Суперблок . . . . .                                    | 14        |
|           | dentry . . . . .                                       | 16        |
|           | inode . . . . .  | 17        |
|           | Файл и файловая система . . . . .                      | 19        |
| 6.1       | Загружаемые модули ядра . . . . .                      | 21        |
| <b>7</b>  | <b>proc</b>  | <b>22</b> |
|           | Файлы последовательностей . . . . .                    | 24        |
|           | single_show . . . . .                                  | 26        |
| <b>8</b>  | <b>Открытые файлы</b>                                  | <b>27</b> |
|           | FILE . . . . .   | 28        |
| <b>9</b>  | <b>Слабы</b>   | <b>30</b> |
| <b>10</b> | <b>Аппаратные прерывания</b>                           | <b>31</b> |
|           | softirq . . . . .                                      | 32        |
|           | Тасклеты . . . . .                                     | 33        |
|           | Очереди работ . . . . .                                | 34        |
| <b>11</b> | <b>Специальные файлы устройств</b>                     | <b>36</b> |
| <b>12</b> | <b>USB</b>   | <b>38</b> |

# 1 Ликбез по 1 семестру

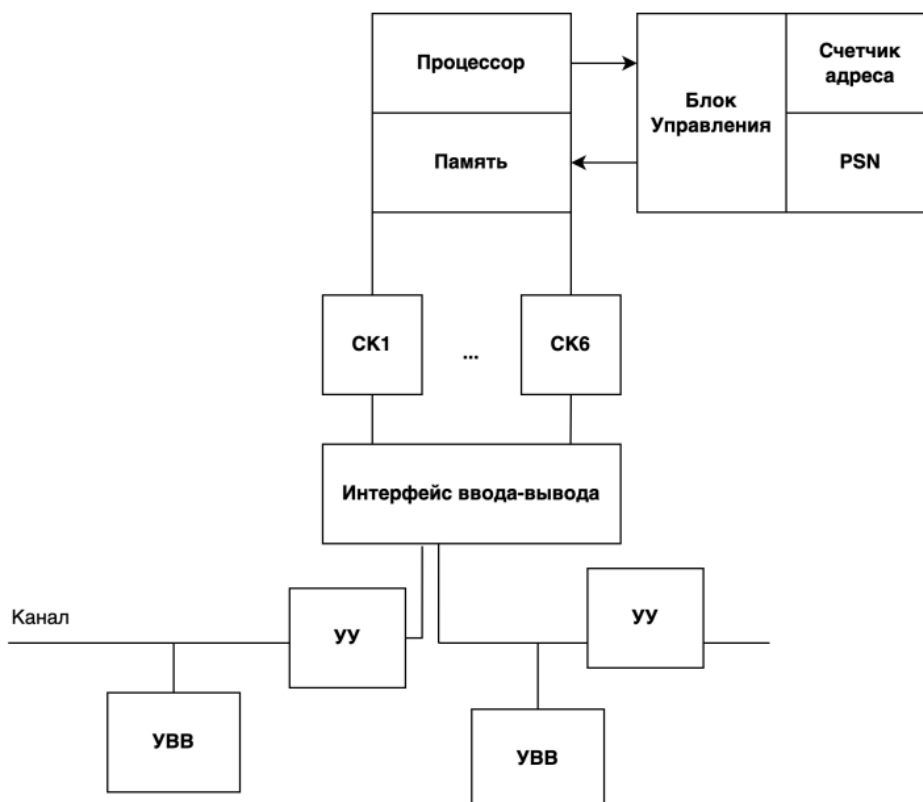
## 1.1 Концептуальная трехшинная архитектура

В компьютерах передаются сигналы трех типов: данные (данные и команды), адреса, сигналы управления.



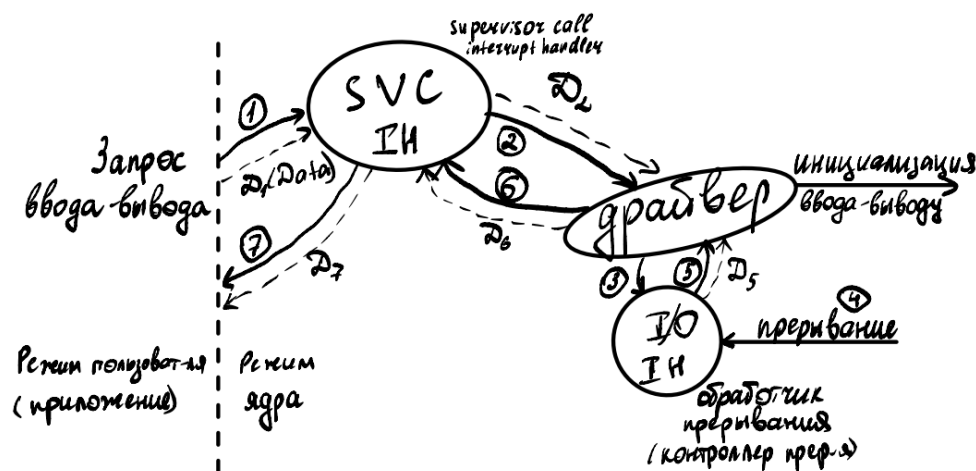
На линию IRQ приходит сигнал. Когда на контроллер приходит прерывание, он формирует сигнал INT который по шине управления поступает на выделенную ножку процессора. В конце цикла выполнения команды процессор проверяет наличие сигнала на выделенной ножке, если пришел сигнал прерывания то контроллер посылает INTA. Получив INTA контроллер по шине данных выставляет вектор прерывания (базовый вектор + номер линии IRQ). В реальном режиме процессор считывает вектор прерывания и переходит на его выполнение, переведя его смещение по таблице и считав адрес обработчика.

## 1.2 Канальная архитектура



Канал – программно-управляемое устройство – получает каналную программу от процессора. Быстрые устройства подключаются к селекторным каналам. Медленные устройства к мультиплексорным.

### 1.3 Прерывания в последовательности ввода-вывода



В современных системах для отдельно стоящей машины блокирующий ввод-вывод — основная модель ИО (запросив ИО, приложение блокируется, в система переходит в режим ядра и запрос обслуживается).

В этом виде ввода/вывода в полной мере реализуется идея распараллеливания функций (3-е поколение ЭВМ), которая базируется на стремлении освободить процессор от непроизводительных действий управления медленными внешними устройствами. Пока медленное внешнее устройство выполняет задачу ввода данных, процессор может перейти на выполнение другой работы.

В результате выполнения функций ядра будет вызван драйвер устройства, он формирует команды и данные и передает их по шине данных УВВ. Далее управление передается каналу или контроллеру. Контроллер сохраняет данные и возвращает их приложению.

### 1.4 Классификация прерываний

- Системные вызовы (программные прерывания) – процесс переключается в режим ядра и выполняет реентерабельный код ядра, после переключается в режим пользователя
- Исключения. Синхронные. 2 типа
  - Исправимые. страничное прерывание.
  - Неисправимые
- Аппаратные. Асинхронные. 3 типа
  - От системного таймера. Является единственным периодическим событием в системе. прерывание.
  - От устройств ввода/вывода.
  - От действий оператора

## 2 Возможно стоит знать

### 2.1 Архитектуры

Сетевое взаимодействие связано с внешним устройством – сетевой картой. Внешние устройства адресуются портами ввода-вывода. АП данных и УВВ начинаются с нуля, возникает наложение. Начиная с PowerPC введены два моста – северный, южный и bus mastering (захват шины). Шина PCI – между внешними устройства (вытеснило ISA шину). Мост определяет кому направить данные – внешнему устройству или памяти. Адресация УВВ и памяти контролируется с помощью захвата шины.

DMA (direct memory access) – передача данных не через регистры процессора. Контроллер прямого доступа к памяти позволяет передавать данные минуя процессор, не загружая данные из ОП и наоборот.

В SMP архитектуре все процессоры равноправны. Процессоры подключаются к общей памяти. NUMA альтернатива к SMP. NUMA (non-uniform memory access) – все процессоры имеют прямой доступ к памяти определенного объема, и могут обращаться к памяти других процессоров. Это позволяет процессорам работать с памятью параллельно.

В шинах PCI не используются физические линии прерываний, прерывания доставляют как MSI

### 2.2 Мультиплексирование

- select/poll. Устаревшее API
- epoll

Существует два режима наблюдения

- Level-triggered – файловый дескриптор возвращается, если остались непрочитанные/записанные данные
- Edge-triggered – файловый дескриптор с событием возвращается, только если с момента последнего возврата произошли новые события (например, пришли новые данные)

epoll instance – ключевой концепт epoll API, структура данных содержащая список файловых дескрипторов для мониторинга и список готовых дескрипторов для работы.

*set\_nonblocking(fd)* – переводит сокет в неблокирующий режим.

*epoll\_create1(0)* – создаёт структуру данных (epoll instance), с которой в дальнейшем идёт работа. Структура одна для всех файловых дескрипторов, за которыми идёт наблюдение. Функция возвращает файловый дескриптор, который в дальнейшем передаётся во все остальные вызовы epoll API.

*epoll\_ctl(epollfd, op, fd, event)* – используется для управления epoll instance, в частности, позволяет выполнять операции EPOLL\_CTL\_ADD (добавление файлового дескриптора к наблюдению), EPOLL\_CTL\_DEL (удаление файлового дескриптора из наблюдения), EPOLL\_CTL\_MOD (изменение параметров наблюдения), EPOLL\_CTL\_DISABLE – для безопасного отключения наблюдения за файловым дескриптором в многопоточных приложениях. EPOLLET устанавливает edge-triggered режим.

*epoll\_wait(epollfd, events, maxevents, timeout)* – возвращает количество (один или более) файловых дескрипторов из списка наблюдения, у которых поменялось состояние (которые готовы к вводу-выводу).

События, за которыми можно наблюдать с помощью epoll:

- EPOLLIN – новые данные (для чтения) в файловом дескрипторе
- EPOLLOUT – файловый дескриптор готов продолжить принимать данные (для записи)
- EPOLLERR – в файловом дескрипторе произошла ошибка
- EPOLLHUP – закрытие файлового дескриптора

### 2.3 Монитор Хоара

2 типа процессов – читатели, могут читать параллельно; писатели, работают в режиме монопольного доступа. 5 семафора: 3 считающих: активные читатели, очередь ждущих писателей, очередь ждущих читателей и 2 бинарный: активный писатель и захват критической секции.

В мониторе определены 4 процедуры: start\_read, stop\_read, start\_write, stop\_write.

Для того, чтобы читатель начал читать, он вызывает start\_read. Читатель может читать, если нет активного писателя и очередь ждущих писателей пуста. В противном случае процесс блокирован на

переменной типа условие `may_read`. Увеличивается число активных читателей и посылается сигнал, что можно читать – (каждый читатель активизирует следующего). Когда читатель завершает чтение он вызывает `stop_read`. Уменьшается число активных читателей. Если нет активных читателей, то посылается сигнал «можно писать».

Для того, чтобы писатель начал писать, он вызывает `start_write`. Писатель может писать, если нет активного писателя и активных читателей. В противном случае процесс блокирован на переменной типа условие `may_write`. Когда процесс заканчивает писать – вызывает `stop_write`. Есть в `stop` есть ждущие читатели, дальше читают они, иначе писатель начинает писать.

В мониторе исключается бесконечное откладывание читателей/писателей.

## 2.4 Семафоры

Набор семафоров, чтобы одной неделимой операцией изменять состояние над частью семафоров набора

Семафоры предложил Дейкстра для того, чтобы убрать активное ожидание (на проверку в цикле тратится процессорное время). Семафор – не отрицательная защищенная переменная, на которой определены 2 неделимые операции: `P`(декремент) и `V`(инкремент). Если семафор принимает только знач. 0,1, то он называется бинарным. Если может принимать неотрицательные целые - считающий. Операция `P(S)` при `S = 0` заблокирует процесс, чтобы остаться в области целых неотрицательных чисел. Операция `V(S)` при `S = 0` разблокирует процесс, ожидающий семафор. За экономией времени платим переходом в режим ядра (`P` и `V` - системные вызовы). Считающие семафоры контролирует значения величин, бинарные – монополярный доступ.

`semop` – операции на семафоре. Для `semop` определена структура `struct sembuf` – номер семафора (наборы семафоров представлены как массивы семафоров), в наборе каждый семафор достается по номеру; `semop` – операция; `semflg` – флаги. `semop` получает `semid`, указатель на массив структур и `intlen` – количество семафоров, на которых будет выполнен `semop` (неделимая операция). `semctl` производит операции управления.

## 2.5 Системные вызовы

`exec` – переводит процесс на новое адресное пространство программы, путь к которой передаётся в качестве аргумента.

`fork` – создает новый процесс потомок, который является копией процесса предка, в том смысле, что он наследует от предка дескрипторы открытых файлов, сигнальную маску и переменные окружения.

## 2.6 Способы взаимного исключения

- программный
- аппаратный
- с помощью семафоров
- с использованием мониторов

### 3 struct task\_struct

Каждый процесс в системе описывается структурой struct task\_struct. В загружаемом модуле можно напрямую получить доступ к полям структуры struct task\_struct, которая содержит информацию о процессе. В ядре определен специальный символ current, который является указателем на структуру struct task\_struct для текущего процесса.

```
1 static int __init my_init(void)
2 {
3     struct task_struct *task;
4     do {
5         printk(KERN_INFO "%d %s %d %d", task->pid, task->comm, task->ppid,
6         task->prio);
7     } while(task = next_task(task) != &init_task);
8     return 0;
9 }
```

#### Флаги

- PF\_NOFREEZE - не может быть блокирован.
- PF\_FORKNOEXEC - не может вызвать exec(), то есть не может сменить адресное пространство по адресное пространство программы, передаваемой в качестве фактического параметра exec(). (не позволяет внедрить вредоносный код в прогу).
- PF\_NOSETAFFINITY - запрещает пользователю менять cpu\_mask
- PF\_KTHREAD - поток ядра

Приоритет: чем меньше значение, тем выше приоритет. Приоритет 0 является максимальным. Для RT-процессов (реального времени) диапазон: 1–99. Для обычных процессов: 100–139.

- SCHED\_NORMAL - используется для обычных процессов в системе. Выделяется процессорное время на основе приоритета - процесс с более высоким приоритетом будет выполняться до тех пор, пока не истечет его квота или пока система не переключится на другой процесс.
- SCHED\_FIFO - процессы организуются в FIFO-очередь. Он будет выполняться, пока не завершит свою задачу или не будет прерван. Является политикой планирования без переключения
- SCHED\_RR - процессы организуются в очередь. Если процесс исчерпал квант, но не успел завершиться, то он помещается в конец очереди. Является политикой планирования с переключением:
- SCHED\_BATCH - предназначена для выполнения «пакетных» задач, которые могут длиться долго и не требуют быстрых откликов. Пакетные процессы выполняются, когда в системе нет активных высокоприоритетных задач.
- SCHED\_IDLE - фоновый (холостой) процесс. Предназначена для процессов, которые выполняются в фоновом режиме и имеют самый низкий приоритет.

Справедливость SCHED\_NORMAL в том, чтобы в итоге все процессы получили процессорное время и по возможности в равных количествах.

#### Состояния

- TASK\_RUNNING - Задача выполняется или готова к выполнению
- TASK\_INTERRUPTIBLE - Задача спит, можно разбудить сигналом
- TASK\_UNINTERRUPTIBLE - Задача спит, нельзя разбудить сигналом

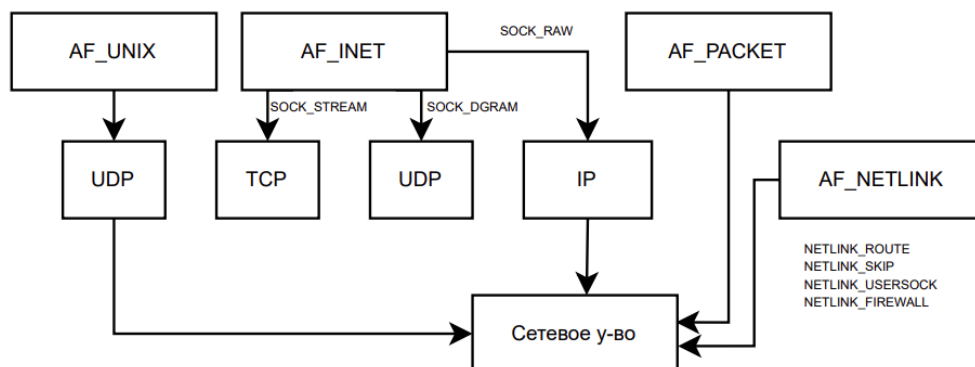
Процесс миграции, например, migration/0, является специальным процессом, который помогает системе управлять процессами, которые необходимо перенести с одного ядра на другое.

## 4 Сокеты

Средства взаимодействия параллельных процессов. Универсальность сокетов – позволяют организовать взаимодействие и на отдельной машине и в распределенной системе.

Сокет – абстракция конечной точки взаимодействия. Создание сокета: *int socket(int family, int type, int protocol)*, family – семейство (сокеты могут принадлежать различным семействам), тип и протокол.

Если сокеты используются для взаимодействия процессов на отдельной машине, то по сути адреса – специальные файлы. В распределенных системах используются сокеты AF\_INET.



Сетевое устройство имеет драйвер и когда на карту приходят пакеты приходит аппаратное прерывание.

Ядро Linux предоставляет для работы с сокетами системный вызов

```
1 asmlinkage long sys_socketcall(int call, unsigned long *args) {
2     unsigned long a[6];
3     int err;
4     if (copy_from_user(a, args, nargs[call]))
5         return -EFAULT;
6     unsigned long a0=a[0];
7     unsigned long a1=a[1];
8     switch(call)
9     {
10         case SYS_SOCKET: err = sys_socket(a0,a1,a[2]); break;
11         case SYS_BIND: err = sys_bind(a0, (struct sockaddr *)a1, a[2]); break;
12         case SYS_CONNECT: err = sys_connect(a0, (struct sockaddr *)a1, a[2]);
13         break
14         ...
15         default: err = -EINVAL; break;
16     }
17 }
```

В switch перечисляют средства все функции сетевого стека. Взаимодействие на сокетах всегда клиент-сервер, даже на отдельных машинах.

```
1 asmlinkage long sys_socket(int family, int typ, int protocol) {
2     struct socket *sock;
3     ...
4     int retval = sock_create(family, type, protocol, &sock);
5     ...
6     return retval;
7 }
```

```
1 struct socket {
2     socket_state state;
3     short type;
4     unsigned long flags;
5     const struct proto_ops *ops; // протокол
6     struct fasync_struct *fasync_list; // список асинхронного запуска
7     struct file *file; // дескриптор файла
8 }
```

```

8     struct sock *sk;
9     wait_queue_head_t wait;
10 };

```

У сокета различают 5 состояний

- SS\_FREE - свободный сокет, с которым можно соединяться;
- SS\_UNCONNECTED - несоединенный сокет;
- SS\_CONNECTING - сокет находится в состоянии соединения;
- SS\_CONNECTED - соединенный сокет;
- SS\_DISCONNECTING - сокет разъединяется в данный момент

Дизайн сокетов следует парадигме UNIX: в идеале отобразить все объекты которые осуществляют чтение/запись на файлы, чтобы с ними можно было работать используя обычные функции чтения/записи из/в файлы.

Семейства

- AF\_UNIX – для взаимодействия на локальной машине;
- AF\_INET – сокеты семейства протоколов TCP, IP, UDP, основанные на протоколе IPv4;
- AF\_INET6 – IPv6;
- AF\_IPX – протокол IPX;
- AF\_UNSPEC – неопределенный домен.

Тип

- SOCK\_STREAM – потоковые сокеты, ориентированные на потоки. Надежное, упорядоченное, полнодуплексное логическое соединение;
- SOCK\_DGRAM – — определяет ненадежную службу без установки логического соединения, где пакеты могут передаваться без сохранения порядка;
- SOCK\_RAW – — низкоуровневый сокет, обеспечивает прямой доступ к сетевому протоколу.

Протокол

- 0 – если AF\_INET и SOCK\_STREAM – по умолчанию TCP, AF\_UNIX и SOCK\_DGRAM – по умолчанию UDP;
- явное указание IPPROTO\_\* (например, IPPROTO\_TCP)

Системный вызов socket не определяет, что с чем соединять, то есть не определяет адрес

```

1 struct sockaddr
2 {
3     sa_family_t sa_family;
4     char sa_data[14];
5 }

```

Парные сокеты *int socketpair(int domain, int type, int protocol, int sv[2]);* - альтернатива программным каналом, обеспечивают дуплексную связь, но не реализовано взаимоисключающие. sockaddr – базовый адрес в работе с сокетами. Для сокетов AF\_UNIX в массиве имя файла и семейство. AF\_UNIX реализуется в пространстве файловых имен. Для сетевых сокетов необходима детализация и используется sockaddr\_in.

На адресацию в сети 14 байт не хватит, введена структура для адресации по сети:

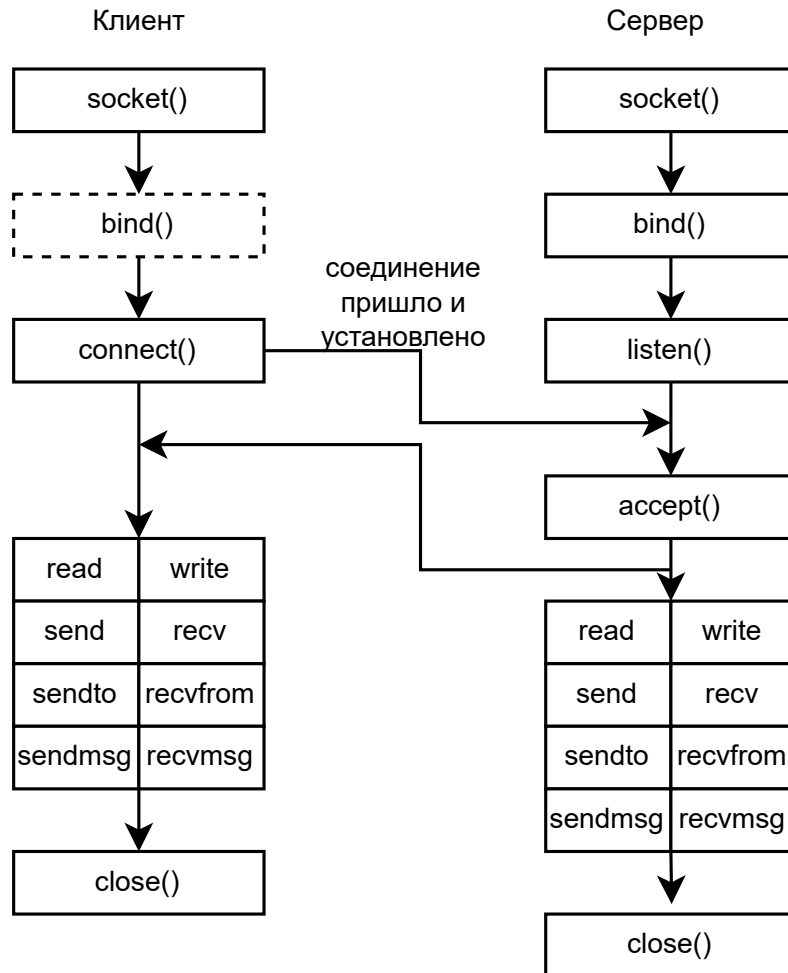
```

1 struct in_addr {
2     _u32 s_addr;
3 };
4
5 struct sockaddr_in
6 {
7     short int sin_family; // Семейство
8     unsigned short int sin_port; // порт
9     struct in_addr sin_addr; // IP-адрес
10    unsigned char sin_zero[8]; // выравнивание
11 };

```

Модель клиент-сервер позволяет создать и сервер и клиент и настроить взаимодействие на одной машине. На сокетах определены системные вызовы, образующие сетевой стек. Стеком является только потому, что функции вызываются в определенном порядке.





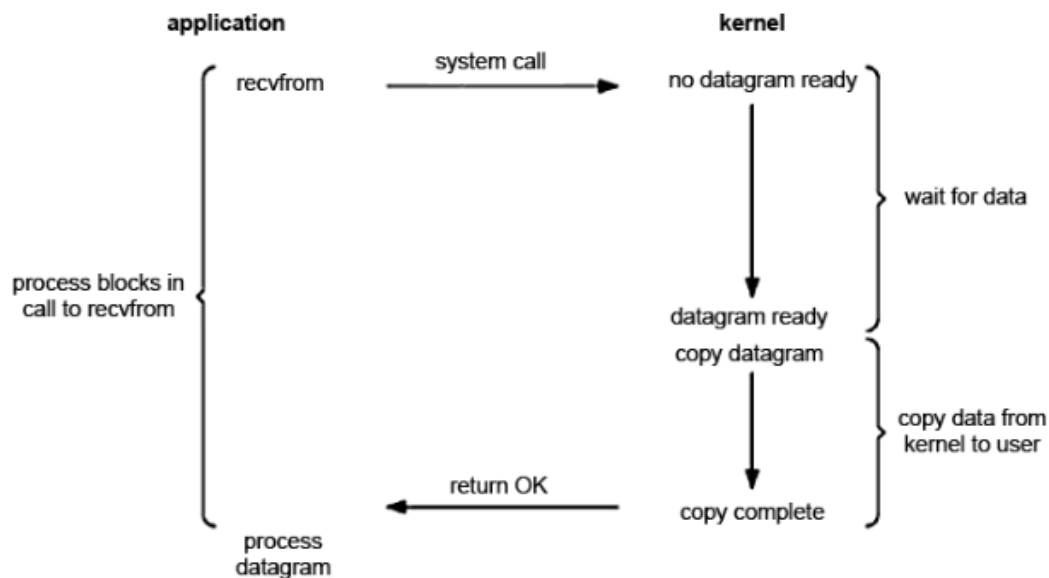
- *int socket(int family, int type, int protocol)* – возвращает файловый дескриптор, который нужно связать с адресом.
- *int bind(int sockfd, struct sockaddr \*addr, int addrlen);*  
Связывает сокет с адресом. На стороне клиента нужен не всегда. После связывания переходит в режим пассивного ожидания соединения.
- *int connect(int sockfd, struct sockaddr \*serv\_addr, int addrlen);*  
Клиенты присоединяются к серверу пытаясь установить связь. Клиент знает адрес сервера
- *int listen(int sockfd, int backlog);*  
Вызывается сервером для получения сообщения от клиентов. Сервер сообщает ядру о переходе в состояние пассивного ожидания запросов клиента.
- *int accept(int sockfd, void \*addr, int \*addrlen);*  
Выполняется на стороне сервера для того, что принять соединение. Если пришел запрос на соединение, то вызывается ассепт. Создает копию исходного сокета, возвращает файловый дескриптор исходного дескриптора, чтобы исходный сокет оставался в режиме прослушивания соединений, а копия находилась в состоянии соединения.

Протокол TCP реализован в виде тройного рукопожатия, чтобы обеспечить надежную связь клиента и сервера.

Два типа порядка байтов: Прямой (host byte order) от старших к младшим разрядам и обратный, сетевой (network byte order) наоборот. Функции для конвертации *htons*, *ntohs* и *htonl*, *ntohl*.

## 5 Модели ввода-вывода

### 5.1 Модель блокирующего ввода-вывода



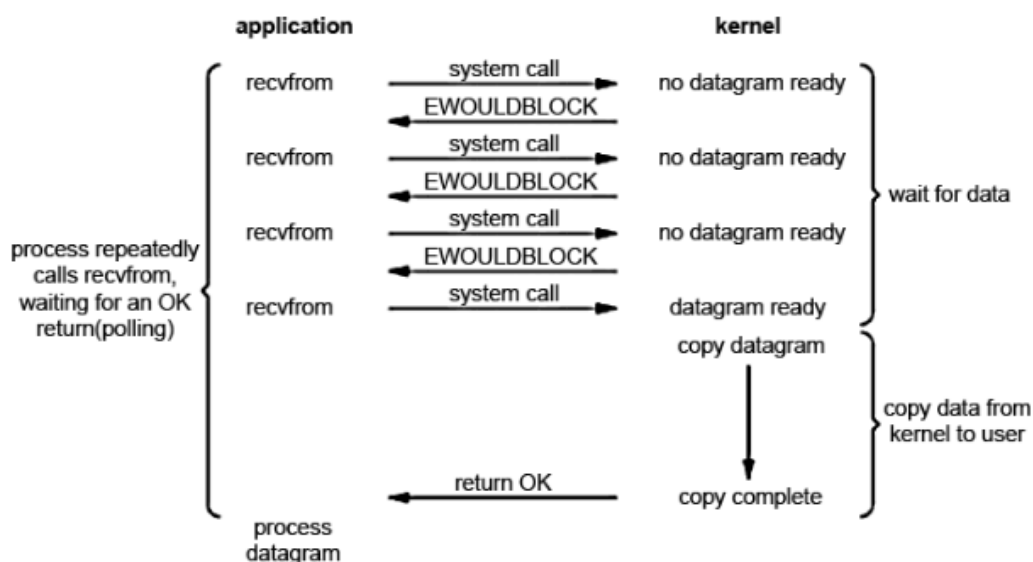
При вводе-выводе, система должна инициировать обращение к внешнему устройству. Получить сразу по запросу приложение данные не может. Обращение к внешним устройствам — длительное действие.

Запросив ввод-вывод, процесс будет блокирован (на диаграмме — `no datagram ready`).

Когда внешнее устройство завершает операцию ввода-вывода, он формирует сигнал прерывания, по которому формируется вектор прерывания. По этому вектору процессор адресует обработчик прерывания.

Данные из буфера устройства записываются в буфер ядра и процесс должен быть разблокирован.

### 5.2 Модель неблокирующего ввода-вывода



Исторически более ранняя модель, процессор тратит время на опрос о готовности данных, нет прерываний.

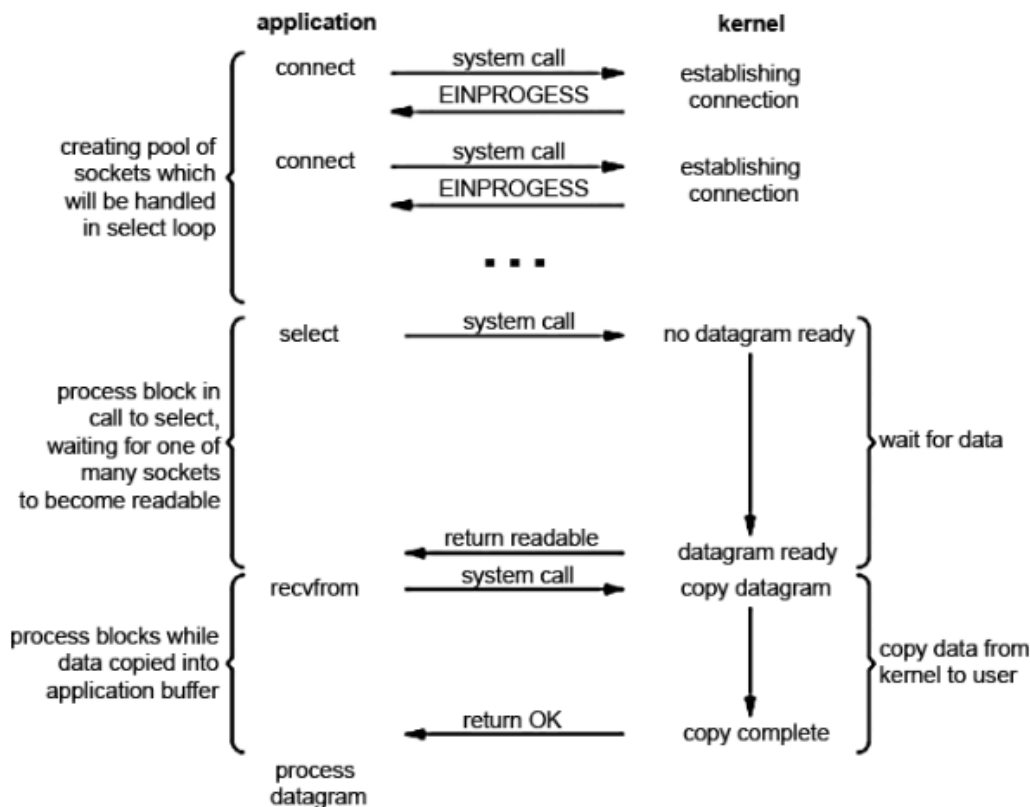
В какой-то момент пакеты данных начинают приходить с большой частотой на сетевую карту — много аппаратных прерываний. Обслуживание аппаратных прерываний затратно. NAPI (new api), переход в режим поллинга (опроса), что быстрее АП.

### 5.3 Мультиплексирование ввода-вывода

Для такой реализации нужно системные вызовы *poll()*, *select()*, *epoll()*, *pselect()*.

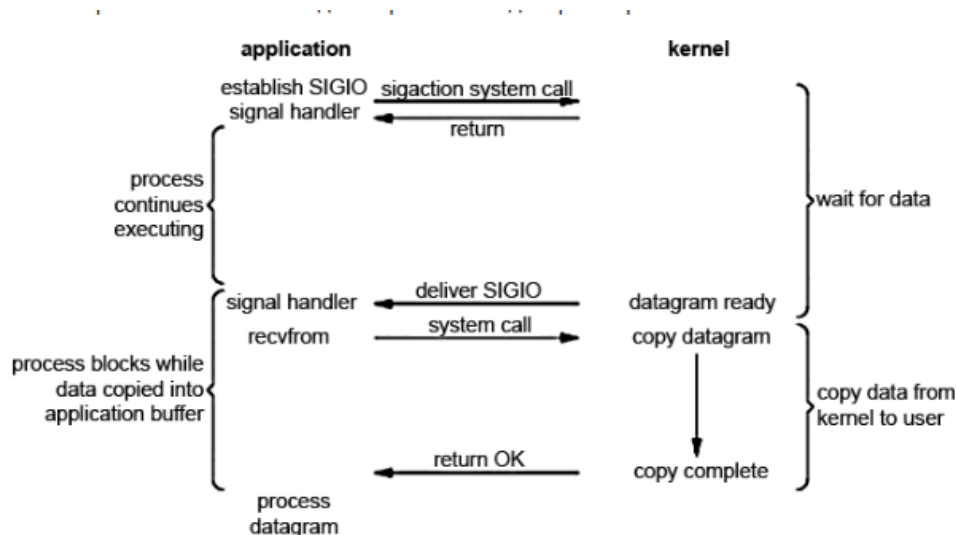
Мультиплексор вызывается до *accept()*, после *listen()*. Ядро выполняет мультиплексирование, начиная опрашивать готовность сокетов. Блокировка по сути разделена на две. Первая блокировка связана с *select()* и закончится когда любой из сокетов станет доступным для взаимодействия. В результате будет установлено соединение и начнется обмен данными. Процесс будет блокироваться так как необходимо получить данные в результате *receive()*

Мультиплексор объединяет информацию из несколько входов и выдает ее по одному входному каналу.



Преимущество мультиплексирования – сокращение времени блокировки.

### 5.4 Управляемый сигналом ввод-вывод

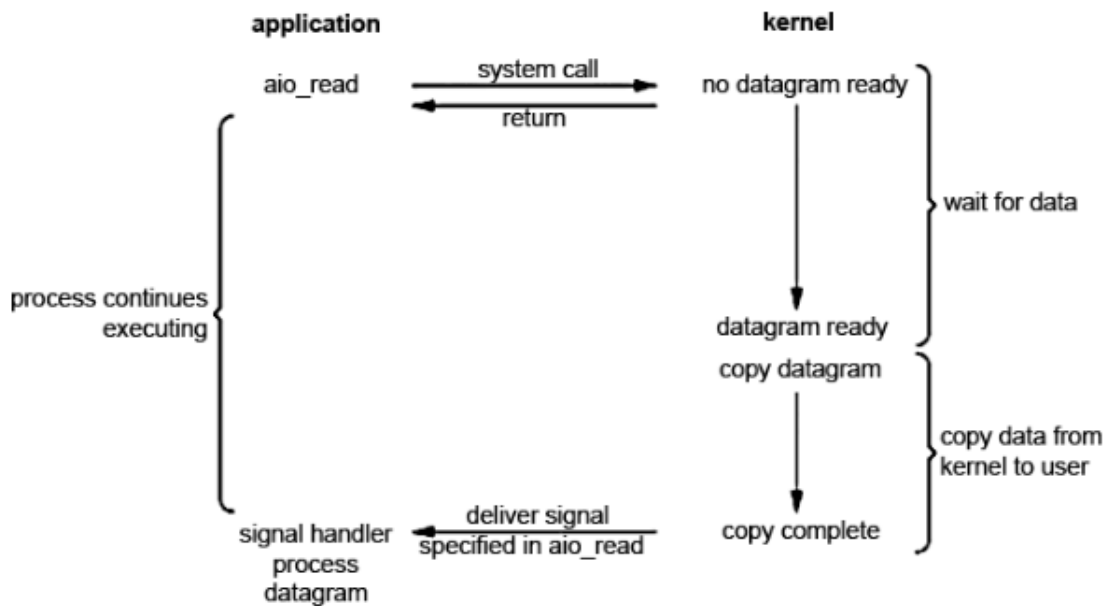


Ядро отслеживает когда данные будут готовы, оно посылает сигнал и вызывается обработчик,

так что приложение не нужно блокировать. Запросив данные может выполняться другое действие не связанное с этим.

## 5.5 Асинхронный ввод-вывод

Ядро должно предоставлять пользователю специальные системные вызовы асинхронного ввода-вывода.



Проблема как и в прошлой модели – найти работу, которую может выполнять процесс пока данные не готовы.

|             | Блокирующий                                     | Неблокирующий            |
|-------------|---|--------------------------|
| Синхронный  | read, write                                     | read, write (O_NONBLOCK) |
| Асинхронный | I/O мультиплексирование (select/poll) + сигналы | AIO                      |

## 6 Файлы

Файл – средство долговременного хранения данных, любая поименованная совокупность данных, может быть бессмысленная, во вторичной памяти. В UNIX все файл.

Файловая система – часть ОС, которая определяет способ организации, хранения и доступа к файлам.

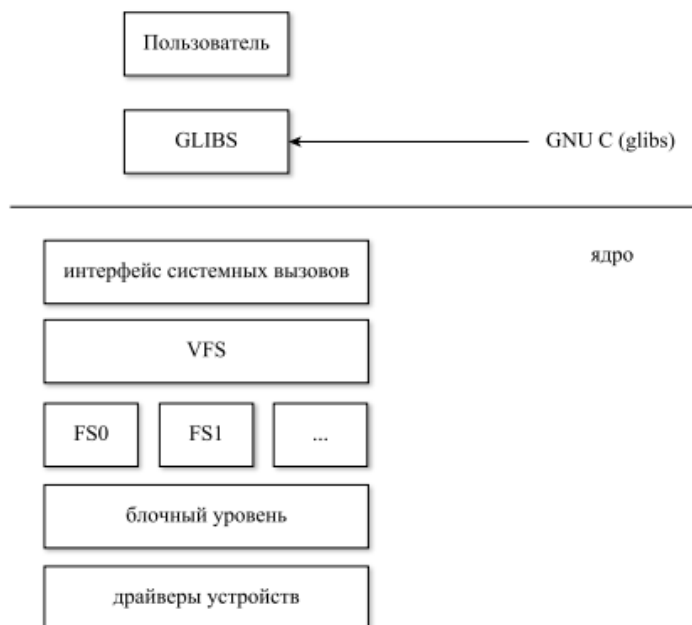
UNIX поддерживает большое количество различных файловых систем. ФС UNIX организована через интерфейс VFS/vnode. В Linux не определена vnode, только VFS. Задача vnode – чтобы ОС без изменений ядра могла поддерживать большое кол-во самых разных ФС.

Ядро UNIX имеет иерархическую структуру по отношению к аппаратному обеспечению. Принято выделять следующие уровни ФС:



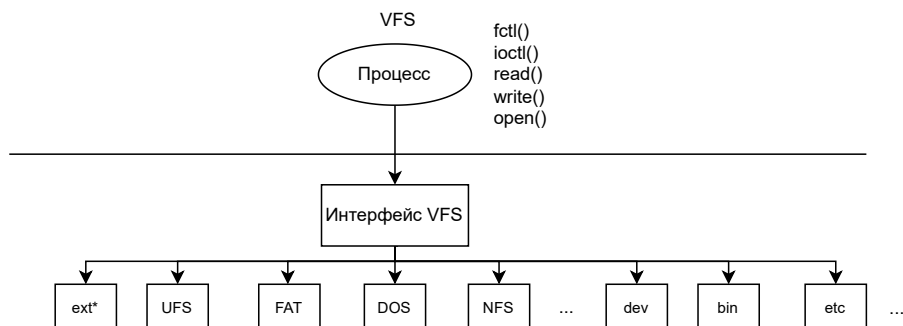
Символьный – именование файлов, строки символов. Не являются идентификатором. Имя файла – это полное имя файла начиная с корневого каталога, содержит все элементы пути и завершается shortname. inode – базовый уровень – уникальный идентификатор файла. Любой файл в логическом представлении начинается с нулевого смещения и его АП непрерывно. Файлы бывают текстовые и бинарные (байтоориентированные). Linux поддерживает FHS (File system standart). FHS определяет структуру и содержимое каталогов. Существует корневая ФС – «\». Корневой каталог и его ветви должны составлять единую ФС, расположенную на одном разделе диска. Существует 2 типа устройств – символьные и блочные.

Структура VFS: четыре структуры для представления VFS – super\_block, inode, dentry и file.

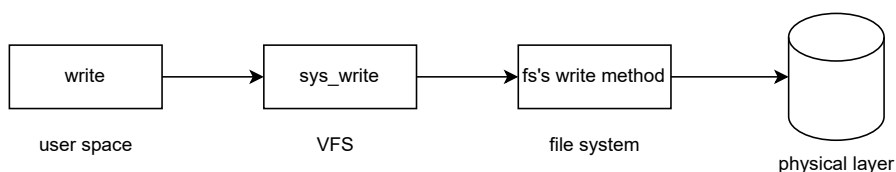


Разделенные на независимые порции потоки данных диска называются секторами. Сектор диска

– участок дорожки магнитного диска представляющий собой наименьший размер порции данных, которая может быть изменена в результате форматирования (номер дорожки + номер сектора). На диске вся информация адресуется, имеет уникальный адрес. Раздел – выделенная на диске АП, характеризуется группами блоков.



Операции над файлами в любой ФС будут подобны



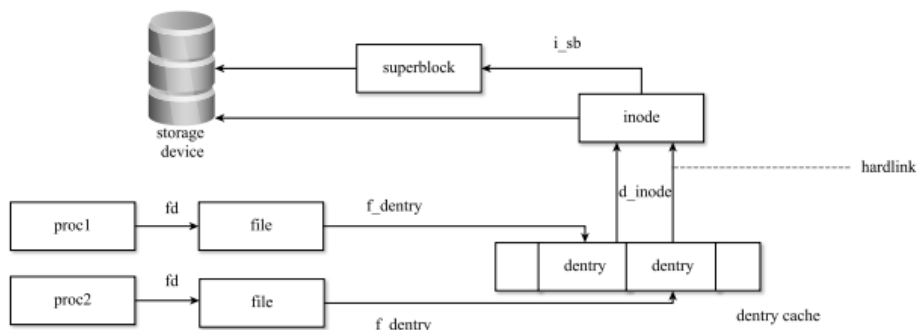
Монтирование ФС – например подключение флешки (у нее своя ФС).

Дерево файлов и каталогов Linux формируется из ветвей которые могут соответствовать различным физическим носителям (дерево формируется из отдельных ФС).

ФС – большое дерево с корневой директорией «\». Монтирование корневой ФС – часть процесса инициализации. Монтирование – последовательность действий в результате которой ФС становится доступной.

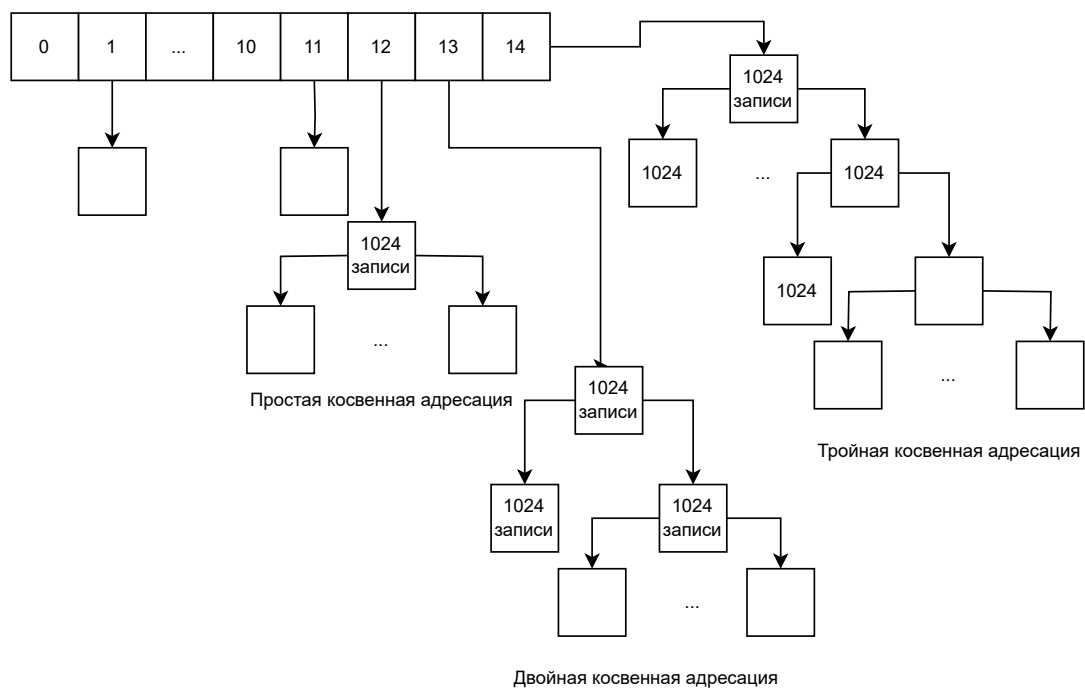
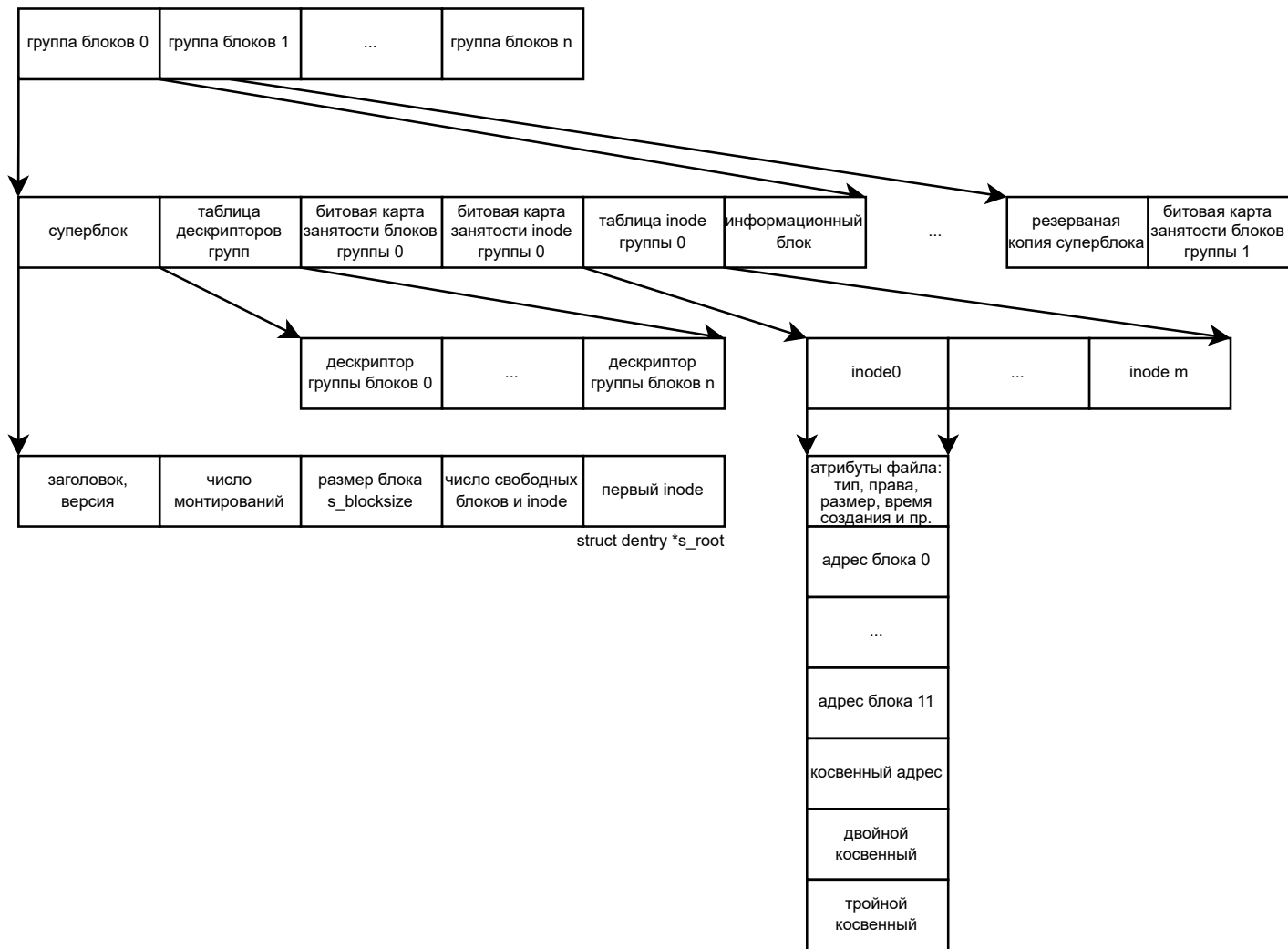
Когда ФС монтируется в существующую директорию все файлы и директорий этой ФС становятся файлами и директориями точки монтирования. Если эти точки содержат какие-то файлы или директории, то они становятся невидимыми.

*mount <ключи> -t <тип ФС> -o <опции ФС> <имя УВВ/Другого устройства содержащего ФС> <точки монтирования>*



Суперблок – структура, содержащая информацию необходимую для монтирования и для доступа к файлам. В каждой ФС ровно 1 суперблок. Инициализированный суперблок располагается в начале раздела. Все АП делится на разделы, в начале суперблок, содержащий всю необходимую информацию для доступа к физическим файлам которые будут храниться в рамках ФС.

struct superblock – структура, описывающая смонтированную ФС. Содержит метаданные для обращения к файлам ФС. АП диска делится на разделы. Один раздел занимает одна ФС.



Таким образом можно хранить огромные файлы, благодаря тому что они хранятся по частям в свободных блоках диска.

```

1 struct superblock {
2     struct list_head s_list;           // обеспечивает организацию связного сп
иска
3     dev_t s_dev; // информация об устройстве, на котором находится файловая система
4     unsigned long s_block_size;        // размер блока в байтах
5     loff_t s_maxbytes;                  // максимально возможный размер файла
6     struct file_system_types *s_type;   // тип файловой системы
7     const struct super_operations *s_op; // все операции, определенные на суперб
локе
8     const struct export_operations *s_export_op; // операции, не учтенные ядром
9     ...
10    unsigned long s_flags;               // флаги монтирования
11    ...
12    unsigned long s_magic;               // магическое число для доступа к супер
блоку
13    struct dentry *s_root;               // точка монтирования
14    ...
15    int s_count;                         // счетчик ссылок на суперблок
16    struct list_head s_mounts;           // список монтирований
17    struct block_device *s_bdev;         // связанное блочное устройство
18    char s_id[32];                       // символьное имя
19    ...
20    fmode_t s_mode;                      // права доступа для монтирования
21    const struct dentry_operations *s_d_op; // в этом суперблоке операции над
dentry
22    struct list_head s_inodes;           // список всех inode'ов
23    ...
24 }

```

```

1 struct super_operations {
2     struct inode *(*alloc_inode)(struct super_block *sb);
3     void (*destroy_inode)(struct inode *);
4     void (*free_inode)(struct inode *);
5     void (*dirty_inode) (struct inode *, int flags); // вызывается VFS, если в
inode были внесены изменения, чтобы не потерять изменения и отразить их в дисковом
inode
6     int (*write_inode) (struct inode *, struct writeback_control *wbc); // записыва
ет указанный inode на диск
7     int (*drop_inode) (struct inode *); // вызывается VFS, когда удаляется последня
я ссылка на inode.
8     void (*put_super) (struct super_block *); // вызывается VFS при размонтировании
файловой системы
9     int (*statfs)(struct dentry *, struct kstatfs *);
10 };

```

В перечне функций super\_operations нет функции создания суперблока, но такая функция есть.

```

1 static struct super_block* alloc_super(struct file_system_type *type, int flags,
    struct user_name_space *user_ns) {
2     struct super_block* s = kzalloc(sizeof(struct super_block), GFP_USER);
3     static const struct super_operations default_op;
4     ...
5     INIT_LIST_HEAD(&s->s_mounts);
6     ...
7     INIT_LIST_HEAD(&s->s_inodes);
8 }

```

Для удаления суперблока используются функции kill\_block\_super, kill\_anon\_super, kill\_litter\_super. struct dentry – определяет элемент пути в имени файла (короткое имя + путь к файлу, то есть все директории с корневого каталога).

```

1 struct dentry {
2     unsigned int d_flags;

```



```

3 struct hlist_bl_node d_hash;
4 struct dentry *d_parent; // Объект dentry родительского каталога
5 struct qstr d_name; // Имя dentry
6 struct inode *d_inode; // Связанный inode
7 unsigned int d_count; //
8 spinlock_t d_lock;
9 unsigned long d_time;
10 void *d_fsdata;
11 unsigned char d_iname[DNAME_INLINE_LEN]; // короткое имя файла
12 const struct dentry_operations *d_op;
13 struct super_block *d_sb;
14 struct list_head d_lru; // указатель на список lru
15 struct list_head d_children;
16 ...
17 struct list_node d_hash; // объекты dentry кэшируются
18 int d_mounted; // Является ли объект точкой монтирования
19 };

```

Состояния dentry: используется (in use), не используется (unused), отрицательный – не связанный с допустимым индексным узлом (d\_inode == NULL). Если in\_use, то данные dentry соответствуют допустимому inode, у такого dentry один или несколько пользователей. Если d\_count == 0, то состояние unused, то есть dentry соответствует допустимому inode, но d\_count == 0. Такой dentry все еще хранится в кэше. Если dentry отрицательный, он все равно сохраняется так как какие-то запросы могут к нему обращаться, но при необходимости можно его и удалить. Для отрицательного dentry фактически нет информации на диске.

Кэш объектов dentry: сам каталог dentry является файлом специального типа. Чтобы каждый раз не обращаться к диску за информацией о dentry, информация заносится в кэш. При проходе по пути к файлу система сначала смотрит на кэш.

3 части кэша: список актуальных dentry, связанные с определенным индексным узлом (может быть несколько dentry с одним и тем же inode). Двусвязный список dentry\_lru – этот список отсортирован по времени обращения к dentry. Каждому объекту dentry присваивается временная метка и список перестраивается при обращении к этому объекту. Если ядру нужно освободить память, элемент удаляется из хвоста списка. Хэш-таблица и хэш-функция для обращения к указанному файлу: таблица в виде массива dentry\_hashtable в котором каждый элемент имеет указатель на dentry. Размер таблицы зависит от объема памяти, значение хэша вычисляется по функции d\_hash(). Поиск по таблице с помощью d\_lookup(), которая возвращает или NULL (если хэш не найден) или указатель если найден.

dentry не соответствует какой-либо структуре данных на диске, он создается на лету ФС. Так как чтение и запись каталога с диска и создание соответствующего объекта dentry требует достаточного времени имеет смысл хранить объекты dentry в памяти, даже если работа с ними закончена, чтобы обратиться к элементу пути повторно при обращении к тому же файлу.

/home/mdir/cf/my\_server.c: VFS при каждом обращении должна пройти по всем dentry данного файла начиная с корневого каталога. Если бы не было кэша то каждый раз нужно было бы обращаться к диску

```

1 struct dentry_operations {
2     int (*d_revalidate)(struct dentry *, unsigned int);
3     int (*d_hash)(const struct dentry *, struct qstr *);
4     int (*d_compare)(const struct dentry *, unsigned int, const
5     int (*d_init)(struct dentry *);
6     void (*d_release)(struct dentry *);
7     char *(*d_name)(struct dentry *, char *, int);
8 };

```

d\_hash создает значения хэша, d\_compare – вызывается когда сравниваются два имени файла, d\_del вызывается при d\_count == 0, d\_release – когда объект dentry теряет связь с inode.

Существует 2 структуры inode: дисковый – описывающий файл и представляющий информацию об адресах блоков в которых находится информация из файла. inode ядра – описывает файл в ядре для ускорения доступа к дисковому inode.

struct inode – описывает конкретный физический файл: в Linux нет понятия vnode в отличие от Unix. Но все еще 2 типа inode: ядра и диска. inode ядра содержит информацию о файле, что делает получение информации быстрее. inode диска содержит адреса блоков диска с данными файла.

```

1 struct inode
2 {
3     umode_t i_mode; //права доступа
4     unsigned short i_opflags;
5     kuid_t i_uid;
6     kgid_t i_gid;
7     unsigned int i_flags;
8     ...
9     const struct inode_operations *i_op; /* указатель на ту таблицу операций с инде
ксом которые определены в конкретной системе для работы с этим индексом.
Такая структура создаётся для конкретной файловой системы */
10    struct super_block *i_sb; /*указатель на связанный суперблок, так как именно он
хранит информацию об inode данной файловой системы. */
11    struct address_space *i_mapping;
12    ...
13    unsigned long i_ino; // индекс inode
14    union
15    {
16        {
17            const unsigned int i_nlink;
18            unsigned int _i_nlink;
19        };
20        ...
21        struct list_head i_lru; // список inode построенный по LRU.
22        ...
23        union{
24            struct hlist_head i_dentry;
25            struct rcu_head i_rcu;
26        };
27        ...
28        atomic_t i_count; /* счётчик ссылок */
29        ...
30        atomic_t i_writecount; /*счётчик использования для записи */
31        ...
32        count struct file_operations *i_top;
33        struct address_space i_data;
34        struct list_head i_devices; /* список блочных устройств */
35        ...
36    }

```

```

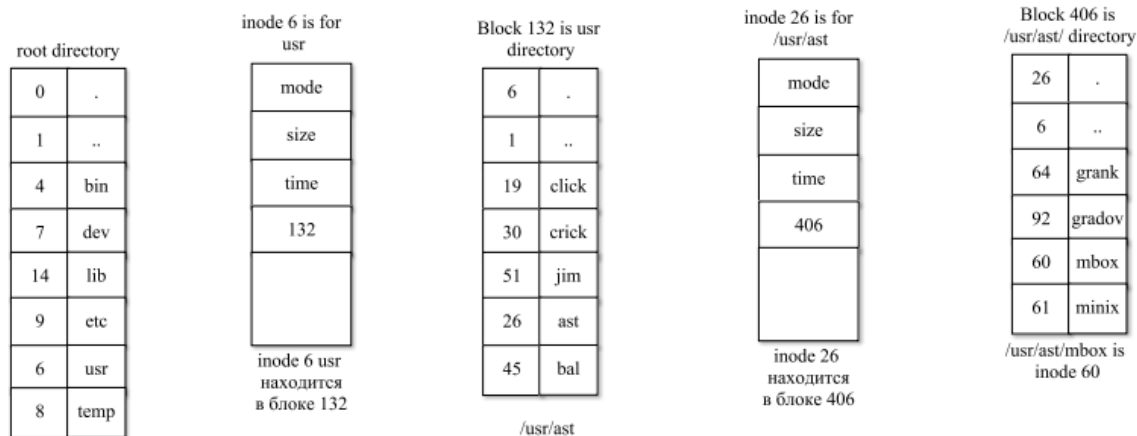
1 struct inode_operations {
2     struct dentry * (*lookup) (struct inode *,struct dentry *, unsigned int); // вы
полняет поиск inode в указанном каталоге
3     ...
4     int (*create) (struct inode *,struct dentry *, umode_t, bool); // создает дескр
иптор файла
5     int (*link) (struct dentry *,struct inode *,struct dentry *);
6     int (*unlink) (struct inode *,struct dentry *);
7     int (*symlink) (struct inode *,struct dentry *,const char *);
8     int (*mkdir) (struct inode *,struct dentry *,umode_t);
9     int (*rmdir) (struct inode *,struct dentry *);
10    int (*mknod) (struct inode *,struct dentry *,umode_t,dev_t);
11    int (*rename) (struct inode *, struct dentry *, struct inode *, struct dentry
*, unsigned int);
12    int (*setattr) (struct dentry *, struct iattr *);
13    int (*getattr) (const struct path *, struct kstat *, u32, unsigned int);
14    ...
15 }

```

Шаги при доступе к файлу */usr/est/mbob*

ФС ищет компонент пути *usr* в корневом каталоге, чтобы найти номер inode файла */usr*. По номеру inode система может найти адрес блока. Получив адрес блока система получает содержимое

| Имя     | Номер inode |
|---------|-------------|
| .       | 3470036     |
| ..      | 3470017     |
| folder1 | 3470031     |
| file1   | 3470043     |
| file2   | 3470023     |
| folder2 | 3470024     |
| file3   | 3470065     |



директории `usr`, у нее есть `inode`. Далее в этой директории система ищет каталог `ast`, по номеру `inode` находит адрес следующего блока. Получив адрес блока система получает содержимое директории `ast`. Из этого `inode` система нашла файл `mbox`. `struct file` – структура описывающая открытый файл (с которым работает процесс).

Структура, описывающая файловую систему

```

1 struct file_system_type {
2     const char* name; /* Название файловой системы */
3     int fs_flags; /* Флаги типа файловой системы */
4     #define FS_REQUIRES_DEV 1
5     #define FS_USERSNS_MOUNT 8
6     struct dentry* (*mount)(struct file_system_type*, int, const char *, void *);
7     void (*kill_sb)(struct super_block*); // используется для прекращения доступа к
    суперблоку
8     struct module* owner; // Модуль, владеющий файловой системой
9     struct file_system_type* next; // Следующая файловая система
10    struct hlist_head fs_supers; // Список объектов типа суперблок
11    ...
12 };

```

Функции монтирования:

```

1 typedef int (*fill_super_t)(struct super_block *, void *, int); // Передаваемая функции
    mount() функция fill_super() заполняет поля структуры struct super_block.
2 struct dentry *mount_bdev(struct file_system_type *fs_type, int flags, const char
    *dev_name, void *data, fill_super_t fill_super); // для монтирования ФС, находящейся
    на блочном устройстве,
3 struct dentry *mount_single(struct file_system_type *fs_type, int flags, void *data,
    fill_super_t fill_super); // для монтирования ФС, точки монтирования которой раздел
    яют один единственный экземпляр ФС.
4 struct dentry *mount_nodev(struct file_system_type *fs_type, int flags, void *data,
    fill_super_t fill_super); // для монтирования ФС, не связанной ни с каким устройств
    ом,

```

В этих функциях вызывается функция `fill_super`, которая выполняет основные действия по инициализации полей `struct super_block` (размер, смещение, магическое число и `super_operations`). Эти функции возвращают объект `dentry` для корневого каталога. Для ФС необходимо создать корневой каталог, что позволяет выполнить монтирование ФС. Чтобы создать корневой каталог, необходимо создать `inode`.

```

1 struct file_system_type my_fs_type =
2 {
3     .owner = THIS_MODULE,
4     .name = "my_fs",
5     .mount = my_mount,
6     .kill_sb = my_kill_sb,
7     .fs_flags = FS_REQUIRES_DEV
8 };
9 static int __init my_init(void) {
10     ...
11     return register_filesystem(&my_fs_type);
12 }
13 static int __init my_init(void) {
14     ...
15     return unregister_filesystem(&my_fs_type);
16 }

```

При инициализации структуры `file_system_type` можно воспользоваться предопределенными структурами ядра

```

1 void generic_shutdown_super(struct super_block *sb) // удалить суперблок
2 int simple_fill_super(struct super_block *s, unsigned long magic,
3 const struct tree_descr *files); // инициализировать поля суперблока

```

Для ФС необходимо определить, какие функции для работы с файлами необходимо реализовать. Для этого существует `simple_dir_operations`.

```

1 struct inode *new_inode(struct super_block *sb)
2 {
3     struct inode *inode;
4     spin_lock_prefetch(&inode_sb_list_lock);
5     inode = new_inode_pseudo(sb);
6     if (inode)
7         inode_sb_list_add(inode);
8     return inode;
9 }
10 EXPORT_SYMBOL(new_inode);
11
12 struct inode *new_inode_pseudo(struct super_block *sb)
13 {
14     struct inode *inode = alloc_inode(sb);
15
16     if (inode) {
17         spin_lock(&inode->i_lock);
18         inode->i_state = 0;
19         spin_unlock(&inode->i_lock);
20         INIT_LIST_HEAD(&inode->i_sb_list);
21     }
22     return inode;
23 }

```

## 6.1 Загружаемые модули ядра

UNIX имеет монолитное ядро – единая программа из структур и функций. Загружаемые модули ядра позволяют вносить в ядро собственную функциональность.

У загружаемого модуля ядра обязательно две точки входа, устанавливаются макросами `MODULE_INIT` и `MODULE_EXIT`.

Точки входа – места с которых начинает выполняться код.

```
1 static int __init my_init(void);  
2 MODULE_INIT(my_init);  
3  
4 static void __exit my_exit(void);  
5 MODULE_EXIT(my_exit);
```

`init` вызывается `insmod`. `exit` вызывается `rmmod`.

В ядре нельзя обращаться к функциям стандартных библиотек. В ядре могут использоваться только функции, определенные в ядре. Нельзя читать и писать сразу в ядро: нарушается безопасность ядра. Для передачи информации из ядра и в ядро используются специальные библиотеки и функции.

Макрос `MODULE_LICENSE` является обязательным. Макрос `MODULE_LICENSE` используется для того, чтобы сообщить ядру, под какой лицензией распространяется исходный код модуля, что влияет на то, к каким функциям и переменным (символы ядра) он может получить доступ в ядре.

Для протоколирования используется функция `printk()`. В функции указывают уровень протоколирования: `KERN_INFO`, `KERN_ERR`, `KERN_DEBUG` (всего 8 уровней). Чем ниже значение макроса, тем выше приоритет. `printk()` пишет в системный журнал, для обращения к нему используют команду `dmesg`.

```
1 printk(KERN_INFO "module loaded");
```

Взаимодействие модулей заключается в том, что один модуль может использовать данные, объявленные и описанные в другом модуле.

В одном модуле переменная объявлена как `public` (доступная другим модулям), в другом, чтобы прошла компиляцию, объявлена как `extern`. Так один модуль может получать доступ к данным другого модуля. Макрос `EXPORT_SYMBOL` – позволяет экспортировать данные модуля и теперь им могут пользоваться другие модули ядра. Модуль может использовать только те имена, которые явно экспортированы.

Модуль связывается с экспортируемым именем по прямому абсолютному адресу. Абсолютный адрес – физический адрес.

Загружаемый модуль становится частью ОС, он загружается в физическую память. Коды ядра работают с физическими адресами. Модуль может обращаться к данным другого модуля только по физическим адресам – имеет значение последовательность загрузки. У модулей есть счетчик ссылок на него, выгрузить его можно когда счетчик равен нулю.

## 7 proc

proc – виртуальная файловая система. Файлы и каталоги такой VFS не хранятся во вторичной памяти и создаются на лету.

proc создавался разработчиками Linux, чтобы предоставить разработчикам ПО в режиме пользователя информацию из режима ядра о разработанном ПО для анализа с точки зрения потребляемых ресурсов во время выполнения.

В proc у каждого процесса есть директория его идентификатора, в ней находятся файлы, символичные ссылки, поддиректории с информацией о каждом процессе `/proc/<pid>`. Можно использовать `/proc/self` для получения информации текущего процесса

| Элемент    | Тип                  | Описание  |
|------------|----------------------|---|
| cmdline    | файл                 | Содержит командную строку процесса, которая была набрана                              |
| cwd        | символическая ссылка | Указывает на директорию процесса  |
| environ    | файл                 | Содержит список переменных окружения процесса   |
| exe        | символическая ссылка | Указывает на исполняемый файл процесса  |
| fd         | директория           | Содержит ссылки на открытые процессом файлы   |
| root       | символическая ссылка | Указывает на корень файловой системы процесса   |
| stat       | файл                 | Содержит информацию о процессе  |
| maps       | файл                 | Содержит текущие регионы ВАП процесса   |
| pagemap    | файл                 | Содержит по одному 64-битному значению на каждую виртуальную страницу памяти процесса |
| mem        | файл                 | информация о страницах памяти процесса  |
| interrupts | файл                 | информация о об источниках прерываний и о частоте возникновения этих прерываний.      |
| task/<tid> | директория           | информация о потоке процесса  |

```
1 * MAPS:
2     address      perms  offset  dev  inode
3 1    5f17aaa8a000-5f17aaa8b000 r--p 00000000 08:03 3932858
4    /home/Desktop/task-03/s.out
5 1    5f17aaa8b000-5f17aaa8c000 r-xp 00001000 08:03 3932858
6    /home/Desktop/task-03/s.out
7 1    5f17aaa8c000-5f17aaa8d000 r--p 00002000 08:03 3932858
8    /home/Desktop/task-03/s.out
9 1    5f17aaa8d000-5f17aaa8e000 r--p 00002000 08:03 3932858
10   /home/Desktop/task-03/s.out
11 1    5f17aaa8e000-5f17aaa8f000 rw-p 00003000 08:03 3932858
12   /home/Desktop/task-03/s.out
13 33   5f17aae6d000-5f17aae8e000 rw-p 00000000 00:00 0          [heap]
14 33   7e67a0000000-7e67a0021000 rw-p 0          00:00 0
15 33   7e67ac000000-7e67ac021000 rw-p 0          00:00 0
16 33   7fffee537000-7fffee558000 rw-p 00000000 00:00 0          [stack]
17 4    7fffee5b0000-7fffee5b4000 r--p 00000000 00:00 0          [vvar]
18 2    7fffee5b4000-7fffee5b6000 r-xp 00000000 00:00 0          [vdso]
19 1    ffffffff600000-ffffffff601000 --xp 00000000 00:00 0 [vsyscall]
```

```
1 * PAGEMAP
2 addr : pfn soft-dirty file/shared swapped present
3 0x5f17aaa8a000 : 59927          1          1          0          1
4 0x5f17aaa8b000 : 2d274          1          1          0          1
5 0x5f17aaa8c000 : 51c00          1          1          0          1
6 0x5f17aaa8d000 : 0          1          0          0          1
7 0x5f17aaa8e000 : 0          1          0          0          1
```

pfn (Page Frame Number) - Номер физической страницы. soft-dirty – страница была модифицирована. file/shared – является файловой или разделяемой, 0 если анонимная страница. swapped – выгружена в swp. present – страница доступна в физической памяти.

```

1 stat
2 pid 17387
3 comm myapp
4 state S,
5 ppid 7980,
6 flags 4194304
7 utime 78074,
8 stime 8573
9 priority 20,
10 nice 0,
11 num_threads 3
12 vsize 2118049792
13 rss 100
14 startcode 104829120306400
15 endcode 104829128962048
16 startstack 140735838567920
17 exit_signal 17
18 processor 4
19 rt_priority 0
20 policy 0
21 start_data 104829129340584
22 end_data 104829129350435

```

```

1 struct proc_dir_entry {
2     atomic_t in_use;
3     const struct inode_operations *proc_iops;
4     union {
5         const struct proc_ops *proc_ops;
6         const struct file_operations *proc_dir_ops;
7     };
8     const struct dentry_operations *proc_dops;
9     union {
10         const struct seq_operations *seq_ops;
11         int (*single_show)(struct seq_file *, void *);
12     };
13     proc_write_t write;
14     void *data;
15     unsigned int low_ino;
16     nlink_t nlink;
17     kuid_t uid;
18     kgid_t gid;
19     loff_t size;
20     struct proc_dir_entry *parent;
21     char *name;
22     u8 flags;
23 };

```

Раньше в `proc_dir_entry` был только `file_operations`. `struct proc_ops` — определяет операции над файлами `proc`, нововведение. Функции `read/write` в `struct file_operations` используются интенсивно, а использование их для ВФС `proc`, которая носит справочный характер, вносит дополнительную нагрузку. Поэтому разработчики Linux ввели `proc_ops`, чтобы разгрузить функции.

```

1 struct proc_ops {
2     unsigned int proc_flags;
3     int (*proc_open)(struct inode *, struct file *);
4     ssize_t (*proc_read)(struct file *, char __user *, size_t, loff_t *);
5     ssize_t (*proc_write)(struct file *, const char __user *, size_t, loff_t *);
6     int (*proc_release)(struct inode *, struct file *);
7 };

```

Для создания файла `proc` используется системный вызов `proc_create`. `proc_create` — обертка `proc_create_data`.

```

1 proc_dir_entry* proc_create_data(const char* name, umode_t mode, struct proc_dir_entry*
   parent, const struct file_operations* proc_fops, void* data);
2
3 static inline struct proc_dir_entry* proc_create(const char* name, umode_t mode,
   struct proc_dir_entry* parent, const struct file_operations* proc_fops) {
4     return proc_create_data(name, mode, parent, proc_fops, NULL);
5 }

```

Создавать каталоги в файловой системе /proc можно используя proc\_mkdir(), а также символические ссылки с proc\_symlink().

```

1 extern struct proc_dir_entry *proc_symlink(const char *, struct proc_dir_entry *,
   const char *);
2 extern struct proc_dir_entry *proc_mkdir(const char *, struct proc_dir_entry *);

```

Для передачи данных из пространства ядра имеется 2 функции

```

1 unsigned long __copy_to_user(void __user* to, const void *
   from, unsigned long n);
2
3
4 unsigned long __copy_from_user(void *to, const void __user
   *from, unsigned long n);

```

Ядро находится в физической памяти, а процесс памяти не имеет, имеет виртуальное адресное пространство. Физические страницы выделяются, когда процесс обращается к данным соответствующей виртуальной страницы. Если при обращении функции ядра к буферу пользователя нужная страница выгружена, то система должна загрузить соответствующие страницы страничным прерыванием. Это прерывание обрабатывается ядром. Происходит переключение между кодами ядра для обработки страничного прерывания и возникнет паника ядра.

Для предотвращения этой ситуации введены эти специальные функции.

```

1 #if LINUX_VERSION_CODE >= KERNEL_VERSION(5,6,0)
2 #define HAVE_PROC_OPS
3 #endif
4
5 #ifdef HAVE_PROC_OPS
6 static struct proc_ops fileops = {
7     .proc_read = fortune_read,
8     .proc_write = fortune_write,
9     .proc_open = fortune_open,
10    .proc_release = fortune_release,
11 };
12 #else
13 static struct file_operations fileops = {
14     .owner = THIS_MODULE,
15     .read = fortune_read,
16     .write = fortune_write,
17     .open = fortune_open,
18     .release = fortune_release,
19 };
20 #endif
21

```

Sequence файлы: второй интерфейс для работы с proc файлами. Для сиквенсов определена структура struct seq\_file и структура struct seq\_operations. В union proc\_dir\_entry два поля: seq\_ops и single\_open. Для seq файлов определены два способа работы: полный и упрощенный.

```

1 struct seq_file {
2     char *buf;
3     size_t size;
4     size_t from;
5     size_t count;
6     size_t pad_until;
7     loff_t index;

```



```

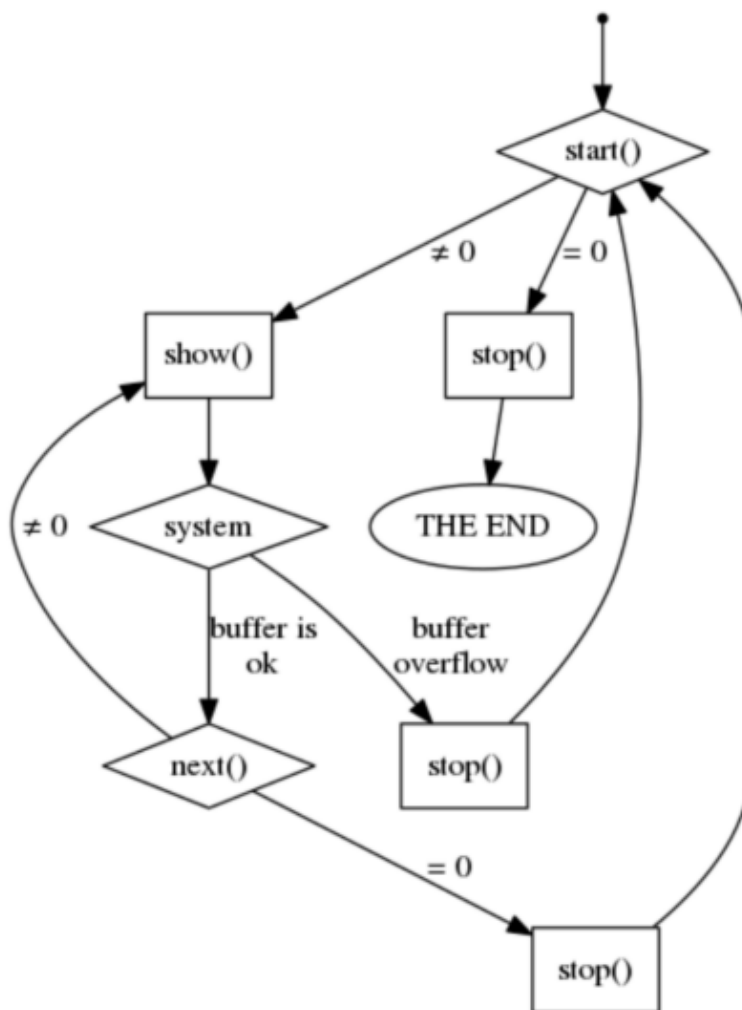
8      loff_t read_pos;
9      struct mutex lock;
10     const struct seq_operations *op;
11     int poll_event;
12     const struct file *file;
13     void *private;
14 };

```

```

1 struct seq_operations {
2     void * (*start) (struct seq_file *m, loff_t *pos);
3     void (*stop) (struct seq_file *m, void *v);
4     void * (*next) (struct seq_file *m, void *v, loff_t *pos);
5     int (*show) (struct seq_file *m, void *v);
6 };

```



Функция stop вызывается в трех случаях:

- Вызов stop сразу после start означает окончание передачи данных
- stop после show означает, что после последней передачи данных буфер, в который передаются данные, заполнен и передача данных завершается.
- вызов после next выполняется, если нужно закончить передачу в буфер данных и завершить работу.

В функции show вызывается seq\_printf. Sequence файлы обеспечивают более надежную передачу информации из пространства ядра в пространство пользователя, за счет того что ядро контролирует указатель на буфер с данными.

```

1 int seq_open(struct file * file, const struct seq_operations *op);
2 ssize_t seq_read(struct file* file, char __user* buf, size_t size, loff_t* ppos);
3 int seq_release ( struct inode * inode, struct file * file);

```

```

4 void seq_printf(struct seq_file *m, const char *fmt, ...);
5 int seq_write(struct seq_file *seq, const void *data, size_t len);
6 loff_t seq_lseek(struct file *, loff_t, int);
7 int seq_release(struct inode *, struct file *);

```

Упрощенный интерфейс single файлов позволяет не писать итератор, так как внутри функции single\_open инициализируются поля структуры seq\_operations.

```

1 int single_open(struct file *file, int (*show)(struct seq_file*, void*), void *data){
2     struct seq_operations *op = kmalloc(sizeof(*op), GFP_KERNEL_ACCOUNT);
3     int res = -ENOMEM;
4     if(op) {
5         op->start = single_start;
6         op->next = single_next;
7         op->stop = single_stop;
8         op->show = show;
9         res = seq_open(file, op);
10        if(!res) {
11            ((struct seq_file*)file->private_data)->private =
12                data;
13        }
14        else {
15            kfree(op);
16        }
17    }
18    return res;
19 }
20 EXPORT_SYMBOL(single_open);

```

Для single\_open размер передаваемых данных ограничен 64 КБ (16 страниц).

## 8 Открытые файлы

Любой процесс может открыть много файлов, причем один файл может быть открыт много раз одним или несколькими процессами.

В дескрипторе процесс есть указатель на открытые процессом файлы и указатель на ФС, которая относится к образу процесса.

```
1 struct task_struct
2 {
3     ...
4     struct fs_struct *fs;
5     struct files_struct *files;
6     ...
7 }
```

struct file – дескриптор открытого файла.

Системный вызов *int open(const char \*pathname, int flags, mode\_t mode);* или *int open(const char \*pathname, int flags);* Флаги:

- O\_CREATE – создать новый файл
- O\_EXCL – (обязательно с O\_CREATE) – заставить систему проверить существование файла, чтобы не потерять информацию ранее записанную в файл
- O\_EXEC
- O\_RDONLY
- O\_RDWR
- O\_APPEND

Информация о файловой системе, к которой относится процесс

```
1 struct fs_struct
2 {
3     atomic_t count;
4     ...
5     struct dentry *root; /* Корневой каталог */
6     struct dentry pwd; /* Текущий рабочий каталог */
7     ...
8     struct vfsmount *rootmnt;
9     ...
10 }
```

В этой структуре хранится вся информация об открытых процессом файлах и файловых дескрипторах.

```
1 struct files_struct {
2     atomic_t count; // счетчик ссылок
3     ...
4     unsigned int next_fd; // число открытых процессом файлов
5     unsigned long close_on_exec_init[1];
6     unsigned long open_fds_init[1];
7     unsigned long full_fds_bits_init[1];
8     struct file __rcu * fd_array[NR_OPEN_DEFAULT]; // массив файловых дескрипторов
9     // файлов, открытых процессом
10 };
```

Каждый процесс имеет собственную таблицу страниц с указателем на системные дескрипторы открытых файлов. Когда создается процесс для него автоматически открываются три файла: stdin, stdout, stderr – 0, 1, 2. Если процесс в результате exec не унаследовал открытые файлы, то первый номер открытого файла будет 3. close\_on\_exec\_init – какие ФД будут закрыты при exec.

```
1 struct file {
2     f_mode_t mode;
3     struct path f_path;
4     struct inode *f_inode;
5     const struct file_operations *f_op;
6     atomic_long_t f_count;
7     ...
8     loff_t f_pos; // отражает как работать с логической структурой файла
```

```

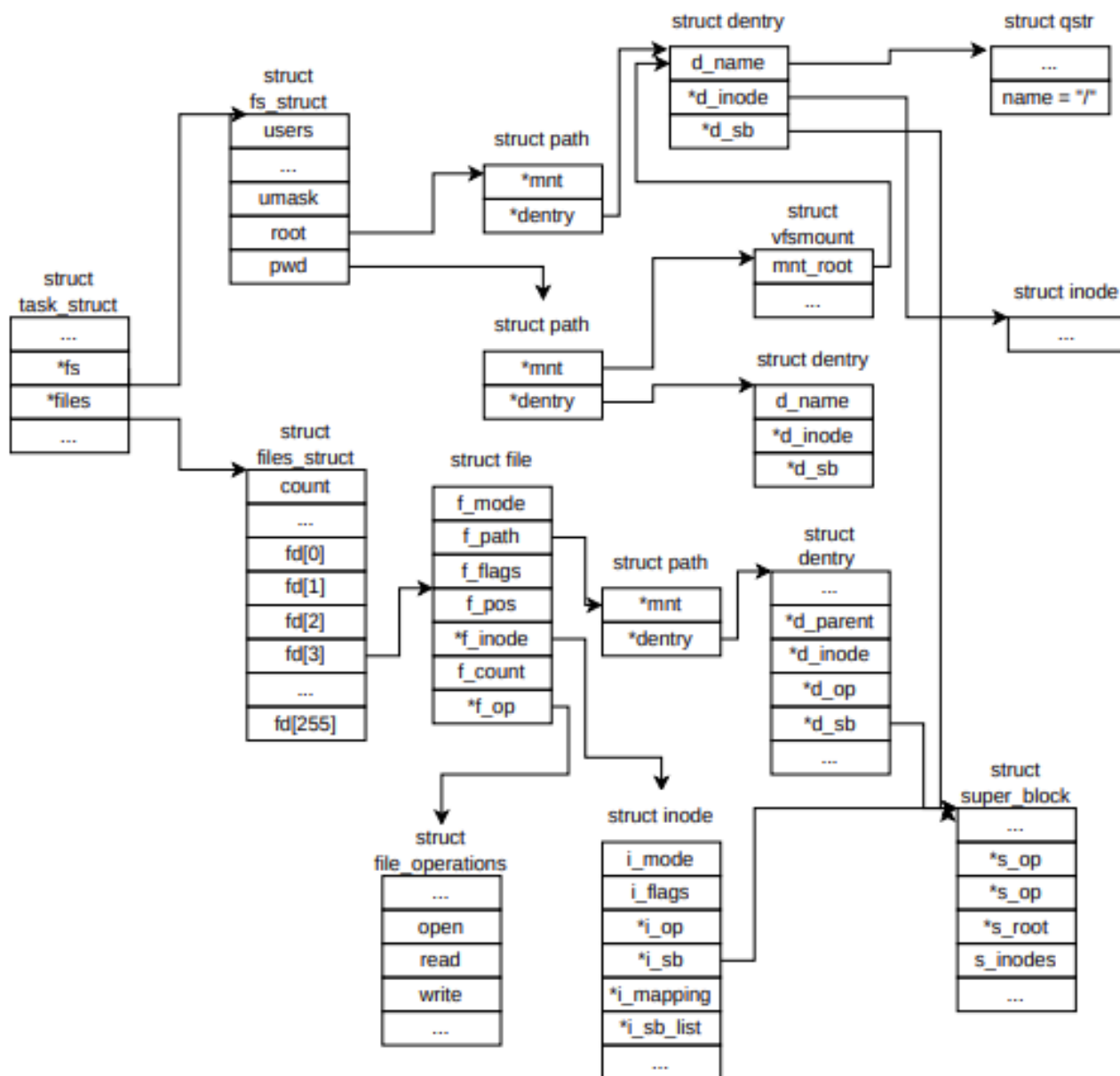
9      ...
10   }
11   struct path {
12       struct vfsmount *mnt;
13       struct dentry *dentry;
14   };

```

struct file работает с логической структурой файла. Любой файл начинается с 0 адреса и любой файл считает, что его смещение начинается с 0.

Каждая операция чтения или записи сдвигает указатель файла на соответствующее количество байт.

Связь структур ядра: stdio – библиотека буферизированного ввода-вывода.



foren() – открывает файл и возвращает связанный с ним указатель. Функция возвращает указатель на структуру FILE, которая называется файловым потоком.

fdopen() принимает дескриптор на открытый файл и возвращает файловый поток.

```

1 struct _IO_FILE {
2     int _flags;

```

/\* High-order word is \_IO\_MAGIC; rest is flags. \*/

```

3  /* The following pointers correspond to the C++ streambuf protocol. */
4  char *_IO_read_ptr;      /* Current read pointer */
5  char *_IO_read_end;      /* End of get area. */
6  char *_IO_read_base;     /* Start of putback+get area. */
7  char *_IO_write_base;    /* Start of put area. */
8  char *_IO_write_ptr;     /* Current put pointer. */
9  char *_IO_write_end;     /* End of put area. */
10 char *_IO_buf_base;      /* Start of reserve area. */
11 char *_IO_buf_end;       /* End of reserve area. */
12
13 /* The following fields are used to support backing up and undo. */
14 char *_IO_save_base; /* Pointer to start of non-current get area. */
15 char *_IO_backup_base; /* Pointer to first valid character of backup area */
16 char *_IO_save_end; /* Pointer to end of non-current get area. */
17
18 struct _IO_marker *_markers;
19
20 struct _IO_FILE *_chain;
21
22 int _fileno;
23 int _flags2;
24 __off_t _old_offset; /* This used to be _offset but it's too small. */
25
26 /* 1+column number of pbase(); 0 is unknown. */
27 unsigned short _cur_column;
28 signed char _vtable_offset;
29 char _shortbuf[1];
30
31 _IO_lock_t *_lock;
32 #ifdef _IO_USE_OLD_IO_FILE
33 };

```

Сначала информация пишется в буфер и оттуда выбирается для обработки. 3 причины перезаписи в файл: буфер заполнен, fflush и fclose().

Пример: файл открывается два раза системным вызовом open() для записи и в него последовательно записывается строка «aaaaaaaaaaaa» по первому дескриптору и затем строка «bbbb» по второму дескриптору, затем файл закрывается два раза. Показать, что будет записано в файл и пояснить результат.

```

1  int main(void)
2  {
3      int fd1 = open("file.txt", O_CREAT | O_EXCL | O_WRONLY);
4      int fd2 = open("file.txt", O_WRONLY);
5
6      char buf1[] = "aaaaaaaaaaaa";
7      char buf2[] = "bbbb";
8      write(fd1, buf1, sizeof(buf1) - 1);
9      write(fd2, buf2, sizeof(buf2) - 1);
10
11     close(fd1);
12     close(fd2);
13
14     return 0;
15 }

```

В ядре есть две таблицы: PFT — process file table и SFT — system file table, в которой находится проинициализированная структура struct file. Процесс может создать только 512 открытых файлов. Базовое значение в ядре: 256. Если надо еще — создается дополнительная таблица еще на 256.

Файл дважды открывается на запись с помощью open(), при этом создаются два экземпляра структуры struct file, которые описывают один физический файл. У каждой из структур инициализируется поле f\_pos = 0. Начиная с нулевой позиции первого открытого файла в файл будет записано аaaaaaaaaа. После чего начиная с нулевой позиции будет записано вvvv. Результатом будет vvvvaaaaaaaaaа.

## 9 Слабы

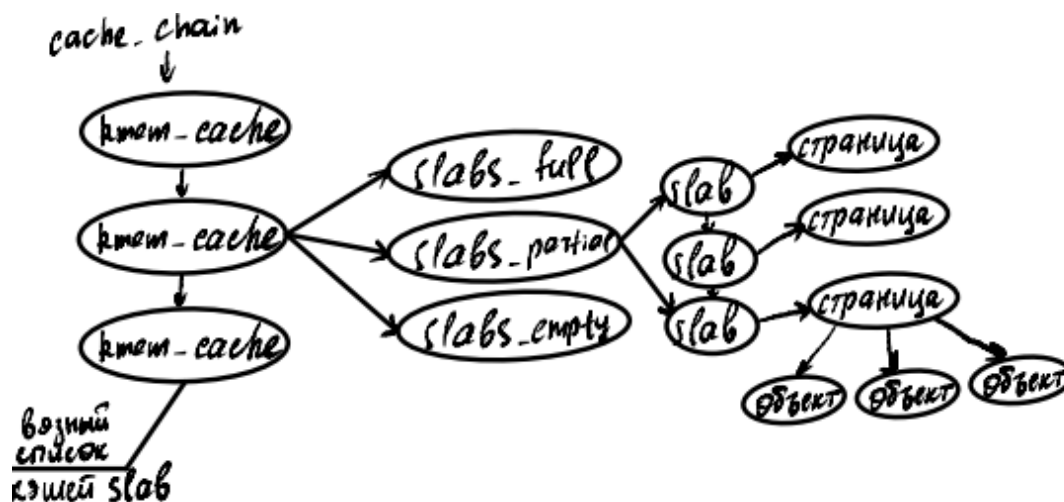
Когда создается объект, инициализируются поля структуры. Если объект удаляется и он понадобится снова его поля нужно будет заново инициализировать. В ядре используются не так много структур, которые описывают базовые объекты. Одни и те же структуры используются многократно для создания одних и тех же объектов. Инициализация полей структур занимает время. Когда нужна в объекте исчезает, объект не удаляется, а остается в слабе для повторного использования.

Побочный эффект слабов – устранение скрытой фрагментации (страница никогда не заполняется полностью). Поскольку память выделяется страницами, скрытая фрагментация присутствует. Но кэш слабов позволяет записывать в страницу несколько слабов одного и того же размера.

```
1 // Создание нового кеша
2 struct kmem_cache* kmem_cache_create(const char* name, unsigned int size, unsigned int
    align, slab_flags_t flags, void (*ctor)(void*));
3 // Для удаления кеша
4 int kmem_cache_destroy(kmem_cache_t *cache);
5 // получение адресного пространства для объекта
6 void *kmem_cache_alloc(kmem_cache_t *cache, int flags);
7 void kmem_cache_free(kmem_cache_t *cache, const void *obj);
```

Слабы работают только со структурами ядра. Каждый slab может находиться в одном из трех состояний:

- slab full — полностью распределенный. В таком слабе нет свободных объектов, поскольку все они распределены и используются.
- slab partial — частично распределенный. В таком слабе есть как выделенные, так и свободные объекты.
- slab empty — пустой, незаполненный. В таком слабе нет ни одного распределенного объекта, все его объекты свободны для последующего использования.



Как только из некоторой части ядра поступает запрос на новый объект, он выделяется из частично распределенного блока, если таковой имеется. В противном случае объект выделяется из пустого блока. Если больше не осталось пустых слабов, они будут созданы по мере необходимости.

Inode кеш имеет следующую структуру:

- Глобальная хеш таблица `inode_hashtable`, в которой каждый `inode` адресуется по указателю на суперблок и 32-битному номеру `inode`.
- Глобальный список `inode_in_use` — `inode` в состоянии `in_use`. В этом списке содержится `inode` с полями `i_count > 0` и `i_nlink > 0`. Только что созданный `inode` добавляется как раз в этот список.
- Глобальный список `inode_unused`, который содержит допустимые `inode` с `i_count = 0`.
- Когда `inode` помечается как "грязный он добавляется к списку `sb->s_dirty` при условии, что он в хэш-таблице.

Через поле `inode->i_list` с `inode` вставляется в список определенного типа, через поле `inode->i_hash`. Каждый `inode` может входить в список `i_hash` и в один и только в один список типа (`in_use`, `unused` или `dirty`).

## 10 Аппаратные прерывания

Аппаратная поддержка аппаратных прерываний развивается. В персональных компьютерах реализована шинная архитектура. Управление внешними устройствами было передано контроллерам (входят в состав внешнего устройства) и адаптерам (на материнской плате).

В 16-разрядных машинах был реализован PIC. В 32-разрядных машинах APIC (advanced PIC), может работать и с 16-разрядными машинами.

В современных системах применяется реализация прерываний от внешних устройств в виде сообщений. MSI — message signaled interrupts — прерывание, сигнализируемое сообщением. Устройства сигнализируют о прерывании, отправляя/записывая сообщение по определенному адресу в памяти.

Рассмотрим не аппаратную составляющую, а ПО позволяющее обрабатывать эти аппаратные прерывания.

Управление внешними устройствами осуществляется драйверами, для каждого устройства пишется драйвер. В составе драйвера входит обработчик прерывания.

В ядре предоставляется возможность зарегистрировать собственный обработчик прерываний.

```
1 typedef irq_return_t(*irq_handler_t)(int, void *);
2 extern int __must_check request_threaded_irq(unsigned int irq, irq_handler_t handler,
3       irq_handler_t thread_fn, unsigned long flags, const char *name, void *dev);
4 static inline int __must_check request_irq(unsigned int irq, irq_handler_t handler,
5       unsigned long flags, const char *name, void *dev) {
6     return request_threaded_irq(irq, handler, NULL, flags, name, dev);
7 }
8 void free_irq(unsigned int irq, void *dev_id); // освободить линию
```

При прерывании по линии IRQ вызывается обработчик прерывания. Когда возникает прерывание, система через соответствующее аппаратное обеспечение переходит на соответствующий обработчик прерывания.

Линии прерываний уже заняты основными обработчиками прерываний. Используются флаги

- IRQF\_SHARED – разделять линии прерывания несколькими обработчиками прерываний;
- IRQF\_PROBE\_SHARED
- IRQF\_NO\_THREAD – не может выполняться на разных ядрах (не может быть распараллелен)
- \_\_IRQF\_TIMER – от системного таймера
- IRQF\_TIMER (\_\_IRQF\_TIMER | IRQF\_NO\_SUSPEND | IRQF\_NO\_THREAD)

```
1 struct irqaction {
2     irq_handler_t      handler;
3     void               *dev_id;
4     void __percpu      *percpu_dev_id;
5     struct irqaction   *next;
6     irq_handler_t      thread_fn;
7     struct task_struct *thread; // система выполняет каждый драйвер как поток
8     struct irqaction   *secondary;
9     unsigned int       irq;
10    unsigned int       flags;
11    unsigned long       thread_flags;
12    unsigned long       thread_mask;
13    const char          *name;
14    struct proc_dir_entry *dir; // /proc/irq/<N>/nameentry
15 } ____cacheline_internodealigned_in_smp;
```

Быстрые и медленные прерывания: в Linux одно быстрое прерывание – от системного таймера.

Проблема: в любом обработчике прерывания выполняется на высочайшем уровне приоритета. В результате никакая другая работа в системе не может выполняться, пока не завершится обработчик. Поэтому все действия в обработчиках – отложенные – через постановку в очередь на выполнение.

Обработчик прерывания выполняется на одном из ядер, при этом на этом ядре запрещаются другие прерывания. А на других ядрах запрещены прерывания по этой линии IRQ.

Медленные делятся на 2 части – top half и bottom half. Верхняя половина нужна, чтобы завершать быстрее в обработчике прерывания минимальный необходимый набор действий. Задача – сохранить данные в буфере ядра

Нижняя половина делится на 3 части: softirq – гибкие прерывания, тасклеты и очереди работ. softirq – отложенные действия, которые определяются статически во время компиляции ядра.

```
1 struct softirq_action {
2     void(*action)(struct softirq_action*); /* ф-я которая должна выполняться */
3 };
```

Кроме этой структуры в ядре определен массив из 10 экземпляров данной структуры:

```
1 static struct softirq_action softirq_vec[NR_SOFTIRQS]
```

| Индекс               | Приоритет | Значение  |
|----------------------|-----------|---|
| HI_SOFTIRQ           | 0         | высоко приоритетные softirq   |
| TIMER_SOFTIRQ        | 1         | программируемые таймеры, которые нужны в ПО для контроллеров              |
| NET_TX_SOFTIRQ       | 2         | отправка сетевых пакетов  |
| NET_RX_SOFTIRQ       | 3         | прием сетевых пакетов   |
| BLOCK_SOFTIRQ        | 4         | блочные устройства  |
| BLOCK_IOPOLL_SOFTIRQ | 5         |   |
| TASKLET_SOFTIRQ      | 6         | тасклеты  |
| SCHED_SOFTIRQ        | 7         | планировщик   |
| HR_TIMER_SOFTIRQ     | 8         | не используются   |
| RCU_SOFTIRQ          | 9         | сохраняется для сохранения нумерации // обязательно последним должен быть |

```
1 const char *const softirq_to_name[NR_SOFTIRQS] = {
2     "HI", "TIMER", "NET_TX", "NET_RX", "BLOCK", "IRQ_POLL",
3     "TASKLET", "SCHED", "HRTIMER", "RCU"
4 };
```

Меньшие номера имеют более высокий приоритет.

```
1 void open_softirq(int nr, void (*action)(struct softirq_action*))
2 {
3     softirq_vec[nr].action = action;
4 }
5 irq_return_t my_interrupt(int irq, void *dev_id) {
6     raise_softirq(MY_SOFT_IRQ);
7     ...
8     return IRQ_HANDLER;
9 }
10 void my_soft_handler(struct softirq_action *) {...}
11 static int __init my_init(void)
12 {
13     /*....*/
14     request_irq(irq, my_interrupt, 0, "my", NULL); //регистрируем свой обработчик
15     open_softirq(MY_SOFT_IRQ, my_soft_handler);
16 }
```

На каждое ядро при загрузке системы создается и инициализируется ksoftirq daemon Задача ksoftirq daemon – обработка отложенных действий, работает с softirq.

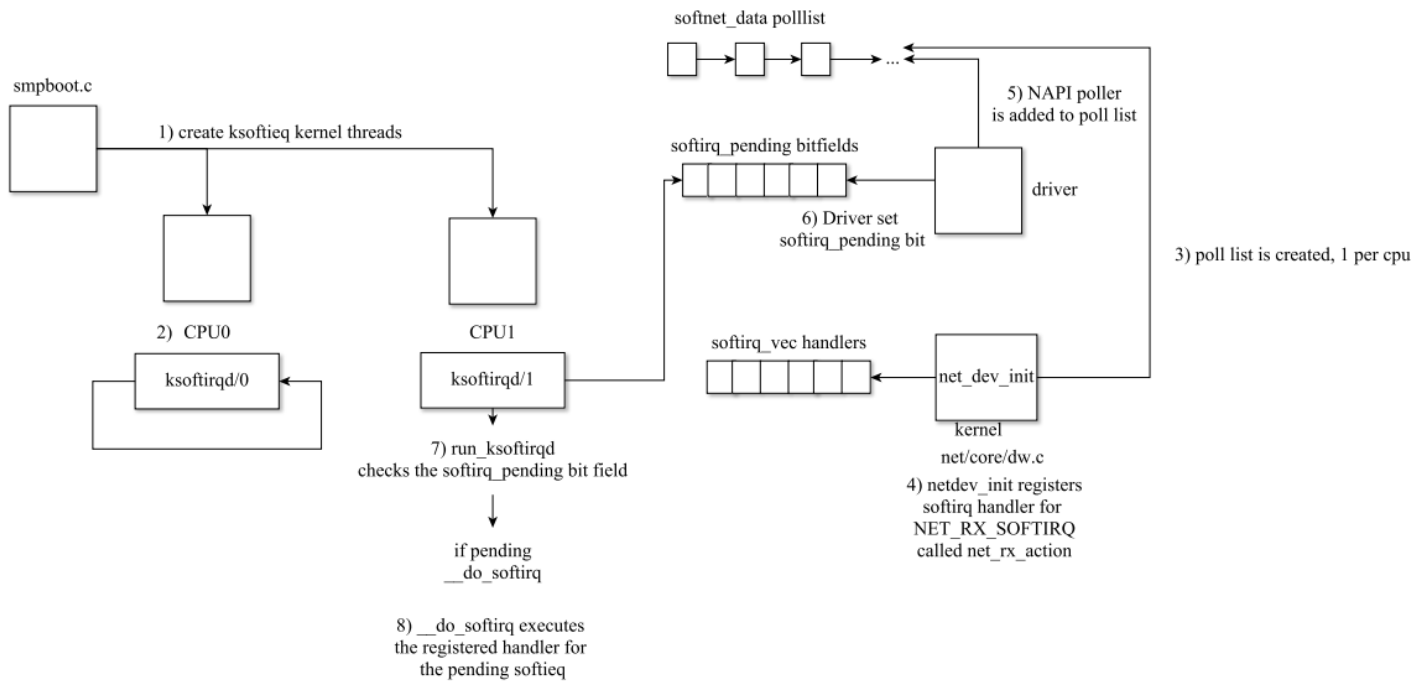
Первое ядро находится в цикле выполнения. run\_ksoftirq проверяет очередь событий для их обработки.

Если каждый входящий пакет приходит на сетевую карту, то каждый пакет нужно обработать – прерывание, что не эффективно.

NAPI – замена прерываний опросом. Основная цель NAPI — сократить количество прерываний, генерируемых при получении пакетов. Создаются poll list, по одному на ядро.

Демон ksoftirq работает с битовой очередью, обозначающая события, требующие обработки. У сетевой карты есть драйвер. Когда пакет приходит, драйвер устанавливает соответствующий бит очереди, демон его анализирует и понимает, что надо обработать пакет.





Если бит выставлен, то пакет пришел и его надо обработать и вызывается функция `do_softirq`. Функция `do_softirq` выполняет зарегистрированный обработчик.

Сетевая подсистема Linux работает по примеру стека BSD. Прием и передача данных с помощью интерфейса сокетов. NAPI – приходит пакет и вызывается `rx_scheduler`, что гарантирует обработку пакетов в дальнейшем. Тасклеты: частный случай `softirq`. Но `softirq` можно сериализовывать (один и тот же тип `softirq` может выполняться параллельно), так что нужны средства взаимоисключающий.

Тасклеты же не сериализуются, один тип тасклета может выполняться только на одном процессоре, так что средства взаимоисключения не требуются. Разные тасклеты могут выполняться параллельно.

В силу этого тасклеты не используются для обработки нагружающих действий. Тасклет – компромисс между производительностью и простотой использования. Тасклеты могут быть зарегистрированы и статически и динамически.

```
1 struct tasklet_struct {
2     struct tasklet_struct *next; /* Указатель на следующий тасклет в списке */
3     unsigned long state;
4     atomic_t count;
5     bool use_callback;
6     void (*func)(unsigned long data); /* Функция-обработчик тасклета */
7     unsigned long data; /* Аргумент функции-обработчика тасклета */
8 };
```

Состояниями тасклета могут быть 0, `TASKLET_STATE_SCHED`, `TASKLET_STATE_RUN`. `TASKLET_STATE_RUN` имеет смысл только в `smp`-архитектуре, так как если тасклет выполняется, то такой тасклет на другом процессоре выполняться не может.

Тасклет можно выполняться только если счетчик ссылок `count == 0`.

Для статического создания тасклета используется макрос. Он создает экземпляр структуры тасклета с именем и функцией.

```
1 DECLARE_TASKLET(name, func);
2 DECLARE_TASKLET_DISABLED(name, func); // count = 1 (тасклет не может выполняться)
```

Динамически тасклет создается с помощью функции ядра

```
1 void tasklet_init(struct tasklet_struct *t, void (*func)(unsigned long), unsigned long data)
```

```

2 {
3     t->next = NULL;
4     t->state = 0;
5     atomic_set(&t->count, 0);
6     t->func = func;
7     t->data = data;
8 }

```

Для того чтобы запланировать тасклет используются функции. В зависимости от функции тасклет помещается в разные списки: или `tasklet_vec` или `tasklet_hi_vec`.

```

1 extern void __tasklet_schedule(struct tasklet_struct *);
2 extern void __tasklet_hi_schedule(struct tasklet_struct *);

```

Особенность тасклетов – не могут блокироваться, можно использовать только спинлоки (`tasklet_trylock()`, `tasklet_unlock()`)

Тасклеты могут использоваться только 1 раз. Они строго сериализуются по отношению к самому себе, но не по отношению к другим. Очереди работ: альтернатива тасклетам. Задачи идентичные, обе предоставляют возможность реализовать отложенные действия. Тасклет выполняется в контексте прерывания, код тасклета неделим. Работа выполняется в контексте специального потока ядра и может быть заблокирована.

Тасклет всегда выполняется на процессоре, на котором завершилась обработка прерывания, поставившая тасклет в очередь на выполнение. Работа по умолчанию тоже на том же ядре, но можно и поменять процессор. Также работу можно отложить на некоторое время.

```

1 struct workqueue_struct *alloc_workqueue(const char fmt, unsigned int flags, int
    max_active);

```

```

1 struct workqueue_struct {
2     struct list_head pwqs;
3     struct list_head list;
4     struct worker *rescuer;
5     char name[WQ_NAME_LEN];
6     struct pool_workqueue __percpu *cpu_pwqs; // pool workqueues на процессор.
7     struct pool_workqueue __rcu *numa_pwq_tbl[]; // для новых систем - numa
8 };

```

Структура представляет некоторую работу, которую решает некий обработчик нижней половины, экземпляр структуры помещается в очередь.

```

1 struct work_struct {
2     atomic_long_t data;
3     struct list_head entry;
4     work_func_t func;
5 };

```

Инициализация полей структуры статически:

```

1 DECLARE_WORK(name, void (*func)(void*));

```

Динамически:

```

1 INIT_WORK(struct work_struct *work, void (*func)(void *), void *data);
2 PREPARE_WORK(struct work_struct *work, void (*func)(void *), void *data); // повторно

```

Для того чтобы отложить работу на выполнение используются функции:

```

1 int queue_work(struct workqueue_struct *wq, struct work_struct *work);
2 // на конкретный cpu
3 int queue_work_on(int cpu, struct workqueue_struct *wq, struct work_struct *work);
4 // с задержкой
5 int queue_delayed_work(struct workqueue_struct *wq, struct work_struct *work, unsigned
    long delay);

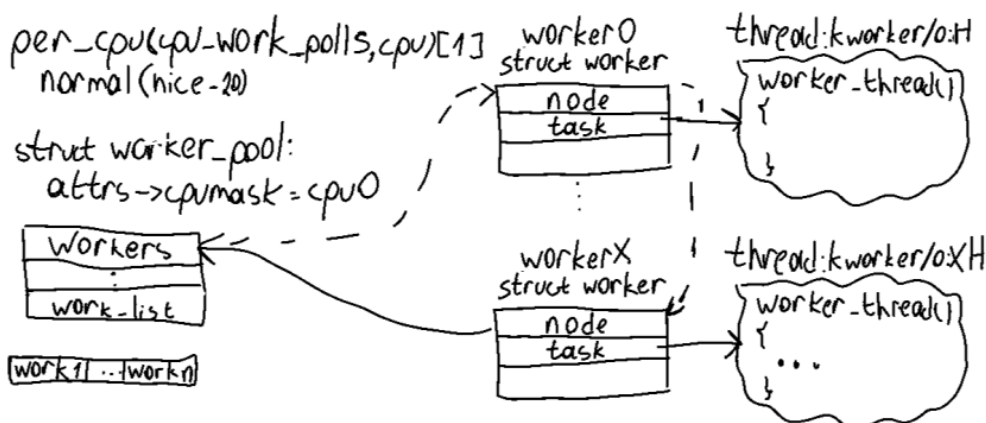
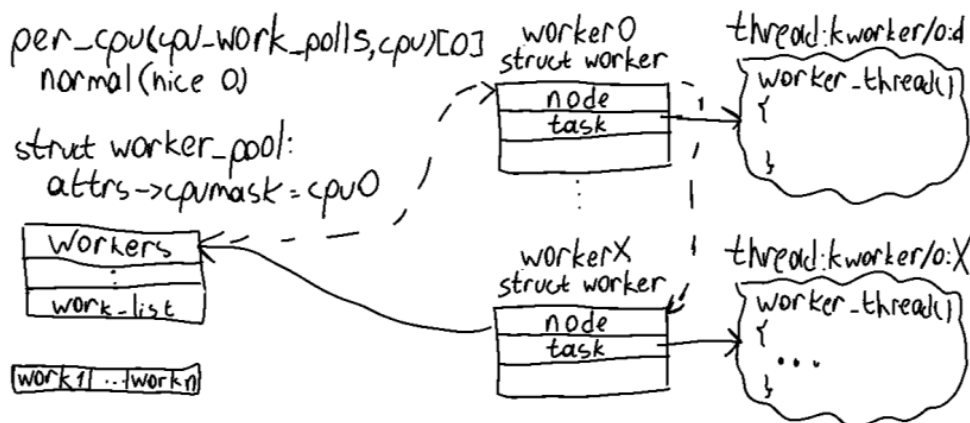
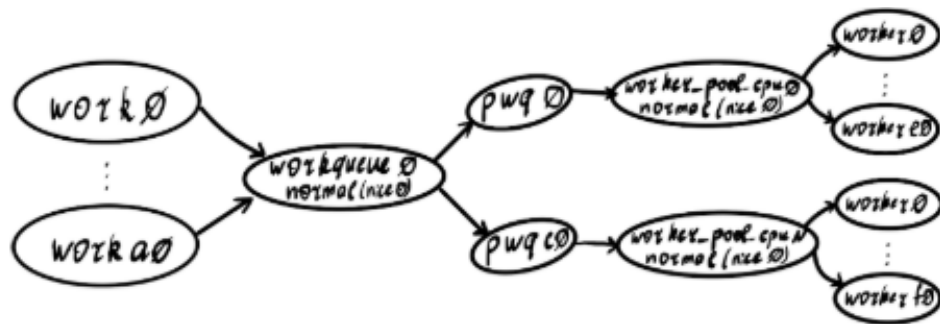
```

```
1 void destroy_workqueue(struct workqueue_struct *wq);
```

Основная идея с которой очереди работ были реализованы: одна параллельная множественная работа, в которой работы не должны блокировать друг друга. Механизм строится с расчетом на экономию потоков, выделенной памяти и планирования.

CMWQ – concurrency managed workqueue.

- Работа – work
- Очереди работ – workqueue. Очередь и работа представлены отношением 1:многие.
- Worker – рабочий. worker соответствует потоку ядра work\_thread.
- PWQ (pool workqueue) – определяется как посредник, отвечающий за установку отношения между work\_pool и work\_queue; work\_queue и pwq представлены отношением 1:многие. work\_pool и pwq представлены отношением 1:многие.



Сначала выполняются потоки worker повышенного приоритета. Тасклет планируется на том же процессе, что и прерывание и как следствие обычно выполняется сразу, поэтому он быстрее очереди работ.

Флаги:

- WQ\_UNBOUND. Очереди работ были разработаны в расчете на выполнение на одном и том же процессоре, что улучшает работу с кэшем.
- WQ\_FREEABLE. Может быть заморожен.

## 11 Специальные файлы устройств

Символьные и блочные устройства. Байт-ориентированные (или символьные) устройства передают данные посимвольно, как непрерывный поток байтов. Блок-ориентированные (или блочные) устройства. Все запоминающие устройства относятся к блочным. Сделаны для обеспечения работы с внешними устройствами как с обычными файлами. При r/w информации с внешнего устройства используются системные вызовы read, write, универсальный механизм.

Файлы устройств могут быть открыты, закрыты, можно читать и писать.

ОС каждому устройству в соответствие ставит 1 специальный файл, обычно в /dev.

Принята система major/minor номеров. Старший номер – идентификатор драйвера устройства. Младший номер – для различия устройства.

В ядре определен тип dev\_t – 32 разрядное число. POSIX определяет существование типа, но не формат полей. Код не должен знать какие-либо предположения об организации номеров внутри числа. Нужно использовать макросы.

```
1 MAJOR(dev_t dev) -- получить старший номер
2 MINOR(dev_t dev) -- получить младший номер
3
4 MKDEV(int major, int minor) -- преобразовать в dev_t
```

Linux поддерживает огромное количество внешних устройств. Старший номер выбирается динамически, функция возвращает его одновременно с первым младшим номером:

```
1 int register_chrdev_region(dev_t from, unsigned count, const char *name);
2 int alloc_chrdev_region(dev_t *dev, unsigned baseminor, unsigned count, const char
   *name) {
3     struct char_device_struct *cd;
4     cd = register_chrdev_region(0, baseminor, count, name);
5     if (IS_ERR(cd)) return PTR_ERR(cd);
6     *dev = MKDEV(cd->major, cd->baseminor);
7     return 0;
8 }
```

Реализация драйвера символьного устройства

```
1 #define IRQ_NO 11
2 static DECLARE_WORK(work, workqueue_fn);
3 static struct file_operations fops = {
4     .owner = THIS_MODULE,
5     .read = ext_read,
6     .write = ext_write,
7     .open = ext_open,
8     .release = ext_release,
9 };
10 static int __init ext_device_init(void) {
11     int ret;
12     if (alloc_chrdev_region(&dev, 0, 1, "ext_dev") < 0) {
13         printk(KERN_ERR "Can't allocate device numbers\n");
14         return -1;
15     }
16     printk(KERN_INFO "major=%d, minor=%d\n", MAJOR(dev), MINOR(dev));
17     /* Creating cdev structure */
18     cdev_init(&ext_dev, &fops);
19     if (cdev_add(&ext_dev, dev, 1) < 0) {
20         printk(KERN_ERR "Can't add cdev\n");
21         unregister_chrdev_region(dev, 1);
22         cdev_del(&ext_cdev);
23         return -1;
24     }
25     /* Creating struct class */
26     dev_class = class_create(THIS_MODULE, "ext_class");
27     if (IS_ERR(dev_class)) {
28         ...
29     }
30     /* Creating device */
31     if (IS_ERR(device_create(dev_class, NULL, dev, NULL, "ext_device"))) {
```

```

32         ...
33     }
34     /* Register IRQ handler */
35     if (request_irq(IRQ_NO, irq_handler, IRQF_SHARED, "ext_device", NULL)) {
36         ...
37     }
38     return 0;
39 }

```

3 типа драйверов устройств

- встроенные в ядро, соответствующие устройства автоматически обнаружатся системой и станут доступными приложениям
- драйверы реализованные загружаемыми модулями ядра.

Обычно драйверы управляют звуковыми и сетевыми картами и SCSI-адаптерами. Обычно при инсталляции системы устанавливаются перечень автоматически подключаемых модулей.

- поделенные между ядром и утилитой

Принтер: система отвечает за взаимодействие с портом, отвечающим за формирование управляющих сигналов

Часто подсистему для работы с внешними устройствами называют подсистемой вводы-вывода, она занимает относительно большой процент кода ядра.

```

1 struct device {
2     struct kobject kobj;
3     struct device *parent;
4     const char *init_name; /* initial name of the device*/
5     const struct device_type *type;
6     const struct bus_type *bus; /* type of bus device is on */
7     struct device_driver *driver; /* which driver has allocated this device */
8     u64 *dma_mask; /* dma mask (if dma'able device) */
9     u64 bus_dma_limit; /* upstream dma constraint */
10    #ifdef CONFIG_NUMA
11    int numa_node; /* NUMA node this device is close to */
12    #endif
13    dev_t devt; /* dev_t, creates the sysfs "dev" */
14 };

```

Выделяют драйвера нижнего уровня, они пишутся разработчиками внешних устройств.

Современные ОС (WIN/Linux) позволяют управлять внешними устройствами с помощью контроллеров.

Общая структура описания драйверов:

```

1 struct device_driver
2 {
3     const char *name;
4     struct bus_type *bus; // тип шины
5     struct module_t *owner;
6     ...
7     // точки входа
8     int (*probe)(struct device *dev);
9     ...
10    int (*remove)(struct device *dev);
11    void (*shutdown)(struct device *dev);
12    int (*suspend)(struct device *dev);
13    int (*resume)(struct device *dev);
14    ...

```

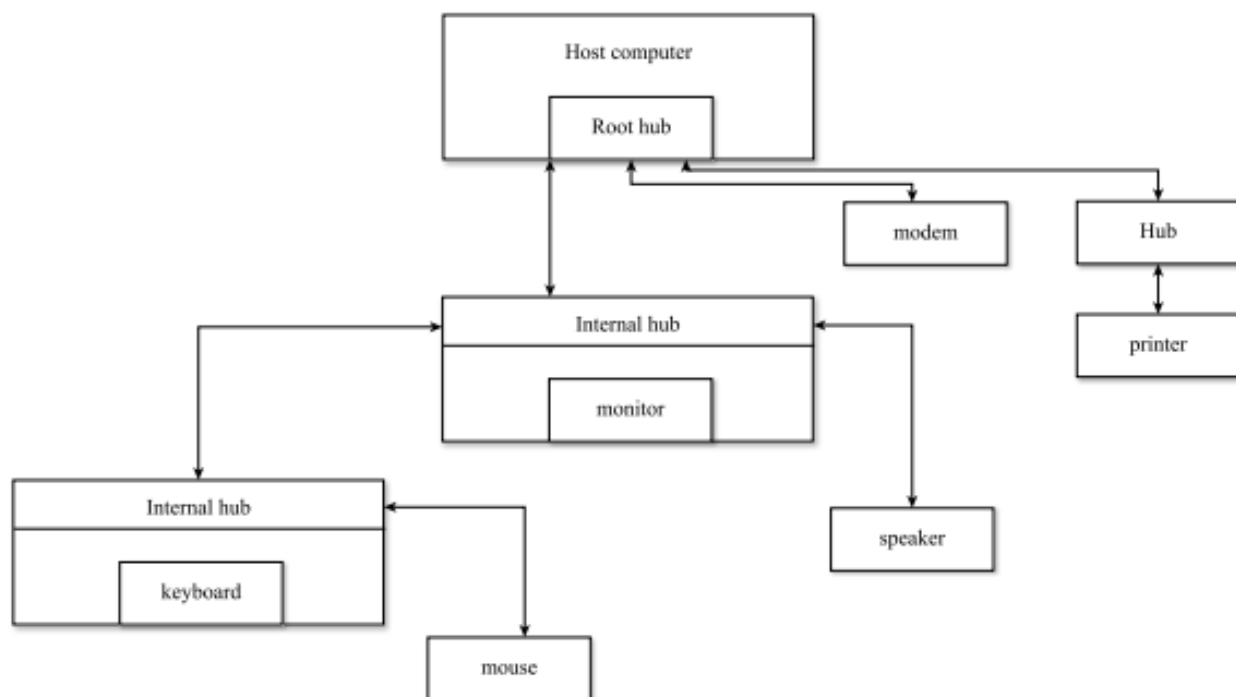
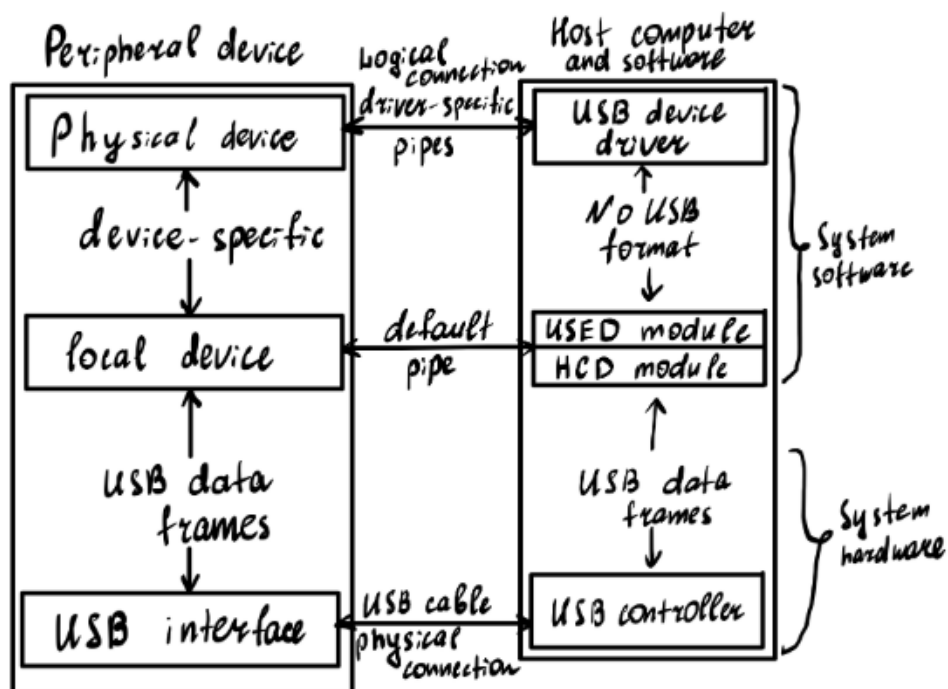


Рис. 12.1: USB топология

Хаб – сетевой концентратор или повторитель, через который соединены узлы сети. Host computer (главный) – корень USB-дерева, который содержит неявные узлы (пока не подключенные). Все USB-устройства работают по принципу горячего подключения: подключил и работает.

Потребность в usb возникла из-за увеличения количества устройств и повышения их быстродействия.

Хаб нужен для распределения электрического сигнала. Такая структура позволяет подключать большое количество устройств ввода/вывода.



В процессе передачи данных

- USB устройство инициирует передачу используя функции интерфейса USB драйвера. Работает по принципу опроса
- USB драйвер отсылает запросы HCD-модулю (host controller driver module)
- HCD делит запросы на отдельные транзакции учитывая возможности шины и характеристики USB-устройства, и планирует транзакции по шине.
- Аппаратная часть выполняет или завершает транзакции.

Транзакции определяются шиной и устройство полностью зависимо.

Базовым понятием при передаче между устройством и шиной является URB (USB request block)

```

1 struct urb {
2     struct kref kref; // reference count of the URB
3     struct list_head urb_list; /* list head for use by the urb's current owner */
4     struct usb_device *dev; /* (in) pointer to associated device */
5     struct usb_host_endpoint *ep; /* (internal) pointer to end point */
6     unsigned int pipe; /* (in) pipe information */
7 }

```

Каждый пакет каждой транзакции содержит номер конечной точки. Драйверы считывают с устройства список конечных точек и создают структуры. Совокупность структур данных в ядре и конечных точек называется каналом (pipe).

4 режима передачи данных

- Control

Двунаправленная передача данных для обмена короткими пакетами вида вопрос/ответ. Обычно осуществляется конечной точкой 0, но могут быть и другие. Позволяет прочитать информацию об устройстве PRODUCT\_ID и VENDOR\_ID

- Isochronous

Гарантирует пропускную способность пакетов за 1 период шины. Также могут быть заданы задержки во время передачи. Данные передаются без подтверждения приёма, используется для устройств реального времени (видео/аудио)

- Interrupt

Короткая, до 64 байт на полной скорости и до 8 байт на низкой. Характерны для ввода символов/координат. Имеют спонтанный характер и должны обслуживаться не медленнее, чем того требует устройство.

- Bulk

Сплошная (поточная) передача данных. Используются устройства, принимающие/отправляющие большое количество данных, но не требующие определенной пропускной способности и времени задержек. Данные имеют самый низкий приоритет и занимают всю свободную полосу пропускания шины.

```

1 struct usb_driver
2 {
3     const char *name;
4     int (*probe)(struct usb_interface *intf, const struct usb_device_id *id);
5     void (*disconnect)(struct usb_interface *intf);
6     ...
7     int (*suspend)(struct usb_interface *intf, pm_message_t message);
8     int (*resume)(struct usb_interface *intf);
9     ...
10    const struct usb_device_id *id_table; // идентификация устройства, обычно
    PRODUCT_ID и VENDOR_ID
11    ...
12    struct device_driver driver;
13 };

```

```

1 struct usb_device_driver {
2     const char *name;
3     bool (*match)(struct usb_device *udev);
4     int (*probe)(struct usb_device *udev);
5     void (*disconnect)(struct usb_device *udev);
6     struct device_driver driver;

```

```
7     const struct usb_device_id *id_table;
```

```
8 };
```

```
1 static int __init my_module_init(void)
```

```
2 {
```

```
3     return usb_register(&my_drivers);
```

```
4 }
```

```
5  
6 static void __exit my_module_exit(void)
```

```
7 {
```

```
8     usb_deregister(&my_driver);
```

```
9 }
```

```
10 static struct usb_driver my_driver =
```

```
11 {
```

```
12     .name = DEV_NAME,
```

```
13     .probe = my_probe,
```

```
14     .disconnect = my_remove,
```

```
15     .id_table = my_table
```

```
16 };
```

```
17 static const struct usb_device_id my_table[] = {{USB_DEVICE(VID, PID)}, {}};
```

```
18  
19 static int my_probe(struct usb_interface *interface, const struct usb_device_id *id)
```

```
20 {
```

```
21     struct usb_endpoint_descriptor *int_in;
```

```
22     ....
```

```
23     device_urb_interrupt = usb_alloc_urb(0, GFP_KERNEL);
```

```
24     int res = usb_register_dev(interface, &my_class);
```

```
25     ...
```

```
26     usb_fill_int_urb(device_urb_interrupt, deviceudev,
```

```
usb_revintpipe(deviceudev, device_endpoint_address), device.data,
```

```
device.buffer_size, process_key_press, NULL, device.time);
```

```
27     ...
```

```
28 }
```

Зачем введены структуры USB? Видимо, вследствие широкого распространения usb и необходимости разработчиков ядра оптимизировать работу с ними. (но это не точно)