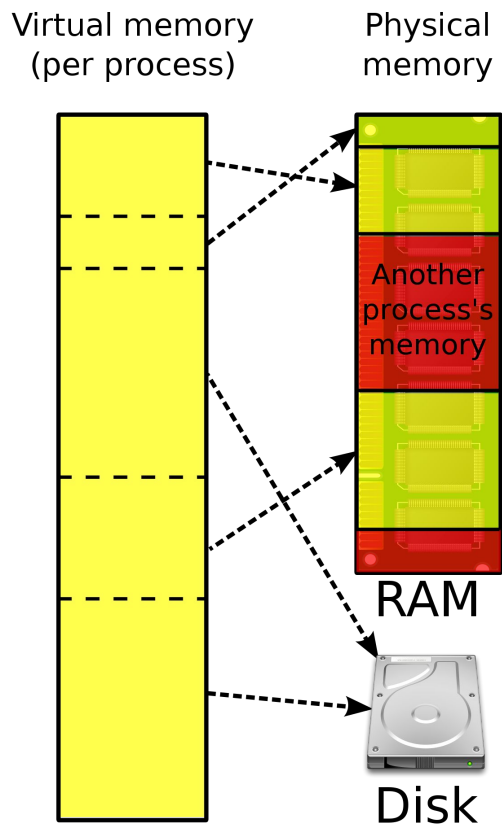


Программирование в Linux

Файлы и IO

Отображение файлов в память



1. Ядро: отображение виртуальные страницы <-> физические — прописывает в спец. регистры CPU (адрес LUT)
2. Ядро: настраивает, что делать, если страница не найдена
 - a. SIGSEGV
 - b. подгрузка
3. Можно настроить отображение виртуальных страниц на подгрузку из определенного файла (mmap)

```
1 void *mmap(void *addr, size_t length, int prot, int flags,  
2           int fd, off_t offset);  
3 int munmap(void *addr, size_t length);
```

Адресное пространство процесса



addr — опционален
(подсказка,
рекомендация по
размещению).
Обычно выровнен по
PAGE_SIZE

размеры должны
быть кратны
PAGE_SIZE



boost::mapped_file — обертка

Анонимные отображения

Отображение файла `/dev/zero` — эквивалентно выделению памяти

Более удобно:

```
1 void* anon_memory = mmap(NULL,  
2                             size,  
3                             PROT_READ | PROT_WRITE,  
4                             MAP_ANONYMOUS | MAP_PRIVATE,  
5                             -1, // fd игнорируется  
6                             0); // смещение тоже
```

Полученные страницы памяти будут заполнены нулями

Настройка прав доступа к отображениям

```
#include <sys/mman.h>

int mprotect(void *addr, size_t len, int prot);

#define _GNU_SOURCE          /* See feature_test_macros(7) */
#include <sys/mman.h>

int pkey_mprotect(void *addr, size_t len, int prot, int pkey);
```

```
33 int main(int argc, char *argv[])
34 {
35     ... configure_signal_handler();
36     ... int pagesize = sysconf(_SC_PAGE_SIZE);
37     ... if (pagesize == -1)
38         ... handle_error("sysconf");
39
40     ... /* Allocate a buffer aligned on a page boundary;
41         ... initial protection is PROT_READ | PROT_WRITE */
42
43     ... char* buffer = memalign(pagesize, 4 * pagesize);
44     ... if (buffer == NULL)
45         ... handle_error("memalign");
46
47     ... printf("Start of region: .....0x%lx\n", (long) buffer);
48
49     ... if (mprotect(buffer + pagesize * 2, pagesize,
50         ... PROT_READ) == -1)
51         ... handle_error("mprotect");
52
53     ... for (char* p = buffer ; ; )
54         ... *(p++) = 'a';
55
56     ... printf("Loop completed\n"); ... /* Should never happen */
57     ... exit(EXIT_SUCCESS);
58 }
```

```
13 #define handle_error(msg) \
14     ... do { perror(msg); exit(EXIT_FAILURE); } while (0)
15
16 static void
17 handler(int sig, siginfo_t* si, void *unused) {
18     ... printf("Got SIGSEGV at address: 0x%lx\n",
19         ... (long) si->si_addr); ...
20     ... exit(EXIT_FAILURE);
21 }
22
23 void configure_signal_handler(void) {
24     ... struct sigaction sa;
25     ... sa.sa_flags = SA_SIGINFO;
26     ... sigemptyset(&sa.sa_mask);
27     ... sa.sa_sigaction = handler;
28     ... if (sigaction(SIGSEGV, &sa, NULL) == -1)
29         ... handle_error("sigaction");
30 }
```

Специальные файлы

/dev/zero: read -> нули; write -> success, discarded

/dev/null: read -> EOF; write -> success discarded

/dev/full: read -> нули; write -> error (no space)

/dev/random, /dev/urandom — генераторы псевдослучайных чисел. random блокируется, если недостаточно энтропии

/dev/uinput — файл для имитации пользовательского ввода (мышь, клавиатура, джойстик...)

```

24  int fd = open("/dev/uinput", O_WRONLY | O_NONBLOCK);
25
26  /* enable mouse button left and relative events */
27  ioctl(fd, UI_SET_EVBIT, EV_KEY);
28  ioctl(fd, UI_SET_KEYBIT, BTN_LEFT);
29
30  ioctl(fd, UI_SET_EVBIT, EV_REL);
31  ioctl(fd, UI_SET_RELBIT, REL_X);
32  ioctl(fd, UI_SET_RELBIT, REL_Y);
33
34  struct uinput_setup usetup;
35  memset(&usetup, 0, sizeof(usetup));
36  usetup.id.bustype = BUS_USB;
37  usetup.id.vendor = 0x1234; /* sample vendor */
38  usetup.id.product = 0x5678; /* sample product */
39  strcpy(usetup.name, "Example device");
40
41  ioctl(fd, UI_DEV_SETUP, &usetup);
42  ioctl(fd, UI_DEV_CREATE);
43
44  /* Give userspace some time to detect new device */
45  sleep(1);
46
47  int i = 50;
48  /* Move the mouse diagonally, 5 units per axis */
49  while (i--) {
50      emit(fd, EV_REL, REL_X, 5);
51      emit(fd, EV_REL, REL_Y, 5);
52      emit(fd, EV_SYN, SYN_REPORT, 0);
53      usleep(15000);
54  }
55  /*
56   * Give userspace some time to read the events before we destroy the
57   * device with UI_DEV_DESTROY.
58   */
59  sleep(1);
60
61  ioctl(fd, UI_DEV_DESTROY);
62  close(fd);
63

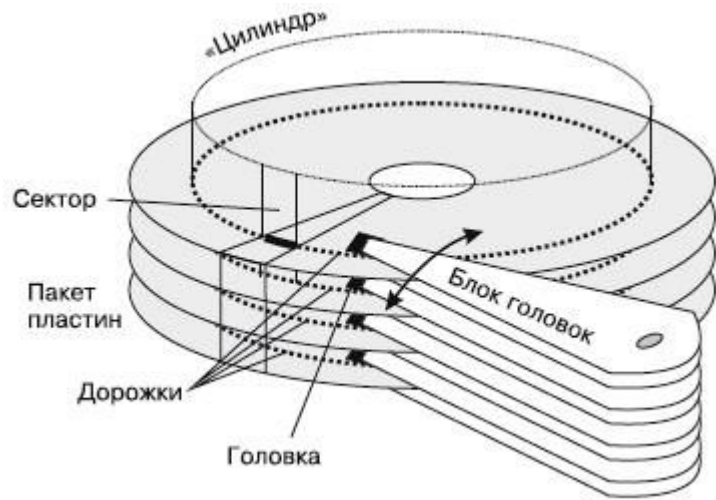
```

```

8  void emit(int fd, int type, int code, int val)
9  {
10     struct input_event ie;
11
12     ie.type = type;
13     ie.code = code;
14     ie.value = val;
15     /* timestamp values below are ignored */
16     ie.time.tv_sec = 0;
17     ie.time.tv_usec = 0;
18
19     write(fd, &ie, sizeof(ie));
20 }

```

Планировщик IO



Скорость вращения 5400 - 15000 rpm

Произвольный доступ → постоянное позиционирование головок

Время позиционирования: 4 - 12 ms!

Основные планировщики

У всех* — sort по позиции на устройстве и объединение смежных запросов

- NOOP
 - общий FIFO
 - не сортирует
 - пытается объединить смежные
 - хорош для SSD
- Deadline
 - Общая отсортированная очередь
 - Timeout устаревания запроса
 - + Две очереди (Read/Write) по времени устаревания
 - Если в голове очереди устаревший запрос — разбирай эту очередь
 - default на сегодняшний день
- CFQ (Completely fair queue)
 - блокирующие запросы — у каждого процесса своя очередь
 - асинхронные запросы — группируем по устройствам
 - Квант времени на очередь
 - очереди переключаются каруселью (Round Robin, RR)
- BFQ (Budget Fair Queue)
 - Модификация CFQ
 - RR учитывает размеры запросов

События файловой системы. inotify

```
4 int fd = inotify_init1(IN_NONBLOCK);
5 int wd = inotify_add_watch(fd, "/path/", IN_ALL_EVENTS);
6 ...
7 char buf[4096] __attribute__((aligned(__alignof__(struct inotify_event))));
8 int len = read(fd, buf, sizeof buf);
9
10 /* Loop over all events in the buffer */
11 for (char* ptr = buf; ptr < buf + len;
12      ptr += sizeof(struct inotify_event) + event->len) {
13     const struct inotify_event * event = (const struct inotify_event *) ptr;
14     // process events
15 }
16 |
```

```
struct inotify_event {
    int wd; /* Watch descriptor */
    uint32_t mask; /* Mask describing event */
    uint32_t cookie; /* Unique cookie associating related
                     events (for rename(2)) */
    uint32_t len; /* Size of name field */
    char name[]; /* Optional null-terminated name */
};
```

```
fd = open("dir/myfile", O_RDWR);
    Generates IN_OPEN events for both dir and dir/myfile.

read(fd, buf, count);
    Generates IN_ACCESS events for both dir and dir/myfile.

write(fd, buf, count);
    Generates IN_MODIFY events for both dir and dir/myfile.

fchmod(fd, mode);
    Generates IN_ATTRIB events for both dir and dir/myfile.

close(fd);
    Generates IN_CLOSE_WRITE events for both dir and dir/myfile.
```