

САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ ПЕТРА ВЕЛИКОГО

КАФЕДРА КОМПЬЮТЕРНЫХ СИСТЕМ И ПРОГРАММНЫХ ТЕХНОЛОГИЙ

Отчёт по лабораторной работе №1

Курс: «Проектирование ОС и компонентов»

Тема: «Системные вызовы»

Выполнил студент:

Бояркин Никита Сергеевич

Группа: 13541/3

Проверил:

Душутина Елена Владимировна

Санкт-Петербург
2018 г.

Содержание

1	Лабораторная работа №1	2
1.1	Цель работы	2
1.2	Программа работы	2
1.3	Индивидуальное задание	2
1.4	Характеристики системы	2
1.5	Системные вызовы	2
1.5.1	Перехват системных вызовов	3
1.5.2	Системный вызов sys_getpid	4
1.5.3	Системный вызов sys_setpid	7
1.5.4	Системный вызов sys_chdir	7
1.5.5	Системный вызов sys_sysinfo	10
1.6	Вывод	14

Лабораторная работа №1

1.1 Цель работы

Ознакомиться с системными вызовами, научиться разрабатывать перехватчики для системных вызовов и вносить изменения в ядро.

1.2 Программа работы

Для каждого системного вызова из индивидуального задания:

- Привести описание.
- Разработать программу с примером использования.
- Осуществить перехват вызова через хук.
- Проанализировать исходный код вызова.

1.3 Индивидуальное задание

Системные вызовы:

- `sys_getpid`
- `sys_setpid`
- `sys_chdir`
- `sys_sysinfo`

1.4 Характеристики системы

```
1 nikita@nikita-VirtualBox:~$ who
2 nikita    tty7          2018-05-07 02:57 (:0)
3 nikita@nikita-VirtualBox:~$ cat /proc/version
4 Linux version 4.4.0-121-generic (buildd@lcy01-amd64-004) (gcc version 5.4.0 20160609 (
    Ubuntu 5.4.0-6ubuntu1~16.04.9) ) #145-Ubuntu SMP Fri Apr 13 13:47:23 UTC 2018
```

1.5 Системные вызовы

В конечном итоге, главная задача операционной системы — это обслуживание потребностей прикладных пользовательских процессов. И обеспечивается это обслуживание через механизм системных вызовов.

В любой (в том числе и микроядерной) операционной системе системный вызов выполняется некоторой выделенной процессорной инструкцией, прерывающей последовательное выполнение команд и передающий управление коду режима супервизора. Это обычно некая команда программного прерывания, в зависимости от архитектуры процессора в разные времена это были команды с мнемониками вида: `svc`, `emt`, `trap`, `int` и им подобными [1]

Команды прерывания для ОС, построенных на архитектуре Intel x86:

OS	desc
MS-DOS	21h
Windows	2Eh
Linux	80h
QNX	21h
MINIX 3	21h

Прикладной процесс вызывает требуемые ему службы посредством библиотечного вызова к множеству библиотек либо вида *.so (динамическое связывание), либо прикомпоновывая к себе фрагмент из библиотеки вида *.a (статическое связывание). Самые известные примеры - это стандартная библиотека языка C libc.so или libpthread.so — библиотека POSIX потоков [1]

Описания системных вызовов (в отличие от библиотечных) отнесены к секции 2 в руководствах man. Все системные вызовы далее преобразуются в вызов ядра функцией syscall(), 1-м параметром которого будет идентификатор выполняемого системного вызова, например __NR_execve.

Найдем описания заданных системных вызовов в таблице векторов системных вызовов syscall_32.tbl:

```
1 # linux/arch/x86/entry/syscalls/syscall_32.tbl
2 #
3 # 32-bit system call numbers and entry vectors
4 #
5 # The format is:
6 # <number> <abi> <name> <entry point> <compat entry point>
7 #
8
9 < ... >
10
11
12 12 i386 chdir sys_chdir __ia32_sys_chdir
13
14 < ... >
15
16 20 i386 getpid sys_getpid __ia32_sys_getpid
17
18 < ... >
19
20 116 i386 sysinfo sys_sysinfo __ia32_compat_sys_sysinfo
21
22 < ... >
```

1.5.1 Перехват системных вызовов

Для перехвата системных вызовов воспользуемся библиотекой kmod_hooking [2] Данная библиотека содержит набор системных функций для регистрации и снятия перехватчиков с заданных функций. Принцип работы основывается на подмене функции внутри ядра. Синтаксис перехватчика иллюстрируется следующим примером перехвата системной функции inode_permission:

```
1 #include <linux/fs.h>
2
3 < ... >
4
5 DECLARE_KHOOK(inode_permission);
6 int khook_inode_permission(struct inode * inode, int mode)
7 {
8     int result;
9
10    KHOOK_USAGE_INC(inode_permission);
11
12    debug("%s(%pK,%08x) [%s]\n", __func__, inode, mode, current->comm);
13
14    result = KHOOK_ORIGIN(inode_permission, inode, mode);
15
16    debug("%s(%pK,%08x) [%s] = %d\n", __func__, inode, mode, current->comm, result);
```

```

17
18 KHOOK_USAGE_DEC(inode_permission);
19
20 return result;
21 }
22
23 < ... >

```

DECLARE_KHOOK указывает на то, какие функции являются перехватчиками. KHOOK_USAGE_INC и KHOOK_USAGE_DEC определяют область, в которой выполняется хук. KHOOK_ORIGIN вызывает непосредственно системный вызов с заданными аргументами.

Добавление собственного кода в ядро производится с помощью модулей. Модули располагаются в директории /lib/modules/(uname -r). Сборка библиотеки производится при помощи следующего makefile:

```

1 NAME      := hooking
2
3 obj-m     := $(NAME).o
4 obj-y     := libudis86/
5
6 ldflags-y := -T$(src)/layout.lds
7
8 $(NAME)-y := module-init.o libudis86/built-in.o
9
10 all:
11     make -C /lib/modules/$(shell uname -r)/build M=$(shell pwd)
12
13 clean:
14     make -C /lib/modules/$(shell uname -r)/build M=$(shell pwd) clean

```

Добавление нового модуля в систему производится через утилиту insmod:

```

1 nikita@nikita-VirtualBox:~/files/student/temp/kmod_hooking$ sudo insmod hooking.ko

```

Удаление модуля из системы производится через утилиту rmmod:

```

1 nikita@nikita-VirtualBox:~/files/student/temp/kmod_hooking$ sudo rmmod hooking

```

1.5.2 Системный вызов sys_getpid

Системный вызов sys_getpid() возвращает идентификатор текущего процесса. Пользовательская функция getpid() предоставляет доступ к данному системному вызову и объявляется в заголовочном файлеunistd.h.

Пример использования

Функция не принимает никаких аргументов и возвращает целочисленный идентификатор процесса:

```

1 #include <iostream>
2 #include <unistd.h>
3
4 int main() {
5     std::cout << "Process PID: " << getpid() << std::endl;
6     return 0;
7 }

```

Результат выполнения программы:

```

1 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ ./executable &
2 [1] 9225
3 Process PID: 9225
4 [1]+  Done                  ./executable

```

Отслеживание системного вызова

Попробуем отследить системный вызов в предыдущей программе через утилиту strace:

```
1 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ strace ./executable
2 execve("./executable", ["/executable"], [/* 74 vars */]) = 0
3
4 < ... >
5
6 brk(NULL)                                = 0x1418000
7 brk(0x144a000)                          = 0x144a000
8 getpid()                                = 9225
9 fstat(1, {st_mode=S_IFCHR|0620, st_rdev=makedev(136, 19), ...}) = 0
10 write(1, "Process PID: 9225\n", 18Process PID: 9225
11 )                                         = 18
12 exit_group(0)                           = ?
13 +++ exited with 0 +++
```

После создания процесса вызывается функция getpid и выводится на экран функцией write, после чего программа завершается.

Перехват системного вызова

Реализуем перехват системного вызова sys_getpid с помощью kmod_hooking:

```
1 #include <linux/syscalls.h>
2
3 < ... >
4
5 DECLARE_KHOOK(sys_getpid);
6 long khook_sys_getpid()
7 {
8     KHOOK_USAGE_INC(sys_getpid);
9
10     long result = KHOOK_ORIGIN(sys_getpid);
11     printk("Syscall 'sys_getpid' hooked. Result: %d.\n", result);
12
13     KHOOK_USAGE_DEC(sys_getpid);
14
15     return result;
16 }
17
18 < ... >
```

Данный код добавляется в linux как модуль ядра, а результат работы журналируется в файл kern.log. Убедимся в правильности работы хука:

```
1 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ ./executable &
2 [1] 9431
3 Process PID: 9431
4 [1]+  Done                  ./executable
5 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ tail /var/log/kern.log
6 May  3 06:10:55 nikita-VirtualBox kernel: [11609.557046] raid6: sse2x2  xor() 8688 MB/s
7 May  3 06:10:56 nikita-VirtualBox kernel: [11609.625101] raid6: sse2x4  gen() 15389 MB/s
8 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692700] raid6: sse2x4  xor() 10589 MB/s
9 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692702] raid6: using algorithm sse2x4
   gen() 15389 MB/s
10 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692703] raid6: .... xor() 10589 MB/s,
   rmw enabled
11 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692704] raid6: using ssse3x2 recovery
   algorithm
12 May  3 06:10:56 nikita-VirtualBox kernel: [11609.693802] xor: automatically using best
   checksumming function:
13 May  3 06:10:56 nikita-VirtualBox kernel: [11609.732719] avx          : 20010.000 MB/sec
14 May  3 06:10:56 nikita-VirtualBox kernel: [11609.749930] Btrfs loaded
15 May  3 08:20:33 nikita-VirtualBox kernel: [19386.834543] Syscall 'sys_getpid' hooked.
   Result: 9431.
```

После запуска программы в kern.log появилась запись о попытке доступа к системному вызову sys_getpid.

Анализ исходного кода

Прототип системного вызова `sys_getpid` находится в `include/linux/syscalls.h` в исходных файлах ядра linux [3]:

```
1 // include/linux/syscalls.h
2
3 asmlinkage long sys_getpid(void);
```

Непосредственно реализация `sys_getpid` находится в `kernel/sys.c` в исходных файлах ядра linux [4]:

```
1 // kernel/sys.c
2
3 SYSCALL_DEFINE0(getpid)
4 {
5     return task_tgid_vnr(current);
6 }
```

Можно заметить, что вызывается функция ядра `task_tgid_vnr`, которая в свою очередь вызывает `__task_pid_nr_ns` в `kernel/pid.c` в исходных файлах ядра linux [5]:

```
1 // kernel/pid.c
2
3 pid_t __task_pid_nr_ns(struct task_struct *task, enum pid_type type,
4                        struct pid_namespace *ns)
5 {
6     pid_t nr = 0;
7
8     rcu_read_lock();
9     if (!ns)
10         ns = task_active_pid_ns(current);
11     if (likely(pid_alive(task))) {
12         if (type != PIDTYPE_PID) {
13             if (type == __PIDTYPE_TGID)
14                 type = PIDTYPE_PID;
15
16             task = task->group_leader;
17         }
18         nr = pid_nr_ns(rcu_dereference(task->pids[type].pid), ns);
19     }
20     rcu_read_unlock();
21
22     return nr;
23 }
```

На время работы функции блокируется чтение, проверяется наличие текущего процесса и наличие `grpId`, после чего вызывается функция `pid_nr_ns`:

```
1 // kernel/pid.c
2
3 pid_t pid_nr_ns(struct pid *pid, struct pid_namespace *ns)
4 {
5     struct upid *upid;
6     pid_t nr = 0;
7
8     if (pid && ns->level <= pid->level) {
9         upid = &pid->numbers[ns->level];
10         if (upid->ns == ns)
11             nr = upid->nr;
12     }
13     return nr;
14 }
```

Если оба аргумента (структуры `pid` и `pid_namespace`) корректны, то возвращается `pid` текущего процесса, в противном случае возвращается 0.

1.5.3 Системный вызов sys_setpid

Такого системного вызова, очевидно, не существует, потому что идентификатор процесса задается единожды при запуске процесса.

Для выделения нового pid используется функция alloc_pid:

```
1 // kernel/pid.c
2
3 struct pid *alloc_pid(struct pid_namespace *ns)
4 {
5     < ... >
6
7     pid = kmem_cache_alloc(ns->pid_cache, GFP_KERNEL);
8     if (!pid)
9         return ERR_PTR(retval);
10
11     tmp = ns;
12     pid->level = ns->level;
13
14     for (i = ns->level; i >= 0; i--) {
15         int pid_min = 1;
16
17         idr_preload(GFP_KERNEL);
18         spin_lock_irq(&pidmap_lock);
19
20         /*
21          * init really needs pid 1, but after reaching the maximum
22          * wrap back to RESERVED_PIDS
23          */
24         if (idr_get_cursor(&tmp->idr) > RESERVED_PIDS)
25             pid_min = RESERVED_PIDS;
26
27         /*
28          * Store a null pointer so find_pid_ns does not find
29          * a partially initialized PID (see below).
30          */
31         nr = idr_alloc_cyclic(&tmp->idr, NULL, pid_min,
32                             pid_max, GFP_ATOMIC);
33         spin_unlock_irq(&pidmap_lock);
34         idr_preload_end();
35
36         if (nr < 0) {
37             retval = nr;
38             goto out_free;
39         }
40
41         pid->numbers[i].nr = nr;
42         pid->numbers[i].ns = tmp;
43         tmp = tmp->parent;
44     }
45
46     < ... >
47 }
```

Функция вызывает kmem_cache_alloc с аргументом в виде набора кэшированных структур для выделения памяти под структуру нового pid [6].

Далее происходит инициализация полей структуры в зависимости от уровня аргумента pid_namespace. Функция idr_alloc_cyclic пытается выделить новый идентификатор процесса. Если получилось, то идентификатор записывается в структуру.

1.5.4 Системный вызов sys_chdir

Системный вызов sys_chdir() задает текущий путь выполнения процесса. Пользовательская функция chdir() предоставляет доступ к данному системному вызову и объявляется в заголовочном файлеunistd.h.

Пример использования

Для примера выведем текущий путь выполнения процесса, поменяем его функцией `chdir()` и выведем путь еще раз. Путь должен измениться на аргумент функции `chdir()`:

```
1 #include <iostream>
2 #include <unistd.h>
3
4 static const int DIRECTORY_PATH_SIZE = 1024;
5 static const char* DIRECTORY_PATH = "/home";
6
7 static char currentDirectory [DIRECTORY_PATH_SIZE];
8
9 const char* getCurrentDirectory ();
10
11 int main() {
12     std::cout << "Current directory: " << getCurrentDirectory() << std::endl;
13     chdir(DIRECTORY_PATH);
14     std::cout << "Current directory after chdir: " << getCurrentDirectory() << std::endl;
15     return 0x0;
16 }
17
18 const char* getCurrentDirectory () {
19     char* result = getcwd(currentDirectory, DIRECTORY_PATH_SIZE);
20
21     if (result == NULL) {
22         return NULL;
23     }
24
25     return currentDirectory;
26 }
```

Результат выполнения программы:

```
1 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ ./executable
2 Current directory: /home/nikita/files/student/ezsem4/sys/1/listings
3 Current directory after chdir: /home
```

Отслеживание системного вызова

Попробуем отследить системный вызов в предыдущей программе через утилиту `strace`:

```
1 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ strace ./executable
2 execve("./executable", ["/home/nikita/files/student/ezsem4/sys/1/listings/./executable"],
3   [/* 74 vars */]) = 0
4
5 < ... >
6
7 getcwd("/home/nikita/files/student/ezsem4/sys/1/listings", 1024) = 49
8 fstat(1, {st_mode=S_IFCHR|0620, st_rdev=makedev(136, 19), ...}) = 0
9 write(1, "Current directory: /home/nikita/"..., 68Current directory: /home/nikita/files/
10   student/ezsem4/sys/1/listings
11   ) = 68
12 chdir("/home")
13   = 0
14 getcwd("/home", 1024)
15   = 6
16 write(1, "Current directory after chdir: /"..., 37Current directory after chdir: /home
17   ) = 37
18 exit_group(0)
19   = ?
20 +++ exited with 0 +++
```

Как и ожидалось, была вызвана цепочка системных вызовов: `getcwd -> write -> chdir -> getcwd -> write`, в результате которой изменилось значение путь выполнения процесса.

Перехват системного вызова

Реализуем перехват системного вызова `sys_chdir` с помощью `kmod_hooking`:

```

1 #include <linux/syscalls.h>
2
3 < ... >
4
5 DECLARE_KHOOK(sys_chdir);
6 long khook_sys_chdir(const char __user *filename)
7 {
8     KHOOK_USAGE_INC(sys_chdir);
9
10    long result = KHOOK_ORIGIN(sys_chdir, filename);
11    printk("Syscall 'sys_chdir' hooked. Arguments: %s. Result: %d.\n", filename, result);
12
13    KHOOK_USAGE_DEC(sys_chdir);
14
15    return result;
16 }
17
18 < ... >

```

Данный код добавляется в linux как модуль ядра, а результат работы журналируется в файл kern.log. Убедимся в правильности работы хука:

```

1 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ ./executable
2 Current directory: /home/nikita/files/student/ezsem4/sys/1/listings
3 Current directory after chdir: /home
4 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ tail /var/log/kern.log
5 May  3 06:10:56 nikita-VirtualBox kernel: [11609.625101] raid6: sse2x4   gen() 15389 MB/s
6 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692700] raid6: sse2x4   xor() 10589 MB/s
7 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692702] raid6: using algorithm sse2x4
8 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692703] raid6: .... xor() 10589 MB/s,
9 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692704] raid6: using ssse3x2 recovery
10 May  3 06:10:56 nikita-VirtualBox kernel: [11609.693802] xor: automatically using best
11 May  3 06:10:56 nikita-VirtualBox kernel: [11609.732719] avx          : 20010.000 MB/sec
12 May  3 06:10:56 nikita-VirtualBox kernel: [11609.749930] Btrfs loaded
13 May  3 08:20:33 nikita-VirtualBox kernel: [19386.834543] Syscall 'sys_getpid' hooked.
14 May  3 08:21:32 nikita-VirtualBox kernel: [19445.923432] Syscall 'sys_chdir' hooked.
    Arguments: /home. Result: 0.

```

После запуска программы в kern.log появилась запись о попытке доступа к системному вызову sys_chdir.

Анализ исходного кода

Прототип системного вызова sys_chdir находится в include/linux/syscalls.h в исходных файлах ядра linux [3]:

```

1 // include/linux/syscalls.h
2
3 asmlinkage long sys_chdir(const char __user *filename);

```

Непосредственно реализация sys_chdir находится в kernel/sys.c в исходных файлах ядра linux [4]:

```

1 // kernel/sys.c
2
3 SYSCALL_DEFINE1(chdir, const char __user *, filename)
4 {
5     return ksys_chdir(filename);
6 }

```

Можно заметить, что вызывается функция ядра ksys_chdir в fs/open.c в исходных файлах ядра linux [7]:

```

1 // fs/open.c
2

```

```

3 int ksys_chdir(const char __user *filename)
4 {
5     struct path path;
6     int error;
7     unsigned int lookup_flags = LOOKUP_FOLLOW | LOOKUP_DIRECTORY;
8 retry:
9     error = user_path_at(AT_FDCWD, filename, lookup_flags, &path);
10    if (error)
11        goto out;
12
13    error = inode_permission(path.dentry->d_inode, MAY_EXEC | MAY_CHDIR);
14    if (error)
15        goto dput_and_out;
16
17    set_fs_pwd(current->fs, &path);
18
19 dput_and_out:
20    path_put(&path);
21    if (retry_estale(error, lookup_flags)) {
22        lookup_flags |= LOOKUP_REVAL;
23        goto retry;
24    }
25 out:
26    return error;
27 }

```

Функция проверяет наличие задаваемого пути и наличие прав на его изменение. После этого вызывается функция `set_fs_pwd`, которая непосредственно задает путь:

```

1 // fs/open.c
2
3 void set_fs_pwd(struct fs_struct *fs, const struct path *path)
4 {
5     struct path old_pwd;
6
7     path_get(path);
8     spin_lock(&fs->lock);
9     write_seqcount_begin(&fs->seq);
10    old_pwd = fs->pwd;
11    fs->pwd = *path;
12    write_seqcount_end(&fs->seq);
13    spin_unlock(&fs->lock);
14
15    if (old_pwd.dentry)
16        path_put(&old_pwd);
17 }

```

1.5.5 Системный вызов `sys_sysinfo`

Системный вызов `sys_sysinfo()` выводит системную информацию, такую как: время с момента запуска системы в секундах, количество памяти и количество свободно памяти в байтах, количество активных процессов, размер свфпа и др. Пользовательская функция `sysinfo()` предоставляет доступ к данному системному вызову и объявляется в заголовочном файле `sys/sysinfo.h`.

Пример использования

Для примера выведем системную информацию в удобном для чтения формате:

```

1 #include <iostream>
2 #include <sys/sysinfo.h>
3
4 static const int MEGABYTE_IN_BYTES = 1024 * 1024;
5 static const int MINUTE_IN_SECONDS = 60;
6 static const int HOUR_IN_SECONDS = MINUTE_IN_SECONDS * 60;
7 static const int DAY_IN_SECONDS = HOUR_IN_SECONDS * 24;
8

```

```

9 int main() {
10     struct sysinfo information;
11     sysinfo(&information);
12
13     std::cout << "Time since boot: days " << information.uptime / DAY_IN_SECONDS << ", "
14     << (information.uptime % DAY_IN_SECONDS) /
15     HOUR_IN_SECONDS << ":"
16     << (information.uptime % HOUR_IN_SECONDS) /
17     MINUTE_IN_SECONDS << ":"
18     << information.uptime % MINUTE_IN_SECONDS <<
19     std::endl;
20
21     std::cout << "Total RAM: " << information.totalram / MEGABYTE_IN_BYTES << " Mb" <<
22     std::endl;
23     std::cout << "Free RAM: " << information.freeram / MEGABYTE_IN_BYTES << " Mb" << std
24     ::endl;
25     std::cout << "Process count: " << information.procs << std::endl;
26     return 0x0;
27 }

```

Результат выполнения программы:

```

1 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ ./executable
2 Time since boot: days 0, 14:42:23
3 Total RAM: 3951 Mb
4 Free RAM: 605 Mb
5 Process count: 631

```

Отслеживание системного вызова

Попробуем отследить системный вызов в предыдущей программе через утилиту strace:

```

1 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ strace ./executable
2 execve("./executable", ["/executable"], [/* 74 vars */]) = 0
3
4 < ... >
5
6 sysinfo({uptime=53030, loads=[31328, 29824, 20160], totalram=4142927872, freeram
7     =628195328, sharedram=67964928, bufferram=390832128, totalswap=4292866048, freeswap
8     =4292866048, procs=633, totalhigh=0, freehigh=0, mem_unit=1}) = 0
9 fstat(1, {st_mode=S_IFCHR|0620, st_rdev=makedev(136, 19), ...}) = 0
10 write(1, "Time since boot: days 0, 14:43:5"... , 34Time since boot: days 0, 14:43:50
11 ) = 34
12 write(1, "Total RAM: 3951 Mb\n", 19Total RAM: 3951 Mb
13 ) = 19
14 write(1, "Free RAM: 599 Mb\n", 17Free RAM: 599 Mb
15 ) = 17
16 write(1, "Process count: 633\n", 19Process count: 633
17 ) = 19
18 exit_group(0) = ?
19 +++ exited with 0 +++

```

Как и ожидалось была получена структура sysinfo через соответствующий системный вызов.

Перехват системного вызова

Реализуем перехват системного вызова sys_sysinfo с помощью kmod_hooking:

```

1 #include <linux/syscalls.h>
2
3 < ... >
4
5 DECLARE_KHOOK(sys_sysinfo);
6 long khook_sys_sysinfo(struct sysinfo __user *info)
7 {
8     KHOOK_USAGE_INC(sys_sysinfo);
9 }

```

```

10 long result = KHOOK_ORIGIN(sys_sysinfo, info);
11 printk("Syscall 'sys_sysinfo' hooked. Arguments: %ld, %ld, %ld, %ld. Result: %d.\n",
12        info->uptime, info->totalram, info->freeram, info->procs, result);
13
14 KHOOK_USAGE_DEC(sys_sysinfo);
15
16 return result;
17 }
18
19 < ... >

```

Данный код добавляется в linux как модуль ядра, а результат работы журналируется в файл kern.log. Убедимся в правильности работы хука:

```

1 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ ./executable
2 Time since boot: days 0, 5:24:33
3 Total RAM: 3951 Mb
4 Free RAM: 208 Mb
5 Process count: 644
6 nikita@nikita-VirtualBox:~/files/student/ezsem4/sys/1/listings$ tail /var/log/kern.log
7 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692700] raid6: sse2x4  xor() 10589 MB/s
8 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692702] raid6: using algorithm sse2x4
9 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692703] raid6: .... xor() 10589 MB/s,
10 May  3 06:10:56 nikita-VirtualBox kernel: [11609.692704] raid6: using ssse3x2 recovery
11 May  3 06:10:56 nikita-VirtualBox kernel: [11609.693802] xor: automatically using best
12 May  3 06:10:56 nikita-VirtualBox kernel: [11609.732719] avx : 20010.000 MB/sec
13 May  3 06:10:56 nikita-VirtualBox kernel: [11609.749930] Btrfs loaded
14 May  3 08:20:33 nikita-VirtualBox kernel: [19386.834543] Syscall 'sys_getpid' hooked.
15 May  3 08:21:32 nikita-VirtualBox kernel: [19445.923432] Syscall 'sys_chdir' hooked.
16 May  3 08:22:01 nikita-VirtualBox kernel: [19473.324429] Syscall 'sys_sysinfo' hooked.

```

После запуска программы в kern.log появилась запись о попытке доступа к системному вызову sys_sysinfo.

Анализ исходного кода

Прототип системного вызова sys_sysinfo находится в include/linux/syscalls.h в исходных файлах ядра linux [3]:

```

1 // include/linux/syscalls.h
2
3 asmlinkage long sys_sysinfo(struct sysinfo __user *info);

```

Непосредственно реализация sys_sysinfo находится в kernel/sys.c в исходных файлах ядра linux [4]:

```

1 // kernel/sys.c
2
3 SYSCALL_DEFINE1(sysinfo, struct sysinfo __user *, info)
4 {
5     struct sysinfo val;
6
7     do_sysinfo(&val);
8
9     if (copy_to_user(info, &val, sizeof(struct sysinfo)))
10         return -EFAULT;
11
12     return 0;
13 }

```

Можно заметить, что вызывается функция ядра do_sysinfo в kernel/sys.c в исходных файлах ядра linux [4]:

```

1 // kernel/sys.c
2
3 static int do_sysinfo(struct sysinfo *info)
4 {
5     unsigned long mem_total, sav_total;
6     unsigned int mem_unit, bitcount;
7     struct timespec tp;
8
9     memset(info, 0, sizeof(struct sysinfo));
10
11     get_monotonic_boottime(&tp);
12     info->uptime = tp.tv_sec + (tp.tv_nsec ? 1 : 0);
13
14     get_avenrun(info->loads, 0, SI_LOAD_SHIFT - FSHIFT);
15
16     info->procs = nr_threads;
17
18     si_meminfo(info);
19     si_swapinfo(info);
20
21     /*
22      * If the sum of all the available memory (i.e. ram + swap)
23      * is less than can be stored in a 32 bit unsigned long then
24      * we can be binary compatible with 2.2.x kernels. If not,
25      * well, in that case 2.2.x was broken anyways...
26      *
27      * -Erik Andersen <andersee@debian.org>
28      */
29
30     mem_total = info->totalram + info->totalswap;
31     if (mem_total < info->totalram || mem_total < info->totalswap)
32         goto out;
33     bitcount = 0;
34     mem_unit = info->mem_unit;
35     while (mem_unit > 1) {
36         bitcount++;
37         mem_unit >>= 1;
38         sav_total = mem_total;
39         mem_total <<= 1;
40         if (mem_total < sav_total)
41             goto out;
42     }
43
44     /*
45      * If mem_total did not overflow, multiply all memory values by
46      * info->mem_unit and set it to 1. This leaves things compatible
47      * with 2.2.x, and also retains compatibility with earlier 2.4.x
48      * kernels...
49      */
50
51     info->mem_unit = 1;
52     info->totalram <<= bitcount;
53     info->freeram <<= bitcount;
54     info->sharedram <<= bitcount;
55     info->bufferram <<= bitcount;
56     info->totalswap <<= bitcount;
57     info->freeswap <<= bitcount;
58     info->totalhigh <<= bitcount;
59     info->freehigh <<= bitcount;
60
61 out:
62     return 0;
63 }

```

В первую очередь структура обнуляется, после чего собирается информация о времени, памяти, свопе и количестве процессов через соответствующие функции `get_monotonic_boottime`, `get_avenrun`, `si_meminfo`,

si_swapinfo. В завершении делается поправка на версии ядра 2.2.x с побитовым сдвигом полей результирующей структуры.

1.6 Вывод

В данной работе были изучены системные вызовы, принципы перехвата системных вызовов, а также добавление модулей в ядро ОС Linux.

В ОС Linux(версия ядра v4.17-rc4) определено около 384 системных вызовов для архитектуры i386. Такое количество системных вызовов позволяет программисту получить удобный и полный доступ ко всем компонентам ОС в рамках прав доступа.

Каждый системный вызов содержит элементы синхронизации там где это необходимо для корректной работы в условиях многопоточности.

Каждый разработчик сообщества может улучшать исходный код системных вызовов и других частей ядра путем создания pull request в GitHub репозитории ядра Linux [8]. На сегодняшний день ядро Linux разрослось до такой степени, что его изучение путем анализа кода малоэффективно.

Перехват системных вызовов достаточно несложно реализуется подменой функции внутри ядра, однако, для этого необходимы права суперпользователя.

Литература

- [1] Часть 5. Системные вызовы [Электронный ресурс], IBM. — URL: https://www.ibm.com/developerworks/ru/library/l-linux_kernel_05/index.html (дата обращения: 03.05.2018).
- [2] Kmod Hooking library, Electronic resource, GitHub. — URL: https://github.com/milabs/kmod_hooking/ (online; accessed: 03.05.2018).
- [3] Linux syscalls.h, Electronic resource, GitHub. — URL: <https://github.com/torvalds/linux/blob/master/include/linux/syscalls.h> (online; accessed: 03.05.2018).
- [4] Linux sys.c, Electronic resource, GitHub. — URL: <https://github.com/torvalds/linux/blob/master/kernel/sys.c> (online; accessed: 03.05.2018).
- [5] Linux pid.c, Electronic resource, GitHub. — URL: <https://github.com/torvalds/linux/blob/master/kernel/pid.c> (online; accessed: 03.05.2018).
- [6] PID Allocation in Linux Kernel, Electronic resource, Medium Corporation. — URL: https://medium.com/@gargi_sharma/pid-allocation-in-linux-kernel-dc0c78d14e77 (online; accessed: 03.05.2018).
- [7] Linux open.c, Electronic resource, GitHub. — URL: <https://github.com/torvalds/linux/blob/master/fs/open.c> (online; accessed: 03.05.2018).
- [8] Linux, Electronic resource, GitHub. — URL: <https://github.com/torvalds/linux> (online; accessed: 03.05.2018).