# HW #4. Interaction between Process and Kernel

## Overview

Modern operating systems such as Windows and Linux are structured into two spaces: user space and kernel space. Most of the operating system functions are implemented in the kernel. Programs in the user space have to use appropriate system calls to invoke the corresponding kernel functions. In this homework, we will take a closer look at the system call mechanism by tracing system calls made by a user process calls. We will then demonstrate how to implement a new system call on Fedora Linux. We will also demonstrate how to copy data from kernel space to user space and vice versa.

## Tasks

**A. Use 'strace' to trace the system calls made by the 'ls' command**

    1. Use 'strace'

```
$ strace ls 2>& strace.txt
```

    2. Open/Cat the output file 'strace.txt' (e.g. Figure 1)

```
1 execve("/bin/ls", ["ls", "2"], [/* 51 vars */]) = 0
2 brk(0)                                = 0x1f93000
3 access("/etc/ld.so.nohwcap", F_OK)    = -1 ENOENT (No such file or directory)
4 mmap(NULL, 8192, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_ANONYMOUS, -1, 0) = 0x7f5d351d6000
5 access("/etc/ld.so.preload", R_OK)    = -1 ENOENT (No such file or directory)
6 open("/etc/ld.so.cache", O_RDONLY)    = 3
7 fstat(3, {st_mode=S_IFREG|0644, st_size=58372, ...}) = 0
8 mmap(NULL, 58372, PROT_READ, MAP_PRIVATE, 3, 0) = 0x7f5d351c7000
9 close(3)                              = 0
10 access("/etc/ld.so.nohwcap", F_OK)   = -1 ENOENT (No such file or directory)
11 open("/lib/librt.so.1", O_RDONLY)    = 3
12 read(3, "\177ELF\2\1\1\0\0\0\0\0\0\0\0\0\3\0>\0\1\0\0\0\220!\0\0\0\0\0\0"..., 832) = 832
13 fstat(3, {st_mode=S_IFREG|0644, st_size=31744, ...}) = 0
14 mmap(NULL, 2128848, PROT_READ|PROT_EXEC, MAP_PRIVATE|MAP_DENYWRITE, 3, 0) = 0x7f5d34db1000
15 mprotect(0x7f5d34db8000, 2093056, PROT_NONE) = 0
16 mmap(0x7f5d34fb7000, 8192, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_FIXED|MAP_DENYWRITE, 3, 0
17 close(3)                             = 0
18 access("/etc/ld.so.nohwcap", F_OK)   = -1 ENOENT (No such file or directory)
19 open("/lib/libselinux.so.1", O_RDONLY) = 3
20 read(3, "\177ELF\2\1\1\0\0\0\0\0\0\0\0\0\3\0>\0\1\0\0\0\20Y\0\0\0\0\0\0"..., 832) = 832
21 fstat(3, {st_mode=S_IFREG|0644, st_size=117592, ...}) = 0
22 mmap(NULL, 2217480, PROT_READ|PROT_EXEC, MAP_PRIVATE|MAP_DENYWRITE, 3, 0) = 0x7f5d34b93000
23 mprotect(0x7f5d34baf000, 2093056, PROT_NONE) = 0
24 mmap(0x7f5d34dae000, 8192, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_FIXED|MAP_DENYWRITE, 3,
25 mmap(0x7f5d34db0000, 1544, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_FIXED|MAP_ANONYMOUS, -1,
26 close(3)                             = 0
```

**Figure 1. screenshot of strace command**

    3. You can see all the system calls made by the ls command in sequential

order. For instance, in Figure 1, we can see that the ls command has invoked the execve, brk, access, and mmap system calls

**B. Add a custom system call**

1. Download the kernel source (same steps as in Homework 2)
2. Add a custom system call to the syscall table (see Figure 2)

```
$ vim [source code directory]/arch/x86/syscall/syscall_64.tbl
```



**Figure 2. add a system call 'sayhello' to syscall table**

3. Add the system call definition to the syscall interface (see Figure 3)

```
$ vim [source code directory]/include/linux/syscalls.h
```



**Figure 3. add the system call 'sayhello' definition to the syscall interface**

4. Implement the custom system call (see Figure 4)

```
$ vim [source code directory]/kernel/sayhello.c
```



**Figure 4. the system call 'sayhello'**

5. Modify the Makefile (e.g. Figure 5)

```
$ vim [source code directory]/kernel/Makefile
```

```
5  obj-y     = fork.o exec_domain.o panic.o printk.o \
6            cpu.o exit.o itimer.o time.o softirq.o resource.o \
7            sysctl.o sysctl_binary.o capability.o ptrace.o timer.o user.o \
8            signal.o sys.o kmod.o workqueue.o pid.o \
9            rcupdate.o extable.o params.o posix-timers.o \
10           kthread.o wait.o kfifo.o sys_ni.o posix-cpu-timers.o mutex.o \
11           hrtimer.o rwsem.o nsproxy.o srcu.o semaphore.o \
12           notifier.o ksysfs.o cred.o \
13           async.o range.o groups.o \
14           sayhello.o
```

**Figure 5. modify the Makefile**

6. Make the new kernel (steps like homework 2)

● For a multi-core PC, you can accelerate the kernel make process with the '-j [number of threads]' option.

```
$ make -j 4
```

**C. Invoke system call by the system all number (see Figure 6)**

1. Include the needed libraries

```
#include <unistd.h>
#include <sys/syscall.h>
```

2. Use function 'syscall'

```
Usage: syscall(int [syscall number], [parameters to syscall])
```

```
1  #include <stdio.h>
2  #include <unistd.h>
3  #include <sys/syscall.h>
4
5  int main() {
6      int ret = syscall(312);
7      printf("ret: %d\n", ret);
8      return 0;
9  }
```

**Figure 6. call a system call in a program**

● For detailed information of syscall, please check Linux man pages

```
$ man syscall
```

3. After running the code, you can use 'dmesg' to see the messages output from printk (e.g. Figure 7)

```
$ dmesg
```



```
oshw4 [/home/ychsu] -ychsu- % dmesg | tail -n 1
[  724.729489] Hello !
```

**Figure 7. the 'printk' messages from 'sayhello' system call**

※ You can download the full source code of the examples in the section B and section C here.

**D. Copy data between user space and kernel space (Figure 8)**

1. Include the header for XXX

```
#include <linux/uaccess.h>
```

2. copy_from_user
   'copy_from_user' is used to copy user space data to a kernel space buffer
   It is defined at '[source code directory]/include/asm-generic/uaccess.h'

```
Usage: copy_from_user (void* dst, void* src, unsigned long len)
```

3. copy_to_user
   'copy_to_user' is used to copy kernel space data to a user space buffer
   It is defined at '[source code directory]/include/asm-generic/uaccess.h'

```
Usage: copy_to_user(void* dst, void* src, unsigned long len)
```

```
1  #include <asm/uaccess.h>
2  #include <linux/kernel.h>
3  #include <linux/string.h>
4
5  #define BUF_SIZE 1000
6
7  asmlinkage long sys_sayhello2(char *ptr) {
8
9      char name[BUF_SIZE];
10     char buf[BUF_SIZE];
11     unsigned long len;
12     memset(name, 0, BUF_SIZE);
13     memset(buf, 0, BUF_SIZE);
14
15     len = strlen(ptr);
16     if(copy_from_user(name, ptr, len)) {
17         return -EFAULT;
18     }
19     printk(KERN_DEBUG "syscall get name: %s(len: %lu)\n", name, len);
20
21     snprintf(buf, BUF_SIZE, "Hello, %s !", name);
22     if(copy_to_user(ptr, buf, strlen(buf))) {
23         return -EFAULT;
24     }
25     printk(KERN_DEBUG "%s\n", buf);
26
27     return 0;
28 }
```

**Figure 8. the 'sayhello2' system call uses 'copy_from_user' to get a name and uses 'copy_to_user' to store the hello message to user space buffer**

※ You can download the full source code of the examples in the section D here.

E.  **Manipulate the task_struct (program control block) of a process**

In Linux, each process has a data structure 'task_struct' to store its information(process id, process state, page table, etc.), and there is a global variable 'current' which points to the current process' task_struct'. As a result, you can get the current process's task_struct information easily through the *current* pointer.(e.g. see Figure 9 to get the current process's state)

```
9      long state = current->state;
10     printk(KERN_DEBUG "get state: %ld\n", state);
11
12     if(copy_to_user(dst, &state, sizeof(state))) {
13         return -EFAULT;
14     }
```

**Figure 9. get process state**

The 'task_struct' is defined in '[source code directory]/include/linux/sched.h' (e.g. Figure 10), you can trace the structure and know more about the process.

**Figure 10. the definition of 'task_struct' in file 'sched.h'**

※ You can download the full source code of the examples in the section E here and the other similar example source code here.

**F.** **Send a signal from kernel space to user space**

- The kernel space

1. Include the required headers

```
#include <asm/siginfo.h>
#include <linux/sched.h>
```

4. Declare and initialize a signal structure

```
// declare a signal structure
struct siginfo info;
// initialization
memset(&info, 0, sizeof(struct siginfo));
info.si_signo = SIGUSR1;
info.si_code = SI_KERNEL;
```

The variable 'si_signo' is the signal number and 'si_code' presents how to send the signal, we set it as 'SI_KERNEL' to indicate that the signal is sent from the kernel.

5. Get the 'task_struct' of a process by process id

```
struct task_struct* task;
task = find_task_by_vpid(pid);
```

The 'find_task_by_vpid' function will return the process task structure with a given process id.

6. Send a signal to the process by its 'task_struct'

```
int ret = send_sig_info(SIGUSR1, &info, task);
```

The first parameter of send_sig_info is the signal number, the second parameter is a pointer to the signal structure, and the last parameter is a pointer to the specified task structure.

- The user space program
1. Include the required header

```
#include <signal.h>
```

2. Declare and initialize a signal handler structure

```
struct sigaction sig;
sig.sa_sigaction = receiveData;
sig.sa_flags = SA_SIGINFO;
```

The 'sa_sigaction' variable is the signal handler function and 'sa_flag' is the signal flags which modify the behavior of the signal.

3. Regist the signal handler when receive the specific signal

```
sigaction(SIGUSR1, &sig, NULL);
```

The first parameter is the signal number for the signal that the program is interested in, the second is the new sigaction structure pointer, and the last is the old sigaction structure pointer.

4. Define the signal handler function(e.g. function receiveData)

```
void receiveData(int signo, siginfo_t *info) {
    // do something when receive the signal
}
```

The first parameter is the received signal number, the second is the siginfo structure from the sender, it is an optional parameter.

※ You can download the full source code of the examples in the section F here.

● For more detailed information, you can use command 'man 7 signal'

## Homework Submission

Usually, we use the command 'ps aux' to get a list of the running processes on a system are. By using the 'ps aux' command, we can see the command line strings for the running processes (e.g. the './a.out' in Figure 11). Internally, the command line string for each process is stored in a field in the process's task_struct in the kernel.



**Figure 11. program a.out is running**

The command line string is located at memory address 'arg_start' and the tail of the string is located at the address 'arg_end'. These two variables are kept in the 'mm_struct' structure (Figure 12), and the mm_struct structure appears in 'task_struct' as the variable 'mm' (Figure 13).

'task_struct' is defined at '[source code directory]/include/linux/sched.h', and 'mm_struct' is defined at '[source code directory]/include/linux/mm_types.h'.

**Figure 12. arg_start, arg_end in mm_struct**



**Figure 13. mm in task_struct**

From the section E, we know the usage of the 'current' variable. Now you can simply use the global variable 'current' and access the inner variable 'mm'. You

can then get the command line string's address, which is stored in the 'arg_start' variable After you locate the string, you can use copy_to_user function to copy the string to user space memory.

In this homework, you need to complete two tasks.

## Task 1

**You have to implement a custom system call to support two features: the first is to return the command line string of a chosen process; the second is to to modify the command line string (i.e. to change the command line string of a chosen process to your student id). You also have to implement a user space program as the front-end to demonstrate that the custom system call is working.**

At the beginning of your user space program, you may delay it and use 'ps aux' to show the original comand line string like Figure 14.

And after calling your custom system call, your user space program may print out the original command line string and you have to use 'ps aux' again to confirm that the modified command line string like Figure 15 and Figure 16, and then capture the screenshots for homework submission.

For example, you implement a system call 'sys_change_cmdline' and a user space program 'change_cmdline'.

At first, execute 'change_cmdline', and use 'ps aux' before calling the system call 'sys_change_cmdline' to show the original command line string './change_cmdline'(e.g. Figure 14)

```
oshw4 [/home/ychsu] -ychsu- % ps aux | grep change_cmdline
ychsu      953  0.0  0.0   4108    308 pts/0    S+   01:40   0:00 ./change_cmdline
ychsu      955  0.0  0.0 109400    904 pts/1    S+   01:40   0:00 grep --color change_cmdline
```

**Figure 14. before the syscall, the command line of the 'change_cmdline' is './change_cmdline'**

After the system call 'sys_change_cmdline' is called, program 'change_cmdline' will print out the original command line string is './change_cmdline' (e.g. Figure 15)

```
oshw4 [/home/ychsu] -ychsu- % ./change_cmdline
call syscall 'sys_change_cmdline'
origin cmdline: ./change_cmdline
```

**Figure 15. after the syscall, 'change_cmdline' will get and print out the original command line string is './change_cmdline'**

And then use 'ps aux' to show the current command line string is 'Yen-Chun Hsu 0056021' (e.g. Figure 16)



```
oshw4 [/home/ychsu] -ychsu- % ps aux | grep 953
ychsu      953  0.0  0.0   4112    308 pts/0    S+   01:40   0:00 Yen-Chun Hsu 0056021
ychsu      966  0.0  0.0 109400    904 pts/1    S+   01:43   0:00 grep --color 953
```

**Figure 16. after the syscall, the command line of 'change_cmdline' will become 'Yen-Chun Hsu 0056021'**

Try to implement the custom system call, and then show the 'ps aux' result to demonstrate your work. Also, please briefly describe how you implement it.

※ For convenience, you can download the user space program prototype of this task here, and the system call program prototype here.

## Task 2

**Think about why we cannot just use memcpy, strcpy, etc. function to copy data from kernel space memory to user space memory?**

Please archive your system call program, your user space program and the assignment document in PDF in an RAR file. Submit this RAR file to E3.

You can download all the example source code from github.

You can search and trace the kernel source code at lxr website.

It will take a long time to build the kernel, so you should start working on the homework early on.