

Spring 2023
Lab Assignment #4: Linux Module Programming
Assigned: May 2
Due: May 22, 11:59:59 PM

1 Introduction

This assignment will help you understand *Linux Kernel Module Programming* based on *Debug File System* interface. It has two parts, the one is tracing process tree from specific process id, and the other is finding physical address using virtual address.

Note: The purpose of this lab is to learn *Linux Kernel Programming* and understand the difference between kernel-level programming and user-level programming. By this lab, you can use kernel data structures and understand the processing mechanism of them.

2 Background

2.1 Loadable Kernel Module

Loadable Kernel Module is an object file which contains kernel code to extend the functionality of kernel. There are two programmatic ways to access kernel space. The one is *System Call* which a program requests a service from the kernel of operating system. For adding *System Call* to the system, the kernel needs to be recompiled after the new *System Call* is enrolled to the system. It spends a lot of time to compile the whole kernel. The second is *Loadable Kernel Module* which can load and unload easily without the kernel compile. The developers can easily add the interface to access the kernel using *Loadable Kernel Module*.

2.2 Debug File System

Debug File System(debugfs) is special file system available in the Linux Kernel. *Debugfs* is a simple-to-use RAM-based file system specially designed for debugging purpose. It exists as a simple way for kernel developers to make information available to user space. Because *debugfs* has no rules at all, the developers can put any information they want. *Debugfs* also supports simple user-to-kernel interfaces in *Linux Kernel Module*. So, the user space developers can access kernel information easily using *debugfs*.

2.3 Linux Kernel Module Convention

The kernel developers have to follow the convention for *Linux Kernel Module*. See the source code below:

```
1 #include <linux/module.h>
2
3 MODULE_LICENSE("GPL");
4
5 static int __init init_my_module(void)
6 {
7     // Running when this module is inserted to system
8 }
9
10 static void __exit exit_my_module(void)
11 {
12     // Running when this module is removed from system
13 }
14
15 module_init(init_my_module);
16 module_exit(exit_my_module);
```

This is a basic frame of code for *Linux Kernel Module*. There are two functions, the one is called when the kernel module is inserted to system and the other is called when the kernel module is removed from system. The two functions are enrolled to the kernel using `module_init` and `module_exit` functions. `MODULE_LICENSE` macro declares which license the module uses.

2.4 Debugfs APIs

Linux Kernel offers some APIs for developers to use debugfs easily. Before seeing the *debugfs* APIs, we have to know how to connect the functions to the file operations interfaces in *Linux Kernel*.

```
1 #include <linux/fs.h>
2
3 static int open_op(struct inode *node, struct file *fp)
4 {
5     // Running when open file operation is called
6 }
7
8 static int release_op(struct inode *node, struct file *fp)
9 {
10    // Running when close file operation is called
11 }
12
13 static ssize_t write_op(struct file *fp,
14                        const char __user *user_buffer,
15                        size_t length,
16                        loff_t *position)
```

```

17 {
18     // Running when write file operation is called
19 }
20
21 static ssize_t read_op(struct file *fp,
22                       char __user *user_buffer,
23                       size_t length,
24                       loff_t *position)
25 {
26     // Running when read file operation is called
27 }
28
29 static const struct file_operations my_fops = {
30     .open = open_op,
31     .release = release_op,
32     .write = write_op,
33     .read = read_op,
34 };

```

file_operation structure offers file operations interfaces to developer. Developer implements file operations following the convention of the functions.

Now, let's see the *debugfs* APIs. There are some *debugfs* APIs based on file operations like below:

```

struct dentry *debugfs_create_dir(const char *name, struct dentry *parent)
struct dentry *debugfs_create_file(const char *name, umode_t mode,
                                   struct dentry *parent, void *data,
                                   const struct file_operations *fops)

struct dentry *debugfs_create_u32(const char *name, umode_t mode,
                                  struct dentry *parent, u32 *value)
struct dentry *debugfs_create_u64(const char *name, umode_t mode,
                                  struct dentry *parent, u64 *value)

struct dentry *debugfs_create_x32(const char *name, umode_t mode,
                                  struct dentry *parent, u32 *value)
struct dentry *debugfs_create_x64(const char *name, umode_t mode,
                                  struct dentry *parent, u64 *value)

struct debugfs_blob_wrapper {
    void *data,
    unsigned long size;
};
struct dentry *debugfs_create_blob(const char *name, umode_t mode,
                                   struct dentry *parent,
                                   struct debugfs_blob_wrapper *blob)

```

3 Hand Out Instructions

3.1 Prepare Environment

In this lab, you will prepare your own development environment using virtual machine.

We recommend to use VirtualBox for this lab. You can get VirtualBox binaries from the VirtualBox official web site.

<https://www.virtualbox.org/>

You should use Ubuntu 20.04 LTS version for this lab.

3.2 Skeleton Code

Start by downloading the skeleton code for the Kernel Lab `kernellab-handout.tar` from eTL to your own environment. Then give the command: `tar xvf kernellab-handout.tar`. This will cause a number of files to be unpacked into the directory.

After unpack the tarball, you can see the project structure under `KernelLab-handout` like this:

```
KernelLab
├── ptree
│   ├── dbfs_ptree.c
│   └── Makefile
└── paddr
    ├── app.c
    ├── dbfs_paddr.c
    └── Makefile
```

`ptree` directory is for part 1: *Process Tree Tracing*. It has skeleton C code(`dbfs_ptree.c`) and build script(`Makefile`). `paddr` directory is for part 2: *Find Physical Address*.

It also has skeleton C code(`dbfs_paddr.c`) and build script(`Makefile`).

4 Hand In Instructions

4.1 Preparation for the Lab

You should submit the capture images of your preparation on eTL. You should capture images of your development environment with following commands are executed.

- `# lsb_release -a` will check your linux version.
- `# uname -ar` will check your kernel version.

- `KernelLab-handout/ptree/# make` and `KernelLab-handout/ptree/# make clean` will check the initial code of part #1 with your environment. If the compilation and module loading are succeed with your environment, the kernel log will give you the following message `dbfs_ptree module initialize done`. When you remove the module from the kernel, the kernel log will give you the following message `dbfs_ptree module exit`. You should capture with kernel log message using `dmesg -w`.
- `KernelLab-handout/paddr/# make` and `KernelLab-handout/paddr/# make clean` will check the initial code of part #2 with your environment. If the compilation and module loading are succeed with your environment, the kernel log will give you the following message `dbfs_paddr module initialize done`. When you remove the module from the kernel, the kernel log will give you the following message `dbfs_paddr module exit`. You should capture with kernel log message using `dmesg -w`.

4.2 Your Implementation

You should submit below files for submissions:

1. A tarball includes your implementation both part 1 and 2
2. A report

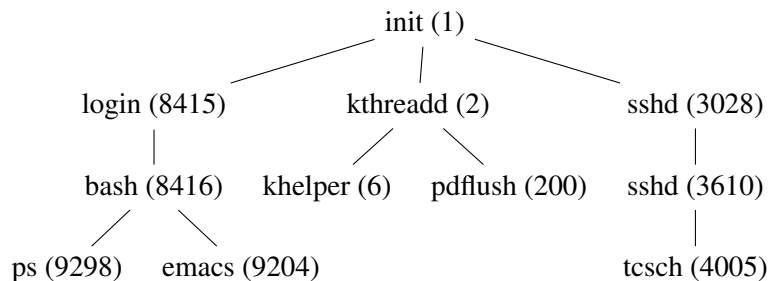
In the report, you describe the goal of Kernel Lab and how to implement for each part. In conclusion section, you explain what you learn in this lab, what was difficult and what was surprising, and so on. The project structure after submission is like below:

```
KernelLab
├── ptree
│   ├── dbfs_ptree.c
│   └── Makefile
└── paddr
    ├── app.c
    ├── dbfs_paddr.c
    └── Makefile
```

Of course, you can add more files for your implementation(e.g. header files, additional C files). However, the project should be built by the command `make` and not allowed any additional command to build the project. And never change the `app.c` and `debugfs` file name (If you change it, the grading scripts can't work correctly).

5 Part #1: Process Tree Tracing

The purpose of this Assignment is tracing process from the leaf to `init` process and logging it using `debugfs`. By this Assignment, you can understand the `task_struct` which has information about the process and manage it. In Linux System, there is process tree to manage the processes in the system. The tree has `init` process whose pid is 1, as a root node and every process except `init` has one parent process.



In user space, process can get its pid using `getpid()` function and also its parent pid using `getppid()`. But, process can't know all of ancestor processes pid. In order to get all ancestor pid, you have to access kernel space. Every process in Linux system has `task_struct` which has the whole information for the process. In kernel space, you can manage `task_struct` to get the information about the process (e.g. pid, parent process `task_struct` pointer, etc.). You can access the parent process `task_struct` recursively and finally you can access `init` process.

You should follow the rules to build and `debugfs` file name. The build command is `make` using `Makefile`. The command `make` should include `insmod` command. After `insmod` the module, you input the specific pid to `/sys/kernel/debug/ptree/input`. Then, you can see the process tree branch from input pid to `init` process using command `cat /sys/kernel/debug/ptree/ptree`. You can see the running processes using command `ps`.

```
unix> ps
  PID TTY          TIME CMD
 2881 pts/0    00:00:00 sudo
 2882 pts/0    00:00:00 su
 2885 pts/0    00:00:00 bash
 2889 pts/0    00:00:00 ps
```

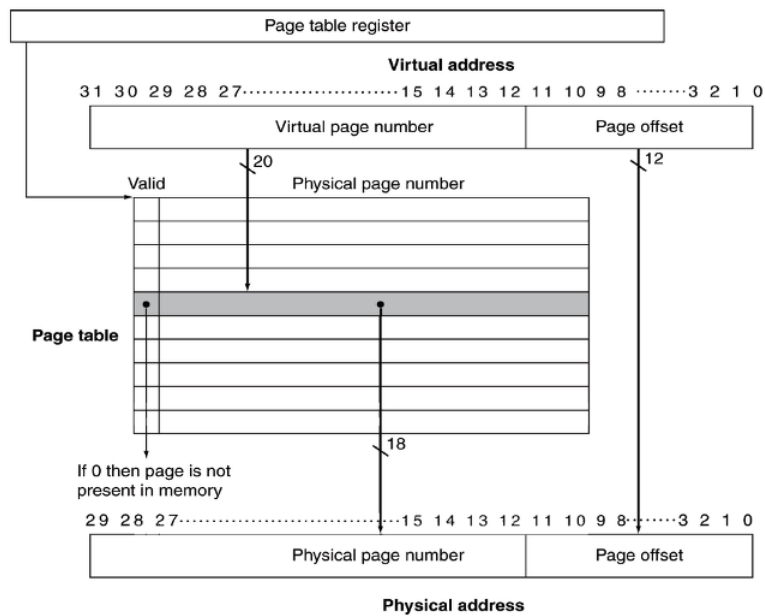
Then, you can choose one of them and write the pid to the input file. The whole stream of project running is like below:

```
unix> make
...
unix> sudo su
unix> cd /sys/kernel/debug/ptree
unix> echo 2881 >> input
unix> cat ptree
init (1)
xfce4-panel (2306)
xfce4-terminal (2408)
bash (2413)
sudo (2881)
```

The sequence of printing should be from `init` to the leaf process in the tree.

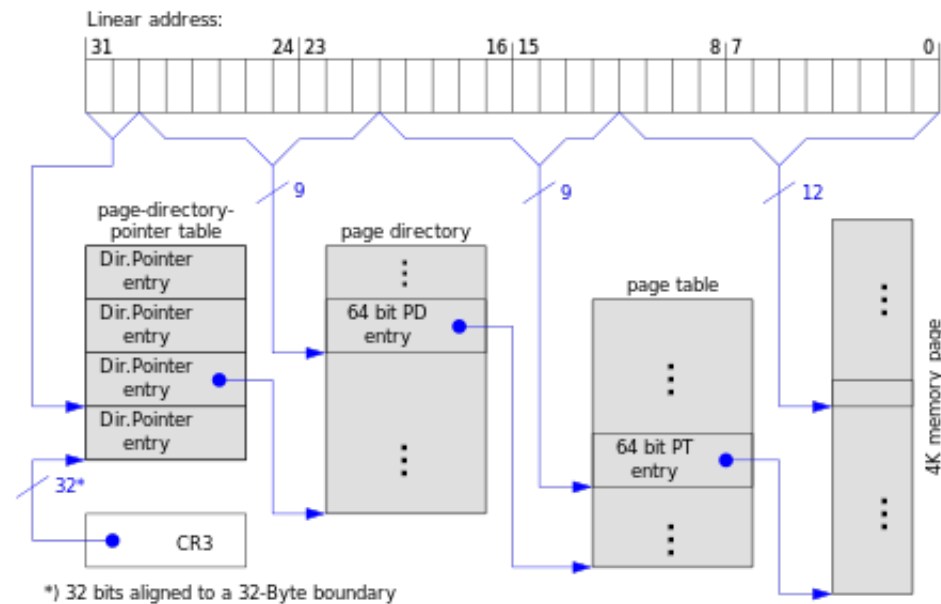
6 Part #2: Find Physical Address

Physical address is a memory address that is represented in the form of a binary number. Historically, every process in 32-bit Linux system has 4 GB memory space which is represented in the virtual address. But, this virtual address is not an actual address in physical memory. So, computer system has to translate the virtual address to physical address to access physical memory data. Every process gets the physical address from virtual address using *page table* which has the physical page number mapped to the virtual page number.



However, modern processors and Linux systems using 64-bit architecture support up to 256 TB memory space. 48 bits are required to represent the virtual address. Therefore, 36 bits are used for virtual page number in 64-bit Linux systems.

The Ubuntu OS used in this Lab has 5-level page table. By using multi-level page table, the process has more compact page table area in memory than single page table system process. Below is an example of a multi-level page table.



`task_struct` in each process has `mm_struct` to manage memory area. `mm_struct` has the pointer of the top level page table entry (`pgd`). You can obtain the pointer of next level page table entry (`p4d`) decoding `pgd` entry. Repeat this process, finally, you can obtain page table entry (`pte`) and find the physical address.

You should follow the rules to build and `debugfs` file name. The build command is make using `Makefile`. The command `make` should include `insmod` same as part 1. After `insmod` the module, you run the application `app`, then you can see the result like this:

```
unix> sudo ./app
[TEST CASE]      PASS
```

The application `app` includes some tests. First, the application makes a *virtual address* which is mapped to predefined *physical address*. Next, The kernel module you made gets the *pid* of `app` and the *virtual address*. After that, the kernel module returns the *physical address*. At last, `app` compares the return value from your kernel module and the *physical address* which is predefined. Your source code should have the process of page walk through multi-level page tables. Don't use the other way to find physical address without page walk process. And you never change the file `app.c` and `debugfs` file name in the skeleton code.