# CS353: Linux Kernel Project 3 Report

Weichen Li

5120309662

eizo.lee@sjtu.edu.cn

Department of Computer
Science and Engineering
School of Electronic Information
and Electrical Engineering
Shanghai Jiao Tong University

## CONTENTS

## I. OBJECT

- Write a module that is called *mtest*.
- When module loaded, module will create a proc fs entry */proc/mtest*.
- */proc/mtest* will accept 3 kinds of input:
  - *"listvma"* will print all vma of current process in the format of

    | start-addr | end-addr | permission |
    | --- | --- | --- |
    | 0x10000 | 0x20000 | rwx |
    | 0x30000 | 0x40000 | r– |
    | . . . | . . . | . . . |

  - *"findpage addr"* will find $va \rightarrow pa$ translation of address in current process's mm context and print it. If there is not $va \rightarrow pa$ translation, printk "translation not found".
  - *"writeval addr val"* will change an unsigned long size content in current process's virtual address into val. Note module should write to identity mapping address of addr and verify it from userspace address addr.
  - Implement a program, in which there is an integer variable *v*. The initial value of *v* is 0. The program should recognize commands such as *"write int"* and *"print"*. The former can write to the file *"/proc/mtest"*, the latter can output the value of *v*.
- All the print can be done with printk and check result with *dmesg*.

## II. COMPILE STEPS

### A. Preprocess working

At the very beginning, to save some trivial operations on authority verification, we need get the root authority first. (Need to input the key)

```
su
```

*B. Create relevant file*

To start our work, we need to create a new directory to store our work files.

```
mkdir /usr/src/Project3
cd /usr/src/Project3
```

*1) mtest.c:* Then we create a *.c* file.

```
vi mtest.c
```

Every time the "vi" command is called, the terminal will jump into the command window of VIM. We can simply input ":wq" and press "Enter" button to get back to the terminal.

Next we need to create a Makefile file:

```
vi Makefile
```

Do not quit the VIM in a hurry, we need to add some content in the Makefile file, press the "i" button on the keyboard first, then input:

```
obj-m := mtest.o

KDIR := /lib/modules/$(shell uname -r)/build

PWD := $(shell pwd)

all:

make -C $(KDIR) M=$(PWD) modules

clean:

rm *.o *.ko *.mod.c Module.symvers modules.order -f
```

Then press the "ESC" button, input ":wq", press "Enter" button, the Makefile file is built. (This is the typical way to modify and save text files using VIM, "wq" means save and quit)

Next, input:

```
vi mtest.c
```

Add content below , save and quit.

```
#include <linux/module.h>
#include <linux/kernel.h>
#include <linux/proc_fs.h>
#include <linux/string.h>
#include <linux/vmalloc.h>
#include <linux/sched.h>
#include <linux/init.h>
#include <linux/slab.h>
#include <linux/mm.h>
#include <linux/vmalloc.h>
#include <linux/highmem.h>
#include <asm/uaccess.h>
#include <linux/errno.h>
#include <linux/fs.h>

static void mtest_dump_vma_list(void)

{

struct task_struct *task = current;              //get the task_struct of the current process

struct mm_struct *mm = task->mm;

struct vm_area_struct *vma;                //get the vma area of the current process

int count = 0;      //the number of vma

down_read(&mm->mmap_sem);
```

```c
for(vma = mm->mmap; vma; vma = vma->vm_next)
{
count++;

printk("%d:  0x%lx 0x%lx ", count, vma->vm_start, vma->vm_end);

if (vma->vm_flags & VM_READ)
printk("r");
else
printk("-");

if (vma->vm_flags & VM_WRITE)
printk("w");
else
printk("-");

if (vma->vm_flags & VM_WRITE)
printk("x");
else
printk("-");

printk("\n");
}

up_read(&mm->mmap_sem);

}
static struct page *
my_follow_page(struct vm_area_struct *vma, unsigned long addr)
{
pgd_t *pgd;
pmd_t *pmd;
pud_t *pud;
pte_t *pte;

spinlock_t *ptl;
```

```c
    struct page *page = NULL;
    struct mm_struct *mm = vma->vm_mm;



    pgd = pgd_offset(mm, addr);      //get pgd
    if (pgd_none(*pgd) || unlikely(pgd_bad(*pgd)))
    goto out;


    pud = pud_offset(pgd, addr);    //get pud
    if (pud_none(*pud) || unlikely(pud_bad(*pud)))
    goto out;


    pmd = pmd_offset(pud, addr);    //get pmd
    if (pmd_none(*pmd) || unlikely(pmd_bad(*pmd)))
    goto out;


    pte = pte_offset_map_lock(mm, pmd, addr, &ptl); //get pte


    if (!pte)
    goto out;


    if (!pte_present(*pte))    //pte not in memory
    goto unlock;


    page = pfn_to_page(pte_pfn(*pte));


    if (!page)
    goto unlock;
    get_page(page);


unlock:
    pte_unmap_unlock(pte, ptl);


out:
    return page;


}
```

```c
static void mtest_find_page(unsigned long addr)

{

struct vm_area_struct *vma;

struct task_struct *task = current;

struct mm_struct *mm = task->mm;

unsigned long kernel_addr;

struct page *page;


down_read(&mm->mmap_sem);

vma = find_vma(mm, addr);

page = my_follow_page(vma, addr);


if (!page)

{

printk("translation failed.\n");

goto out;

}


kernel_addr = (unsigned long) page_address(page);


kernel_addr += (addr & ~PAGE_MASK);


printk("vma 0x%lx -> pma 0x%lx\n", addr, kernel_addr);


out:

up_read(&mm->mmap_sem);

}

static void

mtest_write_val(unsigned long addr, unsigned long val)

{

struct vm_area_struct *vma;

struct task_struct *task = current;

struct mm_struct *mm = task->mm;

struct page *page;

unsigned long kernel_addr;
```

```
down_read(&mm->mmap_sem);

vma = find_vma(mm, addr);


//test if it is a legal vma

if (vma && addr >= vma->vm_start && (addr + sizeof(val)) < vma->vm_end)

{

if (!(vma->vm_flags & VM_WRITE))   //test if we have rights to write

{

printk("cannot write to 0x%lx\n", addr);

goto out;

}


page = my_follow_page(vma, addr);

if (!page)

{

printk("page not found 0x%lx\n", addr);

goto out;

}


kernel_addr = (unsigned long) page_address(page);

kernel_addr += (addr &~ PAGE_MASK);

printk("write 0x%lx to address 0x%lx\n", val, kernel_addr);

*(unsigned long *)kernel_addr = val;

put_page(page);

}

else

{

printk("no vma found for %lx\n", addr);

}

out:

up_read(&mm->mmap_sem);

}

static ssize_t

mtest_write(struct file *file, const char __user *buffer, size_t count, loff_t *data)

{

char buf[128];

unsigned long val, val2;
```

```c
if (count > sizeof(buf))

return -EINVAL;


if (copy_from_user(buf, buffer, count))    //get the command from shell

return -EINVAL;


if (memcmp(buf, "listvma", 7) == 0)

mtest_dump_vma_list();

else if (memcmp(buf, "findpage", 8) == 0)

{

if (sscanf(buf+8, "%lx", &val) == 1)

mtest_find_page(val);

}

else if (memcmp(buf, "writeval", 8) == 0)

{

if (sscanf(buf+8, "%lx %lx", &val, &val2) == 2)

{

mtest_write_val(val, val2);

}

}

return count;

}


static struct

file_operations proc_mtest_operation = {

write: mtest_write,

};


static int __init

mtest_init(void)

{

proc_create("mtest", 0, NULL, &proc_mtest_operation);

printk("Create mtest...\n");

return 0;

}
```

```
        static void __exit

        mtest_exit(void)

        {

        remove_proc_entry("mtest", NULL);

        }


        MODULE_LICENSE("GPL");

        MODULE_DESCRIPTION("memory management task");

        module_init(mtest_init);

        module_exit(mtest_exit);
```

*2) program.c:* Then we create a *.c* file.

```
vi program.c
```

Every time the "vi" command is called, the terminal will jump into the command window of VIM. We can simply input ":wq" and press "Enter" button to get back to the terminal. To edit the file, we can press "i".

Add content below , save and quit.

```
#include<stdio.h>

int main()
{
        int n = 0;
        int i = 0;
        int v = 0;
        char s[80];
        FILE *tt;
        printf("Virtual Memory address of v: %p\n", &v);

        while(1){
                tt = fopen("/proc/mtest", "w+");
                fgets(s, 80, stdin);
                if(memcmp(s, "exit", 4) == 0) break;
                else if(memcmp(s, "print", 5) == 0)
                        printf("%d\n", v);
                else if(memcmp(s, "write", 5) == 0){
                        sscanf(s, "write %d", &i);
                        n = fprintf(tt, "writeval%lx %lx", (unsigned long)(void*)&v, (unsigned long)i);
                }
                fclose(tt);
        }
        return 0;
}
```

Then press the "ESC" button, input ":wq", press "Enter" button, the Makefile file is built. (This is the typical way to modify and save text files using VIM, "wq" means save and quit)

*C. Compile and insert the module*

So far we have completed all the preparation. Next, input:

```
        make
```

Then input:

```
        insmod mtest.ko
```

*D. Test the module and program*

*1) listvma:* Input

```
echo "listvma" > /proc/mtest
dmesg
```

Then we will get the output shown as Fig 1.



Fig. 1.   listvma

*2) findpage:* Input

```
echo "findpage 0x..." > /proc/mtest
dmesg
```
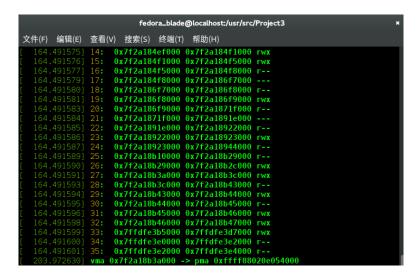
Then we will get the output shown as Fig 2.



Fig. 2.   findpage

*3) writeval:* Input

```
echo "writeval 0x... ..." > /proc/mtest
dmesg
```

Then we will get the output shown as Fig 3.

At this moment, we have successfully completed all the tasks for experiment 3, congratulations!

Fig. 3. writeval

Tips: you can use "rmmod xxx.ko" command to remove module xxx from the kernel, or use "make clean" to delete all the files produced when compiling, then only the initial files left.

*4) Program Execution:* We can execute the program using codes below:

```
./program
```

then we can test the program as shown in Fig 4.



Fig. 4. Program Execution