# CS353 Linux Kernel Project Report

Bingyu Shen

5130309283

sby2013@sjtu.edu.cn

June 3, 2016

# Contents

**6  References**                                                                         **28**

### Abstract

In the projects of this lecture, I finished four parts of the project. The first project requires to upgrade the linux kernel to the newest version; the second project part A focus on writing module program and install it in the kernel program, part B focus on process management. which requires us to add a member in the $task\_struct$ and record how many times it has been executed. The third project focus on memory management, which requires to write a module $mtest$ and can write and find certain memory address of all running process. The fourth project focus on file system, in which we can learn a lot about the romfs file system. All programs test runs on Unbuntu 15.10 and results are checked by TA. The source and modified program are attached.

# 1  Compile the Linux Kernel

## 1.1  Requirements

Upgrade your linux kernel to newest version.

## 1.2  Experiment

I installed Ubuntu 15.10 in vmware workstation, in case the operating system is broken. The version of linux kernel on Ubuntu 15.10 is 4.2.0. I download the newest version of linux kernel on $www.kernel.org$, which is 4.4.3. As you know, it is quite difficult to upgrade the linux kernel from 3.x to 4.x, so I choose the newest version of Ubuntu.

**A.** Preparation
First download and extract kernel to $/usr/src$, so at this moment the folder looks like. Go to the kernel folder. If you are not the first time to compile, you may do the following



Figure 1: Environment preparation

command.

```
1  # make clean
2  # make mrproper
```

**B.** Compile
Then before you do make command, you should do make menu command. I strongly suggest you do the make menuconfig command.

```
1  # make manuconfig
```

Then you can customize your kernel. Since I just want the version to upgrade, I just use the default setting.
Next step is to compile and install the kernel.

```
1  # make
2  # make modules_install
3  # make install
```

The $make$ step will take a long time, since I use the virtual machine, it usually takes 4-5 hours.

**C.** Update GRUB
Use the following command.

```
1  # cd /boot
2  # mkinitrd -o initrd.img-4.4.3
```
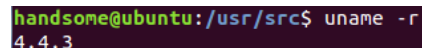
At this step, it fails to change the grub. I looked it up on the internet, and use the following command instead.

```
1  # sudo update-grub
```

It works!

## 1.3 Results

The kernel version is successfully changed.



Figure 2: Kernel version

# 2 Module Program and Process Management

## 2.1 Part A: Module Program

### 2.1.1 Requirements

Compile a kernel and run it in the system.

**A.** Module 1: Load/unload the module can output some info

**B.** Module 2: Module accepts a parameter (an integer) and load the module, output the parameters value

**C.** Module 3: Module creates a proc file, reading the proc file returns some info

### 2.1.2 Experiment

**A.** Module 1
Load and unload the module, it can print some info. We know that to write a module, you need to create the init and exit function. In the init and exit function, you can use $printk$ to print the kernel message. So the first module of program Module_1.c is as follows:

```
1  #include <linux/init.h>
2  #include <linux/module.h>
3  #include <linux/kernel.h>
4
5  static int __init hello_init(void){
6          printk("<3>Hello a linux kernel module 2A_1.\n");
7          return 0;
8  }
9
10 static void __exit hello_exit(void){
11         printk("<3>Bye 2A_1.\n");
12 }
13
14 module_init(hello_init);
15 module_exit(hello_exit);
16 MODULE_LICENSE("GPL");
```

**B.** Module 2

The module can accept parameter when insmod. This requires you to use `module_param_array`.

```
1  module_param_array(name, type, num, perm)
```

name is the parameter's name, it is usually a static variable in the module. The parameter can be set after install the module.

type is the type of parameter, it can be int, short, bool, etc.

num is a $int*$ type variable, and it stores the num of variables.

perm is the authority level of module. It is defined in $stat.h$.

Based on this information, we can write the second module, Module_2.c.

```
1  #include <linux/kernel.h>
2  #include <linux/module.h>
3  #include <linux/init.h>
4  #include <linux/moduleparam.h>
5  static int test[3];
6  int num;
7  module_param_array(test,int,&num,0644);
8  void hello_foo(void){
9          printk("Hello\n");
10 }
11 EXPORT_SYMBOL(hello_foo);
12 static int __init hello_init(void){
13         printk(KERN_INFO"Hello world\n");
14         printk(KERN_INFO"Params:test:%d,%d,%d;\n",test[0],
               test[1],test[2]);
15         return 0;
16 }
17 static void __exit hello_exit(void){
18         printk(KERN_INFO"Goodbye world\n");
19 }
```

```
20 MODULE_LICENSE("GPL");
21 MODULE_DESCRIPTION("Test");
22 module_init(hello_init);
23 module_exit(hello_exit);
```

Notice that we defined a `static int test[3]` in the 5th line. It can be used to accept parameter when we install the module.

The framework of the program is just like the Module1. Just modify the *printk* to print the parameters in `test`.

**C.** Module 3

This module is much more harder than the previous ones.

```
1  #include <linux/module.h>
2  #include <linux/proc_fs.h>
3  #include <linux/seq_file.h>
4
5  static int hello_proc_show(struct seq_file *m, void *v) {
6          seq_printf(m, "I am Bingyu Shen!\n");
7          return 0;
8  }
9  static int hello_proc_open(struct inode *inode, struct
      file *file) {
10          return single_open(file, hello_proc_show, NULL);
11 }
12 static const struct file_operations hello_proc_fops = {
13          .owner = THIS_MODULE,
14          .open = hello_proc_open,
15          .read = seq_read,
16          .llseek = seq_lseek,
17          .release = single_release,
18 };
19 static int __init hello_proc_init(void) {
20          proc_create("hello_proc", 0, NULL, &
              hello_proc_fops);
21          return 0;
22 }
23 static void __exit hello_proc_exit(void) {
24          remove_proc_entry("hello_proc", NULL);
25 }
26 MODULE_LICENSE("GPL");
27 module_init(hello_proc_init);
28 module_exit(hello_proc_exit);
```

This part requires the module to create a proc file, reading the proc file returns some info. It uses the function `proc_create` to create a proc file.

```
1  static inline struct proc_dir_entry *proc_create(const
      char *name, mode_t mode, \\
2   struct proc_dir_entry *parent, const struct
       file_operations *proc_fops);
```

name is the proc name, mode_t pis the type of umask, it is unsigned octal number. It has three bits , each bit stands for the permission of a user. First bit for owner, second bit for group, third and last bit for others.

- Octal value : Permission
- 0 : read, write and execute
- 1 : read and write
- 2 : read and execute
- 3 : read only
- 4 : write and execute
- 5 : write only
- 6 : execute only
- 7 : no permissions

Now, you can use above table to calculate file permission. For example, if umask is set to 077, the permission can be calculated as follows:

- owner: read, write and execute
- group: No permissions
- other: No permissions

struct proc_dir_entry is a proc directory. It can be created use the following function

```
1 struct proc_dir_entry *proc_mkdir(const char *name, struct
        proc_dir_entry *parent);
```

If you want the parent directory is /proc, just set parent = NULL.

file_operations is defined like

```
1 static const struct file_operations hello_proc_fops = {
2         .owner = THIS_MODULE,
3         .open = hello_proc_open,
4         .read = seq_read,
5         .llseek = seq_lseek,
6         .release = single_release,
7 };
```

Note that parameter begin with seq_ and single_ are defined in kernel.
.open operation is a function to read a specific file.
.write operation is a function to write the proc file.

In the module 3, I just use 0 as the mode so that all users can do any operation on this proc. And directory is set to NULL so that is created in /proc directory. And I define the file_operation only with open operations. And it can print some info when open the file.

The makefile is written as follows.

```
1 obj-m := module_1.o module_2.o module_3.o
2 KDIR := /lib/modules/$(shell uname -r)/build
3 PWD := $(shell pwd)
4
```
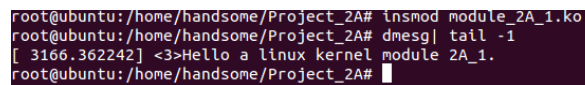
```
5 all:
6         make -C $(KDIR) M=$(PWD) modules
7 clean:
8         rm *.o *.ko *.mod.c Module.symvers modules.order -
           f
```

### 2.1.3 Results
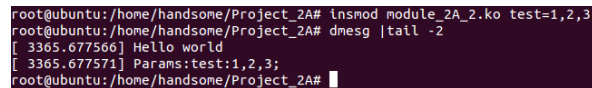
The commands used are as follows.

```
1 # insmod module_2A_1.ko
2 # dmesg
3
4 # insmod module_2A_2.ko test=1,2,3
5 # dmesg
6
7 # insmod module_2A_3.ko
8 # cat /proc/hello_proc
```
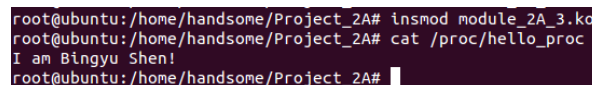
And the results are listed below.



Figure 3: Module 1 output



Figure 4: Module 2 output



Figure 5: Module 3 output

## 2.2 Part B: Process

### 2.2.1 Requirements

Process: schedule in times.

**A.** Add ctx, a new member to $task\_struct$ to record the schedule in times of the process; When a task is scheduled in to run on a cpu, increase ctx of the process;

**B.** Export ctx under /proc/XXX/ctx;

### 2.2.2 Experiment

- First modify the kernel.

  **A.** add `ctx` to `task_struct`.
  `task_struct` is in `/include/linux/sched.h`
  The program after modification is like Figure. 6.



  Figure 6: Add ctx to task_struct

  **B.** Initialize `ctx=0` in the `do_fork` function.
  The function of `do_fork` is in `kernel/fork.c.` After modification it is like Figure.7

  **C.** In kernel file `kernel/sched/core.c`, find the where processes are switched, plus `ctx` by 1.
  The program after modification is like Figure.8

  **D.** Create a new entry in the `pid_entry_tgid_base_stuff[]`. To specify a function to be executed.

  ```
  1  ONE("ctx", S_IRUSR, proc_pid_my_file),
  ```

  Modify it like Figure. 9. And the function `proc_pid_my_file` is defined in the later, like Figure. 10

- CD to the directory of the linux kernel.(`/usr/src/linux-4.4.3`).Then compile the kernel with the following command.

  ```
  1  # make
  2  # make modules_install
  3  # make install
  ```

  The time needed to execute the first command may be very long, as well as the time for the next two commands, so please be patient when the kernel is being compiled.
  After the compilation of our modified kernel, next we can test to show the ctx of each task.

```
/*
 * Do this prior waking up the new thread - the thread pointer
 * might get invalid after that point, if the thread exits quickly.
 */
if (!IS_ERR(p)) {
        struct completion vfork;
        /* bingyu shen*/
        p->ctx = 0;
        /* bingyu shen*/
        struct pid *pid;

        trace_sched_process_fork(current, p);

        pid = get_task_pid(p, PIDTYPE_PID);
        nr = pid_vnr(pid);

        if (clone_flags & CLONE_PARENT_SETTID)
                put_user(nr, parent_tidptr);

        if (clone_flags & CLONE_VFORK) {
                p->vfork_done = &vfork;
                init_completion(&vfork);
                get_task_struct(p);
        }

        wake_up_new_task(p);

        /* forking complete and child started to run, tell ptracer */
        if (unlikely(trace))
                ptrace_event_pid(trace, pid);
                                          1731,14-32    82%
```

Figure 7: Initialize ctx in do_fork function

```
        if (task_on_rq_queued(prev))
                update_rq_clock(rq);

        next = pick_next_task(rq, prev);
        clear_tsk_need_resched(prev);
        clear_preempt_need_resched();
        rq->clock_skip_update = 0;

        if (likely(prev != next)) {
                rq->nr_switches++;
                rq->curr = next;
                ++*switch_count;
                /* bingyu shen*/
                next->ctx++;
                /* bingyu shen*/

                trace_sched_switch(preempt, prev, next);
                rq = context_switch(rq, prev, next); /* unlocks the rq */
                cpu = cpu_of(rq);
        } else {
                lockdep_unpin_lock(&rq->lock);
                raw_spin_unlock_irq(&rq->lock);
        }

        balance_callback(rq);
}

static inline void sched_submit_work(struct task_struct *tsk)
{
                                          3180,18-32    36%
```

Figure 8: Increase ctx by 1

Figure 9: Add function entry



Figure 10: Print seq function

### 2.2.3 Results

To show the `ctx` of each task, we can write a small test program ctx_test.c to test it.

```c
1  #include <stdio.h>
2  int main() {
3      while (1) getchar();
4      return 0;
5  }
```

Use the following command to compile it.

```
1  # gcc ctx_test.c -o ctx_test
```

Then execute by

```
1  # ./ctx_test
```

Fisrt we should get the pid of the process. Open in a new terminal N to execute

```
1  # ps -e | grep ctx_test
```

Then the pid of ctx_test is 2761. Continue to execute

```
1  # cd /proc/2761
2  # cat
```

The result is 1, which is the ctx of process 2761. Then switch to the original terminal to input any char.
Switch to the termial N to execute

```
1  # cat
```

We find the result is 2, increased by 1. The whole process is shown in Figure.11



Figure 11: Experiment of Project B

### 2.2.4 Extention

I think we can use a module to print ctx of all running process. The module's framework is just like the module 3 in the part A. It reads all the running process's name, pid and ctx.

```
1  #include <linux/module.h>
2  #include <linux/list.h>
3  #include <linux/init.h>
4  #include <linux/sched.h>
5  #include <linux/proc_fs.h>
6  #include <linux/seq_file.h>
7  MODULE_LICENSE("Dual BSD/GPL");
8  static int ctx_proc_print(struct seq_file *m, void *v) {
9          struct task_struct *task, *p;
10         struct list_head *pos;
11         int count=0;
12
13         seq_printf(m, "test module init\n");
14
15         task=&init_task;
16
17         list_for_each(pos, &task->tasks)
18         {
19                 p=list_entry(pos, struct task_struct, tasks);
20                 count++;
21                 seq_printf(m, "%-15s\t[pid: %d]\t[ctx: %d]\n",
                           p->comm, p->pid,p->ctx);
22         }
23         seq_printf(m, "Total %d tasks\n", count);
24
25         return 0;
26 }
27
28 static int ctx_proc_open(struct inode *inode, struct file *
      file) {
29         return single_open(file, ctx_proc_print, NULL);
30 }
31
32 static const struct file_operations ctx_proc_fops = {
33         .owner = THIS_MODULE,
34         .open = ctx_proc_open,
35         .read = seq_read,
36         .llseek = seq_lseek,
37         .release = single_release,
38 };
39
40 static int test_init(void)
41 {
42         proc_create("ctx_proc",0,NULL,&ctx_proc_fops);
```

```
43          return 0;
44 }
45 static void test_exit(void)
46 {
47          remove_proc_entry("ctx_proc", NULL);
48 }
49 module_init(test_init);
50 module_exit(test_exit);
```

The makefile is like

```
1 obj-m := ctx_proc.o
2 KDIR := /lib/modules/$(shell uname -r)/build
3 PWD := $(shell pwd)
4 all:
5          make -C $(KDIR) M=$(PWD) modules
6 clean:
7          rm *.o *.ko *.mod.c Module.symvers modules.order -f
```

Then execute

```
1 # make
2 # insmod ctx_proc.ko
3 # cat /proc/ctx_proc
```

The result of the command is shown in Figure. 12.



```
polkit-gnome-au [pid: 2007]    [ctx: 412]
unity-fallback- [pid: 2014]    [ctx: 387]
vmtoolsd        [pid: 2019]    [ctx: 1963]
gconfd-2        [pid: 2031]    [ctx: 46]
gvfs-udisks2-vo [pid: 2037]    [ctx: 46]
udisksd         [pid: 2049]    [ctx: 63]
gvfs-mtp-volume [pid: 2058]    [ctx: 30]
gvfs-gphoto2-vo [pid: 2063]    [ctx: 35]
gvfs-afc-volume [pid: 2068]    [ctx: 33]
evolution-calen [pid: 2079]    [ctx: 627]
gvfsd-burn      [pid: 2097]    [ctx: 34]
gvfsd-trash     [pid: 2110]    [ctx: 88]
evolution-calen [pid: 2127]    [ctx: 145]
evolution-calen [pid: 2136]    [ctx: 144]
evolution-addre [pid: 2139]    [ctx: 74]
evolution-addre [pid: 2160]    [ctx: 137]
telepathy-indic [pid: 2191]    [ctx: 324]
mission-control [pid: 2200]    [ctx: 267]
zeitgeist-datah [pid: 2216]    [ctx: 126]
zeitgeist-daemo [pid: 2222]    [ctx: 151]
zeitgeist-fts   [pid: 2230]    [ctx: 177]
gnome-terminal- [pid: 2299]    [ctx: 6382]
gnome-pty-helpe [pid: 2305]    [ctx: 9]
bash            [pid: 2307]    [ctx: 74]
update-notifier [pid: 2351]    [ctx: 394]
su              [pid: 2364]    [ctx: 24]
bash            [pid: 2365]    [ctx: 203]
deja-dup-monito [pid: 2485]    [ctx: 290]
cat             [pid: 2754]    [ctx: 1]
Total 282 tasks
```

Figure 12: Experiment of module solution

# 3 Memory Management

## 3.1 Requirements

Write a module that is called mtest. When module loaded, module will create a proc fs entry /proc/mtest. /proc/mtest will accept 3 kind of input.

**A.** "listvma" will print all vma of current process in the format of start-addr end-addr permission
e.g
0x10000 0x20000 rwx
0x30000 0x40000 r–

**B.** "findpage addr" will find va-¿pa translation of address in current processs mm context and print it. If there is not va-¿pa translation, prink translation not found

**C.** "writeval addr val" will change an unsigned long size content in current processs virtual address into val. Note module should write to identity mapping address of addr and verify it from userspace address addr.

All the print can be done with printk and check result with dmesg.

## 3.2 Background Knowledge

**A.** Virtual Memory Areas(VMA) and Process
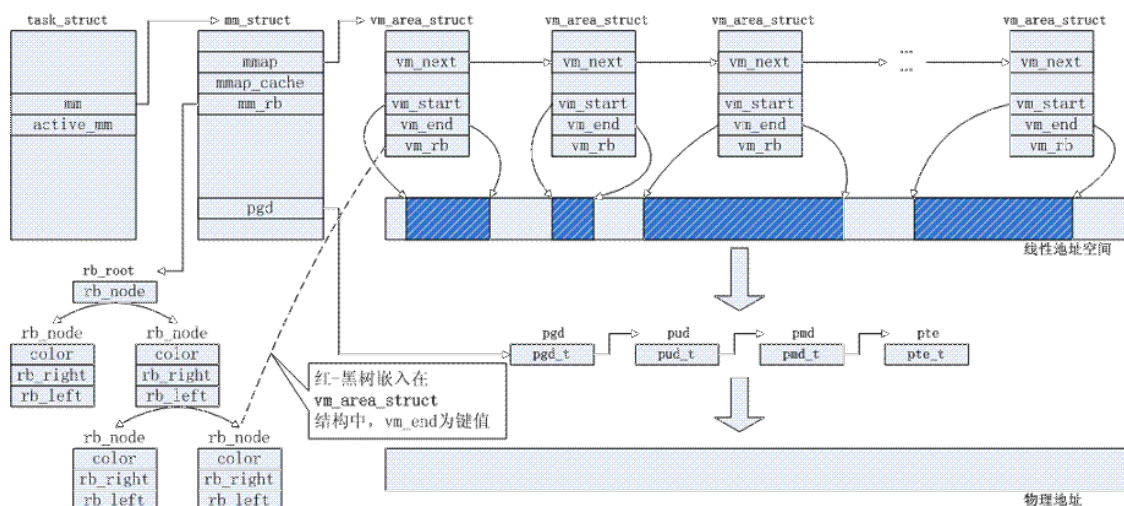


Figure 13: VMA illustration

struct task_struct is the process descriptor in Linux, the mm_struct field records the VMA of the process. It records the memory related parameters. The prototype of mm_struct is as follows(only list the related codes).

```
1  struct mm_struct {
2      ... ...
3      struct vm_area_struct * mmap;/* list of VMAs */
4      struct rb_root mm_rb;
```

```
5       unsigned long mmap_base;/* base of mmap area */
6       unsigned long task_size;/* size of task vm space */
7
8       pgd_t * pgd;
9       /* How many users with user space? */
10      atomic_t mm_users;
11      /* How many references to "struct mm_struct" (users
            count as 1) */
12      atomic_t mm_count;
13      int map_count;        /* number of VMAs */
14      struct rw_semaphore mmap_sem;
15      /* Protects page tables and some counters */
16      spinlock_t page_table_lock;
17      ...
18  };
```

Note that mmap_sem is semaphore to lock the mmap, so that if won't be written in two processes at the same time.

`struct vm_area_struct` is the description of virtual memory area. And use mmap to find all VMAs of this process.

```
1  struct vm_area_struct {
2      /* VM area parameters */
3      struct mm_struct * vm_mm;
4      unsigned long vm_start;
5      unsigned long vm_end;
6
7      unsigned short vm_flags;
8  /* linked list of VM areas per task, sorted by address */
9      struct vm_area_struct * vm_next;
10     ... ...
11  };
```

In which the vm_flags is defined as short, the types we used are as follows.

```
1   /* currently active flags */
2  #define VM_READ          0x0001
3  #define VM_WRITE         0x0002
4  #define VM_EXEC          0x0004
5  #define VM_SHARED        0x0008
```

The structure of VMA is illustrated in Figure. 13.

**B.** Linear Address of Page

The linear page addressing mode is illustrated in Figure. 14.

So from `mm_struct->pgd` we can get the page global directory, which is a physical page frame. This frame contains an array of type pgd_t which is an architecture specific type defined in $< asm/page.h >$. Each active entry in the PGD table points to a page frame containing an array of Page Middle Directory (PMD) entries of type pmd_t which in turn points to page frames containing Page Table Entries (PTE) of type pte_t, which finally points to page frames containing the actual user data.
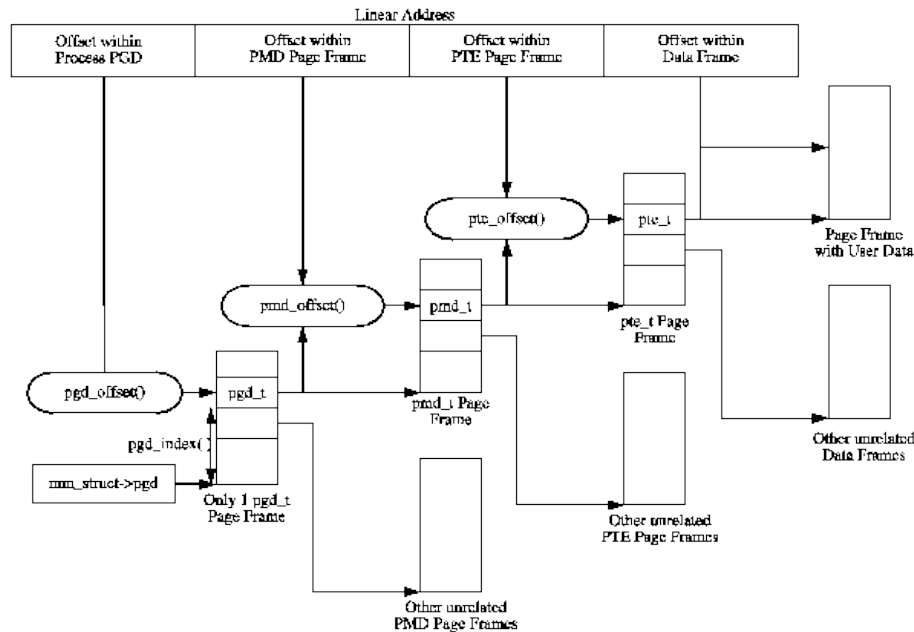
Figure 14: Page Linear Address

The code of get a page address is referenced in $< mm/memory.c >$, the function `follow_page()`. I modify it to make it simpler.

**C.** VMA to PMA

From virtual memory address to physical memory address, we can find the VMA of a given address first. Then get the specific page with the VMA.

First find the page address of the page. Then add the offset to the page address.

```
1 kernel_addr = (unsigned long) page_address(page);
2 kernel_addr += (addr & ~PAGE_MASK);
```

## 3.3 Experiment

The code of mtest is as follows.

```
1  #include <linux/module.h>
2  #include <linux/kernel.h>
3  #include <linux/proc_fs.h>
4  #include <linux/string.h>
5  #include <linux/vmalloc.h>
6  #include <linux/sched.h>
7  #include <linux/init.h>
8  #include <linux/slab.h>
9  #include <linux/mm.h>
10 #include <linux/vmalloc.h>
11 #include <linux/highmem.h>
12 #include <asm/uaccess.h>
13 #include <linux/errno.h>
```

```
14  #include <linux/fs.h>

15

16  static void mtest_dump_vma_list(void)
17  {
18          //get the task_struct of the current process
19      struct task_struct *task = current;
20      struct mm_struct *mm = task->mm;
21          //get the vma area of the current process
22      struct vm_area_struct *vma;

23

24      int count = 0;      //the number of vma

25

26      down_read(&mm->mmap_sem);

27

28      for(vma = mm->mmap; vma; vma = vma->vm_next)
29      {
30          count++;
31          printk("%d:  0x%lx 0x%lx ", count, vma->vm_start, vma
                ->vm_end);
32          if (vma->vm_flags & VM_READ)
33              printk("r");
34          else
35              printk("-");

36

37          if (vma->vm_flags & VM_WRITE)
38              printk("w");
39          else
40              printk("-");

41

42          if (vma->vm_flags & VM_EXEC)
43              printk("x");
44          else
45              printk("-");

46

47          printk("\n");

48

49      }
50      up_read(&mm->mmap_sem);
51  }

52

53  static struct page *
54  my_follow_page(struct vm_area_struct *vma, unsigned long addr)
55  {

56

57      pgd_t *pgd;
58      pmd_t *pmd;
59      pud_t *pud;
```

```
60      pte_t *pte;

62      spinlock_t *ptl;

64      struct page *page = NULL;
65      struct mm_struct *mm = vma->vm_mm;

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

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

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

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

83      if (!pte)
84          goto out;

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

89      page = pfn_to_page(pte_pfn(*pte));

91      if (!page)
92          goto unlock;
93      get_page(page);

95  unlock:
96      pte_unmap_unlock(pte, ptl);
97  out:
98      return page;
99  }


102 static void mtest_find_page(unsigned long addr)
103 {
104     struct vm_area_struct *vma;
105     struct task_struct *task = current;
```

```
106        struct mm_struct *mm = task->mm;
107        unsigned long kernel_addr;
108        struct page *page;
109
110        down_read(&mm->mmap_sem);
111
112        vma = find_vma(mm, addr);
113        page = my_follow_page(vma, addr);
114
115        if (!page)
116        {
117            printk("translation failed.\n");
118            goto out;
119        }
120
121        kernel_addr = (unsigned long) page_address(page);
122        kernel_addr += (addr & ~PAGE_MASK);
123        printk("vma 0x%lx -> pma 0x%lx\n", addr, kernel_addr);
124    out:
125        up_read(&mm->mmap_sem);
126    }
127
128    static void
129    mtest_write_val(unsigned long addr, unsigned long val)
130    {
131        struct vm_area_struct *vma;
132        struct task_struct *task = current;
133        struct mm_struct *mm = task->mm;
134        struct page *page;
135        unsigned long kernel_addr;
136
137        down_read(&mm->mmap_sem);
138        vma = find_vma(mm, addr);
139        //test if it is a legal vma
140        if (vma && addr >= vma->vm_start && (addr + sizeof(val)) <
               vma->vm_end)
141        {
142            if (!(vma->vm_flags & VM_WRITE))   //test if we have
                   rights to write
143            {
144                printk("cannot write to 0x%lx\n", addr);
145                goto out;
146            }
147            page = my_follow_page(vma, addr);
148            if (!page)
149            {
150                printk("page not found 0x%lx\n", addr);
```

```
151              goto out;
152          }
153
154          kernel_addr = (unsigned long) page_address(page);
155          kernel_addr += (addr &~ PAGE_MASK);
156          printk("write 0x%lx to address 0x%lx\n", val,
                  kernel_addr);
157          *(unsigned long *)kernel_addr = val;
158          put_page(page);
159      }
160      else
161      {
162          printk("no vma found for %lx\n", addr);
163          }
164
165      out:
166          up_read(&mm->mmap_sem);
167 }
168
169 static ssize_t
170 mtest_write(struct file *file, const char __user *buffer,
        size_t count, loff_t *data)
171 {
172      char buf[128];
173      unsigned long val, val2;
174      if (count > sizeof(buf))
175          return -EINVAL;
176          //get the command from shell
177      if (copy_from_user(buf, buffer, count))
178          return -EINVAL;
179
180      if (memcmp(buf, "listvma", 7) == 0)
181          mtest_dump_vma_list();
182      else if (memcmp(buf, "findpage", 8) == 0)
183      {
184          if (sscanf(buf+8, "%lx", &val) == 1)
185              mtest_find_page(val);
186      }
187      else if (memcmp(buf, "writeval", 8) == 0)
188      {
189          if (sscanf(buf+8, "%lx %lx", &val, &val2) == 2)
190          {
191              mtest_write_val(val, val2);
192          }
193      }
194
195          return count;
```

```
196  }
197
198
199
200  static struct file_operations proc_mtest_operation = {
201      write: mtest_write,
202  };
203
204
205
206  static int __init mtest_init(void)
207  {
208      proc_create("mtest", 0, NULL, &proc_mtest_operation);
209      printk("Create mtest...\n");
210      return 0;
211  }
212
213
214
215  static void __exit mtest_exit(void)
216  {
217      remove_proc_entry("mtest", NULL);
218  }
219
220
221  MODULE_LICENSE("GPL");
222  MODULE_DESCRIPTION("memory management task");
223  module_init(mtest_init);
224  module_exit(mtest_exit);
```

- write a Makefile

```
1  obj-m := mtest.o
2  KDIR := /lib/modules/$(shell uname -r)/build
3
4  PWD := $(shell pwd)
5  default:
6          $(MAKE) -C $(KDIR) M=$(PWD) modules
7  clean:
8          rm -rf *.o *.ko *.mod.c
```

- type "make" in shell

- type "sudo insmod mtest.ko"

- type "sudo su"

- type: echo "listvma" ¿ /proc/mtest

- type: dmesg (then you will find a lot of vma)

- choose one of them and type : echo "findpage 0x........." ¿ /proc/mtest,
  then type : dmesg
  note: you may find that "translation failed" shows up. But that does not mean you fail
  the test. Please choose an address between the start and the end of the vma listed. e.g.
  0x123 C 0x345, you may want to type: echo findpage 0x300 > /proc/mtest, because you
  cannot be sure if the beginning of the virtual address is used.

- choose an address that you has rights to write.
  Then type: echo writeval 0x... 123 ¿ /proc/mtest
  type: dmesg
  note: 123 can be any unsigned int
  you'd better choose the same address to test in both step 7 and 8, so that you will see
  if the physical address is consistent.



Figure 15: listvma experiment



Figure 16: findpage experiment

Figure 17: listvma experiment

# 4 File System

## 4.1 Requirements

**A.** Source

- Inode.c/Makefile (kernel source of romfs)
- Test.img (a romfs image, you can mount it to a dir with mount Co loop test.img xxx)
- Say test.img is mounted in t, find t output
  - aa
  - bb
  - ft
  - fo
  - fo/aa

**B.** Practice 1

- Change romfs code to hide a file/dir with special name
- Test & result
  - insmod romfs hided_file_name= aa
  - Mount Co loop test.img t
  - then ls t, ls t/fo, no "aa" and "fo/aa". found
  - ls t/aa, or ls fo/aa, no found
  - Without the code change, above two operations can find file aa

**C.** Practice 2

- change the code of romfs to correctly read info of an encrypted romfs
- Test & result
  - insmod romfs encrypted_file_name=bb
  - Mount Co loop test.img t
  - Say bbs original content is bbbbbbb
  - With the change, cat t/bb output cccccccccc

**D.** Practice 3

- change the code of romfs to add x (execution) bit for a specific file
- Test & result
  - insmod romfs exec_file_name=cc

- Mount Co loop test.img t
- Without code changes ls -l t, output is -rw-r–r–
- With the change, output is -rwxr-xr-x
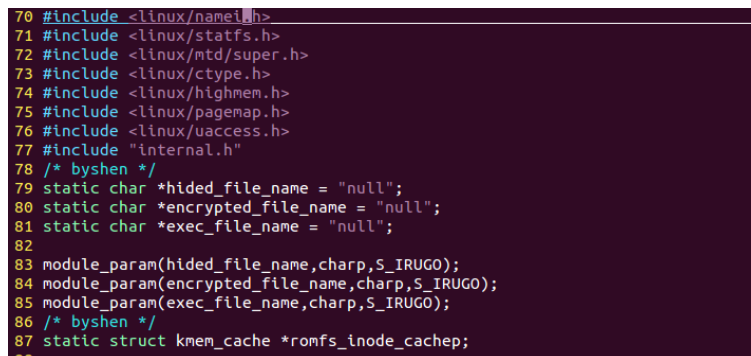
## 4.2  Experiment

- Create and modofy relevant files

```
1  # mkdir Project4
2  # cp -r /usr/src/linux-xxxx/fs/romfs /Project4
3  # cd romfs
```

- Modify Makefile.

```
1  obj-m := romfs.o
2  romfs-y := super.o storage.o
3  KDIR := /lib/modules/$(shell uname -r)/build
4  EXTRA_FLAGS := -I(PWD)
5  PWD := $(shell pwd)
6  all:
7          make -C $(KDIR) M=$(PWD) modules
8  clean:
9          rm *.o *.ko *.mod.c Module.symvers modules.order -
              f
```

- Modify super.c. There are mainly five modifications.

```
70 #include <linux/namei.h>
71 #include <linux/statfs.h>
72 #include <linux/mtd/super.h>
73 #include <linux/ctype.h>
74 #include <linux/highmem.h>
75 #include <linux/pagemap.h>
76 #include <linux/uaccess.h>
77 #include "internal.h"
78 /* byshen */
79 static char *hided_file_name = "null";
80 static char *encrypted_file_name = "null";
81 static char *exec_file_name = "null";
82
83 module_param(hided_file_name,charp,S_IRUGO);
84 module_param(encrypted_file_name,charp,S_IRUGO);
85 module_param(exec_file_name,charp,S_IRUGO);
86 /* byshen */
87 static struct kmem_cache *romfs_inode_cachep;
88
```

Figure 18: Modification 1

- create test.img

```
1  # mkdir test
2  # genromfs -V "vromfs" -f test.img -d test
```

Then add a directory $t$ to mount the image file.

```
1  # mkdir t
```

- compile the program.

```
1  # make
```

```
109 static int romfs_readpage(struct file *file, struct page *page)
110 {
111     struct inode *inode = page->mapping->host;
112     loff_t offset, size;
113     unsigned long fillsize, pos;
114     void *buf;
115     int ret;
116     /* byshen */
117     int i;
118     char *fname = file->f_path.dentry->d_iname;
119     /* byshen */
120
121     buf = kmap(page);
122     if (!buf)
123         return -ENOMEM;
124
125     /* 32 bit warning -- but not for us :) */
126     offset = page_offset(page);
127     size = i_size_read(inode);
128     fillsize = 0;
129     ret = 0;
130     if (offset < size) {
131         size -= offset;
132         fillsize = size > PAGE_SIZE ? PAGE_SIZE : size;
133
134         pos = ROMFS_I(inode)->i_dataoffset + offset;
135
                                                    106,2        15%
```

Figure 19: Modification 2

```
133
134         pos = ROMFS_I(inode)->i_dataoffset + offset;
135
136         ret = romfs_dev_read(inode->i_sb, pos, buf, fillsize);
137         if (ret < 0) {
138             SetPageError(page);
139             fillsize = 0;
140             ret = -EIO;
141         }
142     }
143     /* byshen */
144     if (strcmp(fname, encrypted_file_name)==0) {
145         for (i = 0; i < fillsize; ++i)
146             *((char *)buf + i) = 'c';
147     }
148     /* byshen */
149     if (fillsize < PAGE_SIZE)
150         memset(buf + fillsize, 0, PAGE_SIZE - fillsize);
151     if (ret == 0)
152         SetPageUptodate(page);
153
154     flush_dcache_page(page);
155     kunmap(page);
156     unlock_page(page);
157     return ret;
158 }
159
160 static const struct address_space_operations romfs_aops = {
                                                    155,14-17    19%
```

Figure 20: Modification 3

```
201         j = romfs_dev_strnlen(i->i_sb, offset + ROMFH_SIZE,
202                     sizeof(fsname) - 1);
203         if (j < 0)
204             goto out;
205
206         ret = romfs_dev_read(i->i_sb, offset + ROMFH_SIZE, fsname, j);
207         if (ret < 0)
208             goto out;
209         fsname[j] = '\0';
210
211         ino = offset;
212         nextfh = be32_to_cpu(ri.next);
213         /* byshen */
214         if (strcmp(fsname, hided_file_name)==0) {
215             offset = nextfh & ROMFH_MASK;
216             continue;
217         }
218         /* byshen */
219         if ((nextfh & ROMFH_TYPE) == ROMFH_HRD)
220             ino = be32_to_cpu(ri.spec);
221         if (!dir_emit(ctx, fsname, j, ino,
222                 romfs_dtype_table[nextfh & ROMFH_TYPE]))
223             goto out;
224
225         offset = nextfh & ROMFH_MASK;
226     }
227 out:
228     return 0;
                                            213,6-12     30%
```

Figure 21: Modification 4

```
281
282     /* byshen */
283     if (strcmp(name, exec_file_name)==0) {
284         inode->i_mode |= S_IXUGO;
285     }
286     /* byshen */
287
288     if (IS_ERR(inode)) {
289         ret = PTR_ERR(inode);
290         goto error;
291     }
292     goto outi;
293
294     /*
295      * it's a bit funky, _lookup needs to return an error code
296      * (negative) or a NULL, both as a dentry.  ENOENT should not
297      * be returned, instead we need to create a negative dentry by
                                            282,5-8     40%
```

Figure 22: Modification 5

## 4.3 Results

Before modification.



Figure 23: Practice 1-3 before modification

I write a run.sh to simplify the test.

```
1  echo "--- Practice 1 ---"
2  insmod romfs.ko hided_file_name=aa
3  mount -o loop test.img t
4  find t
5  umount test.img
6  rmmod romfs.ko
7  echo "--- Practice 2 ---"
8  insmod romfs.ko encrypted_file_name=bb
9  mount -o loop test.img t
10 cat t/bb
11 umount test.img
12 rmmod romfs.ko
13 echo " "
14 echo "--- Practice 3 ---"
15 insmod romfs.ko exec_file_name=cc
16 mount -o loop test.img t
17 ls -l t
18 umount test.img
19 rmmod romfs.ko
```

And the test result after modification is correspondingly as Figure.24.

# 5 Conclusion

In the projects of this lecture, I finished four parts of the project. The first project requires to upgrade the linux kernel to the newest version; the second project part A focus on writing

Figure 24: Practice 1-3 after modification

module program and install it in the kernel program, part B focus on process management. which requires us to add a member in the task struct and record how many times it has been executed. The third project focus on memory management, which requires to write a module mtestand can write andfind certain memory address of all running process. Thefourth project focus on file system, in which we can learn a lot about the romfs file system. All programs test runs on Unbuntu 15.10 and results are checked by TA. The source and modified program are attached.

The experiment promote my understanding of the linux proc, memory and file system, and I gain a lot of experiments in reading the source code and kernel programming.

Finally, I appreciate the financial support from Prof. Wu and TA that helped the project.

# 6    References

[1] Linux kernel experiment handbook of SJTU CS353
    [2] CS353 Course materials