

乌有理工大学

《操作系统》课程实验报告

实验题目	实验三：内存管理		
姓 名	Nemo	学 号	2019114514
组 别	—	合 作 者	—
班 级	计科 NaN 班	指导教师	Outis

1 实验概述

1.1 实验介绍

本实验通过在内核态分配内存的任务操作，让学生们了解并掌握操作系统中内存管理的布局，内核态内存分配的实现，了解 openEuler 内核的物理内存管理，熟悉 `struct page` 数据结构。了解 openEuler 内核的物理内存管理，熟悉进程的地址空间的概念以及相关的数据结构。

1.2 任务描述

- 使用 `kmalloc` 分配 1KB、8KB 的内存，打印指针地址。
- 使用 `vmalloc` 分配 8KB、1MB、64MB 的内存，打印指针地址。
- 查看已分配的内存，根据机器是 32 位或 64 位的情况，分析地址落在的区域。
- 正确编写并安装设备驱动程序 `pfstat.ko`，用于读取物理页框属性。
- 正确编写并安装设备驱动程序 `vma.ko`，用于获得每个 VMA 的信息。

1.3 实验目的

- 正确编写满足功能的源文件，正确编译。
- 正常加载、卸载内核模块，且内核模块功能满足任务所述。
- 了解操作系统的内存管理。

2 实验内容

2.1 内核内存分配

2.1.1 实验过程

用 `make` 命令编译 `kmalloc.c` 和 `vmalloc.c`，用 `insmod` 命令加载编译完成的内核模块，用 `rmod` 命令卸载内核模块，用 `dmesg` 命令查看内核模块在内核态的运行结果。

实验结果如下所示 (openEuler 20.03):

```
1 [root@ostest lab3]# cd task1
2 [root@ostest task1]# make
3 make -C /root/kernel M=/root/lab3/task1 modules
4 make[1]: Entering directory '/root/kernel'
5   CC [M]  /root/lab3/task1/kmalloc.o
6   Building modules, stage 2.
7   MODPOST 1 modules
8   CC      /root/lab3/task1/kmalloc.mod.o
9   LD [M]  /root/lab3/task1/kmalloc.ko
10 make[1]: Leaving directory '/root/kernel'
11 [root@ostest task1]# insmod kmalloc.ko
12 [root@ostest task1]# rmmod kmalloc.ko
13 [root@ostest task1]# dmesg | tail -n4
14 [ 208.853053] Start kmalloc!
15 [ 208.853241] kmallocmem1 addr = ffff80016fa4f000
16 [ 208.853543] kmallocmem2 addr = ffff80016fcdc000
17 [ 230.338834] Exit kmalloc!
18 [root@ostest task1]# cd ../task2
19 [root@ostest task2]# make
20 make -C /root/kernel M=/root/lab3/task2 modules
21 make[1]: Entering directory '/root/kernel'
22   CC [M]  /root/lab3/task2/vmalloc.o
23   Building modules, stage 2.
24   MODPOST 1 modules
25   CC      /root/lab3/task2/vmalloc.mod.o
26   LD [M]  /root/lab3/task2/vmalloc.ko
27 make[1]: Leaving directory '/root/kernel'
28 [root@ostest task2]# insmod vmalloc.ko
29 [root@ostest task2]# rmmod vmalloc.ko
30 [root@ostest task2]# dmesg | tail -n5
31 [ 250.533425] Start vmalloc!
32 [ 250.533651] vmallocmem1 addr = ffff00000b940000
33 [ 250.534012] vmallocmem2 addr = ffff0000204c0000
34 [ 250.534508] vmallocmem3 addr = ffff000023f00000
35 [ 261.096602] Exit vmalloc!
```

使用 vi kernel/arch/arm64/configs/openeuler_defconfig 命令, 结果如下:

```
1 #
2 # Automatically generated file; DO NOT EDIT.
3 # Linux/arm64 4.19.93 Kernel Configuration
4 #
5 ...
6 CONFIG_ARM64=y
7 CONFIG_64BIT=y
8 CONFIG_MMU=y
9 CONFIG_ARM64_PAGE_SHIFT=16
10 ...
11 # CONFIG_ARM64_4K_PAGES is not set
12 # CONFIG_ARM64_16K_PAGES is not set
13 CONFIG_ARM64_64K_PAGES=y
14 # CONFIG_ARM64_VA_BITS_42 is not set
15 CONFIG_ARM64_VA_BITS_48=y
16 CONFIG_ARM64_VA_BITS=48
17 ...
18 CONFIG_PGTABLE_LEVELS=3
19 ...
```

可知目前的配置是虚拟地址位数为 48 位, 页表的大小是 64K, 页表转化是 3 级。
kmalloc 和 vmalloc 分配的内存地址, 都位于内核空间。

2.1.2 思考题

内存泄露是指程序中已动态分配的内存未释放或无法释放；内存溢出是指程序在申请内存时，没有足够的内存空间供其使用；内存越界是指程序向系统申请一块内存后，使用时超出申请范围。

分析：下面两个程序是否会产生内存泄露、内存溢出或内存越界？

```
1 void function_which_allocates() { float *a = malloc(sizeof(float) * 45); }
2
3 int main() { function_which_allocates(); }
```

在 `function_which_allocates` 函数中的指针 `a` 在离开函数作用域后就销毁了，所以 `malloc` 给它的大小为 `sizeof(float) * 45 = 180 B` 的内存泄露了。如欲不泄露内存，需要在函数结束之前对指针进行 `free` 处理。

```
1 void func(char* input) {
2     char buffer[16];
3     strcpy(buffer, input);
4 }
5
6 int main() {
7     char longstring[256];
8     int i;
9     for (i = 0; i < 255; i++) {
10         longstring[i] = 'B';
11     }
12     func(longstring);
13 }
```

使用字符串复制函数 `strcpy` 需要考虑来源字符串和目标字符串的长度。`strcpy` 在 Linux 0.11 中的实现方式如下：

```
1 extern inline char *strcpy(char *dest, const char *src) {
2     __asm__(
3         "cld\n"
4         "1:\tlodsb\n\t"
5         "stosb\n\t"
6         "testb %%al,%%al\n\t"
7         "jne 1b"
8         :: "S"(src), "D"(dest));
9     return dest;
10 }
```

观察内嵌的汇编代码可知，如果 `src` 的长度大于 `dest` 的长度，当复制到目的缓冲器时，它会改写到连接目的缓冲器后方的存储器，导致无法预期的结果。这里 `src` 的长度是 256 而 `dest` 长度是 16，所以会产生内存越界。本程序在本机运行的结果是 `illegal hardware instruction`，在有些计算机上结果会是 `segmentation fault`。

2.2 物理页框属性

本实验期望编写一个内核模块，统计当前系统中的页情况，通过设备文件的 open 接口返回给用户态并打印。

本实验使用 `get_num_physpages()` (`linux/mm.h`) 来获取物理页数目,使用 `pfn_to_page` (`asm-generic/memory_model.h`) 宏来获得当前页的 `struct page` 指针。

查阅资料可知, `struct page` 拥有一个 32 位的 `flags` 域,用来存放页的状态。`flag` 的每一位单独表示一种状态,所以总共可以有 32 种不同的状态,这些状态在 `linux/page-flags.h` 中定义。`linux/page-flags.h` 中还定义了众多内联函数,用于测试页是否有某种状态,这些函数的原型是 `static __always_inline int Page##uname(struct page *page)`,如果测试页有状态 `uname` 则返回 1, 否则返回 0。

核心代码和实验结果如下所示 (openEuler 20.03)。

```
1 static ssize_t pfstat_read(struct file *file, char __user *out, size_t size,
2                           loff_t *off) {
3     struct page_info *buf = file->private_data;
4     memset(buf, 0, sizeof(struct page_info));
5
6     buf->num_physpages = get_num_physpages();
7     printk(KERN_INFO "OK: %lu phys page(s) found. \n", buf->num_physpages);
8
9     for (int cur_pfn = 0; cur_pfn < buf->num_physpages; cur_pfn++) {
10         struct page *hoc = pfn_to_page(cur_pfn + ARCH_PFN_OFFSET);
11         buf->valid += pfn_valid(cur_pfn + ARCH_PFN_OFFSET);
12         buf->free += hoc->_refcount.counter == 0;
13         buf->locked += PageLocked(hoc);
14         buf->reserved += PageReserved(hoc);
15         buf->swapcache += PageSwapCache(hoc);
16         buf->referenced += PageReferenced(hoc);
17         buf->slab += PageSlab(hoc);
18         buf->private += PagePrivate(hoc);
19         buf->uptodate += PageUptodate(hoc);
20         buf->dirty += PageDirty(hoc);
21         buf->active += PageActive(hoc);
22         buf->writeback += PageWriteback(hoc);
23         buf->mappedtodisk += PageMappedToDisk(hoc);
24     }
25     printk(KERN_INFO "OK: flags unwrapped in %lu page(s).\n", buf->num_physpages);
26
27     copy_to_user(out, buf, sizeof(struct page_info));
28     printk(KERN_INFO "OK: copy_to_user done.\n");
29     return size;
30 }
```

```
1 [root@ostest pfstat]# ./run.sh
2 make -C /lib/modules/4.19.188/build M=/root/lab3/pfstat modules
3 make[1]: Entering directory '/root/kernel'
4   Building modules, stage 2.
5   MODPOST 1 modules
6 make[1]: Leaving directory '/root/kernel'
7 131072 physpages checked, total memory = 8192 MB
8 130978 valid pages, valid memory = 8186 MB
9 free: 96698
10 ...↓
```

```

11 ...↑
12 locked: 0
13 reserved: 21989
14 swapcache: 0
15 referenced: 4624
16 slab: 1123
17 private: 888
18 uptodate: 9490
19 dirty: 255
20 active: 6500
21 writeback: 0
22 mappedtodisk: 4614

```

使用 vi /proc/meminfo 命令, 结果如下:

```

1 MemTotal:          6975296 kB
2 MemFree:           6177152 kB
3 MemAvailable:      6348672 kB
4 Buffers:           42752 kB
5 Cached:            453312 kB
6 SwapCached:        0 kB
7 Active:            416128 kB
8 Inactive:          230080 kB
9 Active(anon):       155264 kB
10 Inactive(anon):     7936 kB
11 Active(file):       260864 kB
12 Inactive(file):     222144 kB
13 Unevictable:        0 kB
14 Mlocked:           0 kB
15 SwapTotal:         0 kB
16 SwapFree:          0 kB
17 Dirty:             640 kB
18 Writeback:         0 kB
19 AnonPages:         150528 kB
20 Mapped:            104320 kB
21 Shmem:             13120 kB
22 KReclaimable:       28160 kB
23 Slab:              72000 kB
24 SReclaimable:       28160 kB
25 SUnreclaim:        43840 kB
26 KernelStack:       10048 kB
27 PageTables:        7936 kB
28 NFS_Unstable:      0 kB
29 Bounce:            0 kB
30 WritebackTmp:      0 kB
31 CommitLimit:       3487616 kB
32 Committed_AS:      517056 kB
33 VmallocTotal:      133009637312 kB
34 VmallocUsed:        0 kB
35 VmallocChunk:       0 kB
36 Percpu:            2304 kB
37 HardwareCorrupted: 0 kB
38 AnonHugePages:      0 kB
39 ShmemHugePages:     0 kB
40 ShmemPmdMapped:     0 kB
41 CmaTotal:           524288 kB
42 CmaFree:            523776 kB
43 HugePages_Total:    0
44 HugePages_Free:     0
45 HugePages_Rsvd:     0
46 HugePages_Surp:     0
47 Hugepagesize:       524288 kB
48 Hugetlb:            0 kB

```

发现部分结果符合内核模块的运行结果, 内存分配是动态变化的。

2.3 进程地址空间

本实验期望编写一个内核模块，获取当前进程的地址空间信息以及每个 VMA 的信息，通过设备文件的 open 接口返回给用户态并打印。

本内核模块先在 Ubuntu 下编译运行。在 Ubuntu 的 mm_struct 中 locked_vm 是 unsigned long 型的，而 pinned_vm 是 atomic64_t 型的。第一次在 openEuler 运行时，编译报错：

```
1 /root/lab3/vma/vma.c: In function 'vma_read':
2 /root/lab3/vma/vma.c:42:20: error: incompatible types when assigning to type 'long unsigned
3 int' from type 'atomic_long_t {aka struct <anonymous>}'
4     buf->locked_vm = cur_mm->locked_vm;
5         ^
6 /root/lab3/vma/vma.c:43:39: error: request for member 'counter' in something not a structure
7 or union
8     buf->pinned_vm = cur_mm->pinned_vm.counter;
```

发现在 openEuler 中反之，locked_vm 是 atomic_long_t 型的。修改代码以便能在 openEuler 下运行，核心代码和实验结果如下所示。

```
1 static ssize_t vma_read(struct file *file, char __user *out, size_t size,
2                         loff_t *off) {
3     struct mm_struct *cur_mm = current->mm;
4     struct buffer *buf = file->private_data;
5     if (size == sizeof(struct buffer)) {
6         buf->map_count = cur_mm->map_count;
7         printk(KERN_INFO "OK: %d vma(s) found.\n", buf->map_count);
8     } else {
9         buf->task_size = current->mm->task_size;
10        buf->highest_vm_end = cur_mm->highest_vm_end;
11        buf->mm_users = cur_mm->mm_users.counter;
12        buf->mm_count = cur_mm->mm_count.counter;
13        buf->pgtables_bytes = cur_mm->pgtables_bytes.counter;
14        buf->map_count = cur_mm->map_count;
15        buf->total_vm = cur_mm->total_vm;
16        buf->locked_vm = cur_mm->locked_vm.counter;
17        buf->pinned_vm = cur_mm->pinned_vm;
18        buf->data_vm = cur_mm->data_vm;
19        buf->exec_vm = cur_mm->exec_vm;
20        buf->stack_vm = cur_mm->stack_vm;
21        buf->start_code = cur_mm->start_code;
22        buf->end_code = cur_mm->end_code;
23        buf->start_data = cur_mm->start_data;
24        buf->end_data = cur_mm->end_data;
25        buf->start_brk = cur_mm->start_brk;
26        buf->brk = cur_mm->brk;
27        buf->start_stack = cur_mm->start_stack;
28        buf->arg_start = cur_mm->arg_start;
29        buf->arg_end = cur_mm->arg_end;
30        buf->env_start = cur_mm->env_start;
31        buf->env_end = cur_mm->env_end;
32        printk(KERN_INFO "OK: information of mm_struct got.\n");
33
34        int nvma = (size - sizeof(struct buffer)) / sizeof(struct vma_struct);
35        struct vma_struct *ret_vma_list = (struct vma_struct *) (buf + 1);
36        struct vm_area_struct *vma_list = cur_mm->mmap;
37        for (int i = 0; i < nvma; i++) {
38            ...↓
```

```

39 ...↑
40     ret_vma_list[i].vm_start = vma_list[i].vm_start;
41     ret_vma_list[i].vm_end = vma_list[i].vm_end;
42     ret_vma_list[i].vm_flags = vma_list[i].vm_flags;
43     if (vma_list[i].vm_file && vma_list[i].vm_file->f_path.dentry) {
44         strncpy(ret_vma_list[i].filename,
45                 vma_list[i].vm_file->f_path.dentry->d_iname, 40);
46     } else {
47         strncpy(ret_vma_list[i].filename, "NULL", 5);
48     }
49 }
50 printk(KERN_INFO "OK: information of %d vma(s) got.\n", nvma);
51 }
52 copy_to_user(out, buf, size);
53 printk(KERN_INFO "OK: copy_to_user done.\n");
54 return size;
55 }

```

```

1 [root@ostest vma]# ./run.sh
2 make -C /lib/modules/4.19.188/build M=/root/lab3/vma modules
3 make[1]: Entering directory '/root/kernel'
4 CC [M] /root/lab3/vma/vma.o
5 Building modules, stage 2.
6 MODPOST 1 modules
7 CC /root/lab3/vma/vma.mod.o
8 LD [M] /root/lab3/vma/vma.ko
9 make[1]: Leaving directory '/root/kernel'
10 file opened.
11 t = 176, map_count: 20
12 datasize = 4656
13 task_size: 281474976710656, highest_vm_end: 281474958884864, mm_users: 1, mm_count: 1
14 map_count: 20, total_vm: 58, locked_vm: 0, pinned_vm: 0
15 data_vm: 8, exec_vm: 40, stack_vm: 3
16 start_code: 0x400000, end_code: 0x401070, start_data: 0x41fdc8, end_data: 0x420070
17 start_brk: 0x3c900000, brk: 0x3c930000, start_stack: 0xfffffeeef0d50
18 arg_start: 0xfffffeeefd68, arg_end: 0xfffffeeefd71, env_start: 0xfffffeeefd71, env_end: 0xf...
19 addr_space print ok.
20 start: 0x400000, end: 0x410000, flags: 0x875, filename: test.o
21 start: 0xffff80016f3602b8, end: 0xffff4e230000, flags: 0x100073, filename: libbfd-2.33.1.so
22 start: 0xffff04550000, end: 0xffff04570000, flags: 0x75, filename: libpthread-2.28.so
23 start: 0xffff03620000, end: 0xffff03630000, flags: 0x100073, filename: libargon2.so.0
24 start: 0xffff80016f363ae8, end: 0xffffc8e610000, flags: 0x75, filename: timer.h
25 start: 0xaaadc67d0000, end: 0xaaadc67e0000, flags: 0x100871, filename: systemd-logind
26 start: 0xffff80016f3674e8, end: 0xaaadb9480000, flags: 0x100871, filename: gcc
27 start: 0xffffc474f0000, end: 0xffffc47500000, flags: 0xfb, filename: vma.o
28 start: 0xffff03d10000, end: 0xffff03d20000, flags: 0x100073, filename: libselinux.so.1
29 start: 0xffff80016f3601d0, end: 0xffffea07b0000, flags: 0x100073, filename: libm-2.28.so
30 start: 0xffff80016f3645c8, end: 0xffffc08980000, flags: 0x75, filename: rwlock.h
31 start: 0xffff035c0000, end: 0xffff035d0000, flags: 0x100073, filename: libjson-c.so.4.0.0
32 start: 0xffff034a0000, end: 0xffff034b0000, flags: 0x100073, filename: libpcr2-8.so.0.8.0
33 start: 0xffff03be0000, end: 0xffff03c20000, flags: 0x75, filename: liblz4.so.1.9.2
34 start: 0xffff80016f3677a0, end: 0xffffea0f20000, flags: 0x100073, filename: libpthread-2.28.so
35 start: 0xffff04010000, end: 0xffff04020000, flags: 0x100073, filename: libcap.so.2.27
36 start: 0xffff80016f36bab0, end: 0xaaabf2930000, flags: 0x100073, filename: NULL
37 start: 0xffff03e10000, end: 0xffff03e20000, flags: 0x100073, filename: libmount.so.1.1.0
38 start: 0xffff03f50000, end: 0xffff03f60000, flags: 0x100071, filename: libgcrypt.so.20.2.3
39 start: 0xffff80016f36e290, end: 0xffffea08d0000, flags: 0x100071, filename: libgmp.so.10.3.2
40 address of i is 0x0xfffffeeef0bbc
41 start of print_mm: 0x400854
42 global val: 0x420068
43 local val: 0xfffffeeef0bec
44 static val: 0x42006c
45 malloc allocated mem: 0x3c9007e0

```

使用命令 `vi /proc/vmallocinfo`, 发现部分结果符合内核模块的运行结果。

3 小结

通过在内核态分配内存，了解并掌握操作系统中内存管理的布局 and 内核态内存分配的实现，了解 openEuler 内核的物理内存管理。通过编写内核模块熟悉了 `struct page` 数据结构，了解了 openEuler 内核的物理内存管理，熟悉了进程的地址空间的概念以及相关的数据结构。

本实验在华为弹性云服务器 root@121.36.*.* (openEuler 20.03) 完成。