

CSE 5433 OS Laboratory

Lab4 Report

Group 8
Sixiang Ma, Xiaokuan Zhang

11/01/2017

1 Objectives

In lab4, our goal is to implement partial functions of process checkpointing at the system level. In particular, this lab focuses on checkpointing memory data for a process. It includes three components:

- (a) implement the basic function of checkpointing a memory range specified by a user program;
- (b) design an incremental checkpointing scheme;
- (c) implement the incremental checkpointing scheme.

2 Basic Checkpointing Scheme

2.1 High-level Design

On a high level, we need to implement a system call like previous labs, which will be called by a userspace process for checkpointing its memory within certain range (begin, end). In that system call, there are mainly two steps we should do:

- (a) Traverse the virtual memory area of the process, page by page. If a page is fully within the range, then we need to output it to a file.
- (b) Output corresponding pages to a file.

2.2 Implementation

2.2.1 Creating A System Call

This step is similar to what we did in previous labs. we defined two global variables *cp_start*, *cp_end* to store the requested range for checkpointing. Also, we defined a *file_count* variable to store the current file

counter, and a *filename* char array to store the current filename. After updating the *cp_start*, *cp_end*, we call the *traverse()* function to do the real checkpointing.

```

1 unsigned long cp_start;
2 unsigned long cp_end;
3 int file_count = 0;
4 char filename[16];
5
6 asmlinkage long sys_cp_range(unsigned long *start_addr, unsigned long *end_addr){
7     unsigned long startnow;
8     unsigned long endnow;
9     printk("enter cp_range\n");
10    if (copy_from_user(&startnow, start_addr, sizeof(*start_addr)))
11        return -EFAULT;
12    if (copy_from_user(&endnow, end_addr, sizeof(*end_addr)))
13        return -EFAULT;
14    printk("\ncp_range: %lx, %lx\n\n", startnow, endnow);
15    cp_start = startnow;
16    cp_end = endnow;
17    traverse();
18    printk("exit cp_range\n");
19    ++file_count;
20    return 0;
21 }

```

2.2.2 Traversing the Memory Area of a Process

We can use `area = current->mm->mmap` to get a `vm_area_struct*`, which points to the start of the vm area of the current process. Since it is a linked list, we can use `area = area->vm_next` to traverse this list to visit each vm area.

For each vm area, we visit each page to determine if it is fully within the requested range. If it is, we use `vfs_write()` to write this page to a file. To ensure the safety, we also do a page walk to retrieve the pte of an virtual address. If `pte && pte_read(*pte)`, we write this page to a file. The page walk is performed in `lookup_address_userspace()` function.

```

1 pte_t *lookup_address_userspace(struct mm_struct* mm, unsigned long address)
2 {
3     pgd_t *pgd;
4     // pud_t *pud;
5     pmd_t *pmd;
6     pte_t *pte;
7
8     pgd = pgd_offset(mm, address);
9
10    if (pgd_none(*pgd) || pgd_bad(*pgd))
11        return NULL;
12
13    pmd = pmd_offset(pgd, address);
14    if (pmd_none(*pmd) || pmd_bad(*pmd) || !pmd_present(*pmd))
15        return NULL;
16
17    pte = pte_offset_kernel(pmd, address);

```

```

18     if (pte_none(*pte) || !pte_present(*pte))
19         return NULL;
20     return pte;
21 }

1 void traverse(void){
2     struct task_struct *tsk;
3     struct vm_area_struct *area;
4     struct mm_struct *mm;
5     int ret;
6     unsigned long addr;
7     pte_t *pte;
8     printk("enter traverse\n");
9
10    ret = open_file();
11    if (ret > -1) {
12        tsk = current;
13        if (tsk->mm) {
14            mm = tsk->mm;
15            area = tsk->mm->mmap;
16            while (area) {
17                printk("vm_start: %lx, vm_end: %lx\n", area->vm_start, area->vm_end);
18                addr = area->vm_start;
19                while (addr + PAGE_SIZE <= area->vm_end) {
20                    if (addr >= cp_start && addr + PAGE_SIZE <= cp_end) {
21                        pte = lookup_address_userspace(mm, addr);
22
23                        if (pte && pte_read(*pte)) {
24                            vfs_write(filp, (void*)addr, PAGE_SIZE, &pos);
25                            vfs_write(filp, "\n", 1, &pos);
26                            printk("addr written: %lx\n", addr);
27                        }
28                    }
29                    addr += PAGE_SIZE;
30                }
31                area = area->vm_next;
32            }
33        }
34        printk("exit traverse\n");
35    }
36    close_file();
37 }

```

2.2.3 Writing to a File

We use `vfs_write()` function to write to a file. We defined three global variables `pos`, `old_fs`, `filp` and two helper functions, `open_file()` and `close_file()`, which are called in `traverse()`.

```

1 loff_t pos = 0;
2 mm_segment_t old_fs;
3 struct file *filp = NULL;
4 int open_file(void){
5     printk("enter open_file\n");
6     snprintf(filename, sizeof(filename), "%s%d", "cp-out-", file_count);

```

```

7
8     old_fs = get_fs();
9     set_fs(KERNEL_DS);
10    filp = filp_open(filename, O_WRONLY|O_CREAT, 0644);
11    if (IS_ERR(filp))
12    {
13        printk("OPEN ERROR!\n");
14        return -1;
15    }
16    printk("exit open_file\n");
17    return 0;
18
19 }
20
21 void close_file(void){
22     printk("enter close_file\n");
23     filp_close(filp, NULL);
24     set_fs(old_fs);
25     pos = 0;
26     printk("exit close_file\n");
27 }

```

3 Incremental Checkpointing Scheme

3.1 High-level Design

explain the high-level idea

3.2 Implementation

detailed implementation