# EE516 : Project 2

Gaurav Kalra

(Student ID: 2016 45 93)

gvkalra@kaist.ac.kr

```c
66    static int __init init_procmon (void)
67    {
68
69        struct proc_dir_entry *procmon_proc;
70
71
72        procmon_p        proc_create(FILE_NAME, 644, NULL, &fops);
73        if (!procmon_proc) {
74            printk(KERN_ERR "==Cannot create procmon proc entry \n");
75            return -1;
76        }
77        printk(KERN_INFO "== init procmon\n");
78        return 0;
79    }
80
81    static void __exit exit_procmon(void)
82    {
83        remove_proc_entry(FILE_NAME, NULL);
84        printk(KERN_INFO "== exit procmon\n");
85    }
86
87
88
89    module_init(init_procmon);
90    module_exit(exit_procmon);
91
92    MODULE_LICENSE("GPL");
93    MODULE_AUTHOR("Dong-Jae Shin");
94    MODULE_DESCRIPTION("EE516 Project2 Process Monitoring Module");
```

**proc_create()**: creates a file entry in /proc directory
**parameters**:
1. **name** – filename which will show up in proc directory
2. **mode** – file creation mode (w,r,x for user,group,others)
3. **parent directory** – pointer to parent directory in /proc where this file will be created (we used NULL to indicate /proc/ directory)
4. **file operations** – pointer to file operations structure

**remove_proc_entry()**: remove the file entry from /proc directory
**parameters:**
1. **name** – filename as in /proc
2. **parent** – parent dir pointer (same NULL here)

## Figure 1: Creating / Removing a file from procfs

```c
56    struct file_operations fops = {
57        .owner = THIS_MODULE,
          .open = procmon_proc_open,
          .read = seq_read,
          .llseek = seq_lseek,
          .write = procmon_proc_write,
          .release = single_release,
```

**seq_read()**: read method for sequential files
**seq_lseek()**: llseek method for sequential files
**single_release()**: free the structures associated with sequential files

The parameters for these APIs are compatible with respective VFS file operations. And thus can be used as drop-in replacement so that we do not need to write our own implementation.

## Figure 2: Sequence File APIs for VFS file operations

```c
20    static int procmon_proc_show(struct seq_file *m, void *v)
2
          seq_printf(m, "======= Contents ====== \n");
          return 0;
24    }
25
26
27
28    static int procmon_proc_open(struct inode *inode, struct file *file)
29    {
30        printk(KERN_INFO "proc called open\n");
31        return single_open(file, procmon_proc_show, NULL);
32    }
```

**seq_printf()**: It works like printk() for proc files. The first argument must be the sequence file.
**single_open()**: It creates and opens a new file descriptor for sequence files.
**parameters:**
1. **file** – The file structure from VFS
2. **show** – Call-back for printing data in proc filesystem
3. **data** – Call-back data (NULL in our case)

```c
36    static ssize_t procmon_proc_write(struct file *file, const char __user *buf, size_t count, loff_t *ppos)
37    {
38
39        memset(mybuf, 0, sizeof(mybuf));
40
41        if (count > BUF_SIZE) {
42            count = BUF_SIZE;
43        }
44
45        if (copy_from_user(mybuf, buf, count)) {
46            return -EFAULT;
47        }
48
49        printk(KERN_INFO "proc write : %s\n", mybuf);
50        return (ssize_t)count;
51    }
```

## Figure 3: Using Sequence Files for displaying contents of procfs

**Task 2: Traverse Process - tasklist**
Print every task's **PID** and **Process Name** in your proc file system

```
57    /* file operations */
58    static struct file_operations fops = {                          1
59         .owner = THIS_MODULE,
60         .open = pl_open,
61
62         /* read method for sequential files */
63         .read = seq_read,
64
65         /* llseek method for sequential files */
66         .llseek = seq_lseek,
67
68         /* free the structures associated with sequential file */
69         .release = single_release,
70    };
71
72    static void
73    _pl_module_exit(void)
74    {
75         dbg("");
76
77         if (pl != NULL)
78              proc_remove(pl);          2
79    }
80
81    static int __init
82    pl_module_init(void)
83    {
84         dbg("");
85
86         /* create /proc/proc_list */
87         pl = proc_create(PROC_NAME, 0, NULL, &fops);    3
88         if (pl == NULL) {
89              err("Failed to create proc_list");
90              goto error;
91         }
92         return 0;
93
94    error:
95         _pl_module_exit();
96         return -1;
97    }
98
99    static void __exit
100   pl_module_exit(void)
101   {
102        _pl_module_exit();
103   }
104
105   module_init(pl_module_init);
106   module_exit(pl_module_exit);
107
108   MODULE_AUTHOR("Gaurav Kalra");
109   MODULE_DESCRIPTION("PR02 Traverse Process - tasklist");
110   MODULE_LICENSE("GPL");
```

**Figure 4**: Module Initialization / Clean-up (refer task01 for details)

```
12    /*
13    Organization of task information in kernel:
14
15        struct task_struct {
16            ...
17            pid_t pid;
18            ...
19            struct list_head tasks;
20            ...
21            char comm[TASK_COMM_LEN];
22        };
23
24        struct list_head {
25            |  struct list_head *next, *prev;
26        };
27    */
28
29    static int
30    pl_show(struct seq_file *m, void *v)
31    {
32        struct task_struct *tsk;
33        char name[TASK_COMM_LEN];
34
35        /* header */
36        seq_printf(m, "PID        ProcessName        \n");
37
38        /* Print pid & name of each process */
39        for_each_process(tsk) {
40
41            /* print information */
42            seq_printf(m, "%-10u%-20s\n",
43                task_pid_nr(tsk),
44                get_task_comm(name, tsk));
45        }
46
47        return 0;
48    }
49
50    static int
51    pl_open(struct inode *inode, struct file *file)
52    {
53        dbg("");
54        return single_open(file, pl_show, NULL);
55    }
56
```

**Figure 5**: Logic for printing PID & Process Name

```
1    MODNAME := proc_list
2    obj-m := ${MODNAME}.o
3    ${MODNAME}-objs := main.o
4
5    KDIR := /lib/modules/$(shell uname -r)/build
6
7    all:
8        $(MAKE) -C $(KDIR) M=$(PWD) modules
9
10   clean:
11        $(MAKE) -C $(KDIR) M=$(PWD) clean
```

```
gvkalra@gvkalra-desktop ~/Desktop/EE516/PR02/task02 (master) $ make
make -C /lib/modules/4.4.0-38-generic/build M=/home/gvkalra/Desktop/EE516/PR02/task02 modules
make[1]: Entering directory '/usr/src/linux-headers-4.4.0-38-generic'
  CC [M]  /home/gvkalra/Desktop/EE516/PR02/task02/main.o
  LD [M]  /home/gvkalra/Desktop/EE516/PR02/task02/proc_list.o
  Building modules, stage 2.
  MODPOST 1 modules
  CC       /home/gvkalra/Desktop/EE516/PR02/task02/proc_list.mod.o
  LD [M]  /home/gvkalra/Desktop/EE516/PR02/task02/proc_list.ko
make[1]: Leaving directory '/usr/src/linux-headers-4.4.0-38-generic'
```

**Figure 6**: Makefile (no warnings)

```
gvkalra@gvkalra-desktop ~/Desktop/EE516/PR02/task02 (master) $ cat /proc/proc_list | wc -l
290
gvkalra@gvkalra-desktop ~/Desktop/EE516/PR02/task02 (master) $ ps -ef | wc -l
290
```

**Figure 7**: Verification (**proc_list & ps -ef both output 290 processes**)

**Task 3: per Process Memory Usage**
Print every task's **VIRT** and **RSS Memory** Size

```
13    /*
14    Organization of task information in kernel:
15
16       struct task_struct {
17           ...
18           pid_t pid;
19           ...
20           struct list_head tasks;
21           ...
22           char comm[TASK_COMM_LEN];
23           ...
24           struct mm_struct *active_mm;
25       };
26
27       struct list_head {
28           struct list_head *next, *prev;
29       };
30
31       struct mm_struct {
32           ...
33           unsigned long total_vm;
34           ...
35           struct mm_rss_stat rss_stat;
36       };
37
38       struct mm_rss_stat {
39           atomic_long_t count[NR_MM_COUNTERS];
40       };
41
42        enum {
43           MM_FILEPAGES,   //Resident file mapping pages (Type of File Mapped Page)
44           MM_ANONPAGES,   //Resident anonymous pages (Type of Anonymous Page - Stack, Heap)
45           MM_SWAPENTS,    //Anonymous swap entries
46           MM_SHMEMPAGES, //Resident shared memory pages
47           NR_MM_COUNTERS
48       };
49
50       MM_FILEPAGES + MM_ANONPAGES = RSS Memory
51       total_vm = VIRT Memory
52    */
53
```

**Figure 8**: Organization of **VIRT** & **RSS** in task_struct

```
 54    static int
 55    pl_show(struct seq_file *m, void *v)
 56    {
 57        struct task_struct *tsk;
 58        char name[TASK_COMM_LEN];
 59        unsigned long long virt = 0;
 60        long long rss = 0;
 61
 62        /* header */
 63        seq_printf(m, "PID          ProcessName              "
 64                "VIRT(KB)             RSS Mem(KB)             \n");
 65
 66        /* Print name, PID, VIRT & RSS of each process */
 67        for_each_process(tsk) {
 68            virt = rss = 0;
 69
 70            /* It is possible for active_mm to be NULL.
 71             * In this case, we simply assume VIRT and RSS to be 0.
 72             * Ref: https://www.kernel.org/doc/Documentation/vm/active_mm.txt
 73             */
 74            if (tsk->active_mm != NULL) {
 75                /* VIRT memory */
 76     ①        virt = tsk->active_mm->total_vm;
 77                virt *= (PAGE_SIZE >> 10); /* convert pages to KB */
 78
 79                /* RSS Memory */
 80     ②        rss += atomic_long_read(&tsk->active_mm->rss_stat.count[MM_FILEPAGES]);
 81                rss += atomic_long_read(&tsk->active_mm->rss_stat.count[MM_ANONPAGES]);
 82                rss *= (PAGE_SIZE >> 10); /* convert pages to KB */
 83            }
 84
 85            /* print information */
 86            seq_printf(m, "%-10u%-20s%-20llu%-20llu\n",
 87                task_pid_nr(tsk),                              ③
 88                get_task_comm(name, tsk),
 89                virt, rss);
 90        }
 91
 92        return 0;
 93    }
```

**Figure 9**: Logic for calculating and printing VIRT, RSS

```
1   MODNAME := proc_memory
2   obj-m := ${MODNAME}.o
3   ${MODNAME}-objs := main.o
4
5   KDIR := /lib/modules/$(shell uname -r)/build
6
7   all:
8       $(MAKE) -C $(KDIR) M=$(PWD) modules
9
10  clean:
11      $(MAKE) -C $(KDIR) M=$(PWD) clean
```

**Figure 10**: Makefile (no warnings)

```
gvkalra@gvkalra-desktop ~/Desktop/EE516/PR02/task03 (master) $ make
make -C /lib/modules/4.4.0-38-generic/build M=/home/gvkalra/Desktop/EE516/PR02/task03 modules
make[1]: Entering directory '/usr/src/linux-headers-4.4.0-38-generic'
  CC [M]  /home/gvkalra/Desktop/EE516/PR02/task03/main.o
  LD [M]  /home/gvkalra/Desktop/EE516/PR02/task03/proc_memory.o
  Building modules, stage 2.
  MODPOST 1 modules
  CC      /home/gvkalra/Desktop/EE516/PR02/task03/proc_memory.mod.o
  LD [M]  /home/gvkalra/Desktop/EE516/PR02/task03/proc_memory.ko
make[1]: Leaving directory '/usr/src/linux-headers-4.4.0-38-generic'
```

```
gvkalra@gvkalra-desktop ~/Desktop/EE516/PR02/task03 (master) $ cat /proc/proc_memory | grep sublime
16199     sublime_text        1096892              67768        1
gvkalra@gvkalra-de    2  p ~/Desktop/EE516/PR02/task03 (master) $ top -b -n 1 | grep sublime
16199 gvkalra   20    1096892  67768  52368 S   0.0  0.8   0:00.44 sublime_text
```

**Figure 11**: Verification (sublime_text)

```
gvkalra@gvkalra-desktop ~/Desktop/EE516/PR02/task03 (master) $ top -b -n 1 | grep slack
 2529 gvkalra   20    0 1686468 146920   85420 S   0.0  1.8   3:36.89 slack
 2600 gvkalra   20      316428   30296   27184 S   0.0  0.4   0:00.01 slack
 2669 gvkalra   20    0 1102600 104920   61124 S   0.0  1.3   0:08.41 slack
 2736 gvkalra   20    0 1496056 346040  224816 S   0.0  4.3   1:57.91 slack
 2780 gvkalra   20    0 1379216 265472  136304 S   0.0  3.3   0:45.19 slack
gvkalra@gvkalra-desktop ~/Desktop/EE516/PR02/task03 (master) $ cat /proc/proc_memory | grep slack
2529      slack        1686468              146920
2600      slack         316428               30296
2669      slack        1102600              104920        2
2736      slack        1496056              346040
2780      slack        1379216              265472
```

**Figure 12**: Verification (slack)

**Task 4: per Process I/O Usage**
Print every task's I/O usage

```
13    /*
14    Organization of task information in kernel:
15
16        struct task_struct {
17            ...
18            pid_t pid;
19            ...
20            struct list_head tasks;
21            ...
22            char comm[TASK_COMM_LEN];
23            ...
24            struct task_io_accounting ioac;
25        };
26
27        struct list_head {
28            struct list_head *next, *prev;
29        };
30
31        struct task_io_accounting {
32            #ifdef CONFIG_TASK_XACCT
33                u64 rchar; //bytes read
34                u64 wchar; //bytes written
35                u64 syscr; //# of read syscalls
36                u64 syscw; //# of write syscalls
37            #endif
38
39            #ifdef CONFIG_TASK_IO_ACCOUNTING
40                //The number of bytes which this task has caused to be read from storage.
41                u64 read_bytes;
42
43                //The number of bytes which this task has caused, or shall cause to be written to disk.
44                u64 write_bytes;
45
46                //A task can cause "negative" IO too.  If this task truncates some
47                //dirty pagecache, some IO which another task has been accounted for
48                //(in its write_bytes) will not be happening.  We _could_ just
49                //subtract that from the truncating task's write_bytes, but there is
50                //information loss in doing that.
51                u64 cancelled_write_bytes;
52            #endif
53        };
54    */
```

**Figure 13**: Organization of **I/O data** in task_struct

```
56    static int
57    pl_show(struct seq_file *m, void *v)
58    {
59        struct task_struct *tsk, *t;
60        char name[TASK_COMM_LEN];
61        struct task_io_accounting acct;
62
63        /* header */
64        seq_printf(m, "PID          ProcessName              "
65              "rchar(B)              wchar(B)                 "
66              "syscr(#)              syscw(#)                 "
67              "read_bytes(B)          write_bytes(B)         "
68              "cancelled_write_bytes(B)\n");
69
70        /* Print name, PID, I/O stats of each process */
71        for_each_process(tsk) {
72            acct = tsk->ioac; /* initialize accounting data */      (1)
73
74            /* account each thread
75             * Ref: https://github.com/torvalds/Linux/blob/master/fs/proc/base.c
76             * Function: do_io_accounting()
77             */
78            t = tsk;
79            task_io_accounting_add(&acct, &tsk->signal->ioac);
80            while_each_thread(tsk, t)                              (2)
81                task_io_accounting_add(&acct, &t->ioac);
82
83            /* print information */
84            seq_printf(m, "%-10u%-20s%-20llu%-20llu%-20llu%-20llu%-20llu%-20llu%-24llu\n",
85                task_pid_nr(tsk),
86                get_task_comm(name, tsk),
87        (3)    (unsigned long long)acct.rchar,
88                (unsigned long long)acct.wchar,
89                (unsigned long long)acct.syscr,
90                (unsigned long long)acct.syscw,
91                (unsigned long long)acct.read_bytes,
92                (unsigned long long)acct.write_bytes,
93                (unsigned long long)acct.cancelled_write_bytes);
94        }
95
96        return 0;
97    }
```

**Figure 14**: Logic for calculating and printing I/O information

```
 1    MODNAME := proc_io
 2    obj-m := ${MODNAME}.o
 3    ${MODNAME}-objs := main.o
 4
 5    KDIR := /lib/modules/$(shell uname -r)/build
 6
 7    all:
 8        $(MAKE) -C $(KDIR) M=$(PWD) modules
 9
10    clean:
11        $(MAKE) -C $(KDIR) M=$(PWD) clean
```

**Figure 15**: Makefile (no warnings)

```
gvkalra@gvkalra-desktop ~/Desktop/EE516/PR02/task04 (master) $ make
make -C /lib/modules/4.4.0-38-generic/build M=/home/gvkalra/Desktop/EE516/PR02/task04 modules
make[1]: Entering directory '/usr/src/linux-headers-4.4.0-38-generic'
  CC [M]  /home/gvkalra/Desktop/EE516/PR02/task04/main.o
  LD [M]  /home/gvkalra/Desktop/EE516/PR02/task04/proc_io.o
  Building modules, stage 2.
  MODPOST 1 modules
  CC      /home/gvkalra/Desktop/EE516/PR02/task04/proc_io.mod.o
  LD [M]  /home/gvkalra/Desktop/EE516/PR02/task04/proc_io.ko
make[1]: Leaving directory '/usr/src/linux-headers-4.4.0-38-generic'
```

```
gvkalra@gvkalra-desktop ~/Desktop/EE516/PR02/task04 (master) $ cat /proc/proc_io | grep colord
1015      colord          859003        11075         12020         1470          10805248        0               0
gvkalra@gvkalra-desktop ~/Desktop/EE516/PR02/task04 (master) $ sudo cat /proc/1015/io
rchar: 859003
wchar: 11075
syscr: 12020
syscw: 1470
read_bytes: 10805248
write_bytes: 0
cancelled_write_bytes: 0
```

**Figure 16**: Verification (colord)

task05
main.c
Makefile
pm_list.c
pm_list.h
sort.c
sort.h
utils.h

**Task 5: Sorting Features**
Create a process monitoring tool with sorting feature

The task is divided into 2 modules:
1. **Sorting Module** – It is used to view or update current sorting order
2. **PM (procmon) Module** – It is used to create process monitoring

```c
1   #pragma once
2
3   enum {
4       SORT_ORDER_PID = 0,
5       SORT_ORDER_VIRT,
6       SORT_ORDER_RSS,
7       SORT_ORDER_IO
8   };
9
10  /* Initializes sorting module
11   * Return:
12   *       < 0 on error
13   */
14  int
15  sort_module_init(void);
16
17  /* De-initializes sorting module
18   */
19  void
20  sort_module_exit(void);
21
22  /* Returns current sorting order
23   * Refer SORT_ORDER_<xyz> enumeration for return values
24   */
25  inline int
26  sort_get_order(void);
```

```c
88   /* file operations */
89   static struct file_operations sort_ops = {
90       .owner = THIS_MODULE,
91       .open = sort_open,
92
93       /* read method for sequential files */
94       .read = seq_read,
95
96       /* llseek method for sequential files */
97       .llseek = seq_lseek,
98
99       /* write() system call on VFS */
100      .write = sort_write,
101
102      /* free the structures associated with sequential file */
103      .release = single_release,
104  };
105
106  void
107  sort_module_exit(void)
108  {
109      dbg("");
110
111      if (ps != NULL)
112          proc_remove(ps);
113  }
114
115  int
116  sort_module_init(void)
117  {
118      dbg("");
119
120      /* create /proc/procmon_sorting */
121      ps = proc_create(PROC_SORT, 0, NULL, &sort_ops);
122      if (ps == NULL) {
123          err("Failed to create procmon_sorting");
124          goto error;
125      }
126      return 0;
127
128  error:
129      sort_module_exit();
130      return -1;
131  }
132
133  inline int
134  sort_get_order(void)
135  {
136      return sort_order;
137  }
```

**Figure 17**: Sorting Module APIs exposed to PM module

```
48    static ssize_t
49    sort_write(struct file *file, const char __user *user_buf, size_t length, loff_t *offset)
50    {
51        char buf[BUF_SIZE];
52        dbg("");
53
54        /* clear buffer */
55        memset(buf, 0x00, sizeof(buf));
56
57        /* resize */
58        if (length > BUF_SIZE)
59            length = BUF_SIZE;
60
61        /* copy data from user space */
62        if (copy_from_user(buf, user_buf, length))   (1)
63            return -EFAULT; /* Bad address */
64
65        /* null terminate buffer */   (2)
66        buf[length - 1] = '\0';
67
68        dbg("buf: [%s]", buf);
69
70        /* set sorting order */
71        if (strcmp(buf, "pid") == 0)
72            sort_order = SORT_ORDER_PID;
73        else if (strcmp(buf, "virt") == 0)
74            sort_order = SORT_ORDER_VIRT;
75        else if (strcmp(buf, "rss") == 0)
76            sort_order = SORT_ORDER_RSS;
77        else if (strcmp(buf, "io") == 0)
78            sort_order = SORT_ORDER_IO;
79        else
80            info("Invalid value! sort_order is not changed");
81
82        dbg("sort_order: [%d]", sort_order);   (3)
83
84        /* copied from user space, return bytes */
85        return length;
86    }
```

**Figure 18**: Sorting Module write() file operation to update sort_order

```
16  static int
17  sort_show(struct seq_file *m, void *v)
18  {
19      switch (sort_order) {
20      case SORT_ORDER_VIRT:
21          seq_printf(m, "PID \t [VIRT] \t RSS \t I/O\n");
22          break;
23
24      case SORT_ORDER_RSS:
25          seq_printf(m, "PID \t VIRT \t [RSS] \t I/O\n");
26          break;
27
28      case SORT_ORDER_IO:
29          seq_printf(m, "PID \t VIRT \t RSS \t [I/O]\n");
30          break;
31
32      case SORT_ORDER_PID:
33      default: /* fallthrough */
34          seq_printf(m, "[PID] \t VIRT \t RSS \t I/O\n");
35          break;
36      }
37
38      return 0;
39  }
40
41  static int
42  sort_open(struct inode *inode, struct file *file)
43  {
44      dbg("");
45      return single_open(file, sort_show, NULL);
46  }
```

**Figure 19**: Sorting Module open() file operation

```
52    /* file operations */
53    static struct file_operations fops_pm = {
54          .owner = THIS_MODULE,
55          .open = pm_open,
56
57          /* read method for sequential files */
58          .read = seq_read,
59
60          /* llseek method for sequential files */
61          .llseek = seq_lseek,
62
63          /* free the structures associated with sequential file */
64          .release = pm_release,
65    };
66
67    static void
68    _pm_module_exit(void)
69    {
70          dbg("");
71
72          /* de-initialize sorting module */
73          sort_module_exit();
74
75          /* remove procmon */
76          if (pm != NULL)
77                proc_remove(pm);
78    }
79
80    static int __init
81    pm_module_init(void)
82    {
83          int ret;
84          dbg("");
85
86          /* create /proc/procmon */
87          pm = proc_create(PROC_NAME, 0, NULL, &fops_pm);
88          if (pm == NULL) {
89                err("Failed to create procmon");
90                goto error;
91          }
92
93          /* initialize sorting module */
94          ret = sort_module_init();
95          if (ret < 0) {
96                err("Failed to initialize sorting module");
97                goto error;
98          }
99          return 0;
100
101   error:
102         _pm_module_exit();
103         return -1;
104   }
105
106   static void __exit
107   pm_module_exit(void)
108   {
109         _pm_module_exit();
110   }
```

**Figure 20**: PM Module init / exit
(Sorting module is init/exit by PM module)

```c
12    static int
13    pm_show(struct seq_file *m, void *v)
14    {
15        dbg("");
16
17        /* show pm_list */
18        return pm_list_show(m);
19    }
20
21    static int
22    pm_open(struct inode *inode, struct file *file)
23    {
24        int ret;
25        dbg("");
26
27        /* initialize pm_list
28         * it means to parse all processes & save them in
29         * linked list owned by this module
30         */
31        ret = pm_list_init();
32        if (ret < 0) {
33            err("Failed to initialize pm_list: %d", ret);
34            return ret;
35        }
36
37        return single_open(file, pm_show, NULL);
38    }
39
40    static int
41    pm_release(struct inode *inode, struct file *file)
42    {
43        dbg("");
44
45        /* clean-up pm_list */
46        pm_list_deinit();
47
48        /* release sequence file */
49        return single_release(inode, file);
50    }
51
```

**Figure 21**: PM Module open() and release()
(**pm_list** is PM Module's data structure discussed in next page)

```
 7   /* pm_list node entry */
 8   struct pm_list_entry {
 9       struct list_head entries;
10
11       pid_t pid;
12       char name[TASK_COMM_LEN];
13       unsigned long long virt; /* in KB */
14       long long rss; /* in KB */
15       unsigned long long disk_read; /* in KB */
16       unsigned long long disk_write; /* in KB */
17   };
18
19   /* first entry of pm_list */
20   extern struct pm_list_entry pm_list_init_entry;
21
22   /* macro to find next pm_list element */
23   #define pm_list_next_entry(e) \
24       container_of((e)->entries.next, struct pm_list_entry, entries)
25
26   /* initializes pm_list
27    * Return:
28    *    < 0 on error
29    */
30   int pm_list_init(void);
31
32   /* deinits pm_list
33    */
34   void pm_list_deinit(void);
35
36   /* shows pm_list after sorting in
37    * currently set order
38    */
39   int pm_list_show(struct seq_file *m);
```

**Figure 22**: **pm_list** data structure & APIs.
(Each API is discussed next)


Notes:
1. The design of this data structure closely resembles
task_struct from the Linux Kernel.

2. Reference:
https://github.com/torvalds/linux/blob/master/init/init_task.c

3. "struct list_head" is the kernel way of providing linked lists, which
has been referred from
https://github.com/torvalds/linux/blob/master/include/linux/list.h

```
125    int
126    pm_list_init(void)
127    {
128        struct task_struct *tsk, *t;
129        struct pm_list_entry *entry;
130        char name[TASK_COMM_LEN];
131        unsigned long long virt;
132        long long rss;
133        struct task_io_accounting acct;
134
135        dbg("");
136
137        /* Add to linked list */
138        for_each_process(tsk) {                                    (1)
139            dbg("[ADD] name: [%s] pid: [%d]", get_task_comm(name, tsk),
140                task_pid_nr(tsk));
141
142            /* Calculate memory */
143            virt = 0;
144            rss = 0;
145            if (tsk->active_mm != NULL) {
      (2)      virt = tsk->active_mm->total_vm;
147            virt *= (PAGE_SIZE >> 10); /* convert pages to KB */
148
149            rss += atomic_long_read(&tsk->active_mm->rss_stat.count[MM_FILEPAGES]);
150    (3)    rss += atomic_long_read(&tsk->active_mm->rss_stat.count[MM_ANONPAGES]);
151            rss *= (PAGE_SIZE >> 10); /* convert pages to KB */
152        }
153
154        /* I/O */
155        acct = tsk->ioac;
156        t = tsk;
157                                                                   (4)
158        /* Account for each thread */
159        task_io_accounting_add(&acct, &tsk->signal->ioac);
160        while_each_thread(tsk, t)
161            task_io_accounting_add(&acct, &t->ioac);
162
163        /* Allocate linked list node */
164        entry = alloc_pm_list_node(task_pid_nr(tsk),
165                    get_task_comm(name, tsk),
166                    virt, rss,
167                    (acct.read_bytes >> 10), /* divide by 1024 to convert into KB */
168                    (acct.write_bytes >> 10));
169        if (entry == NULL) {
170            info("Skipping PID: [%d]", task_pid_nr(tsk));
171            continue;
172        }
173
174        /* Add to list */                                          (5)
175        list_add(&entry->entries, &pm_list_init_entry.entries);
176    }
177
178        return 0;
179    }
```

**Figure 23**: pm_list_init() – API for creating & initializing a PM list

```
56   int
57   pm_list_show(struct seq_file *m)
58   {
59       struct list_head *cursor, *temp;
60       struct pm_list_entry *entry;
61       dbg("");
62
63       /* Sort list */
64       list_sort(NULL,                          ①
65           &pm_list_init_entry.entries,
66           sort_pm_list_entries);
67
68       /* title */
69       seq_printf(m, "================ Process Monitoring Manager for EE516 ================\n");
70
71       /* header */
72       seq_printf(m, "PID        ProcessName           VIRT(KB)          "
73           "RSS Mem(KB)        DiskRead(KB)        DiskWrite(KB)       "
74           "Total I/O(KB)      \n");
75
76       list_for_each_safe(cursor, temp, &pm_list_init_entry.entries) {
77           /* typecast */
78           entry = list_entry(cursor, struct pm_list_entry, entries);
79
80           /* print information */                                        ②
81           seq_printf(m, "%-10u%-20s%-20llu%-20llu%-20llu%-20llu%-20llu\n",
82           entry->pid,
83           entry->name,
84           entry->virt,
85           entry->rss,
86           entry->disk_read,
87           entry->disk_write,
88           (entry->disk_read + entry->disk_write));
89       }
90
91       return 0;
92   }
```

```
23   /* sort function for pm_list */
24   static int
25   sort_pm_list_entries(void *priv, struct list_head *a, struct list_head *b)
26   {
27       struct pm_list_entry *entry_a, *entry_b;
28       int sort_order;
29
30       /* typecast */
31       entry_a = list_entry(a, struct pm_list_entry, entries);
32       entry_b = list_entry(b, struct pm_list_entry, entries);
33
34    ① /* find current sort order */
35       sort_order = sort_get_order();
36
37       /* sort */
38       switch (sort_order) {
39       case SORT_ORDER_VIRT:                                      ②
40           return (entry_b->virt - entry_a->virt);
41
42       case SORT_ORDER_RSS:
43           return (entry_b->rss - entry_a->rss);
44
45       case SORT_ORDER_IO:
46           /* subtract & then add to prevent overflow */
47           return ((entry_b->disk_write - entry_a->disk_write) + \
48               (entry_b->disk_read - entry_a->disk_read));
49
50       case SORT_ORDER_PID:
51       default:
52           return (entry_a->pid - entry_b->pid);
53       }
54   }
```

**Figure 24**: pm_list_show() – API for displaying PM list in sorted order

```
181    void
182    pm_list_deinit(void)
183    {
184        struct list_head *cursor, *temp;
185        struct pm_list_entry *entry;
186        dbg("");
187
188    [1] list_for_each_safe(cursor, temp, &pm_list_init_entry.entries) {
189            entry = list_entry(cursor, struct pm_list_entry, entries);
190            dbg("[REMOVE] name: [%s] pid: [%d]", entry->name, entry->pid);
191
192    [2] list_del(cursor);
193        free_pm_list_node(entry);
194        }
195
196        dbg("Is empty? %d", list_empty(&pm_list_init_entry.entries));
197    }
```

**Figure 25**: pm_list_deinit() – API for freeing a PM list

```
1    MODNAME := procmon
2    obj-m := ${MODNAME}.o
3    ${MODNAME}-objs := main.o pm_list.o sort.o
4
5    KDIR := /lib/modules/$(shell uname -r)/build
6
7    all:
8        $(MAKE) -C $(KDIR) M=$(PWD) modules
9
10   clean:
11       $(MAKE) -C $(KDIR) M=$(PWD) clean
```

**Figure 26**: Makefile (no warnings)

```
gvkalra@gvkalra-desktop ~/Desktop/EE516/PR02/task05 (master) $ make
make -C /lib/modules/4.4.0-38-generic/build M=/home/gvkalra/Desktop/EE516/PR02/task05 modules
make[1]: Entering directory '/usr/src/linux-headers-4.4.0-38-generic'
  CC [M]  /home/gvkalra/Desktop/EE516/PR02/task05/main.o
  CC [M]  /home/gvkalra/Desktop/EE516/PR02/task05/pm_list.o
  CC [M]  /home/gvkalra/Desktop/EE516/PR02/task05/sort.o
  LD [M]  /home/gvkalra/Desktop/EE516/PR02/task05/procmon.o
  Building modules, stage 2.
  MODPOST 1 modules
  CC      /home/gvkalra/Desktop/EE516/PR02/task05/procmon.mod.o
  LD [M]  /home/gvkalra/Desktop/EE516/PR02/task05/procmon.ko
make[1]: Leaving directory '/usr/src/linux-headers-4.4.0-38-generic'
```

```
root@gvkalra-desktop:~# cat /proc/procmon_sorting
[PID]    VIRT    RSS    I/O
root@gvkalra-desktop:~# echo virt > /proc/procmon_sorting
root@gvkalra-desktop:~# cat /proc/procmon
=============== Process Monitoring Manager for EE516 ===============
PID    ProcessName      VIRT(KB)    RSS Mem(KB)    DiskRead(KB)    DiskWrite(KB)    Total I/O(KB)
3454   synergys         3214860     21156          1252            0                1252
2372   evolution-calen  2196228     75640          600             256              856
2324   compiz           1849444     525492         24892           128              25020
2529   slack            1685444     146980         49980           12604            62584
21950  chrome           1496680     190064         0               0                0
2736   slack            1496056     345404         13428           84               13512
```

**Figure 27**: Verification (VIRT)

```
root@gvkalra-desktop:~# echo rss > /proc/procmon_sorting
root@gvkalra-desktop:~# cat /proc/procmon
=============== Process Monitoring Manager for EE516 ===============
PID    ProcessName      VIRT(KB)    RSS Mem(KB)    DiskRead(KB)    DiskWrite(KB)    Total I/O(KB)
2324   compiz           1849444     525492         24892           128              25020
2736   slack            1496056     345404         13428           84               13512
2780   slack            1379216     264996         4               20               24
21721  chrome           1188368     218420         1180            603128           604308
22011  chrome           980040      192536         0               0                0
21950  chrome           1496680     190180         0               0                0
```

**Figure 28**: Verification (RSS)

In Linux, Processes & Threads are all represented by task_struct. The pid field of task_struct is unique for every process & thread.

However, since POSIX mandates that all threads of a process should have the same ProcessID (as seen by the user), Linux uses "Thread Group ID" (TGID) to satisfy POSIX mandate.

A TGID is the PID of the thread that started the whole process.

```
1474
1475    struct task_struct {
1591
1592            pid_t pid;
1593            pid_t tgid;
1594
```

**Figure 29**: pid/tgid in task_struct.
Source: https://github.com/torvalds/linux/blob/master/include/linux/sched.h

In Linux, both processes & threads are created using clone() system call. The 'flags' passed to clone() specify the degree of data sharing. Threads in Linux have it's own stack but shares Heap, BSS, Data and Text with other threads having same TGID.

The CFS (Completely Fair Scheduler) as used in Kernel doesn't discriminate between a thread & a process. It works on "tasks" and is not aware of thread & process abstractions.

In Linux, user-space is separated from kernel-space in all aspects.
To arbitrate data between user / kernel space, Linux provides two
functions:
1. **copy_to_user**: Copy data from kernel space to user space
2. **copy_from_user**: Copy data from user space to kernel space

A proc file system (procfs) is just another type of file system on
top of Virtual File System (VFS) & thus applies the same rule of
data exchange between user & kernel space.