

COL331 A1

Sarthak Panda (2022CS51217)

March 2025

1 Introduction

The current implementation in total required modification of 10 files. Which we have added to patch files and are discussed in following sections.

File	Status
Kbuild	Modified
arch/x86/entry/syscalls/syscall_64.tbl	Modified
arch/x86/kernel/sys_x86_64.c	Modified
arch/x86/um/syscalls_64.c	Modified
fs/open.c	Modified
hello/Makefile	New File
hello/hello.c	New File
include/linux/hello.h	New File
include/linux/syscalls.h	Modified
mm/mmap.c	Modified

2 Data structure of the Monitoring List

```
1 #ifndef HELLO_H
2 #define HELLO_H
3 #include <linux/list.h>
4 #include <linux/types.h>
5 struct pid_node {
6     struct per_proc_resource* proc_resource;
7     struct list_head next_prev_list;
8 };
9 struct per_proc_resource {
10     pid_t pid;
11     unsigned long heapsize;
12     unsigned long openfile_count;
13     long heap_quota;
14     long file_quota;
15     bool quotas_defined;
16 };
17 extern struct list_head monitored_pids_head_node;
18 extern struct mutex monitored_list_mutex;
19 #endif
```

Listing 1: Kernel code snippet from hello/hello.h

I have updated the per_proc_resource node structure to handle both the resource usage tracker and limiter and added three more attributes (heap_quota, file_quota, and quotas_defined). Further, in the hello.h file, I kept the declaration of the monitored_pids_head_node, which is a global variable representing the sentinel node of the list to be monitored. It is subsequently initialized by the macro LIST_HEAD(monitored_pids_head_node); in hello.c. Similarly, monitored_list_mutex is a mutex lock declared in hello.h, which will be used while accessing the monitoring list (since multiple threads can access the list, it becomes essential to use locks). Again, it is initialized and defined in 'hello.c' by the macro DEFINE_MUTEX(monitored_list_mutex);.

3 Resource Usage Tracker

3.1 sys_register

```
1 SYSCALL_DEFINE1(register, pid_t, pid) /* in hello/hello.c */
2 {
3     struct task_struct *task;
4     struct pid_node *node;
5     struct list_head *pos;
6     printk(KERN_INFO "Syscall sys_register called with pid: %d\n", pid);
7
8     if (pid < 1) {
9         printk(KERN_INFO "Syscall sys_register: PID %d is less than 1\n", pid);
10        return -22;
11    }
12
13    rcu_read_lock();
14    task = find_task_by_vpid(pid);
15    rcu_read_unlock();
16    if (!task) {
17        printk(KERN_INFO "Syscall sys_register: PID %d does not exist\n", pid);
18        return -3;
19    }
20
21    mutex_lock(&monitored_list_mutex);
22    list_for_each(pos, &monitored_pids_head_node) {
23        node = list_entry(pos, struct pid_node, next_prev_list);
24        if (node->proc_resource->pid == pid) {
25            mutex_unlock(&monitored_list_mutex);
26            printk(KERN_INFO "Syscall sys_register: PID %d is already monitored\n", pid);
27            return -23;
28        }
29    }
30
31    node = kmalloc(sizeof(struct pid_node), GFP_KERNEL);
32    if (!node) {
33        mutex_unlock(&monitored_list_mutex);
34        printk(KERN_ERR "Syscall sys_register: Memory allocation failed for PID %d\n", pid);
35        return -ENOMEM;
36    }
37
38    node->proc_resource = kmalloc(sizeof(struct per_proc_resource), GFP_KERNEL);
39    if (!node->proc_resource) {
40        kfree(node);
41        mutex_unlock(&monitored_list_mutex);
42        printk(KERN_ERR "Syscall sys_register: Memory allocation failed for PID %d\n", pid);
43        return -ENOMEM;
44    }
45
46    node->proc_resource->pid = pid;
47    unsigned long heapsize = 0;
48    node->proc_resource->heapsize = 0;
49    node->proc_resource->openfile_count = 0;
50    node->proc_resource->heap_quota = -1;
51    node->proc_resource->file_quota = -1;
52    node->proc_resource->quotas_defined = false;
53
54    list_add(&node->next_prev_list, &monitored_pids_head_node);
55    mutex_unlock(&monitored_list_mutex);
56    return 0;
57 }
```

Listing 2: Kernel code snippet for `sys_register` in `hello/hello.c`

The first important check that `sys_register` needs to perform is to check if the `pid` is valid and there is an `task_struct` corresponding to it. For which we use `find_task_by_vpid` (resource link) .This function is essentially is a wrapper that finds the process by its virtual PID (`vpid`) within the PID namespace of the current process. It calls `find_task_by_pid_ns()` with the provided `vpid` and the active PID namespace (obtained via `task_active_pid_ns(current)`), ensuring that the lookup is within the current namespace boundaries.

```

1 struct task_struct *find_task_by_vpid(pid_t vnr)/*in kernel/pid.c*/
2 {
3     return find_task_by_pid_ns(vnr, task_active_pid_ns(current));
4 }

```

Listing 3: Kernel code snippet for `find_task_by_vpid`

Further, it is important to note that `rcu_read_lock()` is essential before accessing the `task_struct`, because otherwise it can possibly be concurrently updated or deleted by writer threads. Unlike mutexes, which serialize access and block other threads (or CPUs) from entering the critical section, `rcu_read_lock()` is very lightweight and allows multiple readers to run concurrently. It doesn't enforce exclusive access; instead, it ensures that any memory being read won't be reclaimed until all readers exit their critical sections. In read-mostly scenarios (which are common in kernel code), `rcu_read_lock()` offers far lower overhead than a mutex.

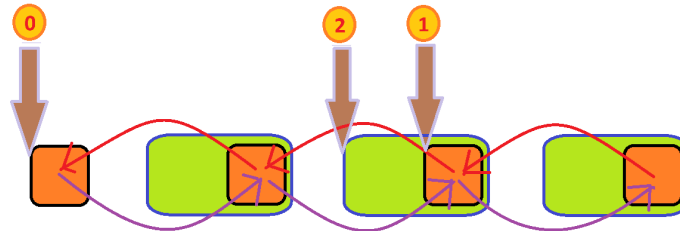
```

1 /**
2  * list_for_each - iterate over a list
3  * @pos: the &struct list_head to use as a loop cursor.
4  * @head: the head for your list.
5  */
6 #define list_for_each(pos, head) \
7     for (pos = (head)->next; !list_is_head(pos, (head)); pos = pos->next)

```

Listing 4: Kernel code snippet for `list_for_each` in `/include/linux/list.h`

Now next important thing is before we iterate over the monitoring list, we must lock the earlier defined lock as multiple threads simultaneously updating the list may lead to inconsistent state. Further, one essential thing to note is that `list_for_each` in `list.h` actually starts from actual head and not sentinel node [arrow 0 denotes it in below figure] that it takes as argument (this is because if we see in Listing 3, `pos` starts from `(head)->next`).



Next `list_entry` is used to get the pointer to starting location of actual node/container of list head (i.e. it shifts the pointer from 1 to 2).

```

1 /**
2  * list_entry - get the struct for this entry
3  * @ptr: the &struct list_head pointer.
4  * @type: the type of the struct this is embedded in.
5  * @member: the name of the list_head within the struct.
6  */
7 #define list_entry(ptr, type, member) \
8     container_of(ptr, type, member)

```

Listing 5: Kernel code snippet for `list_entry` in `/include/linux/list.h`

Now if `pid` already exists then return -23, else we move on allocate space in kernel space using `kmallocc`, The `GFP_KERNEL` flag is used to get free pages in kernel space.(As mentioned in link here it should be our first choice when trying to allocate in kernel space.). Next if allocations are successful we set the node parameters and add it to list using `list_add`(in `list.h`) then unlock mutex and return.

3.2 sys_fetch

```
1 SYSCALL_DEFINE2(fetch, struct per_proc_resource __user *, stats, pid_t, pid)
2 {
3     struct pid_node *node;
4     struct list_head *pos;
5     struct per_proc_resource kstats;
6     if (!stats || !access_ok(stats, sizeof(struct per_proc_resource))) {
7         return -22;
8     }
9     mutex_lock(&monitored_list_mutex);
10    list_for_each(pos, &monitored_pids_head_node) {
11        node = list_entry(pos, struct pid_node, next_prev_list);
12        if (node->proc_resource->pid == pid) {
13            kstats = *node->proc_resource;
14            kstats.heapsize = kstats.heapsize / (1024 * 1024);
15            mutex_unlock(&monitored_list_mutex);
16            if (copy_to_user(stats, &kstats, sizeof(struct per_proc_resource))) {
17                printk(KERN_ERR "sys_fetch: copy_to_user failed for pid %d\n", pid);
18                return -22;
19            }
20            return 0;
21        }
22    }
23    mutex_unlock(&monitored_list_mutex);
24    return -22;
25 }
```

Listing 6: Kernel code snippet for sys_fetch in /hello/hello.c

In the start of sys_fetch we sanity checks the "stats" which is pointer to the user space if that is not null and can be accessed. Further similar to the sys_register we lock the list and then traverse if we find the provided pid in the monitoring list we simply copy that to user space using copy_to_user (in /include/linux/uaccess.h) function (and internally since we are storing the heap-size in bytes we convert it into MB before copying it to location pointed by stats). If fetch went successful we return 0 otherwise in all other cases we return -22.

3.3 sys_deregister

```
1 SYSCALL_DEFINE1(deregister, pid_t, pid)
2 {
3     bool found = false;
4     struct pid_node *node;
5     struct list_head *pos, *n;
6     if (pid < 1) {
7         return -22;
8     }
9     mutex_lock(&monitored_list_mutex);
10    list_for_each_safe(pos, n, &monitored_pids_head_node) {
11        node = list_entry(pos, struct pid_node, next_prev_list);
12        if (node->proc_resource->pid == pid) {
13            list_del(&node->next_prev_list);
14            kfree(node->proc_resource);
15            kfree(node);
16            found = true;
17            break;
18        }
19    }
20    mutex_unlock(&monitored_list_mutex);
21    return found ? 0 : -3;
22 }
```

Listing 7: Kernel code snippet for sys_deregister in /hello/hello.c

Most of things are similar previous syscalls, the important thing is to iterate the list using `list_for_each_safe` which enables us to delete a node in between of a traversal as mentioned in kernel code of `list.h`, further we remove the node from monitoring list using `list_del` and before returning we free up allocated space for the previously registered node using `kfree`.

```

1 /**
2  * list_for_each_safe - iterate over a list safe against removal of list entry
3  * @pos:      the &struct list_head to use as a loop cursor.
4  * @n:       another &struct list_head to use as temporary storage
5  * @head:    the head for your list.
6  */
7 #define list_for_each_safe(pos, n, head) \
8     for (pos = (head)->next, n = pos->next; \
9         !list_is_head(pos, (head)); \
10         pos = n, n = pos->next)

```

Listing 8: Kernel code snippet for `list_for_each_safe` in `/include/linux/list.h`

4 Resource Usage Limiter

4.1 `sys_resource_cap`

```

1 SYSCALL_DEFINE3(resource_cap, pid_t, pid, long, heap_quota, long, file_quota) {
2     struct task_struct *task;
3     struct pid_node *node = NULL;
4     struct list_head *pos, *n;
5     bool should_kill = false;
6     int found = 0;
7     rcu_read_lock();
8     task = find_task_by_vpid(pid);
9     rcu_read_unlock();
10    if (!task) {
11        printk(KERN_INFO "Syscall sys_register: PID %d does not exist\n", pid);
12        return -3;
13    }
14    mutex_lock(&monitored_list_mutex);
15    list_for_each_safe(pos, n, &monitored_pids_head_node) {
16        node = list_entry(pos, struct pid_node, next_prev_list);
17        if (node->proc_resource->pid == pid) {
18            found = 1;
19            break;
20        }
21    }
22    if (!found) {
23        mutex_unlock(&monitored_list_mutex);
24        return -22;
25    }
26    if (node->proc_resource->quotas_defined) {
27        mutex_unlock(&monitored_list_mutex);
28        return -23;
29    }
30    node->proc_resource->heap_quota = heap_quota*1024*1024;
31    node->proc_resource->file_quota = file_quota;
32    node->proc_resource->quotas_defined = true;
33    if ((node->proc_resource->file_quota >= 0 && node->proc_resource->openfile_count > node->
34        proc_resource->file_quota)
35        || (node->proc_resource->heap_quota >= 0 && node->proc_resource->heapsize > node->proc_resource
36            ->heap_quota))
37    {
38        printk(KERN_INFO "File/Heap quota exceeded for PID %d. Sending SIGKILL.\n", pid);
39        list_del(&node->next_prev_list);
40        kfree(node->proc_resource);

```

```

39     kfree(node);
40     should_kill = true;
41 }
42 mutex_unlock(&monitored_list_mutex);
43 if(should_kill){
44     printk(KERN_INFO "RES_CAP killing the process...\n");
45     force_sig(SIGKILL);
46 }
47 return 0;
48 }

```

Listing 9: Kernel code snippet for sys_resource_cap in /hello/hello.c

Most of the strategies used to implement it are discussed in previous sections, only important thing to note if quotas/cap that we are setting up for the current process are already exceeded then we free the node corresponding to it and send the SIGKILL to current process using force_sig (in /kernel/signal.c)

4.2 sys_resource_reset

```

1 SYSCALL_DEFINE1(resource_reset, pid_t, pid)
2 {
3     struct task_struct *task;
4     struct pid_node *node = NULL;
5     struct list_head *pos;
6     rcu_read_lock();
7     task = find_task_by_vpid(pid);
8     rcu_read_unlock();
9     if (!task) {
10         printk(KERN_INFO "sys_resource_reset: PID %d does not exist\n", pid);
11         return -3;
12     }
13     mutex_lock(&monitored_list_mutex);
14     list_for_each(pos, &monitored_pids_head_node) {
15         node = list_entry(pos, struct pid_node, next_prev_list);
16         if (node->proc_resource->pid == pid) {
17             node->proc_resource->heap_quota = -1;
18             node->proc_resource->file_quota = -1;
19             node->proc_resource->quotas_defined = false;
20             mutex_unlock(&monitored_list_mutex);
21             printk(KERN_INFO "sys_resource_reset: Reset quotas for pid %d\n", pid);
22             return 0;
23         }
24     }
25     mutex_unlock(&monitored_list_mutex);
26     return -22;
27 }

```

Listing 10: Kernel code snippet for sys_resource_reset in /hello/hello.c

5 Management of File Count and Heap-size

5.1 Updation of File Count

```
1 static bool update_filecount(void) {
2     bool should_kill = false;
3     pid_t curr_pid = current->pid;
4     struct pid_node *node;
5     struct list_head *pos, *n;
6     mutex_lock(&monitored_list_mutex);
7     list_for_each_safe(pos, n, &monitored_pids_head_node) {
8         node = list_entry(pos, struct pid_node, next_prev_list);
9         if (node->proc_resource->pid == curr_pid) {
10             node->proc_resource->openfile_count++;
11             if (node->proc_resource->quotas_defined &&
12                 node->proc_resource->file_quota >= 0 &&
13                 node->proc_resource->openfile_count > node->proc_resource->file_quota) {
14                 printk(KERN_INFO "update_filecount: File quota exceeded for PID %d. Sending SIGKILL.\n",
15                     curr_pid);
16                 list_del(&node->next_prev_list);
17                 kfree(node->proc_resource);
18                 kfree(node);
19                 should_kill = true;
20             }
21             break;
22         }
23     }
24     mutex_unlock(&monitored_list_mutex);
25     return should_kill;
26 }
```

Listing 11: Kernel code snippet for file count in /fs/open.c

```
1 SYSCALL_DEFINE3(open, const char __user *,
2     filename, int, flags, umode_t, mode)
3 {
4     if (force_o_largefile())
5         flags |= O_LARGEFILE;
6     return do_sys_open(AT_FDCWD, filename, flags,
7         mode);
8 }
```

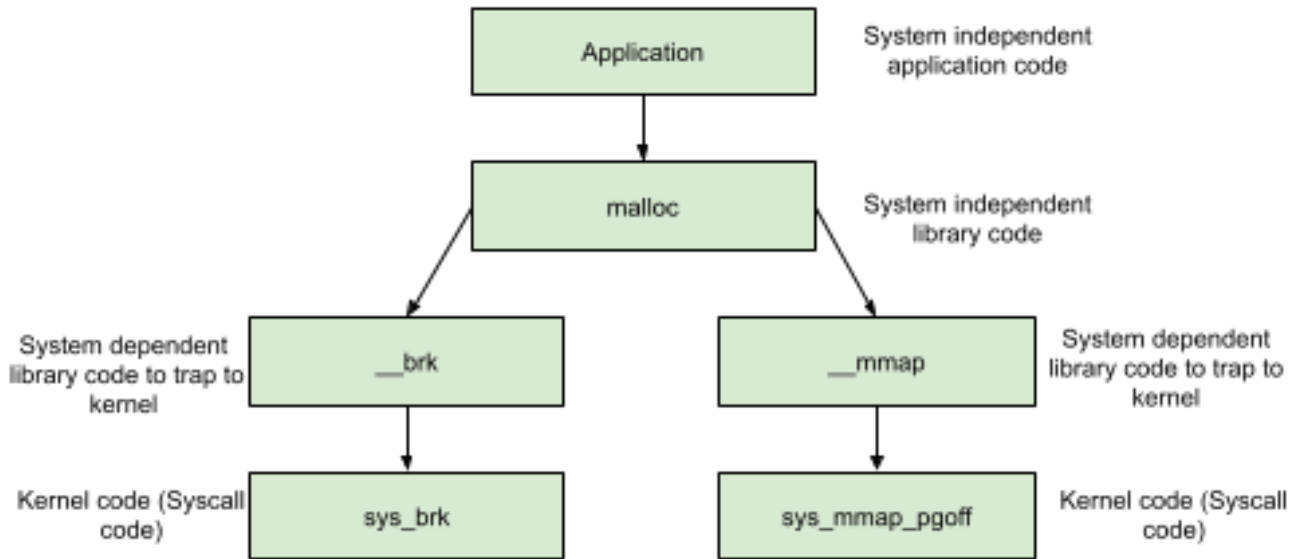
Listing 12: Old Kernel code for open syscall in fs/open.c

```
1 SYSCALL_DEFINE3(open, const char __user *,
2     filename, int, flags, umode_t, mode)
3 {
4     int fd;
5     if (force_o_largefile())
6         flags |= O_LARGEFILE;
7     fd = do_sys_open(AT_FDCWD, filename, flags,
8         mode);
9     if (fd >= 0) {
10         if (update_filecount()) {
11             printk(KERN_INFO "OPEN killing the
12                 process...\n");
13             force_sig(SIGKILL);
14         }
15     }
16     return fd;
17 }
```

Listing 13: New Kernel code for open syscall in fs/open.c

Here to account for the files opened by processes we call `update_filecount` function whenever `open/openat/openat2` syscalls are invoked (In above listings we have given the modification made to `open` syscall a similar modification has been made to `openat` & `openat2`), then the `update_filecount` goes and check if current process which invoked then `open` syscall is in monitoring list, if yes then goes and updates the file count for it, after which it checks if the current file count is within the bounds of quotas/caps set, if not in bounds then it kills the process after returning from the function call (indicated by setting `should_kill` to True)

5.2 Updation of Heap-size



```

1 static bool update_heapsize(unsigned long delta_bytes)
2 {
3     bool should_kill = false;
4     struct pid_node *node;
5     struct list_head *pos, *n;
6     pid_t curr_pid = current->pid;
7     unsigned long delta = delta_bytes;
8     if (delta <= 0)
9     {
10         return should_kill;
11     }
12     mutex_lock(&monitored_list_mutex);
13     list_for_each_safe(pos, n, &monitored_pids_head_node) {
14         node = list_entry(pos, struct pid_node, next_prev_list);
15         if (node->proc_resource->pid == curr_pid) {
16             printk(KERN_INFO "Updating Heap Size (Using MMAP) (For PID: %d): %lu \n", curr_pid,
17                 delta_bytes);
18             node->proc_resource->heapsize += delta;
19             if (node->proc_resource->quotas_defined &&
20                 node->proc_resource->heap_quota >= 0 &&
21                 node->proc_resource->heapsize > node->proc_resource->heap_quota) {
22                 printk(KERN_INFO "Heap quota exceeded for PID %d. Sending SIGKILL.\n", curr_pid);
23                 list_del(&node->next_prev_list);
24                 kfree(node->proc_resource);
25                 kfree(node);
26                 should_kill = true;
27             }
28             break;
29         }
30     }
31     mutex_unlock(&monitored_list_mutex);
32     return should_kill;
33 }

```

Listing 14: Kernel code snippet for accounting heap-size in /mm/mmap.c & /arch/x86/kernel/sys_x86_64.c


```

1 SYSCALL_DEFINE6(mmap, unsigned long, addr,
2   unsigned long, len,
3   unsigned long, prot, unsigned long, flags
4   ,
5   unsigned long, fd, unsigned long, off)
6 {
7     if (off & ~PAGE_MASK)
8         return -EINVAL;
9
10    return ksys_mmap_pgoff(addr, len, prot, flags
11        , fd, off >> PAGE_SHIFT);
12 }

```

Listing 15: Old Kernel code for mmap syscall in /arch/x86/kernel/sys_x86_64.c

```

1 SYSCALL_DEFINE1(brk, unsigned long, brk)
2 {
3     unsigned long newbrk, oldbrk, origbrk;
4     struct mm_struct *mm = current->mm;
5     struct vm_area_struct *brkvma, *next =
6         NULL;
7     unsigned long min_brk;
8     bool populate;
9     bool downgraded = false;
10    ...
11    success:
12    populate = newbrk > oldbrk && (mm->
13        def_flags & VM_LOCKED) != 0;
14    if (downgraded)
15        mmap_read_unlock(mm);
16    else
17        mmap_write_unlock(mm);
18    userfaultfd_unmap_complete(mm, &uf);
19    if (populate)
20        mm_populate(oldbrk, newbrk - oldbrk);
21    return brk;
22 out:
23    mmap_write_unlock(mm);
24    return origbrk;
25 }

```

Listing 17: Old Kernel code for brk syscall in /mm/mmap.c

```

1 SYSCALL_DEFINE6(mmap, unsigned long, addr,
2   unsigned long, len,
3   unsigned long, prot, unsigned long, flags
4   ,
5   unsigned long, fd, unsigned long, off)
6 {
7     unsigned long result;
8     if (off & ~PAGE_MASK)
9         return -EINVAL;
10    result = ksys_mmap_pgoff(addr, len, prot,
11        flags, fd, off >> PAGE_SHIFT);
12    if (!IS_ERR_VALUE(result) && ((flags & (
13        MAP_ANONYMOUS | MAP_PRIVATE)) == (
14        MAP_ANONYMOUS | MAP_PRIVATE)))
15    {
16        unsigned long aligned_len = PAGE_ALIGN(
17            len);
18        if (update_heapsize((long)aligned_len))
19        {
20            force_sig(SIGKILL);
21        }
22    }
23    return result;
24 }

```

Listing 16: New Kernel code for mmap syscall in /arch/x86/kernel/sys_x86_64.c

```

1 SYSCALL_DEFINE1(brk, unsigned long, brk)
2 {
3     unsigned long newbrk, oldbrk, origbrk;
4     struct mm_struct *mm = current->mm;
5     struct vm_area_struct *brkvma, *next =
6         NULL;
7     unsigned long min_brk;
8     unsigned long delta;
9     bool populate;
10    bool downgraded = false;
11    ...
12    success:
13    populate = newbrk > oldbrk && (mm->
14        def_flags & VM_LOCKED) != 0;
15    delta = newbrk - oldbrk;
16    if (downgraded)
17        mmap_read_unlock(mm);
18    else
19        mmap_write_unlock(mm);
20    userfaultfd_unmap_complete(mm, &uf);
21    if (populate)
22        mm_populate(oldbrk, newbrk - oldbrk);
23    if (update_heapsize(delta))
24    {
25        force_sig(SIGKILL);
26    }
27    return brk;
28 out:
29    mmap_write_unlock(mm);
30    return origbrk;
31 }

```

Listing 18: New Kernel code for brk syscall in /mm/mmap.c

Before starting a discussion on implementation, we can look into how malloc allocates the heap. Here is statement from the documentation

Malloc Details

Normally, `malloc()` allocates memory from the heap, and adjusts the size of the heap as required, using `sbrk(2)`. When allocating blocks of memory larger than `MMAP_THRESHOLD` bytes, the glibc `malloc()` implementation allocates the memory as a private anonymous mapping using `mmap(2)`. `MMAP_THRESHOLD` is 128 kB by default, but is adjustable using `mallopt(3)`.

This explains why we need to filter private anonymous mappings while accounting for heap regions that get allocated using `mmap`. If we see `update_heapsize` is mostly similar to the `update_filecount` only, and there invocation in either `brk/mmap` syscall is mostly similar to previous section, only thing to note is we are passing `PAGE_ALIGN` lengths to increment our heap-size.

6 Other Details of Implementation

6.1 Modifications in 'syscall_64.tbl' & 'syscalls.h'

```
375 451 common hello sys_hello
376 452 common register sys_register
377 453 common fetch sys_fetch
378 454 common deregister sys_deregister
379 455 common resource_cap sys_resource_cap
380 456 common resource_reset sys_resource_reset
381
```

Figure 1: Modifications to syscall_64.tbl

```
298 /*
299  * These syscall function prototypes are kept in the same order as
300  * include/uapi/asm-generic/unistd.h. Architecture specific entries go below,
301  * followed by deprecated or obsolete system calls.
302  *
303  * Please note that these prototypes here are only provided for information
304  * purposes, for static analysis, and for linking from the syscall table.
305  * These functions should not be called elsewhere from kernel code.
306  *
307  * As the syscall calling convention may be different from the default
308  * for architectures overriding the syscall calling convention, do not
309  * include the prototypes if CONFIG_ARCH_HAS_SYSCALL_WRAPPER is enabled.
310  */
311 #ifndef CONFIG_ARCH_HAS_SYSCALL_WRAPPER
312 asmlinkage long sys_hello(void);
313 asmlinkage int sys_register(pid_t pid);
314 asmlinkage int sys_fetch(struct per_proc_resource __user *stats, pid_t pid);
315 asmlinkage int sys_deregister(pid_t pid);
316 asmlinkage int sys_resource_cap(pid_t pid, long heap_quota, long file_quota);
317 asmlinkage int sys_resource_reset(pid_t pid);
```

Figure 2: Modifications to syscalls.h

7 Extending the implementation to handle Process and Thread Group

The submission we have made is 'thread specific implementation' here we have discussed some ideas to extend them to 'process specific implementation'. For simplification assume that we are given input only pid of main thread of thread group whose 'tgid' is equal to 'pid'. So when we will register/fetch/deregister/set resource cap we will do it for the entire group of threads in that process/thread group.

7.1 Modification needed to Node structure

```
1 struct per_proc_resource {
2     pid_t pid;
3     unsigned long heapsize;
4     unsigned long openfile_count;
5     long heap_quota;
6     long file_quota;
7     bool quotas_defined;
8     pid_t tgid;
9 };
```

Listing 19: Suggested Update For Node Structure

keeping track of tgid will help us enable faster deletion when pid of main thread of thread group is provided.

7.2 Modifications needed to sys_register

```
1 SYSCALL_DEFINE1(register, pid_t, pid)
2 {
3     struct task_struct *task, *thread, *leader;
4     struct pid_node *node;
5     struct list_head *pos;
6     int ret = 0;
7     printk(KERN_INFO "Syscall sys_register called with pid: %d\n", pid);
8     if (pid < 1) {
9         printk(KERN_INFO "sys_register: PID %d is less than 1\n", pid);
10        return -22;
11    }
12    rcu_read_lock();
13    task = find_task_by_vpid(pid);
14    rcu_read_unlock();
15    if (!task) {
16        printk(KERN_INFO "sys_register: PID %d does not exist\n", pid);
17        return -3;
18    }
19    leader = task->group_leader;
20    mutex_lock(&monitored_list_mutex);
21    for_each_thread(leader, thread) {
22        list_for_each(pos, &monitored_pids_head_node) {
23            node = list_entry(pos, struct pid_node, next_prev_list);
24            if (node->proc_resource->pid == thread->pid) {
25                mutex_unlock(&monitored_list_mutex);
26                printk(KERN_INFO "sys_register: Thread %d (in group %d) is already monitored\n",
27                    thread->pid, leader->pid);
28                return -23;
29            }
30        }
31    }
32    for_each_thread(leader, thread) {
33        node = kmalloc(sizeof(struct pid_node), GFP_KERNEL);
34        if (!node) {
35            ret = -ENOMEM;
36            goto rollback;
```

```

37     }
38     node->proc_resource = kmalloc(sizeof(struct per_proc_resource), GFP_KERNEL);
39     if (!node->proc_resource) {
40         kfree(node);
41         ret = -ENOMEM;
42         goto rollback;
43     }
44     node->proc_resource->pid = thread->pid;
45     node->proc_resource->tgid = leader->pid;
46     node->proc_resource->heapsize = 0;
47     node->proc_resource->openfile_count = 0;
48     node->proc_resource->heap_quota = -1;
49     node->proc_resource->file_quota = -1;
50     node->proc_resource->quotas_defined = false;
51     list_add(&node->next_prev_list, &monitored_pids_head_node);
52     printk(KERN_INFO "sys_register: Registered thread %d in group %d\n",
53            thread->pid, leader->pid);
54 }
55 mutex_unlock(&monitored_list_mutex);
56 return 0;
57 rollback:
58 {
59     struct list_head *q, *tmp;
60     list_for_each_safe(q, tmp, &monitored_pids_head_node) {
61         node = list_entry(q, struct pid_node, next_prev_list);
62         if (node->proc_resource->tgid == leader->pid) {
63             list_del(q);
64             kfree(node->proc_resource);
65             kfree(node);
66         }
67     }
68 }
69 mutex_unlock(&monitored_list_mutex);
70 printk(KERN_ERR "sys_register: Memory allocation failed for thread group %d\n",
71        leader->pid);
72 return ret;
73 }

```

Listing 20: Suggested Update For sys_register

Only main differences includes getting task struct to leader thread, since we assume pid of leader thread will be given as input then we can expect leader=task holds further we can use for `_each_thread` (defined at `/include/linux/sched/signal.h`) to go through all the threads of that thread group and add them to the monitoring list. For the current convention we assume if some thread already we discard entire thread group and free up any allocation made for this group (here tgid stored in node helps in deletion).

7.3 Modifications needed to sys_fetch

```

1  SYSCALL_DEFINE2(fetch, struct per_proc_resource __user *, stats, pid_t, pid)
2  {
3      struct pid_node *node;
4      struct list_head *pos;
5      struct per_proc_resource kstats = {0};
6      bool found = false;
7      if (!stats || !access_ok(stats, sizeof(struct per_proc_resource))) {
8          return -22;
9      }
10     kstats.pid = pid;
11     kstats.tgid = pid;
12     kstats.heapsize = 0;
13     kstats.openfile_count = 0;
14     kstats.heap_quota = -1;
15     kstats.file_quota = -1;

```

```

16 kstats.quotas_defined = false;
17 mutex_lock(&monitored_list_mutex);
18 list_for_each(pos, &monitored_pids_head_node) {
19     node = list_entry(pos, struct pid_node, next_prev_list);
20     if (node->proc_resource->tgid == pid) {
21         kstats.heapsize += node->proc_resource->heapsize;
22         kstats.openfile_count += node->proc_resource->openfile_count;
23         if (!found) {
24             kstats.heap_quota = node->proc_resource->heap_quota;
25             kstats.file_quota = node->proc_resource->file_quota;
26             kstats.quotas_defined = node->proc_resource->quotas_defined;
27             found = true;
28         }
29     }
30 }
31 mutex_unlock(&monitored_list_mutex);
32 if (!found) {
33     printk(KERN_INFO "sys_fetch: No monitored threads found for tgid %d\n", pid);
34     return -22;
35 }
36 kstats.heapsize = kstats.heapsize / (1024 * 1024);
37 if (copy_to_user(stats, &kstats, sizeof(struct per_proc_resource))) {
38     printk(KERN_ERR "sys_fetch: copy_to_user failed for tgid %d\n", pid);
39     return -22;
40 }
41 return 0;
42 }

```

Listing 21: Suggested Update For sys_fetch

Now again there is a similar update that may be needed for sys_fetch. We go through the monitoring list and check nodes where the tgid matches input pid (assuming input pid is leader thread's pid) and then accumulate the file opened and heap allocated and then copy the accumulation to the location pointed by pointer stats.

7.4 Other Concerns related to new implementation

A more or less similar approach to complete the remaining syscalls, and handling updates to file-count and heap-size is mostly same, only thing that will change is killing condition, for which we will need to do an accumulation similar to in previous section and check against cap/quotas. But more important thing that we will need to account for is that some threads may get added to thread group beyond the point when the process got registered and some threads in thread group may finish the execution before the entire process/thread group has actually come to termination. For which we will need to modify fork(), execv() and exit() system calls, when fork() and execv() are used to create a new thread (and not a new process) we can look at it's leader's pid, if that is part of monitoring list we need to add this newly created thread to the list while when thread exits if that thread in monitoring list we need to remove that from monitoring list. So there significant part that needs to taken into account before we extend the current implementation into 'process specific implementation'.