

Container Manager

In this section, we had to implement a container manager data structure which will be in charge of all the containers and will be modified by system calls to create, read, update and delete the containers respectively.

Container manager has the following data structure for generality and storing logs:

```
1 containerStruct containerManager = {  
2     .containerIDs = { 0 },  
3     .numActive = 0,  
4     .procIDs = {{ 0 }},  
5     .notAllowed = {{ 0 }},  
6     .systemCalls = {{ 0 }}  
7 };
```

Algorithm 1: Container Manager

Where container id's is the id's of the current containers, procId's are the processes running in the respective containers, not-allowed are the system calls not allowed in that container, because every process's system call has to pass verification before initiation, and systemCalls are the system calls executed by processes inside. (To ensure that container manager can keep track of system calls and other data related to container).

System calls were created with the following signature:-

- `uint create_container(void)`
It creates a new container with the default parameters from the container manager.
- `uint destroy_container(uint container_id)`
It destroys the specified container and handles clean-up actions like deleting the files and clearing out memory and processes of that container.
- `uint join_container(uint container_id)`
It is used by a process to join a particular existing container. Implemented using fork command. Scheduled by the virtual scheduler later
- `uint leave_container(void)`
It is used by a process to leave the container and go back to the host container (by default) or be cleared.

```
1 // MOD-3 : Container syscalls  
2 #define SYS_create_container    37  
3 #define SYS_destroy_container  38  
4 #define SYS_join_container     39
```

```
5 #define SYS_leave_container      40
6 #define SYS_proc_stat_container 41
7 #define SYS_scheduler_log_on    42
8 #define SYS_scheduler_log_off   43
9 #define SYS_cid                 44
10 #define SYS_memory_log_on       45
11 #define SYS_memory_log_off      46
12 #define SYS_memory_gva          47
```

Algorithm 2: System Calls added

Virtual Scheduler

As we are mimicking the Container behaviour, we have the following mechanism for correctly scheduling the processes of a container. We have edited the Process Control block to include the container information in the PCB itself.

The default scheduler in xv6 is modified to take into account the PCB's container id field and accordingly schedule the processes. We are still following the Round Robin policy.

The new process Control block is :

```
1 struct proc {
2     uint sz;                // Size of process memory (bytes)
3     pde_t* pgdir;           // Page table
4     char *kstack;           // Bottom of kernel stack for this process
5     enum procstate state;   // Process state
6     int pid;                // Process ID
7     struct proc *parent;    // Parent process
8     struct trapframe *tf;   // Trap frame for current syscall
9     struct trapframe *Oldtf; // MOD-1 : Old trap frame
10    struct context *context; // swtch() here to run process
11    void *chan;              // If non-zero, sleeping on chan
12    int killed;              // If non-zero, have been killed
13    struct file *ofile[NOFILE]; // Open files
14    struct inode *cwd;       // Current directory
15    char name[16];           // Process name (debugging)
16    sig_handler sig_handler; // MOD-1 : Signal handler for process
17    char msg[message_size];  // MOD-1 : Signal msg
18    int disableSignals;      // MOD-1 : Disable signals when currently ←
                             // processing one
19    int interrupt;
20    int containerID;         // MOD-3 : Container ID
21    enum v_procstate v_state; // MOD-3 : Virtual process state
22 };
```

Algorithm 3: Process Control Block

And it is scheduled in the following manner :-

```
1 p->v_state = V_RUNNING; // MOD-3 : Make this process virtually running
2 // MOD-3 : Print schedule if log on
3 if(container.containerIDs[p->containerID] == 2)
4     cprintf("Container:%d\tScheduling process:%d\n", p->containerID, p->pid)←
5     ;
6 // Other process which was v_running in this container to shift to ←
7 v_runnable
8 for(p1 = ptable.proc; p1 < &ptable.proc[NPROC]; p1++){
9     if(p1->containerID != -1 && p1->containerID == p->containerID && p1->←
10     v_state == V_RUNNING && p1->pid != p->pid){
11         // checking the schedule status here and virtual scheduling
12         p1->v_state = V_RUNNABLE;
13     }
14 }
```

Algorithm 4: Additions in Scheduler

System Calls Handling

There were many ways to handle the invocation of system calls and their execution such that they are logically separate in different containers and were logged to the containers and the respective container could keep track of the system call invoked from its virtualized processes.

The second challenge here was to keep certain system calls being allowed and certain not allowed from a particular container. Verification of system calls as mentioned in the document.

We managed both the expectations by maintaining a table of not-allowed system calls in the container itself and changed the code of the default syscall(id) function itself to check for the appropriate parameters.

The code change's required was in the container manager for storing the not-allowed calls which is shown in Container Manager data structure. Other was in syscall invocation function, the checking for allowed and not allowed system calls:

```
1 for(int j = 0; j < 100; j++){
2     if(container.notAllowed[curproc->containerID][j] == num)
3         cprintf("Syscall %d not allowed in container with ID %d\n",
4             num, curproc->containerID); // syscall not allowed from this ←
5         container
6 }
```

Algorithm 5: Syscalls Changes

Resource Isolation

We've implemented both parts of this, the page table address translation and the virtual file system.

For the first part, we are keeping a lookup of the addresses used by different processes in the container and thus calculating the GVA, HVA. Calling this function from the changed malloc function. Here are the necessary code changes. :

```
1 int sys_memory_gva(int size, int hva)
2 {
3     argint(0, &size);
4     argint(1, &hva);
5     int id = myproc()->containerID;
6     int addr = container.gva[id];
7     container.gva[id] = addr + size;
8     if(container.memorylog[id] == 1)
9         cprintf("GVA : %d, HVA : %x\n", addr, hva);
10    return addr;
11 }
```

Algorithm 6: Address calculation and translation

For the 2nd part in the virtual file system our model is to create directories for different containers and keep them in isolation, which also provides independence for copy on write.

Our ls prints all the host's files and not the container's directories (by design). When ls through container, due to COW, it shows host's file but on opening/read/write it deals with the local file. This is done in this way because we want host file to remain visible. The relevant code is in split into different functions in the file and too big to include the core part here. When we're opening file for read, we check if the file is in container, if not then pick host file. For write we check if it's in container directory, if it's there, then write at the end of that file. If not, then check it's in the host, then copy the host file to local directory of container and start writing to it.

```
1 ...
2 if(fd < 0 && mode & O_RDWR){ // Not in container and write -> copy to my ←
3     container
4     fd = open(name, O_CREATE|O_RDWR);
5     int fd_host = open(filename, O_RDONLY);
6     while((n = read(fd_host, buf, sizeof(buf))) > 0){
7         write(fd, buf, n);
8     }
9     ...
```

Algorithm 7: Copy on write mechanism

Testing

We tested with two testing scripts in the user level program, `container_test.c` and `cts.c`.

In the script `cts`, a container is joined and a file is there, `file1`, and when you leave the container, the `file1` which opens the host version. This among other function like COW are tested in this script. The script and output are given the the folder.

These scripts and their outputs are given in the root directory.

- `ls()` : Displaying files correctly, detailed output given in folder.
- `ps()` : Running from `container_test` script so showing it's name in the process name column.

```
1 Init 4
2 Child proc 4, num 1
3 Init 5
4 Init 6
5 Child proInit 7
6 Init 8
7 Pid      Name      MemSize Killed  State
8 1        init       12288   0       SLEEPING
9 2        sh         16384   0       SLEEPING
10 3        container_test 16384   0       RUNNING
11 4        container_test 16384   0       SLEEPING
12 5        container_test 16384   0       RUNNING
13 6        container_test 16384   0       RUNNABLE
14 7        container_test 16384   0       RUNNABLE
15 8        container_test 16384   0       RUNNABLE
16 c 6, num 3
17 Child proc 7, num 4
18 Child proc 8, num 5
19 CParent proc 3, num 0
20 Child proc 5, num 2
21 pid:4    name:container_test    state:WAITING
22 pid:5    name:container_test    state:WAITING
23 pid:6    name:container_test    state:RUNNING
24 pid:7    name:container_test    state:RUNNING
25 pid:8    name:container_test    state:RUNNING
26 Container:0    Scheduling process:5
27 ... (Other schedules)
28 Container:0    Scheduling process:4
```

Algorithm 8: `ps` output on `xv6`