

COL331 Assignment 2

DMS/RMS scheduling and PCP

Pranay Shah

Compiled on: August 19, 2023

1 Introduction

2 Deadline Monotonic Scheduling

In DMS, the task τ_i arrives periodically with time period P_i , execution time T_i and deadline D_i . The deadline monotonic scheduling algorithm is a static scheduling algorithm, in which the priorities are assigned to tasks based on their relative deadlines.

2.1 Implementation

We simulate the DMS scheduling in the following way: First add all the registered processes to the real time scheduling class. This will prevent other CFS user processes to intervene. Now, keep all the processes except the one with the highest priority in the "TASK_STOPPED" state. Now, the kernel scheduler won't switch out this process assuming that there are no processes which are not registered for DMS in a higher kernel scheduling class.

We maintain two main queues implemented using kernel `rb_tree` structs: one for all tasks registered under dms scheduling and another for all the tasks which are now in the "ready" state. The ready queue is a subset of the registered queue which corresponds to the list of tasks which have arrived and waiting to execute. All the processes in the registered queue are in the "TASK_STOPPED" state except for one (which is the currently running process).

2.2 The PCB

We create a Process control block for each task registered under dms. The fields of which are as shown in the following snippet.

```
1 struct sched_dms_entity{
2     struct rb_node ready_node;
3     struct rb_node registered_node;
4     pid_t task_pid;
5
6     unsigned int dms_exec_time;
7     unsigned int dms_deadline;
8     unsigned int dms_period;
9     bool ready; // false -> yielded, true -> ready
10    bool isfirst;
11
12    struct timer_list timer;
13};
```

Where `ready_node` and `registered_node` are `rb_nodes` in our registered and ready run-queues.

2.3 Timer Callbacks

Each sched entity has a timer struct associated with it. We can use them to execute a call back function periodically - each execution corresponds to the arrival of the next instance of a task. The call back is as shown in the following snippet:

```
1 void timer_callback(struct timer_list *data)
2 {
3     struct sched_dms_entity *sdmse;
4     unsigned long flags;
5     sdmse = container_of(data, struct sched_dms_entity, timer);
6     if(!find_task_by_vpid(sdmse->task_pid))
7         return;
8
9     spin_lock_irqsave(&dms_rq.rq_lock, flags);
10    if(!sdmse->ready){ // if true, not completed last task; handle this
11        dms_enqueue_ready(sdmse);
12        sdmse->ready = true;
13    }
14    spin_unlock_irqrestore(&dms_rq.rq_lock, flags);
15
16    dms_schedule();
17    mod_timer(data, jiffies + msecs_to_jiffies(sdmse->dms_period));
18 }
```

Here, we change the state of the process to the ready state and insert it to the ready queue. Then we call the `dms_schedule` method which context switches (if necessary) to the new highest priority process. We also set the next time when the call back will be run using `mod_timer`.

2.4 The Scheduling function

The `dms_schedule` method is as shown in the following snippet:

```
1 void dms_schedule(void)
2 {
3     struct rb_node *left;
4     struct sched_dms_entity *nxt_se;
5     struct task_struct *nxt_task, *cur_task;
6     unsigned long flags;
7
8     spin_lock_irqsave(&dms_rq.rq_lock, flags);
9     left = rb_first_cached(&dms_rq.ready_tree);
10    if(!left)
11        goto unlock;
12    nxt_se = container_of(left, struct sched_dms_entity, ready_node);
13    if(nxt_se->task_pid == dms_rq.curr_pid)
14        goto unlock;
```

```

15     nxt_task = find_task_by_vpid(nxt_se->task_pid);
16
17     cur_task = find_task_by_vpid(dms_rq.curr_pid);
18     if(cur_task)
19         send_sig(SIGSTOP, cur_task, 1); // check if task has been killed
20
21     dms_rq.curr_pid = nxt_se->task_pid;
22     send_sig(SIGCONT, nxt_task, 1);
23
24 unlock:
25     spin_unlock_irqrestore(&dms_rq.rq_lock, flags);
26     /* printk("Done Scheduling\n"); */
27 }

```

There are multiple things going on here, but the main gist is that we fetch the highest priority ready task from the ready queue, check if its the currently running process. If not, we stop the running process and start the next running process. The stopping and starting is done by sending the SIGSTOP and SIGCONT signals respectively.

2.5 Yielding

When a process is done with its task, it calls the sys_yield syscall, the Implementation of which is as shown:

```

1 SYSCALL_DEFINE1(yield, pid_t, pid)
2 {
3     ...
4     send_sig(SIGSTOP, cur_task, 1);
5
6     if(sdmse->isfirst){
7         mod_timer(&sdmse->timer, jiffies + msecs_to_jiffies(sdmse->dms_period));
8         sdmse->isfirst = false;
9         return 0;
10    }

```

When a process yields, we stop it and dequeue it from the ready queue (if it was present in it). If a dequeue occurred then rescheduling is required and hence we call `dms_schedule` in the end.

Thus the cycle is complete. Tasks arrive using timer callbacks, get scheduled according to their Deadlines, finish their execution and yield.

2.6 Schedulability test

We perform the Schedulability test according to the algorithm provided by Audsley. For this, we go through all the items in the registered runqueue in the sorted order. The following snippet shows the implementation.

```

1   t = 0;
2   for(i = 0; i < dms_rq.nr_registered; i++){
3       t += exec_times[i];
4       go = true;
5       while(go){
6           interference = 0;
7           for(j = 0; j < i; j++){
8               interference += ((t + periods[j] - 1) / periods[j]) *
                               exec_times[j];
9               if(interference + exec_times[i] <= t)
10                  go = false;
11               else
12                  t = interference + exec_times[i];
13
14               if(t > deadlines[i]){
15                   return false;
16               }
17           }
18       }

```

2.7 RMS

We can clearly see that RMS is a subset of DMS and hence, we route `sys_rm_register` to `sys_dm_register` which takes care of all processes scheduled in the same way.

3 Priority Ceiling Protocol

3.1 Global variables

We maintain a global list of all resources that have been mapped by at least one DMS registered task. Corresponding to each resource we store some information regarding which processes have mapped to it its ceiling priority.

3.2 Integrating with DMS

To integrate this protocol with our DMS Implementation, we first notice that the priority of a task may change over time (inheritance) and there can be two tasks with the same priority. To maintain the ready and registered runqueues in DMS, we need to modify the comparison functions given to the corresponding `rb_trees`. To this extent, we introduce a new struct `sched_prio` as shown in the snippet below.

```

1 struct sdmse_prio{
2 #ifdef CONFIG_PRIO_CEIL_PROTOCOL_DMS
3     struct list_head inherited_priorities_entry;
4     // history stuff

```

```

5     bool inherited;
6     unsigned int inherited_while_holding_rid;
7 #endif
8     pid_t pid;
9     unsigned int dms_deadline;
10 };

```

When the boolean value `inherited` is false, the comparison of two `prio` structs works as in the DMS case: compare deadlines and use pids to break ties. The problem arises when both pid and deadline are equal for the two `prio` structs. This implies that one of them must be inherited and that `prio` has higher priority.

We also add, the following new fields in our PCB for dms:

```

1 #ifndef CONFIG_PRIO_CEIL_PROTOCOL_DMS
2     struct list_head locked_resources_head;
3     struct list_head inherited_priorities_head;
4 #endif

```

The `inherited_priorities_head` corresponds to a stack of inherited priorities. each entry has a corresponding resource id that the task held when performing the inheritance. This is used later on when releasing a resource. The sched prio of a sched dms entity is found in the following way:

1. Check if the `inherited_priorities_head` list is empty or not.
2. If empty then set `inherited` false and return `sdmse prio` with original pid and deadline.
3. otherwise, The current prio is the top entry in the list. return that.

Corresponding function is:

```

1 inline struct sdmse_prio get_sdmse_prio(struct sched_dms_entity *a){
2     struct sdmse_prio cur_prio = {.inherited = false, .pid = a->task_pid,
3     .dms_deadline = a->dms_deadline};
4     if(!list_empty(&a->inherited_priorities_head)){
5         cur_prio = *list_first_entry(&a->inherited_priorities_head, struct
6         sdmse_prio, inherited_priorities_entry);
7         return cur_prio;
8     }
9     return cur_prio;
10 }

```

Using these methods and making small changes in other functions, we can continue using our dms system and build PCP on top of it.

3.3 PCP lock resource

We maintain a global variable `system_ceiling_setter` which is the pid corresponding to the current system ceiling setter. We also maintain a `sdmse prio` struct which corresponds to

the system ceiling priority. The following snippet then follows the PCP.

```
1 SYSCALL_DEFINE2(pcp_lock, pid_t, pid, unsigned int, rid)
2 {
3     ...
4     if(pid == system_ceiling_setter){
5         // I am the ceiling setter
6         goto giveresource;
7     }
8     struct sdmse_prio cur_prio = get_sdmse_prio(sdmse);
9     if(__sdmse_prio_less(cur_prio, system_ceiling_prio)){
10        // higher priority than ceiling
11        // make this the system ceiling setter and give it the resource.
12        system_ceiling_setter = pid;
13        goto giveresource;
14    }
15    else{
16        struct sched_dms_entity *ceil_setter_sdmse;
17        ceil_setter_sdmse = get_sdmse_from_pid_registered(system_ceiling_setter);
18        struct sdmse_prio csprio = get_sdmse_prio(ceil_setter_sdmse);
19        if(__sdmse_prio_less(csprio, cur_prio)){
20            goto rejectresource;
21        }
22        // inherit priority. dequeue, push inherited prio, enqueue. reschedule.
23        spin_lock(&dms_rq.rq_lock);
24        dms_dequeue_ready(ceil_setter_sdmse);
25        inherit_prio(ceil_setter_sdmse, cur_prio);
26        dms_enqueue_ready(ceil_setter_sdmse);
27        spin_unlock(&dms_rq.rq_lock);
28
29        dms_schedule();
30
31        goto rejectresource;
32    }
33
34    giveresource:
35        pcpr->locked = true;
36        calculate_system_ceiling();
37        return 0;
38
39    rejectresource:
40        return -22;
41 }
```

The first if condition corresponds to the first clause: when the requesting task is the system ceiling setter. The second condition corresponds to the second clause: when the priority of requesting task is higher than the system ceiling. In inherit prio, we push the priority of higher priority task to the stack of the system ceiling setter. This is then used when

comparing it to other tasks in the ready queue.

3.4 PCP Unlock Resource

When unlocking a resource, We need to set the priority back to the one which we had when getting the resource. This is done by popping the priority inheritance stack and comparing the resource ids that the task had locked when making the corresponding inheritance. Once all such entries have been popped, we can then use the `get_sdmse_prio` function to get the priority of this task. This follows from the assumption that all the resource acquisitions are nested and perfectly balanced. The following snippet shows the implementation of the above idea.

```
1 SYSCALL_DEFINE2(pcp_unlock, pid_t, pid, unsigned int, rid){
2     ...
3     dms_dequeue_ready(sdmse);
4     struct sdmse_prio *cur_prio;
5     while(!list_empty(&sdmse->inherited_priorities_head)){
6         cur_prio = list_pop_entry(&sdmse->inherited_priorities_head, struct
7             sdmse_prio, inherited_priorities_entry);
8         if(cur_prio->inherited_while_holding_rid != rid){
9             list_push_entry(cur_prio, &sdmse->inherited_priorities_head,
10                 inherited_priorities_entry);
11             break;
12         }
13         kfree(cur_prio);
14     }
15     dms_enqueue_ready(sdmse);
16
17     struct pcpr_resource *pcpr = get_pcpr_from_rid(rid);
18     pcpr->locked = false;
19
20     calculate_system_ceiling();
21
22     return 0;
23 }
```

4 Novelties

4.1 Runqueue Cleaning

Some tasks might get killed before removing themselves using the `sys_remove` syscall. They will interfere with our Schedulability checks and also cause runtime errors if not cleaned properly. Also, memory leaking can occur. I have gracefully handled such situations by placing safety checks everywhere and cleaning the runqueue whenever possible - Every time a process registers itself, whenever `sys_list` has been called etc.

4.2 System Call to clear pcp global variables

I create a new system call `pcp_clear`, that clears out all the global variables used during pcp. This is useful as we may need to reuse the pcp structures and this requires a way to clear them.

4.3 KConfig Options

Both the implementations are made completely modular such that they can be turned on or off by changing the `.config` file. The pcp config option depends on the dms config option. prambal