

הטכניון – מכון טכנולוגי לישראל
מעבדה במערכות הפעלה 046210
תרגיל בית מס' 3

תאריך הגשה: 21.1.2016, עד 23:55

Introduction	3
Working Environment	3
Compiling the Kernel	3
Detailed Description	4
Notes	4
Background Information	5
Task States	5
Task Scheduling	5
Kernel Timing	7
Modified System Call API	8
Useful Information	9
Testing Your Custom Kernel	9
Tips for the Solution	10
Submission Procedure	11
Emphasis Regarding Grade	11

Introduction

In the previous assignment you have implemented a TODO tasks system for processes. In this assignment you will learn about the scheduling algorithm and how to update its policy.

Your mission in this assignment is to add scheduling mechanism to the TODO system. Each TODO task receives a deadline. Late TODO tasks are TODO tasks whose deadline has passed. Your scheduling mechanism should identify processes that have late TODO tasks, and 'punish' them.

Working Environment

You will be working on the same REDHAT 8.0 Linux virtual machine, as in the previous assignments.

Compiling the Kernel

In this assignment you will apply modifications to the Linux kernel. Errors in the kernel can render the machine unusable. Therefore you will apply the changes to a 'custom' kernel. The sources of the custom kernel can be found in `/usr/src/linux-2.4.18-14custom`. All files that you will work with are in this directory. This kernel has already been configured and compiled. Below are the steps needed for recompiling the kernel after applying your changes:

1. Make your changes to the kernel source file(s).
2. Invoke **`cd /usr/src/linux-2.4.18-14custom`**
3. Invoke **`make bzImage`**. The bzImage is the compressed kernel image created with command **`make bzImage`** during kernel compilation. The name bzImage stands for "Big Zimage". Both zImage and bzImage are compressed with gzip. The kernel includes a mini-gunzip to uncompress the kernel and boot into it.
4. Invoke **`make modules`**
5. Invoke **`make modules_install`**
6. Invoke **`make install`**
7. Invoke **`cd /boot`**
8. Invoke **`mkinitrd -f 2.4.18-14custom.img 2.4.18-14custom`**
9. Invoke **`reboot`**. This command will restart the machine.
10. After rebooting choose "custom kernel" in the Grub menu.

The system should boot properly with your new custom kernel.

Important Note: Steps 4 & 5 are necessary only in case you touched any header files (*.h & *.S) since the last time you compiled the kernel. If you modify only kernel source files (*.c) then you can skip these steps and save compilation time.

Detailed Description

This assignment builds on the previous assignment. For this work you should use the TODO task system you implemented in the previous assignment and implement the scheduling algorithm for it. The scheduling algorithm includes the following components:

1. An algorithm that identifies late TODO tasks.
2. A scheduling policy that punishes processes that have late TODO tasks.
3. A system call interface that sets deadlines to the TODO tasks.

Any time a TODO task is added to a process, a deadline is attached to it. The TODO task's deadline is measured in seconds. The TODO tasks list should be ordered according to the deadlines, from the earliest deadline to the latest. When the deadline of an uncompleted TODO task expires (its deadline is smaller than the current time) the TODO task is declared as late. Once a TODO task is declared as late it is removed from the list of TODO tasks. The process, to which the TODO task belongs, is put to sleep for a predetermined time: 60 seconds.

Below are the modifications to the system calls, required for the deadline mechanism:

1. **sys_add_TODO(pid, TODO_description, description_size, TODO_deadline)**: Add a TODO to the TODO's queue of a process identified by **pid**. **TODO_description** is a string of chars of size **description_size**. **TODO_deadline** is the deadline of the TODO task in seconds. A process can add a TODO only to itself or to its descendants.
2. **sys_read_TODO(pid, TODO_index, TODO_description, TODO_deadline, status)**: Get the description of TODO identified by its position, **TODO_index**, in the TODO's queue of process **pid**. The description is returned in **TODO_description**. The deadline of the TODO task is returned in **TODO_deadline**. The status of the TODO is returned in **status**. A process can query the description of a TODO that belongs to itself or to one of its descendants.

The above system calls replace the corresponding system calls that you implemented in the previous assignment. You are required to implement both the system call and its wrapper function (wrapper functions simplify the invocation of system calls from user space). Detailed description of the updated system call and its wrapper function is given in a later section.

Notes

1. The deadline time is set in seconds relative to the value returned by the **C** function **time()**. For example if you want to set a deadline 60 seconds away, and the value returned by **time()** is 100000, then the value you should use in the function **sys_add_TODO** is 100060.
2. The **C** function **time()** can be called only by User Mode applications. To get the current time from inside the kernel you can use the **CURRENT_TIME** macro defined in "sched.h".
3. The TODO tasks list is ordered according to the deadlines of the TODO tasks. Each time a new TODO task is added to a process it is placed in the list in the order corresponding to its deadline.
4. Only uncompleted TODO tasks (**status** = 0) can be declared as late. The deadline of completed TODO tasks (**status** != 0) is not checked.

Background Information

This assignment requires basic understanding of the task scheduling in the Linux kernel. Below is some information to get you started. More information can be found in the recommended books and the following [link](#) (from the CS Operating System course).

Linux is a multitasking operating system. A multitasking operating system achieves the illusion of concurrent execution of multiple processes, even on systems with single CPU. This is done by switching from one process to another very quickly. Linux uses **Preemptive Multitasking**. This means that the kernel decides when a process is to cease running and a new process is to begin running. Tasks can also intentionally **block** or **sleep** until some event occurs (keyboard press, passage of time and etc.). This enables the kernel to better utilize the resources of the system and give the user a responsive feeling.

Task States

The state field of the process descriptor describes what is currently happening to the process. The process can be in one of the following states:

TASK_RUNNING: The process is either executing on a CPU or waiting to be executed.

TASK_INTERRUPTIBLE: The process is suspended (sleeping) until some condition becomes true. Raising a hardware interrupt, releasing a system resource the process is waiting for, or delivering a signal are examples of conditions that might wake up the process (put its state back to TASK_RUNNING).

TASK_UNINTERRUPTIBLE: Like TASK_INTERRUPTIBLE, except that delivering a signal to the sleeping process leaves its state unchanged. This process state is seldom used.

TASK_STOPPED: Process execution has been stopped; the process enters this state after receiving a SIGSTOP, SIGTSTP, SIGTTIN, or SIGTTOU signal.

TASK_ZOMBIE: Process execution is terminated, but the parent process has not yet issued a **wait4()** or **waitpid()** system call to return information about the dead process.[*] Before the wait()-like call is issued, the kernel cannot discard the data contained in the dead process descriptor because the parent might need it.

The value of the state field can be set using simple assignment, i.e.,

```
p->state = TASK_RUNNING;
```

or using the macros **set_current_state** and **set_task_state**.

Task Scheduling

The scheduling algorithm is implemented in “kernel/sched.c”. Linux scheduling is based on “time sharing” technique. The CPU time is divided into *slices*, one for each runnable process (processes with the **TASK_RUNNING** state). Each CPU runs only one process at a time. The kernel keeps track of time using timer interrupts. When the time slice of the currently running process expires, the kernel scheduler is invoked and another task is set to run for the duration of its time slice. Switching between tasks is done through **context switch**. Switching of the currently running task can also occur before the expiration of its time slice. This can

occur due to interrupts that wake up processes with higher priority or when the currently running process yields execution to the kernel (e.g. **blocks** or **sleeps**).

The next task is selected according to its **priority**. This value plays an important part in the scheduling algorithm. The kernel uses it to distinguish between different types of processes:

- Interactive processes: Processes that interact with the user. An interactive process spends most of its life time in the **TASK_INTERRUPTIBLE** state waiting for user activity (e.g. key press). But when the event for which it is waiting occurs, it should become the current running process quickly or else the operating system will appear unresponsive to the user. Therefore, the priority of these tasks should be high.
- Batch process: These processes don't need the user's interaction, and usually run in the background. Typical batch processes are compilers and scientific applications.
- Real-Time process. RT processes should never be interrupted by a normal process (Interactive and Batch processes). These types of processes are used in time critical applications like video and motor controllers. These processes have the highest priority. As long as there are runnable RT processes normal processes are not allowed to run.

The **priority** of a process is a dynamic value that is determined both by the user and the kernel (by collecting statistics on the activity of the process).

The kernel holds all runnable processes (processes with the **TASK_RUNNING** state) in a data structure called a **runqueue**. The runnable processes are further divided to processes which are yet to exhaust their time slice and those whose time slice has expired. The runnable processes are listed in two lists - **active** and **expired** (according to the mentioned division) - in the **runqueue**. The **runqueue** also points to the current running process (which obviously resides in the **active** list). Each time the **schedule()** function is called it selects the next running process from the **active** list. Once the time slice of a process expires it is moved to the **expired** list. When the **active** is empty, the **expired** and **active** lists are switched and the **active** processes are assigned new time slices.

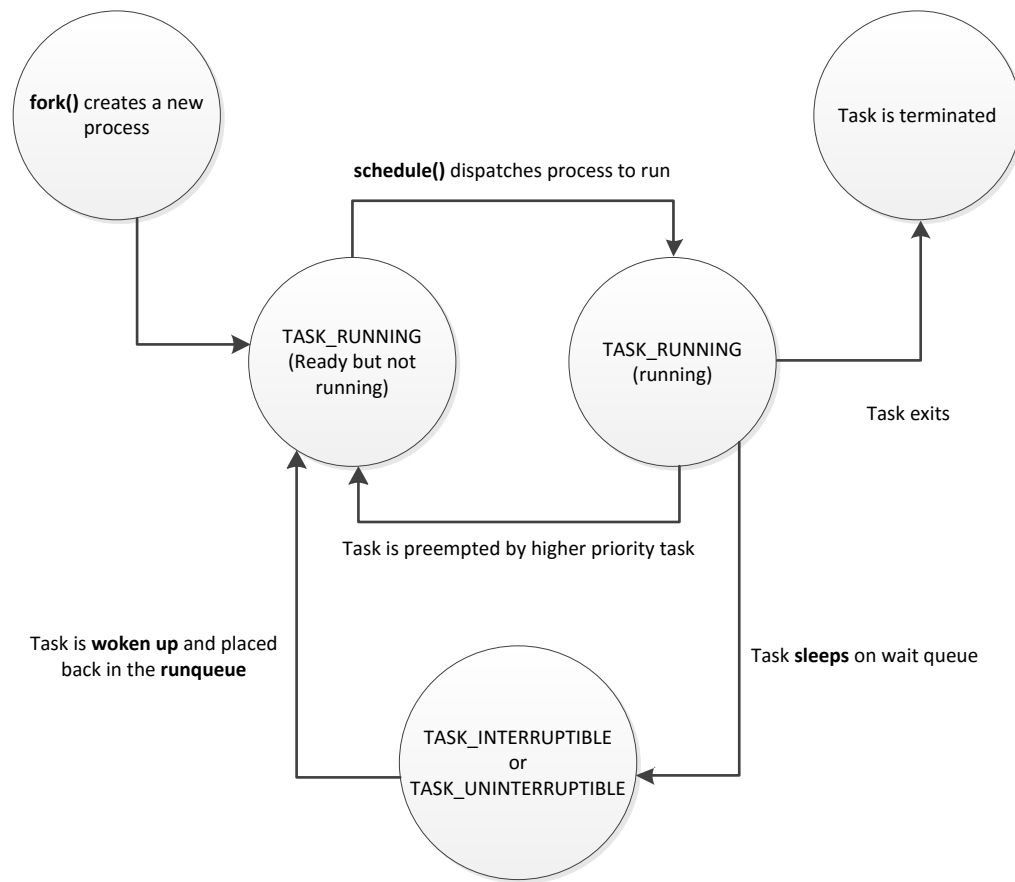
A process can yield its execution to the kernel. This usually happens when the process needs to wait for some event or for a specified period of time. A process can yield its execution by several means. The simplest one is to set its state to **TASK_INTERRUPTIBLE** and then call the **schedule()** function. For example process A can **sleep** by:

```
set_current_state(TASK_INTERRUPTIBLE);
schedule();
//
// After waking up, the process execution continues from here...
//
```

Another process can wake up process A by calling:

```
wake_up_process(process_A);
```

The life cycle of a process is shown in the following diagram:



Kernel Timing

The kernel keeps track of time using the *timer interrupt*. The timer interrupt is issued by the system timer (implemented in hardware). The period of the system timer is called 'tick'. The *timer interrupt* advances the tick counter (called **jiffies**), and initiates time dependent activities in the kernel (decrease the time slice of the current running process, wake up processes that **sleep** waiting for a timer event etc.). The **jiffies** variable (defined in "include/linux/sched.c") counts the system 'tick's event since start up. To 'tick' period duration depends on the specific linux version. The **HZ** variable is use for converting the 'tick's to seconds:

$$time\ in\ seconds = \frac{jiffies}{HZ}$$

Modified System Call API

You should implement the following wrapper function:

1. **int add_TODO(pid_t pid, const char *TODO_description, ssize_t description_size, time_t TODO_deadline)**
 - a. Description:
Add a TODO to the TODO's queue of a process identified by **pid**. **TODO_description** is a string of chars of size **description_size**. **TODO_deadline** holds the deadline time of the TODO task.
 - b. Return value:
 - i. on failure: -1
 - ii. on success: 0
 - c. On failure **errno** should contain one of following values:
 - i. "ENOMEM" (Out of memory): Failure allocating memory.
 - ii. "ESRCH" (No such process): No such process exists or the current process is not allowed to manage the TODO's queue of the process identified by **pid**.
 - iii. "EFAULT" (Bad address): Error copying description from user space.
 - iv. "EINVAL" (Invalid argument) **TODO_description** is NULL or **description_size** < 1 or **TODO_deadline** < **CURRENT_TIME**.
2. **ssize_t read_TODO(pid_t pid, int TODO_index, char *TODO_description, time_t* TODO_deadline, int* status)**
 - a. Description:
Get the description of TODO identified by its position, **TODO_index**, in the TODO's queue of process **pid**. The description is returned in **TODO_description** which is a string of chars. You can assume that the array is large enough for copying the complete TODO's description. The deadline of the TODO is returned in **TODO_deadline**. The status of the TODO is returned in **status**.
 - b. Return value:
 - i. on failure: -1
 - ii. on success: size of the string copied into **TODO_description**.
 - c. On failure **errno** should contain one of following values:
 - i. "ESRCH" (No such process): No such process exists or the current process is not allowed to manage the TODO's queue of the process identified by **pid**.
 - ii. "EFAULT" (Bad address): Error copying description from user space.
 - iii. "EINVAL" (Invalid argument) The TODO's queue size is smaller than **TODO_index** or **TODO_description** is NULL.

Your wrapper functions should follow the example given in the previous assignment. The wrapper functions should be stored in a file called "todo_api.h". The system call should use the following numbering:

System call	Number
sys_add_TODO	243
sys_read_TODO	244
sys_mark_TODO	245
sys_delete_TODO	246

Useful Information

- You can assume that the system is with a single CPU.
- More on system calls and task scheduling can be found in the “Understanding The Linux Kernel” book.
- Use **printk** for debugging (see [link](#)). It is easiest to see **printk**’s output in the textual terminals: Ctrl+Alt+Fn (n=1..6). Note, due to the fact that you are using the VMplayer you might need to press Ctrl+Alt+Space, then release the Space while still holding Ctrl+Alt and then press the required Fn.
- Use **copy_to_user** & **copy_from_user** to copy buffers between User space and Kernel space (see [link](#)).
- You are not allowed to use **syscall** functions to implement code wrappers, or to write the code wrappers for your system calls using the macro **_syscall1**. You should write the code wrappers according to the example of the code wrapper given above.

Testing Your Custom Kernel

You should test your new kernel thoroughly (all functionality and error messages that you can simulate). Note that your code will be tested using an automatic tester. This means that you should pay attention to the exact syntax of the wrapper functions, their names and the header file that defines them. You can use whatever file naming you like for the source/header files that implement the system calls themselves, but they should compile and link using the kernel make file. To do so add your source file in the following line inside the “Makefile” file located in the “kernel” folder:

```
obj-y += sched.o dma.o ... <your_file_name>.o
```

Tips for the Solution

- The deadline mechanism should catch late TODO tasks independently of what the process is doing, i.e. it is not enough to test for a late deadline whenever one of your system calls is used. A possible solution is to check the TODO task during the timer interrupt.
- The timer interrupt occurs every several milli-seconds. It is used for several purposes, e.g. incrementing the jiffies value. It also calls the **scheduler_tick()** function (inside **sched.c**) that updates the time slice of the current process. It is possible to check the status of a TODO_deadline inside the **scheduler_tick()** function.
- Other possibility is to test for late TODO tasks just before the kernel returns from system calls/interrupts to the user mode, (see line 206 in the **entry.S** file).
- If you decide to check the TODO task deadline inside the **scheduler_tick()** function, remember that this function is run inside the interrupt handler. This means that you cannot call **schedule()** or any other function that calls **schedule()** (e.g. **schedule_timeout()** or **interruptible_sleep_on_timeout()**) from inside the **scheduler_tick()** function. You neither can **sleep()** or use any other action that might require interrupts as these are masked inside the interrupt handler. The kernel solves this by using a task flag called **need_resched**. Whenever the **scheduler_tick()** decides that that kernel should reschedule, it turns on this flag. The kernel tests this flag just before it returns to user mode (see line 206 in the **entry.S** file).
- To punish the process with the late TODO task, you need to force the kernel not to schedule it for a period of time. One possible way to do this is to change the status of the process to **TASK_INTERRUPTIBLE** and set a timeout after which it will be woken. Setting a timeout can be done using the **schedule_timeout()** function (inside the **timer.c** file). Note that this function calls the **schedule()** function therefore you cannot call it from inside the **schedule()** function itself. If you want to initiate a timeout from inside the **schedule()** function, you will need to manually set up a timer (see how it is done in the **schedule_timeout()** function).
- When testing the deadline mechanism, you need to create a delay so that some deadline expires. Note that if you use **sleep()**, then the **scheduler_tick()** will not be able to test the TODO tasks of your process as it is never scheduled while it is **sleeping**. Instead, you can create a delay by using a very long loop.

Submission Procedure

1. Submission deadline: 21/01/16 till 23:55.
2. Submissions is allowed in pairs only.
3. You should submit through the moodle website (one submission per pair).
4. You should submit one zip file containing:
 - a. All files you added or modified in your custom kernel. The files should be arranged in folders that preserve their relative path to the root path of the kernel source, i.e:

```
zipfile -+
|
+- submitters.txt
|
+- todo_api.h
|
+- kernel/ -+
|           |
|           +-...
|
+- include/ -+
|           |
|           +-...
...
```

- b. The wrapper functions file "todo_api.h".
- c. A file named "submitters.txt" which lists the name **email** and ID of the participating students. The following format should be used:

```
ploni almoni ploni@t2.technion.ac.il 123456789
john smith john@gmail.com 123456789
```

Note: that you are required to include your email.

Emphasis Regarding Grade

- Your grade for this assignment makes 35% of final grade.
- Pay attention to all the requirements including error values.
- Your submissions will be checked using an automatic checker, pay attention to the submission procedure.
- You are allowed (and encouraged) to consult with your fellow students but you are not allowed to copy their code.
- Your code should be adequately documented and easy to read.
- Delayed submissions will be penalized 2 points for each day (up to 15 points).