

CS345 - Tutorial 4

Implementation of the “Least Slack Time (LST)”
scheduling policy in Linux Kernel

Eva Papadogiannaki
papadogian@csd.uoc.gr

Process Scheduling

- Switching from one process to another in a very short time frame
- Scheduler:
 - When to switch processes
 - Which process to choose next
- Major part of the operating system kernel

Scheduler

- Allows the execution of multiple tasks at the “same” time
- Responsible for:
 - Task coordination among processors
 - Avoiding task starvation and preserving fairness
 - Taking into account system-level tasks (e.g., drivers)

Linux Scheduler

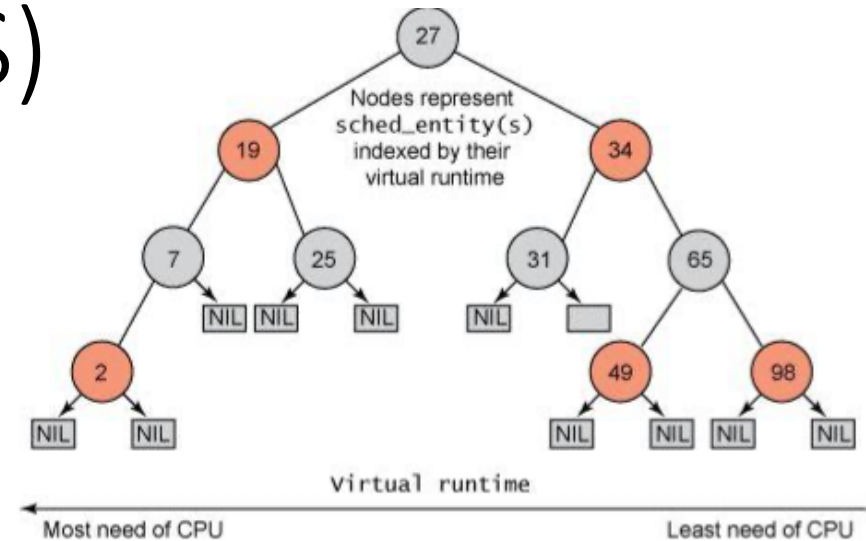
- Executes multiple programs at the “same” time, sharing the CPU with users of varying needs
 - Minimizes response time
 - Maximizes overall CPU utilization
- Ideal Scheduling: n tasks share $100/n$ percentage of CPU effort each
- Preemptive: Higher priority processes evict lower-priority running processes
- Quantum Duration
 - Variable
 - Keep it as long as possible, while keeping good response time

History of schedulers in Linux

- Linux v1.2 – Round Robin
- Linux v2.2 – Scheduling Classes & Policies, Categorizing tasks as non/real-time, non-preemptible
- Linux v2.4 – Division in epochs, goodness of function
- Linux v2.6.0 – v2.6.22 – $O(1)$, Runqueues & priority arrays
- Linux v2.6.23 (and after) – Completely Fair Scheduler (CFS)

Completely Fair Scheduler (CFS)

- Virtual time concept
- Running tasks are sorted using “vruntime”
- Time-ordered red-black tree instead of queue
- Maintains balance in providing processor time to tasks
- At each scheduling invocation:
 - The vruntime of the current task is incremented (time spent in processor)
 - The scheduler chooses the leftmost leaf in the tree (lowest vruntime)
- The leftmost node is cached ($O(1)$)
- Reinsertion of a preempted task takes $O(\log n)$



Scheduling Classes and Policies

Modular design to easily support different scheduling policies

- Each task belongs to a scheduling class
- The scheduling class defines the scheduling policy
- Scheduling policy is set by `sched_setscheduler()`
- Some scheduling policies:
 - `SCHED_NORMAL` – Default linux task policy (CFS, fair)
 - `SCHED_FIFO` – Special time-critical tasks (real-time)
 - `SCHED_RR` – Round-robin scheduling (real-time)

```
/ include / linux / sched.h
```

```
32
33  /*
34   * Scheduling policies
35   */
36  #define SCHED_NORMAL      0
37  #define SCHED_FIFO       1
38  #define SCHED_RR         2
39  #define SCHED_BATCH      3
40  /* SCHED_ISO: reserved but not implemented yet */
41  #define SCHED_IDLE       5
42  /* Can be ORed in to make sure the process is reved
43  #define SCHED_RESET_ON_FORK 0x40000000
```

Linux Kernel source files

- Browse easily through the Linux Kernel source files using this link <https://elixir.bootlin.com/linux/v2.6.38.1/source>
- Actual context switch code, runqueue struct definition, etc.
 - kernel/sched.c <https://elixir.bootlin.com/linux/v2.6.38.1/source/kernel/sched.c>
- Implementation of Completely Fair Scheduling (CFS)
 - kernel/sched_fair.c https://elixir.bootlin.com/linux/v2.6.38.1/source/kernel/sched_fair.c
- Implementation of Real-Time Scheduling (RT)
 - kernel/sched_rt.c https://elixir.bootlin.com/linux/v2.6.38.1/source/kernel/sched_rt.c
- Tasks are abstracted as struct sched_entity and struct sched_rt_entity (for rt class); Also, check struct sched_class
 - include/linux/sched.h <https://elixir.bootlin.com/linux/v2.6.38.1/source/include/linux/sched.h>

Some code snippets

A task's scheduling state (defined in include/linux/sched.h)

```
/ include / linux / sched.h All symb(▼ Search I

170
171 /*
172  * Task state bitmask. NOTE! These bits are also
173  * encoded in fs/proc/array.c: get_task_state().
174  *
175  * We have two separate sets of flags: task->state
176  * is about runnability, while task->exit_state are
177  * about the task exiting. Confusing, but this way
178  * modifying one set can't modify the other one by
179  * mistake.
180  */
181 #define TASK_RUNNING          0
182 #define TASK_INTERRUPTIBLE    1
183 #define TASK_UNINTERRUPTIBLE  2
184 #define  __TASK_STOPPED       4
185 #define  __TASK_TRACED        8
186 /* in tsk->exit_state */
187 #define EXIT_ZOMBIE           16
188 #define EXIT_DEAD             32
189 /* in tsk->state again */
190 #define TASK_DEAD             64
191 #define TASK_WAKEKILL         128
192 #define TASK_WAKING           256
193 #define TASK_STATE_MAX       512
194
195 #define TASK_STATE_TO_CHAR_STR "RSDTtZXxKW"
196
197 extern char  __assert_task_state[1 - 2*!!(
198             sizeof(TASK_STATE_TO_CHAR_STR)-1 != ilog2(TASK_STATE_MAX)+1)];
199
200 /* Convenience macros for the sake of set_task_state */
201 #define TASK_KILLABLE          (TASK_WAKEKILL | TASK_UNINTERRUPTIBLE)
202 #define TASK_STOPPED          (TASK_WAKEKILL |  __TASK_STOPPED)
203 #define TASK_TRACED           (TASK_WAKEKILL |  __TASK_TRACED)
```

Some code snippets

The rq (runqueue) struct (defined in kernel/sched.c)

```
/ kernel / sched.c All symb<v | sea
440
441 /*
442  * This is the main, per-CPU runqueue data structure.
443  *
444  * Locking rule: those places that want to lock multiple runqueues
445  * (such as the load balancing or the thread migration code), lock
446  * acquire operations must be ordered by ascending &runqueue.
447  */
448 struct rq {
449     /* runqueue lock: */
450     raw_spinlock_t lock;
451
452     /*
453      * nr_running and cpu_load should be in the same cacheline because
454      * remote CPUs use both these fields when doing load calculation.
455      */
456     unsigned long nr_running;
457     #define CPU_LOAD_IDX_MAX 5
458     unsigned long cpu_load[CPU_LOAD_IDX_MAX];
459     unsigned long last_load_update_tick;
460 #ifdef CONFIG_NO_HZ
461     u64 nohz_stamp;
462     unsigned char nohz_balance_kick;
463 #endif
464     unsigned int skip_clock_update;
465
466     /* capture load from *all* tasks on this cpu: */
467     struct load_weight load;
468     unsigned long nr_load_updates;
469     u64 nr_switches;
470
471     struct cfs_rq cfs;
472     struct rt_rq rt;
473
474 #ifdef CONFIG_FAIR_GROUP_SCHED
475     /* list of leaf cfs_rq on this cpu: */
476     struct list_head leaf_cfs_rq_list;
477 #endif
478 #ifdef CONFIG_RT_GROUP_SCHED
479     struct list_head leaf_rt_rq_list;
480 #endif
481
482     /*
483      * This is part of a global counter where only the total sum
484      * over all CPUs matters. A task can increase this counter on
485      * one CPU and if it got migrated afterwards it may decrease
```

Some code snippets

The schedule function (in kernel/sched.c)

```
/ kernel / sched.c All symbols Search
3930
3931 /*
3932  * schedule() is the main scheduler function.
3933  */
3934 asmlinkage void __sched schedule(void)
3935 {
3936     struct task_struct *prev, *next;
3937     unsigned long *switch_count;
3938     struct rq *rq;
3939     int cpu;
3940
3941     need_resched:
3942     preempt_disable();
3943     cpu = smp_processor_id();
3944     rq = cpu_rq(cpu);
3945     rcu_note_context_switch(cpu);
3946     prev = rq->curr;
3947
3948     release_kernel_lock(prev);
3949     need_resched_nonpreemptible:
3950
3951     schedule_debug(prev);
3952
3953     if (sched_feat(HRTICK))
3954         hrtick_clear(rq);
3955
3956     raw_spin_lock_irq(&rq->lock);
3957
3958     switch_count = &prev->nivcsw;
3959     if (prev->state && !(preempt_count() & PREEMPT_ACTIVE)) {
3960         if (unlikely(signal_pending_state(prev->state, prev))) {
3961             prev->state = TASK_RUNNING;
3962         } else {
3963             /*
3964              * If a worker is going to sleep, notify and
3965              * ask workqueue whether it wants to wake up a
3966              * task to maintain concurrency. If so, wake
3967              * up the task.
3968              */
3969             if (prev->flags & PF_WQ_WORKER) {
3970                 struct task_struct *to_wakeup;
3971
3972                 to_wakeup = wq_worker_sleeping(prev, cpu);
3973                 if (to_wakeup)
3974                     try_to_wake_up_local(to_wakeup);
3975             }

```

Some code snippets

The pick_next_task function (in kernel/sched.c)

```
/ kernel / sched.c All symbols
3902
3903 /*
3904  * Pick up the highest-prio task:
3905  */
3906 static inline struct task_struct *
3907 pick_next_task(struct rq *rq)
3908 {
3909     const struct sched_class *class;
3910     struct task_struct *p;
3911
3912     /*
3913      * Optimization: we know that if all tasks are in
3914      * the fair class we can call that function directly:
3915      */
3916     if (likely(rq->nr_running == rq->cfs.nr_running)) {
3917         p = fair_sched_class.pick_next_task(rq);
3918         if (likely(p))
3919             return p;
3920     }
3921
3922     for_each_class(class) {
3923         p = class->pick_next_task(rq);
3924         if (p)
3925             return p;
3926     }
3927
3928     BUG(); /* the idle class will always have a runnable task */
3929 }
```

Some code snippets

The sched_class struct (defined in include/linux/sched.h)

```
/ include / linux / sched.h All symbc Search Ident:
1054
1055 struct sched_class {
1056     const struct sched_class *next;
1057
1058     void (*enqueue_task) (struct rq *rq, struct task_struct *p, int flags);
1059     void (*dequeue_task) (struct rq *rq, struct task_struct *p, int flags);
1060     void (*yield_task) (struct rq *rq);
1061
1062     void (*check_preempt_curr) (struct rq *rq, struct task_struct *p, int flags);
1063
1064     struct task_struct * (*pick_next_task) (struct rq *rq);
1065     void (*put_prev_task) (struct rq *rq, struct task_struct *p);
1066
1067 #ifndef CONFIG_SMP
1068     int (*select_task_rq) (struct rq *rq, struct task_struct *p,
1069                          int sd_flag, int flags);
1070
1071     void (*pre_schedule) (struct rq *this_rq, struct task_struct *task);
1072     void (*post_schedule) (struct rq *this_rq);
1073     void (*task_waking) (struct rq *this_rq, struct task_struct *task);
1074     void (*task_woken) (struct rq *this_rq, struct task_struct *task);
1075
1076     void (*set_cpus_allowed) (struct task_struct *p,
1077                             const struct cpumask *newmask);
1078
1079     void (*rq_online) (struct rq *rq);
1080     void (*rq_offline) (struct rq *rq);
1081 #endif
1082
1083     void (*set_curr_task) (struct rq *rq);
1084     void (*task_tick) (struct rq *rq, struct task_struct *p, int queued);
1085     void (*task_fork) (struct task_struct *p);
1086
1087     void (*switched_from) (struct rq *this_rq, struct task_struct *task,
1088                          int running);
1089     void (*switched_to) (struct rq *this_rq, struct task_struct *task,
1090                        int running);
1091     void (*prio_changed) (struct rq *this_rq, struct task_struct *task,
1092                        int oldprio, int running);
1093
1094     unsigned int (*get_rr_interval) (struct rq *rq,
1095                                    struct task_struct *task);
1096
1097 #ifndef CONFIG_FAIR_GROUP_SCHED
1098     void (*task_move_group) (struct task_struct *p, int on_rq);
1099 #endif
1100 };
```

Some code snippets

Handling struct sched_class for fair vs. rt scheduling class

```
/ kernel / sched_fair.c
4166 /*
4167  * All the scheduling class methods:
4168  */
4169 static const struct sched_class fair_sched_class = {
4170     .next = &idle_sched_class,
4171     .enqueue_task = enqueue_task_fair,
4172     .dequeue_task = dequeue_task_fair,
4173     .yield_task = yield_task_fair,
4174
4175     .check_preempt_curr = check_preempt_wakeup,
4176
4177     .pick_next_task = pick_next_task_fair,
4178     .put_prev_task = put_prev_task_fair,
4179
4180 #ifdef CONFIG_SMP
4181     .select_task_rq = select_task_rq_fair,
4182
4183     .rq_online = rq_online_fair,
4184     .rq_offline = rq_offline_fair,
4185
4186     .task_waking = task_waking_fair,
4187 #endif
4188
4189     .set_curr_task = set_curr_task_fair,
4190     .task_tick = task_tick_fair,
4191     .task_fork = task_fork_fair,
4192
4193     .prio_changed = prio_changed_fair,
4194     .switched_to = switched_to_fair,
4195
4196     .get_rr_interval = get_rr_interval_fair,
4197
4198 #ifdef CONFIG_FAIR_GROUP_SCHED
4199     .task_move_group = task_move_group_fair,
4200 #endif
4201 };
```

```
/ kernel / sched_rt.c
1760
1761 static const struct sched_class rt_sched_class = {
1762     .next = &fair_sched_class,
1763     .enqueue_task = enqueue_task_rt,
1764     .dequeue_task = dequeue_task_rt,
1765     .yield_task = yield_task_rt,
1766
1767     .check_preempt_curr = check_preempt_curr_rt,
1768
1769     .pick_next_task = pick_next_task_rt,
1770     .put_prev_task = put_prev_task_rt,
1771
1772 #ifdef CONFIG_SMP
1773     .select_task_rq = select_task_rq_rt,
1774
1775     .set_cpus_allowed = set_cpus_allowed_rt,
1776     .rq_online = rq_online_rt,
1777     .rq_offline = rq_offline_rt,
1778     .pre_schedule = pre_schedule_rt,
1779     .post_schedule = post_schedule_rt,
1780     .task_woken = task_woken_rt,
1781     .switched_from = switched_from_rt,
1782 #endif
1783
1784     .set_curr_task = set_curr_task_rt,
1785     .task_tick = task_tick_rt,
1786
1787     .get_rr_interval = get_rr_interval_rt,
1788
1789     .prio_changed = prio_changed_rt,
1790     .switched_to = switched_to_rt,
1791 };
```

Assignment 4 – LST Scheduling Algorithm

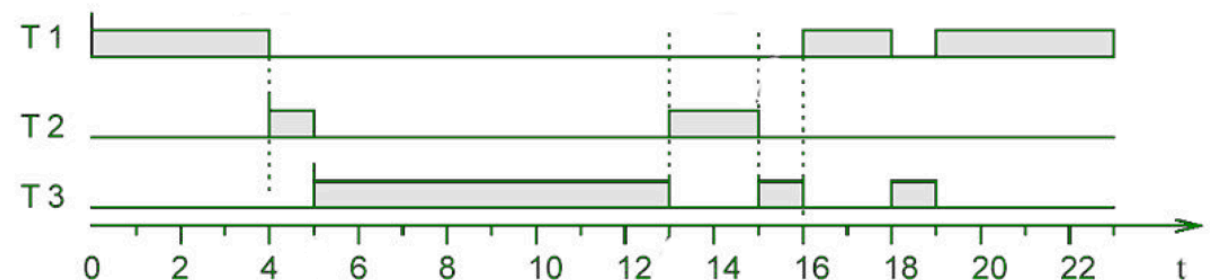
- Implement the Least Slack Time (LST) scheduling algorithm, which assigns priority based on the slack time of a process
- $\text{Slack_time} = \text{deadline} - (\text{computation_time} - \text{elapsed_runtime}) - \text{time}$
- Use your code from Assignment 3
 - You will need the `set_deadlines` system call you implemented
- Use the guidelines from the previous assignment to compile Linux Kernel and run it

Assignment 4 – LST Scheduling Algorithm

	arrival	duration	deadline
T 1	0	10	33
T 2	4	3	28
T 3	5	10	29

- At time $t=0$: Only task T1, has arrived. T1 is executed till time $t=4$.
- At time $t=4$: T2 has arrived.
Slack time of T1: $33-4-6=23$
Slack time of T2: $28-4-3=21$
Hence T2 starts to execute till time $t=5$ when T3 arrives.
- At time $t=5$:
Slack Time of T1: $33-5-6=22$
Slack Time of T2: $28-5-2=21$
Slack Time of T3: $29-5-10=12$
Hence T3 starts to execute till time $t=13$
- At time $t=13$:
Slack Time of T1: $33-13-6=14$
Slack Time of T2: $28-13-2=13$
Slack Time of T3: $29-13-2=14$
Hence T2 starts to execute till time $t=15$

- At time $t=15$:
Slack Time of T1: $33-15-6=12$
Slack Time of T3: $29-15-2=12$
Hence T3 starts to execute till time $t=16$
- At time $t=16$:
Slack Time of T1: $33-16-6=11$
Slack Time of T3: $29-16-2=12$
Hence T1 starts to execute till time $t=18$ and so on..



Assignment 4 – Pre-processing and Filtering

- Before *schedule()* selects the next process
- Scan all processes in the runqueue list and find if there is any process that has a deadline (deadline $\neq -1$)
 - Calculate its slack time
 - If this process has exceeded the given deadline, remove this process from the runqueue list so it'll never be executed
 - If not, iterate the runqueue list *rq*. For each process *p*, check if *p* has less slack value
 - If so execute process *p* first

Assignment 4 – Demonstrate the new scheduler

Create at least 1 demo.

This demo should do the following:

- Create up to 10 child processes
- For each child process, the parent process will set its remaining computation time and its deadline (set deadline to `gettimeofday() + 100 seconds`)
- Each child process will spin for the given computation time (e.g., use `while`, `for`)

Guidelines for Assignment 4

- Browse kernel code with: <https://elixir.bootlin.com/linux/v2.6.38.1/source>
- Another way to map source code is by using ctag: http://www.tutorialspoint.com/unix_commands/ctags.htm
- Understand how the scheduler works
 - For example, you can start with printing inside the *schedule()* function
- Follow the function call path from *schedule* in order to find out how the next task is picked
- Use the *printk()* function often, its syntax is close to *printf* and it's an easy way to observe the kernel's behaviour from the user level (with *dmesg*)
- Reuse existing code snippets within the kernel source code (e.g., to traverse data structures or access members in struct nodes)
- Compile after small changes in the source code (good for easy debugging)
- Submit anything you can that helps you show your effort!