Linux Kernel Programming
**Process Scheduling**

Pierre Olivier

Systems Software Research Group @ Virginia Tech

February 21, 2017

## Outline

1. General information

2. Linux Completely Fair Scheduler

3. CFS implementation

4. Preemption and context switching

5. Real-time scheduling policies

6. Scheduling-related syscalls

Virginia Tech

# Outline

Virginia
Tech

# General information
Scheduling

- ▶ **Scheduler**: OS entity that decide which process should run, when, and for how long
- ▶ Multiplex processes in time on the processor: enables **multitasking**
    - ▶ Gives the user the illusion that processes are executing at the same time
- ▶ Scheduler is responsible for making the best use of the resource that is the CPU time
- ▶ Basic principle:
    - ▶ When in the system there are more ready-to-run processes than the number of cores
        - ▶ **The scheduler decides which process should run**

Virginia Tech

# General information
Multitasking

- ▶ Single core: gives the illusion that multiple processes are running concurrently
- ▶ Multi-cores: enable true parallelism
- ▶ 2 types of multitasking OS:
    - ▶ **Cooperative multitasking**
        - ▶ A process does not stop running until it decides to do so (*yield* the CPU)
        - ▶ The operating system cannot enforce fair scheduling
          - For example in the case of a process that never yields
    - ▶ **Preemptive multitasking**
        - ▶ The OS can interrupt the execution of a process: *preemption*
        - ▶ Generally after the process expires its *timeslice*
        - ▶ And/or based on tasks priorities

Virginia Tech

# General information
A bit of Linux scheduler history

- ▶ From v1.0 to v2.4: simple implementation
  - ▶ **But it did not scale to numerous processes and processors**
- ▶ V2.5 introduced the $O(1)$ scheduler
  - ▶ **Constant time scheduling decisions**
    - ▶ Scalability and execution time determinism
  - ▶ More info in [5]
  - ▶ Issues with latency-sensitive applications (Desktop computers)
- ▶ O(1) scheduler was replaced in 2.6.23 by what is still now the standard Linux scheduler:
  - ▶ **Completely Fair Scheduler (CFS)**
    - ▶ Evolution of the *Rotating Staircase Deadline* scheduler [2, 3]

Virginia Tech

# General information
Scheduling policy - I/O vs compute-bound tasks

- ▶ Scheduling policy are the set of rules determining the choices made by a given model of scheduler
- ▶ **I/O-bound processes**:
  - ▶ Spend most of their time *waiting for I/O*: disk, network, but also keyboard, mouse, etc.
    - ▶ Filesystem, network intensive, GUI applications, etc.
    - ▶ Response time is important
  - ▶ Should run ***often and for a small time frame***
- ▶ **Compute-bound processes**:
  - ▶ *Heavy use of the CPU*
    - ▶ SSH key generation, scientific computations, etc.
    - ▶ Caches stay hot when they run for a long time
  - ▶ Should ***not run often, but for a long time***

Virginia Tech

# General information
Scheduling policy - Priority

- ▶ **Priority**
  - ▶ Order process according to their "importance" from the scheduler standpoint
  - ▶ A process with a higher priority will execute before a process with a lower one
- ▶ **Linux has 2 priority ranges**:
  - ▶ **Nice value**: ranges from -20 to +19, default is 0
    - ▶ High values of nice means lower priority
    - ▶ List process and their nice values with `ps ax -o pid,ni,cmd`
  - ▶ **Real-time priority**: range configurable (default 0 to 99)
    - ▶ Higher values mean higher priority
    - ▶ For processes labeled *real-time*
    - ▶ Real-time processes always executes before standard (nice) processes
    - ▶ List processes and their real-time priority using `ps ax -o pid,rtprio,cmd`

Virginia Tech

# General information
## Scheduling policy - Priority (2)

▶ User space to kernel priorities mapping:

User space view:  **[0                                99][-20            +20]**

| Real-time | Non-real-time |
|:---:|:---:|

Kernel view:  **[0                                                       139]**

# General information
Scheduling policy - Timeslice

- **Timeslice** (quantum):
  - How much time a process should execute before being preempted
  - Defining the default timeslice in an absolute way is tricky:
    - Too long $\rightarrow$ bad interactive performance
    - Too short $\rightarrow$ high context switching overhead
- **Linux CFS does not use an absolute timeslice**
  - The timeslice a process receives is **function of the load of the system**
    - it is *a proportion* of the CPU
  - In addition, that timeslice is **weighted by the process priority**
- When a process $P$ becomes runnable:
  - $P$ will preempt the currently running process $C$ is $P$ consumed a smaller proportion of the CPU than $C$

Virginia Tech

# General information
Scheduling policy - Policy application example

- ▶ 2 tasks in the system:
    - ▶ **Text editor**: I/O-bound, latency sensitive (interactive)
    - ▶ **Video encoder**: CPU-bound, background job
- ▶ Text editor:
    - ▶ A. *Needs a large amount of CPU time*
        - ▶ Does not need to run for long, but needs to have CPU time available whenever it needs to run
    - ▶ B. *When ready to run, needs to preempt the video encoder*
    - ▶ A + B = good interactive performance
- ▶ On a classical UNIX system, needs to set a correct combination of priority and timeslice
- ▶ Different with Linux: *the OS guarantee the text editor a specific proportion of the CPU time*
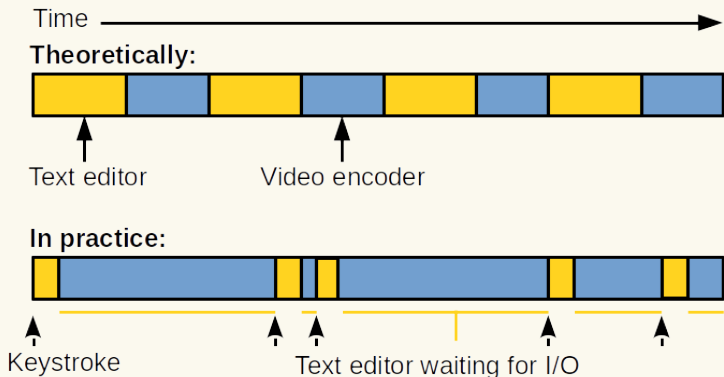
Virginia Tech

# General information
Scheduling policy - Policy application example (2)

▶ Imagine only the two processes are present in the system and run at the same priority
  ▶ Linux gives 50% of CPU time to each
▶ Considering an absolute timeframe:
  ▶ **Text editor does not use fully its 50%** as it often blocks waiting for I/O
    ▶ Keyboard key pressed
  ▶ **CFS keeps track of the actual CPU time used by each program**
  ▶ When the text editor wakes up:
    ▶ CFS sees that **it actually used less CPU time than the video encoder**
    ▶ Text editor preempts the video encoder

# General information
## Scheduling policy - Policy application example (3)



Time

**Theoretically:**

Text editor          Video encoder

**In practice:**

Keystroke          Text editor waiting for I/O

- ▶ Good interactive performance
- ▶ Good background, CPU-bound performance

# Outline

Virginia
Tech

# Linux Completely Fair Scheduler
Scheduling classes

- ▶ **CPU classes**: coexisting CPU algorithms
  - ▶ Each task belongs to a class
- ▶ CFS: SCHED_OTHER, implemented in `kernel/sched/fair.c`
- ▶ Real-time classes: SCHED_RR, SCHED_FIFO, SCHED_DEADLINE
  - ▶ For predictable schedule
- ▶ `sched_class` data structure:

```
1  struct sched_class {
2    void (*enqueue_task) (/* ... */);
3    void (*dequeue_task) (/* ... */);
4    void (*yield_task) (/* ... */);
5    void (*check_preempt_curr) (/* ... */);
6    struct task_struct * (*pick_next_task) (/* ... */);
7    void (*set_cur_task) (/* ...*/);
8    void (*task_tick) (/* ... */);
9    /* ... */
10 }
```

Virginia Tech

# Linux Completely Fair Scheduler

sched_class hooks

- ▶ Functions descriptions:
  - ▶ enqueue_task(...)
    - ▶ Called when a task enters a runnable state
  - ▶ dequeue_task(...)
    - ▶ Called when a task becomes unrunnable
  - ▶ yield_task(...)
    - ▶ Yield the processor (dequeue then enqueue back immediatly)
  - ▶ check_preempt_curr(...)
    - ▶ Checks if a task that entered the runnable state should preempt the currently running task
  - ▶ pick_next_task(...)
    - ▶ Chooses the next task to run
  - ▶ set_curr_task(...)
    - ▶ Called when the currentluy running task changes its scheduling class or task group to the related scheduler
  - ▶ task_tick(...)
    - ▶ Called regularly (default: 10 ms) from the system timer tick handler, might lead to context switch

Virginia Tech

# Linux Completely Fair Scheduler
Unix scheduling

- ► Classical UNIX systems *map priorities (nice values) to absolute timeslices*
- ► Leads to several issues:
    - ► **What is the absolute timeslice that should be mapped to a given nice value?**
        - ► Sub-optimal switching behavior for low priority processes (small timeslices)
    - ► **Relative nice values and their mapping to timeslices**
        - ► Nicing down a process by one can have very different effects according to the tasks priorities
    - ► **Timeslice must be some integer multiple of the timer tick**
        - ► Minimum timeslice and difference between two consecutive timeslices are bounded by the timer tick frequency

Virginia Tech

# Linux Completely Fair Scheduler
Fair scheduling

- ▶ **Perfect multitasking**:
  - ▶ From a single core standpoint
    - ▶ At each moment, each process of the same priority has received an exact amount of the CPU time
    - ▶ What we would get if we could **run n tasks in parallel on the CPU while giving them 1/n of the CPU processing power** → not possible in reality
    - ▶ Or if we could **schedule tasks for infinitely small amounts of time** → context switch overhead issue
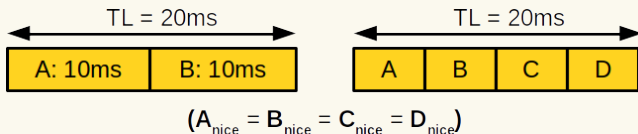- ▶ **3 main (high-level) CFS concepts:**
  1. CFS runs a process for some times, then swaps it for the runnable process that has run the least
  2. No default timeslice, CFS calculates how long a process should run according to the number of runnable processes
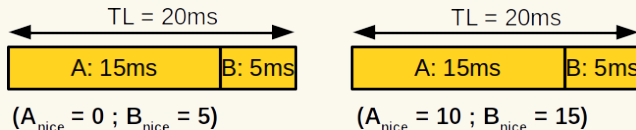  3. That dynamic timeslice is weighted by the process priority (nice)

Virginia Tech

# Linux Completely Fair Scheduler
## Fair scheduling (2)

▶ **Targeted latency**: period during which all runnable processes should be scheduled at least once

▶ Example: processes with the same priority



$$(A_{nice} = B_{nice} = C_{nice} = D_{nice})$$

▶ Example: processes with different priorities



$$(A_{nice} = 0 \; ; \; B_{nice} = 5) \qquad (A_{nice} = 10 \; ; \; B_{nice} = 15)$$

▶ **Minimum granularity**: floor at 1 ms (default)

# Outline

Virginia Tech

# CFS implementation

► **4 main components:**

1. Time accounting
2. Process selection
3. Scheduler entry point (calling the scheduler)
4. Sleeping & waking up

# CFS implementation
Time accounting

▶ `sched_entity` **structure in the** `task_struct` (`se` **field**)

```
1  struct sched_entity
2  {
3    struct load_weight  load;
4    struct rb_node      run_node;
5    struct list_head    group_node;
6    unsigned int        on_rq;
7
8    u64      exec_start;
9    u64      sum_exec_runtime;
10   u64      vruntime;
11   u64      prev_sum_exec_runtime;
12
13   /* additional statistics not shown here */
14 }
```

▶ **Virtual runtime**
  ▶ How much time a process has been executed (ns)

Virginia Tech

# CFS implementation
Time accounting (2)

```
1  static void update_curr(struct cfs_rq *
       cfs_rq)
2  {
3    struct sched_entity *curr =
4      cfs_rq->curr;
5    u64 now = rq_clock_task(rq_of(cfs_rq));
6    u64 delta_exec;
7
8    if (unlikely(!curr))
9      return;
10
11   delta_exec = now - curr->exec_start;
12   if (unlikely((s64)delta_exec <= 0))
13     return;
14
15   curr->exec_start = now;
16
17   schedstat_set(curr->statistics.exec_max,
18           max(delta_exec, curr->statistics
       .exec_max));
```

```
18   curr->sum_exec_runtime += delta_exec;
19   schedstat_add(cfs_rq->exec_clock,
       delta_exec);
20
21   curr->vruntime += calc_delta_fair(
       delta_exec, curr);
22   update_min_vruntime(cfs_rq);
23
24   if (entity_is_task(curr)) {
25     struct task_struct *curtask
26       = task_of(curr);
27
28     trace_sched_stat_runtime(curtask,
       delta_exec, curr->vruntime);
29     cpuacct_charge(curtask, delta_exec);
30     account_group_exec_runtime(curtask,
       delta_exec);
31   }
32
33   account_cfs_rq_runtime(cfs_rq,
       delta_exec);
34 }
```

▶ Invoked regularly by the system timer, and when a process becomes runnable/unrunnable
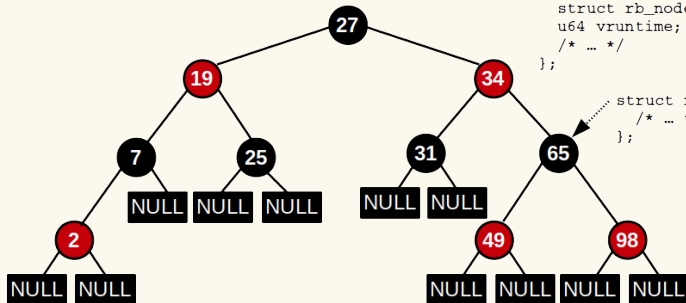
Virginia Tech

# CFS implementation
Process selection



```
struct task_struct {
  /* … */
  struct sched_entity se;
  /* … */
};
```

```
struct sched_entity {
  struct rb_node run_node;
  u64 vruntime; /* 65 here */
  /* … */
};
```

```
struct rb_node {
  /* … */
};
```

Need to execute

▶ Adapted from [1]

# CFS implementation
Process selection (2)

- ▶ When CFS needs to choose which runnable process to run next:
  - ▶ **The process with the smallest `vruntime` is selected**
  - ▶ It is the leftmost node in the tree

```
struct sched_entity *__pick_first_entity(struct cfs_rq *cfs_rq)
{
    struct rb_node *left = cfs_rq->rb_leftmost;

    if (!left)
        return NULL;

    return rb_entry(left, struct sched_entity, run_node);
}
```

Virginia Tech

# CFS implementation
Process selection: adding a process to the tree

- A process is added through `enqueue_entity`:

```
1  static void
2  enqueue_entity(struct cfs_rq *cfs_rq,
        struct sched_entity *se, int flags)
3  {
4    bool renorm = !(flags & ENQUEUE_WAKEUP)
        || (flags & ENQUEUE_MIGRATED);
5    bool curr = cfs_rq->curr == se;
6
7    if (renorm && curr)
8      se->vruntime += cfs_rq->min_vruntime;
9
10   update_curr(cfs_rq);
11
12   if (renorm && !curr)
13     se->vruntime += cfs_rq->min_vruntime;
14
15   update_load_avg(se, UPDATE_TG);
16   enqueue_entity_load_avg(cfs_rq, se);
17   account_entity_enqueue(cfs_rq, se);
18   update_cfs_shares(cfs_rq);
```

```
19   if (flags & ENQUEUE_WAKEUP)
20     place_entity(cfs_rq, se, 0);
21
22   check_schedstat_required();
23   update_stats_enqueue(cfs_rq, se, flags);
24   check_spread(cfs_rq, se);
25   if (!curr)
26     __enqueue_entity(cfs_rq, se);
27   se->on_rq = 1;
28
29   if (cfs_rq->nr_running == 1) {
30     list_add_leaf_cfs_rq(cfs_rq);
31     check_enqueue_throttle(cfs_rq);
32   }
33 }
```

Virginia Tech

# CFS implementation
Process selection: adding a process to the tree (2)

- __enqueue_entity:

```
1  static void __enqueue_entity(struct cfs_rq
        *cfs_rq, struct sched_entity *se)
2  {
3    struct rb_node **link = &cfs_rq->
        tasks_timeline.rb_node;
4    struct rb_node *parent = NULL;
5    struct sched_entity *entry;
6    int leftmost = 1;
7
8    /*
9     * Find the right place in the rbtree:
10    */
11   while (*link) {
12     parent = *link;
13     entry = rb_entry(parent, struct
        sched_entity, run_node);
14     /*
15      * We dont care about collisions.
        Nodes with
16      * the same key stay together.
17      */
```

```
18     if (entity_before(se, entry)) {
19       link = &parent->rb_left;
20     } else {
21       link = &parent->rb_right;
22       leftmost = 0;
23     }
24   }
25   /*
26    * Maintain a cache of leftmost tree
        entries (it is frequently
27    * used):
28    */
29   if (leftmost)
30     cfs_rq->rb_leftmost = &se->run_node;
31
32   rb_link_node(&se->run_node, parent, link
        );
33   rb_insert_color(&se->run_node, &cfs_rq->
        tasks_timeline);
34 }
```

Virginia Tech

# CFS implementation
Process selection: removing a process from the tree

- ▶ dequeue_entity:

```
1  static void
2  dequeue_entity(struct cfs_rq *cfs_rq,
       struct sched_entity *se, int flags)
3  {
4      update_curr(cfs_rq);
5      dequeue_entity_load_avg(cfs_rq, se);
6
7      update_stats_dequeue(cfs_rq, se, flags
         );
8
9      clear_buddies(cfs_rq, se);
10
11     if (se != cfs_rq->curr)
12         __dequeue_entity(cfs_rq, se);
13     se->on_rq = 0;
14     account_entity_dequeue(cfs_rq, se);
```

```
15     if (!(flags & DEQUEUE_SLEEP))
16         se->vruntime -= cfs_rq->
         min_vruntime;
17
18     return_cfs_rq_runtime(cfs_rq);
19
20     update_cfs_shares(cfs_rq);
21
22     if ((flags & (DEQUEUE_SAVE |
       DEQUEUE_MOVE)) == DEQUEUE_SAVE)
23         update_min_vruntime(cfs_rq);
24  }
```

Virginia Tech

# CFS implementation
Process selection: removing a process from the tree (2)

▶ __dequeue_entity:

```
1  static void __dequeue_entity(struct cfs_rq *cfs_rq, struct sched_entity *se)
2  {
3      if (cfs_rq->rb_leftmost == &se->run_node) {
4          struct rb_node *next_node;
5
6          next_node = rb_next(&se->run_node);
7          cfs_rq->rb_leftmost = next_node;
8      }
9
10     rb_erase(&se->run_node, &cfs_rq->tasks_timeline);
11 }
```

Virginia
Tech

# CFS implementation
Entry point: `schedule()`

▶ The kernel calls **schedule()** anytime it wants to invoke the scheduler
  ▶ Calls `pick_next_task()`

```
1  static inline struct task_struct *
2  pick_next_task(struct rq *rq, struct
       task_struct *prev, struct pin_cookie
        cookie)
3  {
4    const struct sched_class *class = &
         fair_sched_class;
5    struct task_struct *p;
6
7    if (likely(prev->sched_class == class &&
8        rq->nr_running == rq->cfs.
       h_nr_running)) {
9      p = fair_sched_class.pick_next_task(rq
       , prev, cookie);
10     if (unlikely(p == RETRY_TASK))
11       goto again;
12
13     if (unlikely(!p))
14       p = idle_sched_class.pick_next_task(
       rq, prev, cookie);
15     return p;
16   }
```

```
17 again:
18   for_each_class(class) {
19     p = class->pick_next_task(rq, prev,
         cookie);
20     if (p) {
21       if (unlikely(p == RETRY_TASK))
22         goto again;
23       return p;
24     }
25   }
26
27   BUG(); /* the idle class will always
         have a runnable task */
28 }
```

Virginia Tech

# CFS implementation
## Sleeping and waking up

▶ Multiple reasons for a task to sleep:
  ▶ Specified amount of time, waiting for I/O, blocking on a mutex, etc.
▶ Going to sleep - steps:
  1. Task marks itself as sleeping
  2. Task enters a *waitqueue*
  3. Task leaves the rbtree of runnable processes
  4. Task calls `schedule()` to select a new process to run
▶ Inverse steps for waking up
▶ Two states associated with sleeping:
  ▶ `TASK_INTERRUPTIBLE`
    ▶ Will be awaken on signal reception
  ▶ `TASK_UNINTERRUPTIBLE`
    ▶ Ignore signals

Virginia Tech

# CFS implementation
Sleeping and waking up: wait queues

- ▶ **Wait queue**:
    - ▶ List of processes waiting for an event to occur

```
1  typedef struct __wait_queue_head
        wait_queue_head_t
2  struct wait_queue_head {
3    spinlock_t lock;
4    struct list_head task_list;}
```

- ▶ Some simple interfaces used to go to sleep have races:
    - ▶ It is possible to go to sleep *after* the event we are waiting for has occurred
    - ▶ Recommended way:

```
1  /* We assume the wait queue we want to wait on is accessible through a variable q */
2
3  DEFINE_WAIT(wait); /* initialize a wait queue entry */
4
5  add_wait_queue(q, &wait);
6  while (!condition) { /* event we are waiting for */
7    prepare_to_wait(&q, &wait, TASK_INTERRUPTIBLE);
8    if(signal_pending(current))
9      /* handle signal */
10    schedule();
11 }
12 finish_wait(&q, &wait);
```

# CFS implementation
Sleeping and waking up: wait queues (2)

▶ **Steps for waiting on a waitqueue:**

1. Create a waitqueue entry (DEFINE_WAIT())
2. Add the calling process to a wait queue (add_wait_queue())
3. Call prepare_to_wait() to change the process state
4. If the state is TASK_INTERRUPTIBLE, a signal can wake the task up → need to check
5. Executes another process with schedule()
6. When the task awakens, check the condition
7. When the condition is true, get out of the wait queue and set the state accordingly using finish_wait()

Virginia Tech

# CFS implementation
Sleeping and waking up: `wake_up()`

- ▶ **Waking up** is taken care of by `wake_up()`
  - ▶ **Awakes all the processes on a waitqueue** by default

```
1  #define wake_up(x) __wake_up(x, TASK_NORMAL, 1, NULL)
2  /* type of x is wait_queue_head_t */
```

- ▶ `__wake_up()` calls `__wake_up_common()`:

```
1  static void __wake_up_common(wait_queue_head_t *q, unsigned int mode,
2          int nr_exclusive, int wake_flags, void *key)
3  {
4    wait_queue_t *curr, *next;
5
6    list_for_each_entry_safe(curr, next, &q->task_list, task_list) {
7      unsigned flags = curr->flags;
8
9      if (curr->func(curr, mode, wake_flags, key) &&
10         (flags & WQ_FLAG_EXCLUSIVE) && !--nr_exclusive)
11       break; /* wakes up only a subset of 'exclusive' tasks */
12    }
13  }
```

- ▶ Exclusive tasks are added through
  `prepare_to_wait_exclusive()`

# CFS implementation
Sleeping and waking up: `wake_up()` (2)

- ► A wait queue entry contains a pointer to a wake-up function
    - ► `include/linux/wait.h`:

```
1  typedef struct __wait_queue wait_queue_t;
2  typedef int (*wait_queue_func_t)(wait_queue_t *wait, unsigned mode, int flags, void *key);
3  int default_wake_function(*wait_queue_func_t)(wait_queue_t *wait, unsigned mode,
4                                                int flags, void *key);
5
6  /* ... */
7
8  struct __wait_queue {
9    /* ... */
10   wait_queue_func_t func;
11   /* ... */
12 }
```

- ► `default_wake_function()` calls `try_to_wake_up()` ...
    - ► ... which calls `ttwu_queue()` ...
    - ► ... which calls `ttwu_do_activate()` (put the task back on runqueue) ...
    - ► ... which calls `ttwu_do_wakeup` ...
    - ► ... which **sets the task state to `TASK_RUNNING`**

Virginia Tech

# CFS implementation
CFS on multicores: brief, high-level overview

- ▶ **Per-CPU runqueues** (rbtrees)
  - ▶ To avoid costly accesses to shared data structures
- ▶ **Runqueues must be kept balanced**
  - ▶ Ex: dual-core with one large runqueue of high-priority processes, and a small one with low-priority processes
    - ▶ High-priority processes get less CPU time than low-priority ones
  - ▶ A load balancing algorithm is run periodically
    - ▶ Balances the queues based on processes priorities and their actual CPU usage
- ▶ More info: [6]

# Outline

Virginia
Tech

# Preemption and context switching
Context switch

- ▶ A **context switch** is the action of swapping the process currently running on the CPU to another one
  - ▶ **Performed by the `context_switch()` function**
    - ▶ Called by `schedule()`

    1. Switch the address space through `switch_mm()`
    2. Switch the CPU state (registers) through `switch_to()`

- ▶ A task can voluntarily relinquish the CPU by calling `schedule()`
  - ▶ **But when does the kernel check if there is a need of preemption?**
    - ▶ `need_resched` flag (per-process, in the `thread_info` of `current`)
  - ▶ `need_resched` is set by:
    1. `scheduler_tick()` when the currently running task needs to be preempted
    2. `try_to_wake_up()` when a process with higher priority wakes up

# Preemption and context switching
need_resched, user preemption

- ▶ The need_resched flag is checked:
    1. Upon returning to user space (from a syscall or an interrupt)
    2. Upon returning from an interrupt
- ▶ If the flag is set, schedule() is called
- ▶ **User preemption happens:**
    1. When returning to user space from a syscall
    2. When returning to user space from an interrupt
- ▶ With Linux, **the kernel is also subject to preemption**

Virginia
Tech

# Preemption and context switching
Kernel preemption

- ▶ In most of Unix-like, kernel code is non-preemptive:
  - ▶ It runs until it finishes
- ▶ **Linux kernel code is preemptive**
  - ▶ A task can be preempted in the kernel as long as execution is **in a safe state**
    - ▶ *Not holding any lock* (kernel is SMP safe)
- ▶ preempt_count in the thread_info structure
  - ▶ Indicates the current lock depth
- ▶ If need_resched && !preempt_count → safe to preempt
  - ▶ Checked when returning to the kernel from interrupt
  - ▶ need_resched is also checked when releasing a lock and preempt_count is 0
- ▶ Kernel code can also call directly schedule()

Virginia Tech

# Preemption and context switching
Kernel preemption (2)

- ▶ **Kernel preemption can occur:**
  1. On return from interrupt to kernel space
  2. When kernel code becomes preemptible again
  3. If a task explicitly calls `schedule()` from the kernel
  4. If a task in the kernel blocks (ex: mutex, result in a call to `schedule()`)

# Outline

Virginia
Tech

# Real-time scheduling policies
SCHED_FIFO **and** SCHED_RR

- *Soft real-time* scheduling classes:
    - Best effort, no guarantees
- **Real-time task of any scheduling class will always run before non-real time ones** (CFS, **SCHED_OTHER)**
    - schedule() → pick_next_task() → for_each_class()
- 2 "classical" RT scheduling policies (kernel/sched/rt.c):
    - **SCHED_FIFO**
        - Tasks run until it blocks/yield, only a higher priority RT task can preempt it
        - Round-robin for tasks of same priority
    - **SCHED_RR**
        - Same as SCHED_FIFO, but with a fixed timeslice

Virginia Tech

# Real-time scheduling policies
Other scheduling policies

- **SCHED_DEADLINE**:
  - Real-time policies mainlined in v3.14 enabling *predictable RT scheduling*
  - EDF implementation based on a period of activation and a worst case execution time (WCET) for each task
  - More info: `Documentation/sched-deadline.txt`, [4], etc.
- **SCHED_BATCH**: non-real-time, low priority background jobs
- **SCHED_IDLE**: non-real-time, *very* low priority background jobs

Virginia Tech

# Outline

Virginia
Tech

Pierre Olivier (SSRG@VT)                    LKP - Process Scheduling                    February 21, 2017      45 / 51

# Scheduling-related syscalls
Scheduling syscalls list

- ▶ sched_getscheduler, sched_setscheduler
- ▶ nice
- ▶ sched_getparam, sched_setparam
- ▶ sched_get_priority_max, sched_get_priority_min
- ▶ sched_getaffinity, sched_setaffinity
- ▶ sched_yield

# Scheduling-related syscalls
## Usage example

```c
1  #define _GNU_SOURCE
2
3  #include <stdio.h>
4  #include <stdlib.h>
5  #include <sys/types.h>
6  #include <unistd.h>
7  #include <sched.h>
8  #include <assert.h>
9
10 void handle_err(int ret, char *func)
11 {
12   perror(func);
13   exit(EXIT_FAILURE);
14 }
15
16 int main(void)
17 {
18   pid_t pid = -1;
19   int ret = -1;
20   struct sched_param sp;
21   int max_rr_prio, min_rr_prio = -42;
22   size_t cpu_set_size = 0;
23   cpu_set_t cs;
```

```c
24 /* Get the PID of the calling process */
25   pid = getpid();
26   printf("My pid is: %d\n", pid);
27
28   /* Get the scheduling class */
29   ret = sched_getscheduler(pid);
30   if(ret == -1)
31     handle_err(ret, "sched_getscheduler");
32   printf("sched_getscheduler returns: "
33     "%d\n", ret);
34   assert(ret == SCHED_OTHER);
35
36   /* Get the priority (nice/RT) */
37   sp.sched_priority = -1;
38   ret = sched_getparam(pid, &sp);
39   if(ret == -1)
40     handle_err(ret, "sched_getparam");
41   printf("My priority is: %d\n",
42     sp.sched_priority);
43
44   /* Set the priority (nice value) */
45   ret = nice(1);
46   if(ret == -1)
47     handle_err(ret, "nice");
```

# Scheduling-related syscalls
## Usage example (2)

```
46     /* Get the priority */
47     sp.sched_priority = -1;
48     ret = sched_getparam(pid, &sp);
49     if(ret == -1)
50       handle_err(ret, "sched_getparam");
51     printf("My priority is: %d\n",
52       sp.sched_priority);
53
54     /* Switch scheduling class to FIFO and
55      * the priority to 99 */
56     sp.sched_priority = 99;
57     ret = sched_setscheduler(pid,
58       SCHED_FIFO, &sp);
59     if(ret == -1)
60       handle_err(ret, "sched_setscheduler");
61
62     /* Get the scheduling class */
63     ret = sched_getscheduler(pid);
64     if(ret == -1)
65       handle_err(ret, "sched_getscheduler");
66     printf("sched_getscheduler returns:"
67       " %d\n", ret);
68     assert(ret == SCHED_FIFO);
```

```
65     /* Get the priority */
66     sp.sched_priority = -1;
67     ret = sched_getparam(pid, &sp);
68     if(ret == -1)
69       handle_err(ret, "sched_getparam");
70     printf("My priority is: %d\n",
71       sp.sched_priority);
72
73     /* Set the RT priority */
74     sp.sched_priority = 42;
75     ret = sched_setparam(pid, &sp);
76     if(ret == -1)
77       handle_err(ret, "sched_setparam");
78     printf("Priority changed to %d\n",
79       sp.sched_priority);
80
81     /* Get the priority */
82     sp.sched_priority = -1;
83     ret = sched_getparam(pid, &sp);
84     if(ret == -1)
85       handle_err(ret, "sched_getparam");
86     printf("My priority is: %d\n",
87       sp.sched_priority);
```

Virginia Tech

# Scheduling-related syscalls
Usage example (2)

```
85    /* Get the max priority value for SCHED_RR */
86    max_rr_prio = sched_get_priority_max(SCHED_RR);
87    if(max_rr_prio == -1)
88      handle_err(max_rr_prio, "sched_get_priority_max");
89    printf("Max RR prio: %d\n", max_rr_prio);
90
91    /* Get the min priority value for SCHED_RR */
92    min_rr_prio = sched_get_priority_min(SCHED_RR);
93    if(min_rr_prio == -1)
94      handle_err(min_rr_prio, "sched_get_priority_min");
95    printf("Min RR prio: %d\n", min_rr_prio);
96
97    cpu_set_size = sizeof(cpu_set_t);
98    CPU_ZERO(&cs);  /* clear the mask */
99    CPU_SET(0, &cs);
100   CPU_SET(1, &cs);
101   /* Set the affinity to CPUs 0 and 1 only */
102   ret = sched_setaffinity(pid, cpu_set_size, &cs);
103   if(ret == -1)
104     handle_err(ret, "sched_setaffinity");
```

```
105   /* Get the CPU affinity */
106   CPU_ZERO(&cs);
107   ret = sched_getaffinity(pid,
108     cpu_set_size, &cs);
109   if(ret == -1)
110     handle_err(ret,
111       "sched_getaffinity");
112   assert(CPU_ISSET(0, &cs));
113   assert(CPU_ISSET(1, &cs));
114   printf("Affinity tests OK\n");
115
116   /* Yield the CPU */
117   ret = sched_yield();
118   if(ret == -1)
119     handle_err(ret,
120       "sched_yield");
121
122   return EXIT_SUCCESS;
123 }
```

Virginia Tech

# Additional documentation

- **CFS:**
    - `http://www.ibm.com/developerworks/library/`
      `l-completely-fair-scheduler/`
    - `http://elinux.org/images/d/dc/Elc2013_Na.pdf`
- **Linux scheduling:**
    - `https://www.cs.columbia.edu/˜smb/classes/`
      `s06-4118/l13.pdf`

# Bibliography I

[1] Inside the linux 2.6 completely fair scheduler.
https://www.ibm.com/developerworks/library/l-completely-fair-scheduler/.
Accessed: 2017-02-01.

[2] The rotating staircase deadline scheduler.
https://lwn.net/Articles/224865/.
Accessed: 2017-01-29.

[3] Rsdl completely fair starvation free interactive cpu scheduler.
https://lwn.net/Articles/224654/.
Accessed: 2017-01-29.

[4] Sched_deadline: a status update.
http://events.linuxfoundation.org/sites/events/files/slides/SCHED_DEADLINE-20160404.pdf.
Accessed: 2017-02-06.

[5] Understanding the linux 2.6.8.1 cpu scheduler.
https://web.archive.org/web/20131231085709/http://joshaas.net/linux/linux_cpu_scheduler.pdf.

Accessed: 2017-01-28.

[6] LOZI, J.-P., LEPERS, B., FUNSTON, J., GAUD, F., QUÉMA, V., AND FEDOROVA, A.
The linux scheduler: A decade of wasted cores.
In *Proceedings of the Eleventh European Conference on Computer Systems* (2016), EuroSys '16, ACM, pp. 1:1–1:16.

Virginia
Tech