Gabriele Proietti Mattia

### Advanced Operating Systems and Virtualization

[11] Scheduling



Department of Computer, Control and Management Engineering "A. Ruberti", Sapienza University of Rome

gpm.name · proiettimattia@diag.uniroma1.it

A.Y. 2020/2021 · V2

#### Outline

- 1. Introduction
- 2. Priorities and Weights
- 3. Scheduler Core
  - 1. Wait Queues
  - 2. Scheduler Entry Point
  - 3. Scheduler Algorithms
- 4. Context Switch

11.1

11. Scheduling

### Introduction

DIAG

Advanced Operating Systems and Virtualization

#### **Process Scheduling**

Like every time sharing system, Linux achieves the effect of an apparent simultaneous execution of multiple processes by switching from one process to another in a very short time frame.

The scheduling policy is concerned with **when** to switch and **which** process to choose. The scheduling algorithm of traditional Unix operating systems must fulfill several conflicting objectives:

- fast process response time
- good throughput for background jobs
- no starvation
- reconciliation of the needs of low- and high- priority processes and so on

### **Linux Scheduling**

Linux Scheduling is based on the time sharing technique: several processes run in "time multiplexing" because the CPU is divided into **slices**, one for each runnable process. Obviously one CPU can run only one process at a given instant, therefore when the currently running process is not terminated when its time slice or **quantum** expires, a process switch may take place.

Time sharing relies on timer interrupts and is thus transparent to process, no additional code needs to be inserted in the programs to ensure CPU time sharing.

The developing history of Linux has seen different scheduling algorithms.

11.2

11. Scheduling

## **Priorities and Weights**



Advanced Operating Systems and Virtualization

### Priority

In Linux, process priority is dynamic. The scheduler keeps track of what processes are doing and adjusts their priorities **periodically**. In this way, processes that have been denied the use of CPU for a long time interval are boosted by dynamically increasing their priority (and vice versa). Process in general are:

- CPU Bound, if require a lot of cpu time
- I/O Bound, if require a lot of I/O operations

Or, according to another classification:

- **interactive**, if they interact a lot with the user, therefore they spend much time waiting for keystrokes and mouse operations; these processes must be very responsive
- **batch**, if they do not need user interaction, since they run in background
- **real-time**, they have to follow strict scheduling requirements, they should be never blocked by a lower priority process (e.g. video, sound applications)

The Linux scheduler allows only to define real-time processes.

#### Nice and RT

Unix demands for priority based scheduling. To each process is associated a "nice" number in [-20, 19]:

- the **higher** the nice, the **lower** the priority
- this tells how nice a process is towards others

There is also the notion of "real time" processes

- **Hard real time**: bound to strict time limits in which a task must be completed (not supported in mainstream Linux)
- **Soft real time**: there are boundaries, but don't make your life depend on it. Examples: burning data to a CD ROM, VoIP

#### **Priorities**

In Linux, real time priorities are in [0, 99]. Here higher value means lower priority. Implemented according to the Real-Time. Extensions of POSIX.

Both nice and rt priorities are mapped to a single value in [0, 139] in the kernel:

- o to 99 are reserved to rt priorities
- 100 to 139 for nice priorities (mapping exactly to [-20, 19])

Priorities are defined in include/linux/sched/prio.h



#### Figure 2-14: Kernel priority scale.

Mauerer, Wolfgang. Professional Linux kernel architecture. John Wiley & Sons, 2010.

#### **Priorities**

#### Macros

```
#define MAX_NICE 19
#define MIN_NICE - 20
#define NICE_WIDTH (MAX_NICE - MIN_NICE + 1)
#define MAX_USER_RT_PRIO 100
#define MAX_RT_PRIO MAX_USER_RT_PRIO
#define MAX_PRIO (MAX_RT_PRIO + NICE_WIDTH)
#define DEFAULT_PRIO (MAX_RT_PRIO + NICE_WIDTH / 2)
```

```
ps -eo pid,rtprio,cmd ('-' = no realtime)
chrt -p pid
chrt -p prio pid
```

#### **Priorities**

```
/*
* Convert user-nice values [ -20 ... 0 ... 19 ]
* to static priority [ MAX RT PRIO..MAX PRIO-1 ],
* and back.
*/
#define NICE TO PRIO(nice) ((nice) + DEFAULT PRIO)
#define PRIO TO NICE(prio) ((prio) - DEFAULT PRIO)
/*
* 'User priority' is the nice value converted to something we
* can work with better when scaling various scheduler parameters,
* it's a [ 0 ... 39 ] range.
*/
```

```
#define USER_PRIO(p) ((p)-MAX_RT_PRIO)
#define TASK_USER_PRIO(p) USER_PRIO((p)->static_prio)
#define MAX_USER_PRIO (USER_PRIO(MAX_PRIO))
```

#### Priorities in task\_struct

There are several fields for representing the priority in the task\_struct:

- **static\_prio** (<u>static</u>): priority given "statically" by a user (and mapped into kernel's representation);
- normal\_priority (<u>dynamic</u>): based on static\_prio and scheduling policy of a process: tasks with the same static priority that belong to different policies will get different normal priorities. Child processes inherit the normal priorities from their parent processes when forked;
- **prio** (<u>dynamic</u>): it is the priority considered by the scheduler, it changes over the process execution keeping track CPU bound processes to penalize, and I/O bound processes to boost;
- **rt\_priority** (static): the realtime priority for realtime tasks in [0, 99].

\*<u>static</u> means that is assigned at process startup but then it can be changed by issuing a system call, <u>dynamic</u> means that the scheduler recomputes it during the process execution

#### Policy is SCHED\_DEADLINE Computing prio V5.11 Policy is SCHED FIFO/RR In kernel/sched/core.c \* Calculate the expected normal priority: i.e. priority \* without taking RT-inheritance into account. Might be \* boosted by interactivity modifiers. Changes upon fork, p->prio = effective prio(p); \* setprio syscalls, and whenever the interactivity 1639 \* estimator recalculates. 1641 \*/ 1642 static inline int normal prio(struct task struct \*p) 1643/ 1655 /\* 1644 int prio; 1645 1656 \* Calculate the current priority, i.e. the priority 1646 if (task has dl policy(p)) \* taken into account by the scheduler. This value might 1657 1647 prio = MAX DL PRIO-1; \* be boosted by RT tasks, or might be boosted by 1658 else if (task has rt policy(p)) prio = MAX RT PRIO-1 - p->rt priority; 1649 1659 \* interactivity modifiers. Will be RT if the task got 1650 else \* RT-boosted. If not then it returns p->normal prio. 1660 1651 prio = \_\_normal\_prio(p); return prio: 1661 \*/ 1662 static int effective prio(struct task struct \*p) https://elixir.bootlin.com/linux/v5.11/source/kernel/sched/core.c#L1642 1663 { 1664 p->**normal prio** = **normal prio**(p); 1665 1666 \* If we are RT tasks or we were boosted to RT priority, 1667 \* keep the priority unchanged. Otherwise, update priority 1668 \* to the normal priority: Check if p->prio > 100 1669 \*/ 1670 if (!rt prio(p->prio)) normal prio - return the priority that is based on the static prio \*/ 1671 return p->normal prio; static inline int \_\_\_\_normal\_prio(struct task\_struct \*p) 1630 1672 return p->**prio**; return p->static prio; 1673

https://elixir.bootlin.com/linux/v5.11/source/kernel/sched/core.c#L1662

### Load weights

The importance of a task is not only described by the priority, but also by a load weight, used to scale the time slice assigned to a scheduled process. The weight can be found in:

```
task_struct->se->load (of type struct load_weight)
```

327	<pre>struct load_weight {</pre>	
328	unsigned long	weight;
329	u32	<pre>inv_weight;</pre>
330	};	
	https://elixir.bootlin.com/linux/v5.11/source/include/linu	JX/sched.h#L327

#### Load Weights

V5.11

9140	/*									
9141	* Nice leve	els are multa	iplicative,	with a gent	tle 10% chai	nge for every				
9142	* nice level changed. I.e. when a CPU-bound task goes from nice 0 to									
9143	* nice 1, i	t will get -	-10% less Cl	PU time thar	n another Cl	PU-bound task				
9144	* that rema	ined on nice	e 0.							
9145	*									
9146	* The "10%	* The "10% effect" is relative and cumulative: from any nice level,								
9147	* if you go	up 1 level,	it's -10%	CPU usage,	if you go	down 1 level				
9148	* it's +10%	S CPU usage.	(to achieve	e that we us	se a multip	lier of 1.25.				
9149	* If a task	goes up by	~10% and a	nother task	goes down	by ~10% then				
9150	* the relat	ive distance	e between ti	hem is ~25%.	)					
9151	*/									
9152	const int sc	hed_prio_to	weight[40]	= {						
9153	/* -20 */	88761,	71755,	56483,	46273,	36291,				
9154	/* -15 */	29154,	23254,	18705,	14949,	11916,				
9155	/* -10 */	9548,	7620,	6100,	4904,	3906,				
9156	/* -5 */	3121,	2501,	1991,	1586,	1277,				
9157	/* 0 */	1024,	820,	655,	526,	423,	nice = 19			
9158	/* 5 */	335,	272,	215,	172,	137,	/			
9159	/* 10 */	110,	87,	70,	56,	45,				
9160	/* 15 */	36,	29,	23,	18,	15, 🖌				
9161	};	https://www.	antin comlinuulus - t	urse (kernel/sched/serre						

https://elixir.bootlin.com/linux/v5.11/source/kernel/sched/core.c

### Load weights

#### Examples

Two tasks running at nice 0 (weight 1024):

→ Both get 50% of time: 1024/(1024+1024) = 0.5

Task 1 is moved to nice -1 (priority boost):

- $\rightarrow$  T1: 1277/(1024+1277)  $\approx$  0.55
- → T2: 1024/(1024+1277) ≈ 0.45 (10% difference)

Task 2 is then moved to nice 1 (priority drop):

- $\rightarrow$  T1: 1277/(820+1277)  $\approx$  0.61
- → T2: 820/(820+1277) ≈ 0.39 (22% difference)

#### Conclusions

Therefore for now we introduced that for a scheduling process we need:

- a set of dynamic and static priorities
- a load weight

11.3

11. Scheduling

### **Scheduler Core**



Advanced Operating Systems and Virtualization

### **Scheduling Classes/Policies**

Every Linux process is always scheduled according to one of the following scheduling classes. The policy is described in task\_struct->policy:

- SCHED\_NORMAL (also called SCHED\_OTHER), conventional time-shared process, in general has a soft priority mechanism over the 'nice' range of -20 to +19 (static priority of 100-139) which decides according to the priority which task goes first, and how much timeslice it gets. This system dynamically alters the priority to allow interactive tasks to go first, and is designed to prevent starvation of lower priority tasks with an expiration policy
- SCHED\_RR, is a fixed real time policy over the static range of o-99 where a lower number (higher priority) task will repeatedly go ahead of \_any\_ tasks lower priority than itself. It is called RR because if multiple tasks are at the same priority it will Round Robin between those tasks
- SCHED\_FIFO, is a fixed real time policy the static range of o-99 where a lower number (higher priority) task will repeatedly go ahead of \_any\_ tasks with lower priority than itself. Unlike RR, if a task does not give up the cpu it will run indefinitely even if other tasks are the same static priority as itself.

### **Scheduling Classes/Policies**

Every Linux process is always scheduled according to one of the following scheduling classes:

- SCHED\_BATCH, does not preempt nearly as often as regular tasks would, thereby allowing tasks to run longer and make better use of caches but at the cost of interactivity. This is well suited for batch jobs.
- SCHED\_IDLE, even weaker than SCHED\_BATCH
- SCHED\_DEADLINE, implementation of the Earliest Deadline First (EDF) scheduling algorithm, augmented with a mechanism (called Constant Bandwidth Server, CBS) that makes it possible to isolate the behavior of tasks between each other. CBS has been replaced with Greedy Reclamation of Unused Bandwidth (GRUB) from kernel 4.13.

### **Scheduling Classes**

For each scheduling class/policy a set of standard function is defined as follows, in this way in order to schedule a task the scheduler core uses always the same set of APIs. We have:

- **enqueue\_task(...)** Called when a task enters a runnable state. It puts the scheduling entity (task) into the red-black tree and increments the nr\_running variable.
- **dequeue\_task(...)** When a task is no longer runnable, this function is called to keep the corresponding scheduling entity out of the red-black tree. It decrements the nr\_running variable.
- yield\_task(...) This function is basically just a dequeue followed by an enqueue, unless the compat\_yield sysctl is turned on; in that case, it places the scheduling entity at the right-most end of the red-black tree.
- **check\_preempt\_curr(...)** This function checks if a task that entered the runnable state should preempt the currently running task.
- **pick\_next\_task(...)** This function chooses the most appropriate task eligible to run next.
- **set\_curr\_task(...)** This function is called when a task changes its scheduling class or changes its task group.
- task\_tick(...) This function is mostly called from time tick functions; it might lead to process switch. This drives the running preemption.

```
Run-queue
       struct sched_class {
1814
       #ifdef CONFIG UCLAMP TASK
1817
               int uclamp_enabled;
1818
       #endif
1820
               void (*enqueue task) (struct rg *rg, struct task struct *p, int flags);
               void (*degueue task) (struct rg *rg, struct task struct *p, int flags);
               void (*yield task)
                                    (struct rq *rq);
               bool (*yield to task)(struct rq *rq, struct task struct *p);
1824
               void (*check preempt curr)(struct rg *rg, struct task struct *p, int flags);
1826
1827
               struct task struct *(*pick next task)(struct rq *rq);
1828
               void (*put prev task)(struct rg *rg, struct task struct *p);
               void (*set next task)(struct rq *rq, struct task struct *p, bool first);
      #ifdef CONFIG SMP
               int (*balance)(struct rq *rq, struct task struct *prev, struct rq flags *rf);
1834
                    (*select task rq)(struct task struct *p, int task cpu, int flags);
               int
1835
               void (*migrate_task_rq)(struct task_struct *p, int new cpu);
1836
               void (*task woken)(struct rg *this rg, struct task struct *task);
               void (*set_cpus_allowed)(struct task_struct *p,
                                        const struct cpumask *newmask,
                                        u32 flags):
1843
               void (*rq online)(struct rq *rq);
               void (*rq offline)(struct rq *rq);
               struct rq *(*find lock_rq)(struct task struct *p, struct rq *rq);
1847
       #endif
               void (*task tick)(struct rq *rq, struct task struct *p, int queued);
               void (*task fork)(struct task struct *p);
               void (*task dead)(struct task struct *p);
```

https://www.kernel.org/doc/html/latest/scheduler/sched-design-CES.html?highlight=sched\_normal#scheduling-classes - https://elixir.bootlin.com/linux/v5.11/source/kernel/sched/sched.h#L1812



Scheduling can be activated in 2 ways: when a task goes to sleep (or yield the CPU) or by a periodic mechanism.



Figure 2-13: Overview of the components of the scheduling subsystem.

Mauerer, Wolfgang. Professional Linux kernel architecture. John Wiley & Sons, 2010.

#### **Scheduler Code Organization**

General code base and specific scheduler classes are found in kernel/sched/

- **core.c**: the common codebase
- fair.c: implementation of the basic scheduler (CFS: Completely Fair Scheduler), it implements SCHED\_NORMAL, SCHED\_BATCH and SCHED\_IDLE
- rt.c: the real-time scheduler implements SCHED\_FIFO and SCHED\_RR
- idle\_task.c: generic entry points for the idle threads and implementation of the idle task scheduling class (<u>not related to SCHED IDLE</u>) for scheduling the idle task (i.e. do\_idle)

#### **Run Queues**

The central data structure of the core scheduler that is used to manage active processes is known as **run queue**. Each CPU has its own run queue, and each active process appears on just one run queue.



#### **Run Queues**

DECLARE\_PER\_CPU\_SHARED\_ALIGNED(struct rq, runqueues);

```
#define cpu_rq(cpu) (&per_cpu(runqueues, (cpu)))
#define this_rq() this_cpu_ptr(&runqueues)
#define task_rq(p) cpu_rq(task_cpu(p))
#define cpu_curr(cpu) (cpu_rq(cpu)->curr)
```

11.3.1

**11. Scheduling** 3. Scheduler Core

#### Wait Queues



Advanced Operating Systems and Virtualization

Defined in include/linux/wait.h Wait Queues implement conditional waits on events: a process wishing to wait for a specific event places itself in the proper wait queue and relinquishes control. Therefore a wait queue represents a set of sleeping processes, which are woken up by the kernel when some condition becomes true.

Wait Queues changed many times in the history of the kernel. In the earlier version they suffered from the "Thundering Herd" performance problem.

#### **Thundering Herd Effect**



Taken from 1999 Mindcraft study on Web and File Server Comparison

#define WQ\_FLAG\_EXCLUSIVE 0x01 ~

```
struct wait_queue_entry {
```

```
unsigned int flags;
```

```
void *private;
```

```
wait_queue_func_t func; 🖛
```

```
struct list_head entry;
```

```
};
```

```
struct wait_queue_head {
    spinlock_t lock;
    struct list_head head;
};
```

typedef struct wait\_queue\_head wait\_queue\_head\_t;

https://elixir.bootlin.com/linux/v5.11/source/include/linux/wait.h

 For solving the Thundering Herd problem the kernel defines two kinds of sleeping processes:

- exclusive, are selectively woken up by the kernel
- *non-exclusive*, are always woken up by the kernel

Specifies how waking up the sleeping process

APIs

- static inline void init\_waitqueue\_entry(struct wait\_queue\_entry \*wq\_entry, struct task\_struct \*p)
- static inline void wait\_event\_interruptible(wq\_head, condition) sleep until a condition gets true
- static inline void wait\_event\_interruptible\_timeout(wq\_head, condition, timeout) - sleep until a condition gets true or a timeout elapses
- static inline void wait\_event\_hrtimeout(wq\_head, condition, timeout)
- static inline void wait\_event\_interruptible\_hrtimeout(wq, condition, timeout)

#### Adding entry to a wait queue

18	<pre>void add_wait_queue(struct wait_queue_head *wq_head, struct wait_queue_entry *wq_entry)</pre>
19	{
20	unsigned long flags;
21	
22	w <b>q_entry</b> ->flags &= ~W <b>Q_FLAG_EXCLUSIVE</b> ;
23	<pre>spin_lock_irqsave(&amp;wq_head-&gt;lock, flags);</pre>
24	add_wait_queue(wq_head, wq_entry);
25	<pre>spin_unlock_irqrestore(&amp;wq_head-&gt;lock, flags);</pre>
26	}

```
29
     void add wait queue exclusive(struct wait queue head *wq head, struct wait queue entry *wq entry)
30
31
             unsigned long flags;
32
33
             wq entry->flags |= WQ FLAG EXCLUSIVE;
             spin_lock_irqsave(&wq_head->lock, flags);
34
               add_wait_queue_entry_tail(wq_head, wq_entry);
35
36
             spin unlock irgrestore(&wg head->lock, flags);
37
38
     EXPORT SYMBOL(add wait queue exclusive);
```

https://elixir.bootlin.com/linux/v5.11/source/kernel/sched/wait.c#L29

#### Removing entry from a wait queue



https://elixir.bootlin.com/linux/v5.11/source/kernel/sched/wait.c#L51

Tasks organization



#### APIs (2)

These functions wake up TASK\_NORMAL = TASK\_INTERRUPTIBLE + TASK\_UNINTERRUPTIBLE

- wake\_up(x)
- wake\_up\_nr(x, nr)
- wake\_up\_all(x)
- wake\_up\_locked(x)
- wake\_up\_all\_locked(x)

These instead only TASK\_INTERRUPTIBLE

- wake\_up\_interruptible(x)
- wake\_up\_interruptible\_nr(x, nr)
- wake\_up\_interruptible\_all(x)
- wake\_up\_interruptible\_sync(x)

11.3.2

**11. Scheduling** 3. Scheduler Core

### Scheduler Entry Point

DIAG

Advanced Operating Systems and Virtualization



Scheduling can be activated in 2 ways: when a task goes to sleep (or yield the CPU) or by a periodic mechanism.



Figure 2-13: Overview of the components of the scheduling subsystem.

Mauerer, Wolfgang. Professional Linux kernel architecture. John Wiley & Sons, 2010.

#### **Scheduler Main Functions**

The main functions used by the scheduler are:

```
scheduler tick()
    Keeps the time slice counter of current up-to-date
try to wake up()
    Awakens a sleeping process
recalc task prio()
    Updates the dynamic priority of a process
schedule()
    Selects a new process to be executed
load balance()
    Keeps the runqueues of a multiprocessor system balanced
     Bovet, Daniel P., and Marco Cesati. Understanding the Linux Kernel: from I/O ports to process management. " O'Reilly Media,
```

```
Inc.", 2005.
```

#### **Scheduler Entry Point**

The entry point for the scheduler is schedule(void) in kernel/sched.c. This is called from several places in the kernel:

- **Direct Invocation**: an explicit call to schedule() is issued
- Lazy Invocation: some hint is given to the kernel indicating that schedule() should be called soon (see need\_resched)

These invocations can be triggered by the **Main Scheduler** or the **Periodic Scheduler**.

In general **schedule()** entails 3 distinct phases, which depend on the scheduler implementation:

- 1. some **checks** on the current process (e.g., with respect to signal processing)
- 2. **selection** of the process to be activated
- 3. context **switch**

#### **Periodic Scheduler**

The function scheduler\_tick() is called from update\_process\_times(), called at every tick of the current CPU (remind the Time Management chapter).

This function has two goals:

- managing scheduling-specific statistics
- calling the scheduling method of the class

```
4528
       /*
4529
        * This function gets called by the timer code, with HZ frequency.
4530
        * We call it with interrupts disabled.
4531
        */
4532
       void scheduler tick(void)
4533
4534
               int cpu = smp processor id();
               struct rq *rq = cpu rq(cpu);
4535
4536
               struct task struct *curr = rg->curr;
4537
               struct rq flags rf;
               unsigned long thermal pressure;
4538
4539
               arch scale freq tick();
4540
4541
               sched clock tick();
4542
4543
               rq lock(rq, &rf);
4544
               update rg clock(rg);
4545
               thermal pressure = arch scale thermal pressure(cpu of(rq));
4546
               update thermal load avg(rq clock thermal(rq), rq, thermal pressure);
4547
4548
              curr->sched class->task tick(rg, curr, 0);
4549
              > calc global load tick(rq);
4550
               psi task tick(rq);
4551
4552
               rq unlock(rq, &rf);
4553
4554
               perf event task tick();
4555
4556
       #ifdef CONFIG SMP
4557
               rq->idle balance = idle cpu(cpu);
4558
               trigger load balance(rq);
4559
       #endif
4560
```

https://elixir.bootlin.com/linux/v5.11/source/kernel/sched/core.c#L4532

#### **Main Scheduler**

The main scheduler function (schedule()) is invoked directly in many points in the kernel to allocate the CPU to a process other than the currently active one. After returning from system calls the kernel also checks whether the flag TIF\_NEED\_RESCHED of the current process is set (by the Periodic Scheduler for example), and if it is checked schedule() is called.



#### **Task States**

The state field in the PCB tracks the current state of the process/thread. Values are defined in include/linux/sched.h:

- **TASK\_RUNNING** the process is either executing on CPU or waiting to be executed
- **TASK\_INTERRUPTIBLE** the process is sleeping until some condition becomes true
- TASK\_UNINTERRUPTIBLE like TASK\_INTERRUPTIBLE but they can be only woken up by the kernel, not by external signals
- **TASK\_STOPPED** process has stopped (after signal SIGSTOP, SIGTSTP)
- TASK\_PARKED
- TASK\_DEAD
- TASK\_WAKEKILL is designed to wake the process on receipt of fatal signals
- TASK\_WAKING
- TASK\_NOLOAD
- TASK\_NEW the task has been just created
- TASK\_STATE\_MAX

Convenience macros for the sake of set\_current\_state:

- #define TASK\_KILLABLE (TASK\_WAKEKILL | TASK\_UNINTERRUPTIBLE)
- #define TASK\_STOPPED (TASK\_WAKEKILL | \_\_TASK\_STOPPED)
- #define TASK\_TRACED (TASK\_WAKEKILL | \_\_TASK\_TRACED)
- #define TASK\_IDLE (TASK\_UNINTERRUPTIBLE | TASK\_NOLOAD)

Convenience macros for the sake of wake\_up():

- #define TASK\_NORMAL (TASK\_INTERRUPTIBLE | TASK\_UNINTERRUPTIBLE)

#### All the PCBs registered in the runqueue are TASK\_RUNNING.

#### **Task State Transition**



#### **TASK\_\*INTERRUPTIBLE**

In case an operation cannot be completed immediately (think of a read()) the task goes to sleep in a wait queue. While doing this, the task enters either the TASK\_INTERRUPTIBLE or TASK\_UNINTERRUPTIBLE state. At this point, the kernel thread calls schedule() to effectively put to sleep the currently-running one and pick the new one to be activated.

Dealing with TASK\_INTERRUPTIBLE can be difficult when the syscall is interrupted for example:

- at kernel level, understand that the task has been resumed due to an interrupt
- clean up all the work that has been done so far
- return to userspace with EINTR
- userspace has to understand that a syscall was interrupted (bugs here!)

Conversely, a TASK\_UNINTERRUPTIBLE might never be woken up again (the dreaded D state in ps). TASK\_KILLABLE is handy for this (since 2.6.25), same as TASK\_UNINTERRUPTIBLE except for fatal sigs.

### Waking up sleeping tasks

The event a task is waiting for calls one of the wake\_up\*() functions on the corresponding wait queue. A task is set to runnable and put back on a runqueue.

It the woken up task has a higher priority than the other tasks on the runqueue, TIF\_NEED\_RESCHED is flagged.

11.3.3

**11. Scheduling** 3. Scheduler Core

## **Scheduling Algorithms**



Advanced Operating Systems and Virtualization

### **Brief History**

- **v1.2**: circular queue for runnable task management that operated with a round-robin scheduling policy. This scheduler was efficient for adding and removing processes (with a lock to protect the structure)
- **v2.2**: introduced the idea of scheduling classes, permitting scheduling policies for real-time tasks, non-preemptible tasks, and non-real-time tasks. The 2.2 scheduler also included support for symmetric multiprocessing (SMP)
- **v2.4**: relatively simple scheduler that operated in O(N) time (as it iterated over every task during a scheduling event), time divided into epochs, inefficient for real-time tasks
- v2.6: O(1) scheduler, was designed to solve many of the problems with the 2.4 scheduler—namely, the scheduler was not required to iterate the entire task list to identify the next task to schedule, very efficient but the code base was gigantic and obscure with magic constants, difficult to maintain. Due to the pressure of these problems and another proposal for a scheduler by Con Kolivas (the Rotating Staircase Deadline Scheduler), the O(1) was replaced by the CFS (in v2.6.23), that is today in the stable branch of the kernel.

#### Up to kernel 2.6.7

The scheduling algorithm used in earlier versions of Linux was quite simple: at every process switch the kernel scanned the entire list of runnable processes (O(n)), computed their priorities and selected the "best" process to run.

For running the algorithm, the time is divided into epochs, at the end of an epoch, every process has run once, using its own quantum if possible. If a process did not use the whole quantum, they have half of the <sup>3</sup> remaining time slice added to the new timeslice.

```
asmlinkage void schedule(void) {
    int this cpu, c; /* weight */
    . . .
    repeat schedule:
    /* Default process to select.. */
    next = idle_task(this_cpu);
    c = -1000; /* weight */
    list_for_each(tmp, &runqueue_head) {
        p = list_entry(tmp, struct task_struct, run_list);
        if (can_schedule(p, this_cpu)) {
          int weight = goodness(p, this_cpu, prev->active_mm);
          if (weight > c)
              c = weight, next = p;
```

#### The Goodness function

The goodness is computed as follows:

```
goodness (p) = 20 - p->nice (base time quantum)
+ p->counter (ticks left in the time quantum)
+ 1 (if page table is shared with the previous process)
+ 15 (if SMP and p was last running of the same CPU)
```

Possible values:

- -1000 never select this process to run
- O out of time slice (p->counter == 0)
- >0 the goodness value, the higher the better
- +1000 real-time process, always select

#### **Epoch management**

```
. . .
/* Do we need to re-calculate counters? */
if (unlikely(!c)) {
    struct task struct *p;
    spin unlock irg(&rungueue lock);
    read lock(&tasklist lock);
    for each task(p)
        p->counter = (p->counter >> 1) + NICE TO TICKS(p->nice);
    read unlock(&tasklist lock);
    spin lock irg(&rungueue lock);
                                                                6 - p->nice/4
    goto repeat schedule;
. . .
```

#### Analysis

Disadvantages:

- a non-runnable task is also searched to determine its goodness
- mixture of runnable/non-runnable tasks into a single runqueue in any epoch
- Performance problems on SMP, as the length of critical sections depends on system load

Advantages:

- perfect Load Sharing
- no CPU underutilization for any workload type
- no (temporary) binding of threads to CPUs

#### Contention on SMP



#### From kernel 2.6.8

The O(1) scheduler has been introduced from version 2.6.8 by Ingo Molnàr. The principal characteristic of the algorithm is that schedules tasks in **constant time**, independently of the number of active processes.

It introduced:

- the global priority scale which we discussed;
- **early preemption**: if a task enters the TASK\_RUNNING state its priority is checked to see whether to call schedule();
- **static** priority for real-time tasks;
- **dynamic** priority for other tasks, recalculated at the end of their timeslice (increases interactivity).

#### Runqueues

}

As we already discussed, each CPU has its own **struct** runqueue, the concept has been introduced with this scheduler. However, this data structure keeps two pointers (to 2 sub-runqueues): one to the list of **active** processes and one to the list of **expired** processes. These are not just linked list, they are pointers to **struct** prio\_array.

```
struct runqueue {
    /* number of runnable tasks */
    unsigned long nr_running;
```

```
struct prio_array *active;
struct prio_array *expired;
struct prio_array arrays[2];
```

The **struct** prio\_array maintains an array of list\_head for each possible value of the priority, so 140 linked lists and a bitmap.

```
struct prio_array {
    int nr_active;
    unsigned long bitmap[BITMAP_SIZE];
    struct list_head queue[MAX_PRIO];
}
```

Runqueues



#### The job of the scheduler

The idea behind the two sub-runqueues, is simple: when a task on the active runqueue uses all of its time slice it's moved to the expired runqueue. During this move the time slice is recalculated (and so its priority). If no task exists on the active runqueue the pointers for the active and the expired runqueues are swapped.

The scheduler always **chooses the task on the highest priority list to execute**. To make this process efficient, a bitmap is used to defined when tasks are on a given priority list. Therefore, on most architectures, the instruction find\_first\_bit\_set() is used to find the highest priority bit set in one of five 32-bit words.

As a consequence, the time it takes to find a task to execute depends not on the number of active tasks but instead on the **fixed number of priorities**.

#### Runqueues: the bitmap

0



13

Х Х Х Х Х Bit 139, priority 139 Х Х Χ Х

#### Prioritization

To prevent tasks from hogging the CPU and thus starving other tasks that need CPU access, the O(1) scheduler can dynamically alter a task's priority. It does so by penalizing tasks that are bound to a CPU and rewarding tasks that are I/O bound. I/O-bound tasks commonly use the CPU to set up an I/O and then sleep awaiting the completion of the I/O. This type of behavior gives other tasks access to the CPU.

Because I/O-bound tasks are viewed as altruistic for CPU access, their priority is decreased (a reward) by a maximum of five priority levels. CPU-bound tasks are punished by having their priority increased by up to five levels. Tasks are determined to be I/O-bound or CPU-bound based on an *interactivity heuristic*. A task's interactiveness metric is calculated based on how much time the task executes compared to how much time it sleeps. Note that because I/O tasks schedule I/O and then wait, an I/O-bound task spends more time sleeping and waiting for I/O completion. This increases its interactive metric.

It's important to note that priority adjustments are performed only on user tasks, not on real-time tasks.

#### **Cross-CPU Scheduling**

Once a task lands on a CPU, it might use up its time slice and get put back on a prioritized queue for rerunning -- but how might it ever end up on another processor?

If all the tasks on CPU<sub>i</sub> exit, the CPU<sub>i</sub> stand idle while another CPU<sub>j</sub> round-robins three, ten or several dozen other tasks. The 2.6 scheduler must, on occasion, see if cross-CPU balancing is needed. Every 200ms a CPU checks to see if any other CPU is out of balance and needs to be balanced with that processor. If the processor is idle, it checks every 1ms so as to get started on a real task earlier.

		1791	/*
10	$c_{\rm D} = c_{\rm D} p_{\rm D} c_{\rm D} c_$	1792	<pre>* idle_balance is called by schedule() if this_cpu is about to become</pre>
45	$cpu = sinp_processor_ru();$	1793	* idle. Attempts to pull tasks from other CPUs.
44	if (unlikely(!rg->nr running)) {	1794	*/
15	idle halance(cpu, rg);	1795	<pre>static inline void idle_balance(int this_cpu, runqueue_t *this_rq)</pre>
45	race_bacance(cpu, rd),	1796	{
46	lt (!rq->nr_runnıng) {	1797	<pre>struct sched_domain *sd;</pre>
47	next = <b>ra</b> -> <b>idle</b> :	1798	
10	$r_{\rm c}$ , even in the state $-0$	1799	for_each_domain(this_cpu, sd) {
48	rq->expired_timestamp = 0;	1800	if (sd->flags & SD_BALANCE_NEWIDLE) {
49	<pre>wake sleeping dependent(cpu, rg);</pre>	1801	if (load_balance_newidle(this_cpu, this_rq, sd)) {
50	doto switch tasks:	1802	/* We've pulled tasks over so stop searching ?
50	goto switch_tasks,	1803	break ;
51	}	1804	}
52	1	1805	}
52	L	1806	}
	https://elixir.bootlin.com/linux/v2.6.8/source/kernel/sched.c#L2245	1807	}
			https://elixir.bootlin.com/linux/v2.6.8/source/kernel/sched.c#L1795

#### **Staircase Scheduler**

The Staircase scheduler was proposed by Con Kolivar, 2004 but none of its schedulers have been merged in the Kernel tree.

The goal of the staircase scheduler is to increase "responsiveness" and reduce the complexity of the O(1) Scheduler. It is mostly based on dropping the priority recalculation, replacing it with a simpler rank-based scheme

It is supposed to work better up to ~10 CPUs (tailored for desktop environments).

#### **Staircase Scheduler**

The expired array is removed and the staircase data structure is used instead. An expired process will be put back into the staircase, but at the next lower rank. It can, thus, continue to run, but at a lower priority. When it exhausts another time slice, it moves down again. And so on. The following little table shows how long the process spends at each priority level:

	Priority rank										
Iteration	Base	-1	-2	-3	-4	-5	-6	-7	-8	-9	•••
1	1	1	1	1	1	1	1	1	1	1	
2		2	1	1	1	1	1	1	1	1	
3			3	1	1	1	1	1	1	1	

When a process reaches the end of the staircase (iteration 2), it gets the previous base priority -1 but with one more timeslice. If a process sleeps (i.e., an interactive process) it gets back up in the staircase.

This approach favors interactive processes rather CPU-bound ones.

https://lwn.net/Articles/87729/

#### From v2.6.23

The Completely Fair Scheduler has been merged in October 2007. This is since then the default Scheduler. The CFS models an "ideal, precise multitasking CPU" on real hardware.

It is based on a red-black tree, where nodes are ordered by process execution time in nanoseconds. A maximum execution time is also calculated for each process.

The main idea behind the CFS is to maintain balance (**fairness**) in providing processor time to tasks. This means processes should be given a fair amount of the processor. When the time for tasks is out of balance (meaning that one or more tasks are not given a fair amount of time relative to others), then those out-of-balance tasks should be given time to execute.

#### The Virtual Runtime

To determine the balance, the CFS maintains the amount of time provided to a given task in what's called the *virtual runtime*. The smaller a task's virtual runtime - meaning the smaller amount of time a task has been permitted access to the processor - the higher its need for the processor. The CFS also includes the concept of sleeper fairness to ensure that tasks that are not currently runnable (for example, waiting for I/O) receive a comparable share of the processor when they eventually need it.

But rather than maintain the tasks in a run queue, as has been done in prior Linux schedulers, the CFS **maintains a time-ordered red-black tree**. A red-black tree is a tree with a couple of interesting and useful properties. First, it's self-balancing, which means that no path in the tree will ever be more than twice as long as any other. Second, operations on the tree occur in  $O(\log n)$  time (where n is the number of nodes in the tree). This means that you can insert or delete a task quickly and efficiently.



With tasks stored in the time-ordered red-black tree, tasks with the gravest need for the processor (lowest virtual runtime) are stored toward the **left side** of the tree, and tasks with the least need of the processor (highest virtual runtimes) are stored toward the **right side** of the tree. The scheduler then, to be fair, picks the left-most node of the red-black tree to schedule next to maintain fairness. The task accounts for its time with the CPU by adding its execution time to the virtual runtime and is then inserted back into the tree if runnable. In this way, tasks on the left side of the tree are given time to execute, and the contents of the tree migrate from the right to the left to maintain fairness. Therefore, each runnable task chases the other to maintain a balance of execution across the set of runnable tasks.

#### Where have the priorities gone?

CFS doesn't use priorities directly but instead **uses them as a decay factor** for the time a task is permitted to execute. Lower-priority tasks have higher factors of decay, where higher-priority tasks have lower factors of delay. This means that the time a task is permitted to execute dissipates more quickly for a lower-priority task than for a higher-priority task. That's an elegant solution to avoid maintaining run queues per priority.



11. Scheduling

### **Context Switch**



Advanced Operating Systems and Virtualization

#### **Context Switch**

Context switch starts with the function in kernel/sched/core.c

The function does some checks on memory (\*active\_mm) and according to the situations performs some operation. Remember that a context switch can happen in the following cases.

4288	/*				
4289	*	kernel	->	kernel	lazy + transfer active
4290	*	user	->	kernel	lazy + mmgrab() active
4291	*				
4292	*	kernel	->	user	<pre>switch + mmdrop() active</pre>
4293	*	user	->	user	switch
4294	*/	/			

https://elixir.bootlin.com/linux/v5.11/source/kernel/sched/core.c#L4276

#### **Context Switch**

Then the function calls switch\_to().

4326		/* Here we just switch the register state and the stack. */
4327		<pre>switch_to(prev, next, prev);</pre>
4328		barrier();
4329		
4330		return finish_task_switch(prev);
4331	}	

https://elixir.bootlin.com/linux/v5.11/source/kernel/sched/core.c#L4276

The switch\_to() is architecture-specific and mainly executes the following two tasks:

- TSS update
- CPU control registers update

https://elixir.bootlin.com/linux/v5.11/source/arch/x86/include/asm/switch\_to.h#L47

#### **Context Switch**

#### switch\_to()

#### As you can see:

- firstly the callee-saved register are pushed in the stack
- then the stack pointer is swapped
- in the end the callee-saved registers are popped from the stack (since they was pushed by the new process to be scheduled)

After these actions the control is passed to \_\_\_switch\_to that performs the switch of the FPU registers and of FS and GS registers.

But where the instruction pointer is set?

```
/*
       * %rdi: prev task
       * %rsi: next task
224
       */
      .pushsection .text, "ax"
226
      SYM FUNC START( switch to asm)
              /*
228
                * Save callee-saved registers
               * This must match the order in inactive task frame
230
                */
               pushq
                      %rbp
232
               pushq
                       %rbx
                       %r12
               pushq
234
                       %r13
               pushq
               pusha
                       %r14
236
               pushq
                       %r15
238
              /* switch stack */
239
                      %rsp, TASK threadsp(%rdi)
               movq
                      TASK threadsp(%rsi), %rsp
               mova
241
242
      #ifdef CONFIG_STACKPROTECTOR
243
               movq
                      TASK stack canary(%rsi), %rbx
244
                      %rbx, PER CPU VAR(fixed percpu data) + stack canary offset
              mova
245
      #endif
246
247
      #ifdef CONFIG RETPOLINE
248
              /*
249
               * When switching from a shallower to a deeper call stack
250
               * the RSB may either underflow or use entries populated
               * with userspace addresses. On CPUs where those concerns
               * exist, overwrite the RSB with entries which capture
               * speculative execution to prevent attack.
254
                */
255
              FILL RETURN BUFFER %r12, RSB_CLEAR_LOOPS, X86_FEATURE_RSB_CTXSW
256
      #endif
257
258
              /* restore callee-saved registers */
259
                       %r15
               popq
260
                       %r14
               popq
                       %r13
               popq
                       %r12
               popq
                       %rbx
               popq
               popq
                       %rbp
                       switch to
              jmp
      SYM FUNC END( switch to asm)
      .popsection
```

# Advanced Operating Systems and Virtualization

[11] Scheduling

LECTURER Gabriele **Proietti Mattia** 

BASED ON WORK BY http://www.ce.uniroma2.it/~pellegrini/



gpm.name · proiettimattia@diag.uniroma1.it

