# Source Reading: CFS

Chao Gao

5142029014

June 13, 2017

# 1   Introduction

In this project, we are required to read the linux kernel source and totally understand a part of linux by the code level. There is no doubt that this will greatly help us to dig into the linux kernel and get to know how it works. And I try to know more about the CFS.

The **Completely Fair Scheduler** is the default process scheduler in the linux kernel from 2.6.23. It handles CPU resource allocation for executing processes, and aims to maximize overall CPU utilization while also maximizing interactive performance.

The original scheduler is of O(1) time complexty based on run queues, and CFS is baesd on red-black tree. In part 3 I will dig into the source code.

My OS environment is Ubuntu 16.04 LTS, the kernel version is **4.9.13**, all the scheduler part code is in kernel/sched folder, and the main part about CFS is in the file **kernel/sched/fair.c**.

# 2   Process scheduler

A good process scheduler should be designed as follows.

- Fairness: every process gets resonable CPU time.

- Efficient: always has a process running on the CPU.

- Response time: make the iteractive process quicker to get response.

- Thoughput: handle more processes in a specific time.

There are two kinds of processes, normal processes and real-time processes. CFS is used for normal processes.

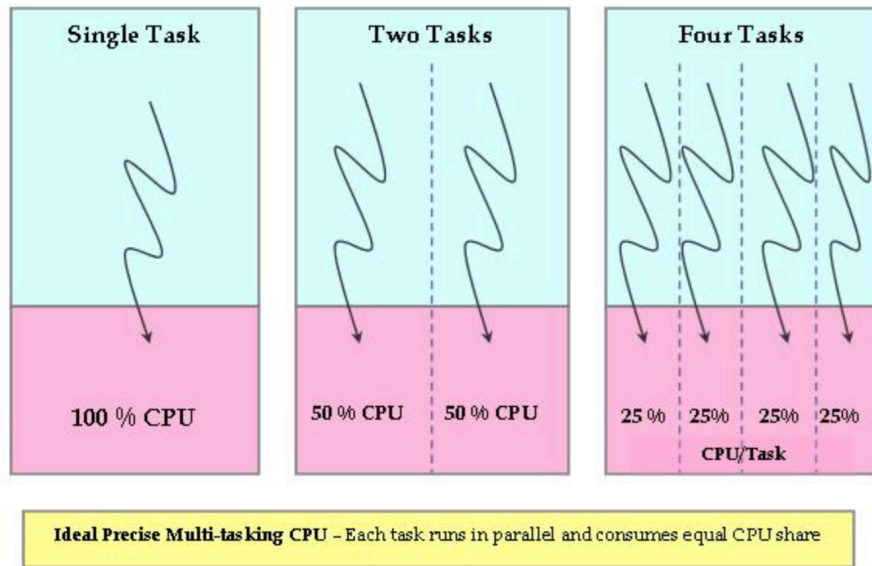The motivation of CFS is getting all tasks executed fairly.

Figure 1: Figure 1.

Simply speaking about CFS algorithm in my own words, put all the tasks as nodes of a **red-black tree**, use **virtual tuntime** to measures the time that the task has used, the task with the least virtual time is the **leftmost node** of the tree, which has been cached. When the picked process is removed, the tree is updated.

# 3 Code detail

## 3.1 sched_entity

In process scheduling, process is not the actual entity. Because in linux, there exists not only process scheduling but also group scheduling. The struct of process is defined in linux/sched.h.

```
1  struct task_struct {
2  #ifdef CONFIG_THREAD_INFO_IN_TASK
3      struct thread_info thread_info;
4  #endif
5      volatile long state; //-1 unrunnable,0 runnable,>0 stopped
```

```
6        void *stack;
7        atomic_t usage;
8        unsigned int flags; //per process flags, defined below
9        unsigned int ptrace;
10
11       ...
12
13       int prio, static_prio, normal_prio;
14       unsigned int rt_priority;
15       const struct sched_class *sched_class;
16       struct sched_entity se;
17       struct sched_rt_entity rt;
18  #ifdef CONFIG_CGROUP_SCHED
19       struct task_group *sched_task_group;
20  #endif
21       struct sched_dl_entity dl;
22
23       ...
```

Here **struct sched_entity se** is the scheduling entity object. And struct sched_entity is also defined in linux/sched.h.

```
1  struct sched_entity {
2      struct load_weight   load;    /* for load−balancing */
3      struct rb_node       run_node;//node on red−black tree
4      struct list_head     group_node;
5      unsigned int         on_rq; //if on the run queue
6
7      u64          exec_start;
8      u64          sum_exec_runtime; //whole running time
9      u64          vruntime;    //virtual run time
10     u64          prev_sum_exec_runtime;
11
12     u64          nr_migrations;
```

```
13        struct  cfs_rq              * cfs_rq ;
14        ...
```

There is a **sched_class** in task_struct which consists of scheduling functions. Because in different scheduling algorithms, the operations are different. Therefore, when we are using CFS, just create the corresponding sched_class and invoke corresponding scheduling functions.

```
1   struct  sched_class {
2       const  struct  sched_class *next ;
3
4       void  (* enqueue_task ) ( struct  rq *rq ,
5                   struct  task_struct *p, int  flags );
6       void  (* dequeue_task ) ( struct  rq *rq ,
7                   struct  task_struct *p, int  flags );
8       void  (* yield_task ) ( struct  rq *rq );
9       bool  (* yield_to_task ) ( struct  rq *rq ,
10                  struct  task_struct *p, bool  preempt );
11
12      void  (* check_preempt_curr ) ( struct  rq *rq ,
13                  struct  task_struct *p, int  flags );
14
15      struct  task_struct * (* pick_next_task ) ( struct  rq *rq ,
16                          struct  task_struct *prev ,
17                          struct  pin_cookie cookie );
18      void  (* put_prev_task ) ( struct  rq *rq ,
19                      struct  task_struct *p);
20      ...
21      void  (* set_curr_task ) ( struct  rq *rq );
22      void  (* task_tick ) ( struct  rq *rq ,
23                      struct  task_struct *p, int  queued );
24      ...
25  };
```

Just take a look at some functions.

- enqueue_task: when a task is about to running, put the sched_entity to the red-black tree.

- dequeue_task: when a task is finished, remove the sched_entity from the red-black tree.

- pick_next_task: choose the next process.

- task_tick: process switch.

## 3.2  Vruntime

Vruntime is the very important concept in the CFS.

Let's recap the basic idea of CFS, in the ideal state, each process can get the same time slice, and run on the CPU at the same time. However, in fact at one time there is only one process occupying CPU. In other words, when a process is running, other processes have to wait. In order to achieve fairness, CFS must punish the currently running process so that the processes that are waiting are scheduled in the next time.

CFS uses vruntime to record the execution time of the process.

Every period counter generates a tick interrupt. At this time, function **scheduler_tick** would be invoked, which is defined in kernel/sched/core.c.

```
1  void scheduler_tick(void)
2  {
3      ...
4      raw_spin_lock(&rq->lock);
5      update_rq_clock(rq);
6
7      //invoke scheduler
8      curr->sched_class->task_tick(rq, curr, 0);
9      cpu_load_update_active(rq);
10     calc_global_load_tick(rq);
11     raw_spin_unlock(&rq->lock);
```

```
12
13        perf_event_task_tick();
14
15        ...
```

Since the architecture is made up of modules. So if we use CFS, at this time program will jump to function **task_tick_fair**, which is defined in kernel/sched/-fair.c.

```
1  static void task_tick_fair(struct rq *rq,
2          struct task_struct *curr, int queued)
3  {
4      ...
5
6      for_each_sched_entity(se) {
7          cfs_rq = cfs_rq_of(se);
8          entity_tick(cfs_rq, se, queued);
9      }
10     ...
11 }
```

Then take a look at function **entity_tick**.

```
1  static void
2  entity_tick(struct cfs_rq *cfs_rq,
3              struct sched_entity *curr, int queued)
4  {
5      /*
6       * Update run-time statistics of the 'current'.
7       */
8      update_curr(cfs_rq);
9      ...
10
```

```
11          if (cfs_rq->nr_running > 1)
12              check_preempt_tick(cfs_rq, curr);
13  }
```

As the source code comments say, function **update_curr** is used to update the statistics of the run-time(vruntime), and kernel will deside how to schedule the processes.

```
1   static void update_curr(struct cfs_rq *cfs_rq)
2   {
3       struct sched_entity *curr = cfs_rq->curr;
4       u64 now = rq_clock_task(rq_of(cfs_rq));
5       u64 delta_exec;
6
7       if (unlikely(!curr))
8           return;
9       //get the CPU time after the latest upate of load
10      delta_exec = now - curr->exec_start;
11      if (unlikely((s64)delta_exec <= 0))
12          return;
13
14      curr->exec_start = now;
15
16      schedstat_set(curr->statistics.exec_max,
17              max(delta_exec, curr->statistics.exec_max));
18
19      //add the delta_exec to the current execution time
20      curr->sum_exec_runtime += delta_exec;
21      schedstat_add(cfs_rq->exec_clock, delta_exec);
22
23      //calculate the vruntime
24      curr->vruntime += calc_delta_fair(delta_exec, curr);
25      update_min_vruntime(cfs_rq);
26
```

```
27          ...

28

29     }
```

Note that function **calc_delta_fair** translates realtime to vruntime. The priority is larger, the vruntime grows more slowly.

```
1   static  inline  u64  calc_delta_fair (u64 delta ,
2               struct  sched_entity *se )
3   {
4       if  ( unlikely ( se−>load . weight != NICE_0_LOAD ))
5         delta = __calc_delta ( delta , NICE_0_LOAD, &se−>load );
6
7       return  delta ;
8   }
```

Take a look at function **__calc_delta**.

```
1   static u64 __calc_delta (u64 delta_exec ,
2               unsigned long weight , struct load_weight *lw )
3   {
4       u64 fact = scale_load_down ( weight );
5       int  shift = WMULT_SHIFT;
6
7       __update_inv_weight (lw );
8
9       if  ( unlikely ( fact >> 32)) {
10          while  ( fact >> 32) {
11              fact >>= 1;
12              shift −−;
13          }
14      }
15
```

```
16        /* hint to use a 32x32->64 mul */
17        fact = (u64)(u32)fact * lw->inv_weight;
18
19        while (fact >> 32) {
20            fact >>= 1;
21            shift --;
22        }
23
24        return mul_u64_u32_shr(delta_exec, fact, shift);
25 }
```

We can know **vruntime = real-execution-time * NICE_0_LOAD/lw.weight**.

## 3.3    CFS schduler class

Let's return back to the struct **sched_entity**.

There is a member called **cfs_rq**, this is the CFS runnning queue, which is defined in kernel/sched/sched.h.

```
1  struct cfs_rq {
2      struct load_weight load;
3      unsigned int nr_running, h_nr_running;
4
5      u64 exec_clock;
6      u64 min_vruntime;
7  #ifndef CONFIG_64BIT
8      u64 min_vruntime_copy;
9  #endif
10     struct rb_root tasks_timeline;
11
12     //next scheduler node(leftmost node of red-black tree)
13     struct rb_node *rb_leftmost;
14
15     struct sched_entity *curr, *next, *last, *skip;
```

```
16    ...

17

18

19    #ifdef CONFIG_FAIR_GROUP_SCHED
20        struct rq *rq;    //cpu runqueue to which this cfs_rq is attached
21    ...

22

23    };
```

Note that **rb_node *rb_leftmost** is the leftmost node of the red-black tree, which
will be the next scheduling node. And sched_entity ***curr** is the current running
process, sched_entity ***next** means some urgent processes which have to run on
the CPU immediately even with complict to CFS. And one important point is that
the root of red-black tree is **rb_node * tasks_timeline**.



Figure 2: Figure 2.

Now we get to the very important stage, we get to know the CFS schedule
class, which consists of all the functions with regard to scheduling algorithms.
**fair_sched_class** in defined in kernel/fair.c.

```
1    const struct sched_class fair_sched_class = {
```

```
2          . next              = &idle_sched_class ,
3          . enqueue_task      = enqueue_task_fair ,
4          . dequeue_task      = dequeue_task_fair ,
5          . yield_task        = yield_task_fair ,
6          . yield_to_task     = yield_to_task_fair ,
7
8          . check_preempt_curr = check_preempt_wakeup ,
9
10         . pick_next_task    = pick_next_task_fair ,
11         . put_prev_task     = put_prev_task_fair ,
12         ...
13
14         . set_curr_task     = set_curr_task_fair ,
15         . task_tick         = task_tick_fair ,
16         . task_fork         = task_fork_fair ,
17
18         . prio_changed      = prio_changed_fair ,
19         . switched_from     = switched_from_fair ,
20         . switched_to       = switched_to_fair ,
21
22         . get_rr_interval   = get_rr_interval_fair ,
23
24         . update_curr       = update_curr_fair ,
25         ...
26    };
```

Give a partial list of the function.

- enqueue_task: called when a task enters a runnable state. It puts the scheduling entity into the red-black tree.

- 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.

- 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 choose 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 derives the running preemption.

## 3.4   Process switch

Since we have known about CFS schduler class and vruntime, now we can turn to the process switch part.

In **pick_next_task** function which is defined in kernel/core.c, scheduler will do the scheduling. And next the CFS scheduler's function will be invoked, which is called **pick_next_task_fair** function. And the task chosen will start executing.

```
1  static struct task_struct *
2  pick_next_task_fair(struct rq *rq,
3          struct task_struct *prev, struct pin_cookie cookie)
4  {
5      struct cfs_rq *cfs_rq = &rq->cfs;
6      struct sched_entity *se;
7      struct task_struct *p;
8      int new_tasks;
9
10     ...
11
12     if (!cfs_rq->nr_running)
13         goto idle;
14
15     put_prev_task(rq, prev);
16
17     do {
```

```
18          se = pick_next_entity(cfs_rq, NULL);
19          set_next_entity(cfs_rq, se);
20          cfs_rq = group_cfs_rq(se);
21      } while (cfs_rq);
22
23      p = task_of(se);
24
25      if (hrtick_enabled(rq))
26          hrtick_start_fair(rq, p);
27
28      return p;
29
30      ...
31  }
```

We can see **pick_next_entity** function to pick the next sched entity, and remove the node in red-black tree and update the red-black tree by using **set_new_entity** function.

Here I dig into pick_next_entity function. And it is implemented by a core function called **__pick_first_entity**.

```
1  struct sched_entity *__pick_first_entity(
2                          struct cfs_rq *cfs_rq)
3  {
4      struct rb_node *left = cfs_rq->rb_leftmost;
5
6      if (!left)
7          return NULL;
8
9      return rb_entry(left, struct sched_entity, run_node);
10 }
```

We can see the leftmost node of red-black tree is chosen, which is the same as our

knowledge learned in class. Our scheduler successfully works.

# 4    Conclusion

CFS's design can be summed up in a single sentence: CFS basically models an "ideal, precise multi-tasking CPU" on real hardware.

Through the code reading, I get to know more about CFS. In this report, I mainly cover the following parts.

- know the sched entity in scheduling.

- understand the vruntime.

- know the CFS scheduler class.

- basic process switch in CFS.

Of course this is not the end, I will get to know more about group scheduler in the future.

In conclusion, in this project, I mainly read the source code about CFS, and since the code is too large, I just get a brief understanding but it is still really helpful for me.

# 5    Summary

Though kernel source reading, I truly understand more about CFS algorithm, and the tough thing is linux kernel changes so rapidly that many resources are not useful at today's time. And it is a hot potato to "get through" the deep function trees, luckily I get over that.

In linux kernel course this semester, I really learn a lot about linux, with pains. But this is only a beginning, linux will be more frequently used in my next learning career.

Finally, thanks a lot to Prof.Chen for teaching and TAs for answering questions.