

Report on Project 2

6/14/2018

516030910306 Hongyi Guo (Henry Guo)

1.include/linux/sched.h

(1) Define WRR scheduler as 6

```
43  #define SCHED_WRR 6 /* added by henry */
```

(2) Define wrinfo, to store information of CPU, for multi-CPU.

```
54  #define MAX_CPUS 8
55  struct wrinfo {
56      int num_cpus;
57      int nr_running[MAX_CPUS];
58      int total_weight[MAX_CPUS];
59  };
60
```

(3) Declare wr_rq, run queue of wr scheduler.

```
162  struct wr_rq; /* added by henry */
```

(4) Define wr scheduler entity. Run_list is the list of tasks of this entity. Time slice depends on whether this process is a foreground/background process.

```
1262  /* added by henry */
1263  struct sched_wrr_entity {
1264      struct list_head run_list;
1265      unsigned int time_slice;
1266
1267      #ifdef CONFIG_SMP
1268          int weight;
1269      #endif
1270
1271      #ifdef CONFIG_WRR_GROUP_SCHED
1272          struct sched_wrr_entity *parent;
1273          /* rq on which this entity is (to be) queued: */
1274          struct wr_rq *wrr_rq;
1275          /* rq "owned" by this entity/group: */
1276          struct wr_rq *my_q;
1277      #endif
1278  };
```

(5) Define the time slice for wrd scheduler, 100 for foreground process and 10 for background processes.

```
1285  /*
1286  * added by henry, used only for SCHED_WRR tasks
1287  * for foreground task, timeslice is 100 msecs
1288  * for background task, timeslice is 10 msecs
1289  */
1290  #define WRR_TIMESLICE_FORE (100 * HZ / 1000)
1291  #define WRR_TIMESLICE_BACK (10 * HZ / 1000)
1292
```

(6) Add wrd_priority and wrd scheduler entity in task_struct

```
1317  unsigned int wrd_priority; /* added by henry */
1320  struct sched_wrd_entity wrd; /* added by henry */
```

2. kernel/sched/sched.h

(1) Define whether to use multi-queues on priority, and whether to print extra debugging information.

```
8  #define CONFIG_WRR_PRIO 1
9  // #define CONFIG_WRR_DEBUG 1
```

(2) If we want to take multi-CPU into consideration, then declare extern struct and spinlock.

```
13  /* by henry */
14  #ifdef CONFIG_SMP
15  extern struct wrd_info my_wrd;
16  extern raw_spinlock_t wrd_info_locks[MAX_CPUS];
17  #endif
```

(3) Define priority array, who has up to MAX_RT_PRIO queues. Its bitmap indicates whether the corresponding queue is empty(0) or not(1).

```

66  /*
67   * by henry
68   *
69   * This is the priority-queue data structure of the WRR scheduling class:
70   */
71  #ifdef CONFIG_WRR_PRIO
72  struct wrr_prio_array {
73      DECLARE_BITMAP(bitmap, MAX_RT_PRIO+1); /* include 1 bit for delimiter */
74      struct list_head queue[MAX_RT_PRIO];
75  };
76  #endif

```

(4) Declare runqueue of wrr scheduler.

```

101 struct wrr_rq; /* added by henry */

```

(5) Define struct wrr_rq. "wrr_nr_running" is the tasks running on all wrr runqueues.

In priority-based multi-queue mode, "active" is an array of wrr_runqueue's.

In single-queue mode, one list is defined by its head "queue".

If SMP is defined, define wrr_nr_total as the total weight of the tasks in this runqueue.

```

301  ////////////////////////////////////////////////// This is the part to study ///////////////////////////////////
302
303  /*
304   * added by henry
305   * Weighted Round-Robin classes' related field in a runqueue:
306   */
307  struct wrr_rq {
308  #ifdef CONFIG_WRR_PRIO
309      struct wrr_prio_array active;
310  #else
311      struct list_head queue;
312  #endif
313      unsigned int wrr_nr_running;
314
315  #ifdef CONFIG_SMP
316      unsigned long wrr_nr_migratory;
317      unsigned long wrr_nr_total;
318      int overloaded;
319      struct plist_head pushable_tasks;
320  #endif
321  };
322

```

(6) Define a wrr_rq named wrr in struct rq.

```

417 struct wrr_rq wrr; /* added by henry */

```

(7) Declare extern wrr_sched_class.

```
896 extern const struct sched_class wrr_sched_class; /* added by henry */
```

(8) Declare extern idle balance function "idle_balance_wrr". Note this is active only when CONFIG_SMP defined.

```
902 extern void idle_balance_wrr(struct rq *this_rq);
```

(9) Declare extern initialization function "init_sched_wrr_class" defined in wrr.c.

```
919 extern void init_sched_wrr_class(void); /* by henry */
```

(10) Declare extern runqueue initialization function "init_wrr_rq" defined in wrr.c.

```
1192 extern void init_wrr_rq(struct wrr_rq *wrr_rq, struct rq *rq); /* by henry */
```

3. kernel/sched/core.c

(1) Define wrr_info struct to store information of wrr scheduler, only used when CONFIG_SMP is defined.

```
90 /* by henry */
91 #ifdef CONFIG_SMP
92 struct wrr_info my_wrr = {
93     .num_cpus = 0,
94     .nr_running = { 0 },
95     .total_weight = { 0 },
96 };
97 #endif
```

(2) Call "INIT_LIST_HEAD" to initialize sched_wrr_entity in "__sched_fork" function.

```
1743 INIT_LIST_HEAD(&p->wrr.run_list); /* by henry */
```

(3) If CONFIG_SMP defined, call idle balance function of wrr scheduler when no wrr task is running on this rq.

```
3388 /* by henry */
3389 #ifdef CONFIG_SMP
3390     if (unlikely(!rq->wrr.wrr_nr_running))
3391         idle_balance_wrr(rq);
3392 #endif
```

(4) Defined in __sched_setscheduler.

```
4255 printk("I AM IN __SCHED_SETSCHEDULER\n");
```

Take SCHED_WRR into consideration.

```
4268         if (policy != SCHED_FIFO && policy != SCHED_RR &&
4269             policy != SCHED_NORMAL && policy != SCHED_BATCH &&
4270             policy != SCHED_IDLE && policy != SCHED_WRR)
4271             return -EINVAL;
4272     }
```

```
4285     if ((rt_policy(policy) || (policy == SCHED_WRR))
4286         != (param->sched_priority != 0)) {
4287         printk("recheck 2 failed\n");
4288         return -EINVAL;
4289     }
```

(5) Defined in “__setscheduler”, which do the actual job of priority change. Set “sched_class” of “p” to “wrr_sched_class” if new policy is “SCHED_WRR”, and set wrr_priority.

```
4191  /* Actually do priority change: must hold rq lock. */
4192  static void
4193  __setscheduler(struct rq *rq, struct task_struct *p, int policy, int prio)
4194  {
4195  #ifdef CONFIG_SCHED_DEBUG
4196      // printk("__setscheduler called\n");
4197      printk("[%d]\tpriority set to %d\n", p->pid, prio);
4198  #endif
4199      p->policy = policy;
4200      p->rt_priority = prio;
4201      p->wrr_priority = prio; /* by henry */
4202      p->normal_prio = normal_prio(p);
4203      /* we are holding p->pi_lock already */
4204      p->prio = rt_mutex_getprio(p);
4205
4206      if (policy == SCHED_WRR) { /* by henry */
4207          p->sched_class = &wrr_sched_class;
4208  #ifdef CONFIG_SCHED_DEBUG
4209          printk("set task scheduler as wrr\n");
4210  #endif
4211      }
4212      else if (rt_prio(p->prio)) {
4213          p->sched_class = &rt_sched_class;
4214  #ifdef CONFIG_SCHED_DEBUG
4215          printk("set task scheduler as rt\n");
4216  #endif
4217      }
4218      else {
4219          p->sched_class = &fair_sched_class;
4220  #ifdef CONFIG_SCHED_DEBUG
4221          printk("set task scheduler as fair\n");
4222  #endif
4223      }
4224      set_load_weight(n);
```

(6) A little bit revise in "void __init sched_init(void)", to initialize wrn runqueue at beginning, and calculate total number of possible CPUs.

```
7200   for_each_possible_cpu(i) {
7201       struct rq *rq;
7202
7203       rq = cpu_rq(i);
7204       raw_spin_lock_init(&rq->lock);
7205       rq->nr_running = 0;
7206       rq->calc_load_active = 0;
7207       rq->calc_load_update = jiffies + LOAD_FREQ;
7208       init_cfs_rq(&rq->cfs);
7209       init_rt_rq(&rq->rt, rq);
7210       init_wrn_rq(&rq->wrn, rq); /*by henry */
7211
7212   #ifdef CONFIG_SMP
7213       my_wrn.num_cpu++;
7214   #endif
```

4. /kernel/sched/rt.c

(1) A little bit revise, add wrn_sched_class into the circular list.

```
2041   const struct sched_class rt_sched_class = {
2042       .next           = &wrn_sched_class,
```

5. kernel/sched/wrn.c

I write this part based on the provided "rt.c".

(1) Overview.

```

500  const struct sched_class wrr_sched_class = {
501      .next                = &fair_sched_class,      /* never */
502
503      /* [required] adding/removing a task to/from a priority array */
504      .enqueue_task        = enqueue_task_wrr,        /* required */
505      .dequeue_task        = dequeue_task_wrr,        /* required */
506
507      .yield_task          = yield_task_wrr,          /* required */
508      .check_preempt_curr  = check_preempt_curr_wrr, /* required */
509
510      .pick_next_task      = pick_next_task_wrr,      /* required */
511      .put_prev_task       = put_prev_task_wrr,       /* required */
512
513      .task_fork           = task_fork_wrr,           /* required */
514
515      #ifdef CONFIG_SMP
516      .select_task_rq      = select_task_rq_wrr,      /* never */
517      .set_cpus_allowed    = set_cpus_allowed_wrr,    /* never */
518      .rq_online           = rq_online_wrr,           /* never */
519      .rq_offline          = rq_offline_wrr,          /* never */
520      .pre_schedule        = pre_schedule_wrr,        /* never */
521      .post_schedule       = post_schedule_wrr,       /* never */
522      .task_woken          = task_woken_wrr,          /* never */
523      #endif
524
525      .switched_from       = switched_from_wrr,       /* never */
526
527      .set_curr_task       = set_curr_task_wrr,       /* required */
528      .task_tick           = task_tick_wrr,           /* required */
529
530      .get_rr_interval     = get_rr_interval_wrr,     /* required */
531
532      .prio_changed        = prio_changed_wrr,        /* never */
533      .switched_to         = switched_to_wrr,         /* required */
534  };
535

```

(2) init_wrr_rq

To initialize wrr runqueue, we need to initialize every runqueue for each priority. Totally MAX_RT_PRIO+1 runqueues. Note that we need to set the last bit of array's bitmap into 1, for it's the delimiter for bitsearch.

```

130 void init_wrr_rq(struct wrr_rq *wrr_rq, struct rq *rq)
131 {
132     #ifdef CONFIG_WRR_PRIO
133         printk("init_wrr_rq called!\n");
134
135         struct wrr_prio_array *array;
136         int i;
137
138         array = &wrr_rq->active;
139         for (i = 0; i < MAX_RT_PRIO; i++) {
140             INIT_LIST_HEAD(array->queue + i);
141             __clear_bit(i, array->bitmap);
142         }
143         /* delimiter for bitsearch: */
144         __set_bit(MAX_RT_PRIO, array->bitmap);
145         wrr_rq->wrr_nr_running = 0;
146
147     #else /* !CONFIG_WRR_PRIO */
148
149         INIT_LIST_HEAD(&wrr_rq->queue);
150         wrr_rq->wrr_nr_running = 0;
151
152     #endif /* CONFIG_WRR_PRIO */
153 }

```

(3) pick_next_task_wrr

Use sched_find_first_bit function to find the first "1" in the bitmap, which is the non-empty queue with the highest priority.

```

174 static struct task_struct *pick_next_task_wrr(struct rq *rq)
175 {
176     struct task_struct *p;
177     struct wrr_rq *wrr_rq;
178     struct sched_wrr_entity *next = NULL;
179
180     wrr_rq = &rq->wrr;
181
182     #ifdef CONFIG_WRR_PRIO
183
184         struct wrr_prio_array *array = &wrr_rq->active;
185         struct list_head *queue;
186         int idx;
187
188         if (!wrr_rq->wrr_nr_running)
189             return NULL;
190
191         idx = sched_find_first_bit(array->bitmap);
192
193         BUG_ON(idx >= MAX_RT_PRIO);
194
195         queue = array->queue + idx;
196         next = list_entry(queue->next, struct sched_wrr_entity, run_list);
197
198         BUG_ON(!next);

```


(4) enqueue_task_wrr

Enqueue the task into the corresponding queue of its priority, and set the same bit to 1 in the bitmap, indicating this queue is non-empty.

```
260     struct sched_wrr_entity *wrr_se;
261     struct wrr_rq *wrr_rq;
262
263     wrr_se = &p->wrr;
264     wrr_rq = &rq->wrr;
265
266     #ifdef CONFIG_WRR_PRIO
267         int prio = wrr_se_prio(wrr_se);
268         struct wrr_prio_array *array = &wrr_rq->active;
269         struct list_head *queue = array->queue + prio;
270
271         if (wrr_se == NULL)
272             return;
273
274         if (flags & ENQUEUE_HEAD) {
275             list_add(&wrr_se->run_list, queue);
276         }
277         else {
278             list_add_tail(&wrr_se->run_list, queue);
279         }
280
281         __set_bit(prio, array->bitmap);
282
283         WARN_ON(!rt_prio(prio));
284
285     #else /* !CONFIG_WRR_PRIO */
```

```
296     wrr_rq->wrr_nr_running++;
297     inc_nr_running(rq);
```

If using multi-CPU, total weight should be updated after enqueueing.

```
299     #ifdef CONFIG_SMP
300         cpu = cpu_of(rq);
301         raw_spin_lock(&wrr_info_locks[cpu]);
302         my_wrr.nr_running[cpu]++;
303         my_wrr.total_weight[cpu] += wrr_se->weight;
304         raw_spin_unlock(&wrr_info_locks[cpu]);
305     #endif
```

(5) dequeue_task_wrr

Basically the same with enqueue, but for checking whether the corresponding queue is empty, and if so, we need to clear the bit in the bitmap.

And when using multi-CPU, we need to decrease the weight of this CPU.

```

222     struct sched_wrr_entity *wrr_se = &p->wrr;
223     struct wrr_rq *wrr_rq = &rq->wrr;
224
225     #ifdef CONFIG_WRR_PRIO
226         int prio = wrr_se_prio(wrr_se);
227         struct wrr_prio_array *array = &wrr_rq->active;
228     #endif
229
230     if (wrr_se == NULL)
231         return;
232
233     list_del_init(&wrr_se->run_list); /* delete and initialize node */
234
235     #ifdef CONFIG_WRR_PRIO
236         if (list_empty(array->queue + prio))
237             __clear_bit(prio, array->bitmap);
238         WARN_ON(!rt_prio(prio));
239     #endif
240
241     WARN_ON(!wrr_rq->wrr_nr_running);
242     wrr_rq->wrr_nr_running--;
243     dec_nr_running(rq);
244
245     #ifdef CONFIG_SMP
246         cpu = cpu_of(rq);
247         raw_spin_lock(&wrr_info_locks[cpu]);
248         my_wrr.nr_running[cpu]--;
249         my_wrr.total_weight[cpu] -= wrr_se->weight;
250         raw_spin_unlock(&wrr_info_locks[cpu]);
251     #endif

```

(6) `requeue_task_wrr`

Move the task into the end of this queue, usually called by "yield_task_wrr".

Use "list_move" and "list_move_tail" function defined in list.h.

```

314     struct sched_wrr_entity *wrr_se;
315     struct wrr_rq *wrr_rq;
316
317     wrr_se = &p->wrr;
318     wrr_rq = &rq->wrr;
319
320     #ifdef CONFIG_WRR_PRIO
321         int prio = wrr_se_prio(wrr_se);
322         struct wrr_prio_array *array = &wrr_rq->active;
323         struct list_head *queue = array->queue + prio;
324
325         if (wrr_se == NULL)
326             return;
327         if (head)
328             list_move(&wrr_se->run_list, queue);
329         else
330             list_move_tail(&wrr_se->run_list, queue);
331     #else

```

(7) yield_task_wrr

To let current task on rq give out its time slice, we just move it to the end of the queue.

```
338 static void yield_task_wrr(struct rq *rq)
339 {
340     requeue_task_wrr(rq, rq->curr, 0);
341 }
```

(8) task_tick_wrr

Decrease the time_slice of task p. If it run out of time slice, use "cgroup_path" to decide which time slice should be given to it.

```
4 /* for group path */
5 #ifndef PATH_MAX
6 #define PATH_MAX 4096
7 #endif
```

```
387 static char group_path[PATH_MAX];
388 static void task_tick_wrr(struct rq *rq, struct task_struct *p, int queued)
389 {
390     struct sched_wrr_entity *wrr_se;
391     wrr_se = &p->wrr;
392
393     if (wrr_se == NULL) {
394         return;
395     }
396     watchdog(rq, p);
397
398     printk("[%d]\tprio is %d\n", p->pid, p->wrr_priority);
399     printk("[%d]\t%d time slice left\n", p->pid, p->wrr.time_slice);
400
401     if (--wrr_se->time_slice) {
402         return;
403     }
404
405     cgroup_path(task_group(p)->css.cgroup, group_path, PATH_MAX);
406
407     if (strcmp(group_path, "/") == 0) {
408         wrr_se->time_slice = WRR_TIMESLICE_FORE;
409 #ifdef CONFIG_WRR_DEBUG
410         printk("[%d]\tttime slice is set as FOREground\n", p->pid);
411 #endif
412     } else {
413         wrr_se->time_slice = WRR_TIMESLICE_BACK;
414 #ifdef CONFIG_WRR_DEBUG
415         printk("[%d]\tttime slice is set as BACKground\n", p->pid);
416 #endif
417     }
418 }
```

```

419 #ifdef CONFIG_WRR_PRIO
420     /* decrease the priority */
421     if (p->wrr_priority < MAX_RT_PRIO-1) {
422         dequeue_task_wrr(rq, p, 1);
423         p->wrr_priority++;
424         enqueue_task_wrr(rq, p, 1);
425     }
426 #endif
427
428     if (wrr_se->run_list.prev != wrr_se->run_list.next) {
429         requeue_task_wrr(rq, p, 0);
430         set_tsk_need_resched(p);
431         return;
432     }
433 }

```

Each time task p run out of time slice, we degrade its priority by plus its wrr_priority by 1, and insert it into new queue.

(9) switched_to_wrr

When a task is switched to wrr, we distribute an initial time slice to it based on it being a foreground or background process.

```

435 static void switched_to_wrr(struct rq *rq, struct task_struct *p)
436 {
437     int foreground;
438     cgroup_path(task_group(p)->css.cgroup, group_path, PATH_MAX);
439     foreground = (strcmp(group_path, "/") == 0);
440
441     printk("group=%s\n", group_path);
442     printk("Switched to a %s WRR entity, pid=%d, proc=%s\n",
443         foreground ? "foreground" : "background", p->pid, p->comm);
444
445     p->wrr.time_slice = foreground ? WRR_TIMESLICE_FORE : WRR_TIMESLICE_BACK;
446
447     if (p->prio < rq->curr->prio)
448         resched_task(rq->curr);
449 }

```

(10) get_rr_interval_wrr

```

462 static unsigned int get_rr_interval_wrr(struct rq *rq, struct task_struct *task)
463 {
464     return task->wrr.time_slice;
465 }

```

(11) some other useful inline function.

```

10 static inline struct task_struct *wrr_task_of(struct sched_wrr_entity *wrr_se)
11 {
12     return container_of(wrr_se, struct task_struct, wrr);
13 }
14
15 static inline int wrr_se_prio(struct sched_wrr_entity *wrr_se)
16 {
17     return wrr_task_of(wrr_se)->wrr_priority;
18 }

```

(12) pick_pullable_task_wrr (when using multi-CPU's)

Check the first non-empty queue, if it has a process not running currently, then we're done, and this task can be pulled and insert to target CPU.

Otherwise, we check the next non-empty queue and on and on and on...

```

20 #ifdef CONFIG_SMP
21 static struct task_struct *pick_pullable_task_wrr(struct rq *src_rq,
22     int this_cpu)
23 {
24     struct task_struct *p;
25     struct list_head *queue;
26     struct wrr_prio_array *array;
27     struct sched_wrr_entity *pos;
28     int idx;
29     int flag; /* done or not */
30
31     array = &src_rq->active;
32     idx = sched_find_first_bit(array->bitmap);
33
34     flag = 0;
35     p = NULL;
36
37     while (idx < MAX_RT_PRIO) {
38         queue = array->queue + idx;
39         list_for_each_entry(pos, queue, run_list) {
40             p = wrr_task_of(pos);
41             if (p == src_rq->curr)
42                 continue;
43             if (cpumask_test_cpu(this_cpu, &p->cpus_allowed))
44                 flag = 1;
45         }
46
47         if (!flag)
48             idx = find_next_bit(array->bitmap, MAX_RT_PRIO, idx+1);
49         else
50             break;
51     }
52
53     return p;
54 }

```

(13) idle_balance_wrr

Find the pullable task on a CPU which has more than 1 task. And move it to the new CPU where rq is on.

```
56 void idle_balance_wrr(struct rq *this_rq)
57 {
58     int this_cpu = this_rq->cpu, cpu;
59     struct task_struct *p;
60     struct rq *src_rq;
61
62     for_each_possible_cpu(cpu) {
63         if (this_cpu == cpu)
64             continue;
65
66         src_rq = cpu_rq(cpu);
67         double_lock_balance(this_rq, src_rq);
68
69         if (src_rq->wrr.wrr_nr_running <= 1)
70             goto skip;
71
72         /* pick a task from src_rq */
73         p = pick_pullable_task_wrr(src_rq, this_cpu);
74
75         if (p) {
76             if (p == src_rq->curr)
77                 goto skip;
78             WARN_ON(!p->on_rq);
79
80             deactivate_task(src_rq, p, 0);
81             set_task_cpu(p, this_cpu);
82             activate_task(this_rq, p, 0);
83             double_unlock_balance(this_rq, src_rq);
84             return;
85         }
86     skip:
87         double_unlock_balance(this_rq, src_rq);
88     }
89 }
```

(14) select_task_rq_wrr

Find the CPU with the minimum weight, so, we can transfer tasks to this CPU.

```

91 static int
92 select_task_rq_wrr(struct task_struct *p, int sd_flag, int flags)
93 {
94     int cpu;
95     int min_cpu = task_cpu(p);
96     int min_weight = INT_MAX;
97
98     if (p->nr_cpus_allowed == 1)
99         goto out;
100
101     /* For anything but wake ups, just return the task_cpu */
102     if (sd_flag != SD_BALANCE_WAKE && sd_flag != SD_BALANCE_FORK)
103         goto out;
104
105     /* search for the cpu with minimized load weight */
106     for_each_possible_cpu(cpu) {
107         if (my_wrr.total_weight[cpu] < min_weight) {
108             min_weight = my_wrr.total_weight[cpu];
109             min_cpu = cpu;
110         }
111     }
112     return min_cpu;
113 out:
114     return min_cpu;
115 }
116
117 #else /* !CONFIG_SMP */

```