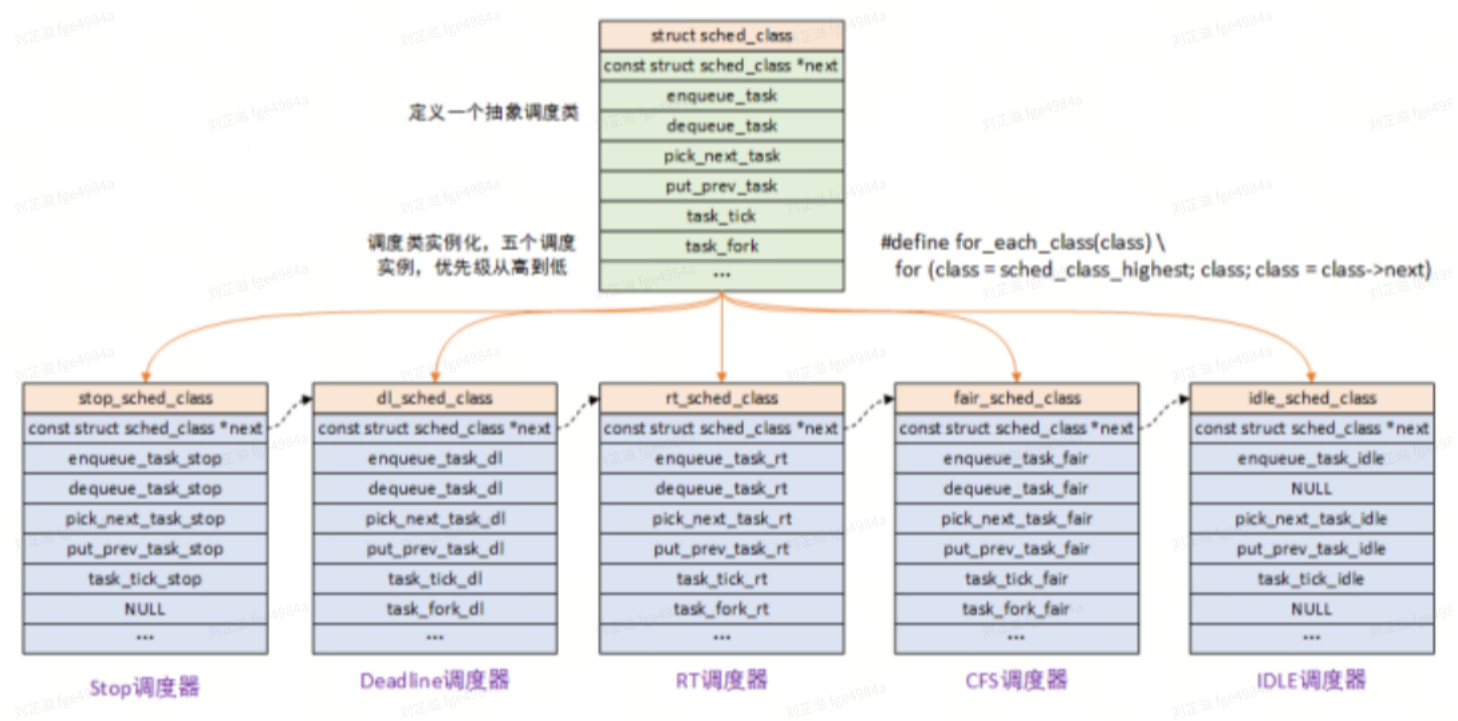


调度器类

```

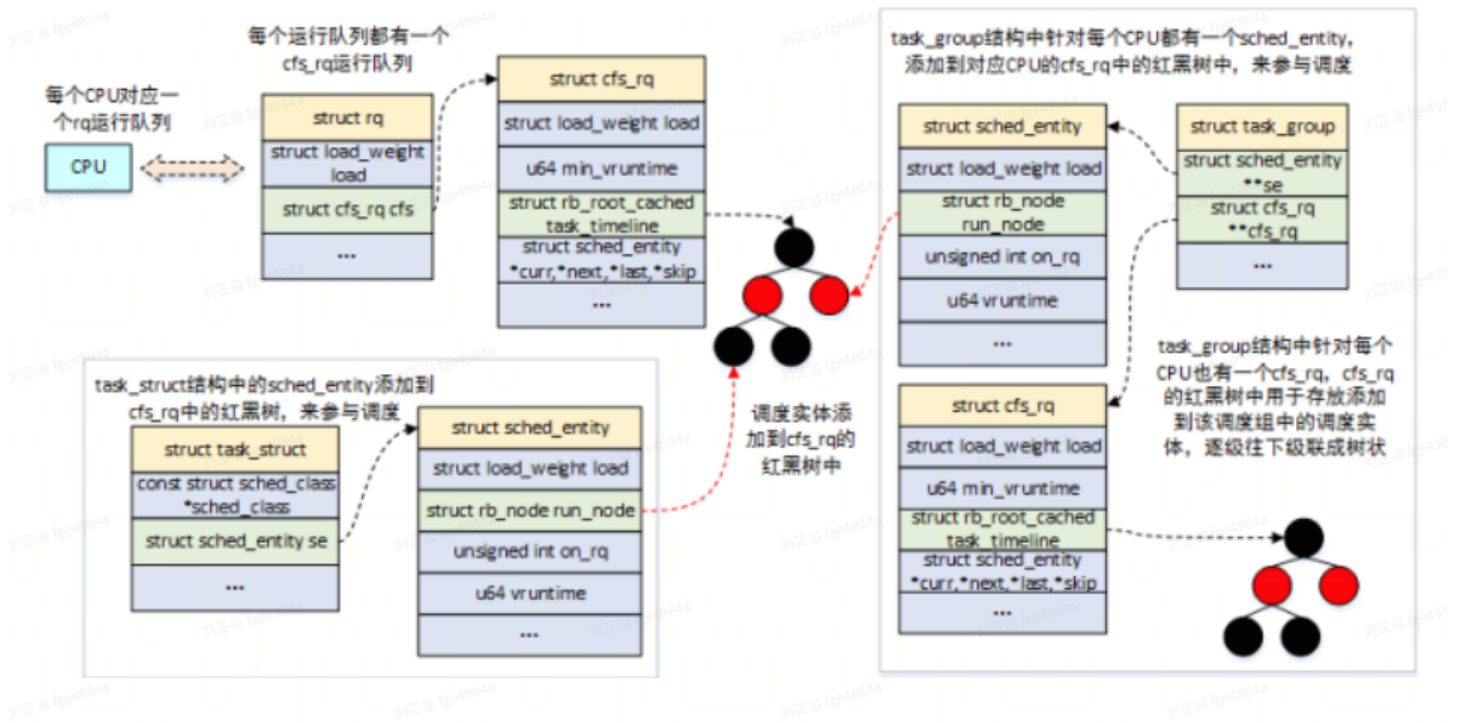
1 struct sched_class {
2     const struct sched_class *next;
3
4     void (*enqueue_task) (struct rq *rq, struct task_struct *p, int wakeup,
5         bool head);
6     void (*dequeue_task) (struct rq *rq, struct task_struct *p, int sleep);
7     void (*yield_task) (struct rq *rq);
8
9     void (*check_preempt_curr) (struct rq *rq, struct task_struct *p, int
10        flags);
11
12     struct task_struct * (*pick_next_task) (struct rq *rq);
13     void (*put_prev_task) (struct rq *rq, struct task_struct *p);
14
15     void (*set_curr_task) (struct rq *rq);
16     void (*task_tick) (struct rq *rq, struct task_struct *p, int queued);
17     void (*task_fork) (struct task_struct *p);
18
19     void (*switched_from) (struct rq *this_rq, struct task_struct *task,
20         int running);
21     void (*switched_to) (struct rq *this_rq, struct task_struct *task,
22         int running);
23     void (*prio_changed) (struct rq *this_rq, struct task_struct *task,
24         int oldprio, int running);
25
26     unsigned int (*get_rr_interval) (struct rq *rq,
27         struct task_struct *task);
28 };

```



rq/cfs_rq/task_group/sched_entity

1. struct rq: 每个 CPU 都有一个对应的运行队列;
2. struct cfs_rq: CFS 运行队列，该结构中包含了 struct rb_root_cached 红黑树，用于链接调度实体 struct sched_entity。rq 运行队列中对应了一个 CFS 运行队列，此外，在 task_group 结构中也会为每个 CPU 再维护一个 CFS 运行队列;
3. struct task_struct: 任务的描述符，包含了进程的所有信息，该结构中的 struct sched_entity，用于参与 CFS 的调度;
4. struct task_group: 组调度，Linux 支持将任务分组来对 CPU 资源进行分配管理，该结构中为系统中的每个 CPU 都分配了 struct sched_entity 调度实体和 struct cfs_rq 运行队列，其中 struct sched_entity 用于参与 CFS 的调度;
5. struct sched_entity: 调度实体，这个也是 CFS 调度管理的对象了;



sched_entity

```

1  struct sched_entity {
2      struct load_weight load; /* for load-balancing */
3      struct rb_node run_node;
4      struct list_head group_node;
5      unsigned int on_rq;
6
7      u64 exec_start;
8      u64 sum_exec_runtime;
9      u64 vruntime;
10     u64 prev_sum_exec_runtime;
11
12     u64 last_wakeup;
13     u64 avg_overlap;
14
15     u64 nr_migrations;
16
17     u64 start_runtime;
18     u64 avg_wakeup;
19     //...
20 #ifdef CONFIG_FAIR_GROUP_SCHED
21     struct sched_entity *parent;
22     /* rq on which this entity is (to be) queued: */
23     struct cfs_rq *cfs_rq;
24     /* rq "owned" by this entity/group: */
25     struct cfs_rq *my_rq;
26 #endif
27 };
  
```

cfs_rq

```
1  /* CFS-related fields in a runqueue */
2  struct cfs_rq {
3      struct load_weight load;           //CFS运行队列的负载权重值
4      unsigned int nr_running, h_nr_running; //nr_running: 运行的调度实体数（参与时
      间片计算）
5
6      u64 exec_clock;    //运行时间
7      u64 min_vruntime; //最少的虚拟运行时间，调度实体入队出队时需要进行增减处理
8  #ifndef CONFIG_64BIT
9      u64 min_vruntime_copy;
10 #endif
11
12     struct rb_root_cached tasks_timeline; //红黑树，用于存放调度实体
13
14     /*
15      * 'curr' points to currently running entity on this cfs_rq.
16      * It is set to NULL otherwise (i.e when none are currently running).
17      */
18     struct sched_entity *curr, *next, *last, *skip; //分别指向当前运行的调度实
      体、下一个调度的调度实体、CFS运行队列中排最后的调度实体、跳过运行的调度实体
19
20     #ifdef CONFIG_SCHED_DEBUG
21         unsigned int nr_spread_over;
22     #endif
23
24     #ifdef CONFIG_SMP
25         /*
26          * CFS load tracking
27          */
28         struct sched_avg avg;           //计算负载相关
29         u64 runnable_load_sum;
30         unsigned long runnable_load_avg; //基于PELT的可运行平均负载
31     #endif CONFIG_FAIR_GROUP_SCHED
32         unsigned long tg_load_avg_contrib; //任务组的负载贡献
33         unsigned long propagate_avg;
34     #endif
35         atomic_long_t removed_load_avg, removed_util_avg;
36     #ifndef CONFIG_64BIT
37         u64 load_last_update_time_copy;
38     #endif
39
40     #ifdef CONFIG_FAIR_GROUP_SCHED
41         /*
42          *  $h\_load = weight * f(tg)$ 
```

```

43      *
44      * Where  $f(tg)$  is the recursive weight fraction assigned to this group.
45      */
46      unsigned long h_load;
47      u64 last_h_load_update;
48      struct sched_entity *h_load_next;
49  #endif /* CONFIG_FAIR_GROUP_SCHED */
50  #endif /* CONFIG_SMP */
51
52  #ifdef CONFIG_FAIR_GROUP_SCHED
53      struct rq *rq; /* cpu runqueue to which this cfs_rq is attached */ //指向CFS运行队列所属的CPU RQ运行队列
54
55      /*
56       * leaf cfs_rqs are those that hold tasks (lowest schedulable entity in
57       * a hierarchy). Non-leaf lrs hold other higher schedulable entities
58       * (like users, containers etc.)
59       *
60       * leaf_cfs_rq_list ties together list of leaf cfs_rq's in a cpu. This
61       * list is used during load balance.
62       */
63      int on_list;
64      struct list_head leaf_cfs_rq_list;
65      struct task_group *tg; /* group that "owns" this runqueue */ //CFS运行队列所属的任务组
66
67  #ifdef CONFIG_CFS_BANDWIDTH
68      int runtime_enabled; /* CFS运行队列中使用CFS带宽控制 */
69      u64 runtime_expires; /* 到期的运行时间 */
70      s64 runtime_remaining; /* 剩余的运行时间 */
71
72      u64 throttled_clock, throttled_clock_task; /* 限流时间相关 */
73      u64 throttled_clock_task_time;
74      int throttled, throttle_count; /* throttled: 限流, throttle_count: CFS运行队列限流次数 */
75      struct list_head throttled_list; /* 运行队列限流链表节点, 用于添加到cfs_bandwidth结构中的cfttle_cfs_rq链表中 */
76  #endif /* CONFIG_CFS_BANDWIDTH */
77  #endif /* CONFIG_FAIR_GROUP_SCHED */
78  };

```

runtime与vruntime

1. Linux 内核默认的 `sysctl_sched_latency` 是 6ms, 这个值用户态可设。`sched_period` 用于保证可运行任务都能至少运行一次的时间间隔;

2. 当可运行任务大于 8 个的时候，`sched_period` 的计算则需要根据任务个数乘以最小调度颗粒值，这个值系统默认为 0.75ms；
3. 每个任务的运行时间计算，是用 `sched_period` 值，去乘以该任务在整个 CFS 运行队列中的权重占比；
4. 虚拟运行的时间 = 实际运行时间 * `NICE_0_LOAD` / 该任务的权重；

函数调用

```
1  /*
2   * delta /= w
3   */
4  static inline unsigned long
5  calc_delta_fair(unsigned long delta, struct sched_entity *se)
6  {
7      if (unlikely(se->load.weight != NICE_0_LOAD))
8          delta = calc_delta_mine(delta, NICE_0_LOAD, &se->load);
9
10     return delta;
11 }
12
13 /*
14  * The idea is to set a period in which each task runs once.
15  *
16  * When there are too many tasks (sysctl_sched_nr_latency) we have to stretch
17  * this period because otherwise the slices get too small.
18  *
19  * p = (nr <= nl) ? l : l*nr/nl
20  */
21 static u64 __sched_period(unsigned long nr_running)
22 {
23     u64 period = sysctl_sched_latency;
24     unsigned long nr_latency = sched_nr_latency;
25
26     if (unlikely(nr_running > nr_latency)) {
27         period = sysctl_sched_min_granularity;
28         period *= nr_running;
29     }
30
31     return period;
32 }
33
34 /*
35  * We calculate the wall-time slice from the period by taking a part
36  * proportional to the weight.
```

```

37  *
38  *  $s = p \cdot P[w/rw]$ 
39  */
40  static u64 sched_slice(struct cfs_rq *cfs_rq, struct sched_entity *se)
41  {
42      u64 slice = __sched_period(cfs_rq->nr_running + !se->on_rq);
43
44      for_each_sched_entity(se) {
45          struct load_weight *load;
46          struct load_weight lw;
47
48          cfs_rq = cfs_rq_of(se);
49          load = &cfs_rq->load;
50
51          if (unlikely(!se->on_rq)) {
52              lw = cfs_rq->load;
53
54              update_load_add(&lw, se->load.weight);
55              load = &lw;
56          }
57          slice = calc_delta_mine(slice, se->load.weight, load);
58      }
59      return slice;
60  }
61
62  /*
63   * We calculate the vruntime slice of a to be inserted task
64   *
65   *  $vs = s/w$ 
66   */
67  static u64 sched_vslice(struct cfs_rq *cfs_rq, struct sched_entity *se)
68  {
69      return calc_delta_fair(sched_slice(cfs_rq, se), se);
70  }

```


默认6ms, 可以用用户层设置

sysctl_sched_latency

任务数<8

获取sched_latency

__sched_period

nr_running > sched_nr_latency

任务数>=8

nr_running *
sysctl_sched_min_granularity

0.75ms * 任务数

计算runtime

sched_slice

for_each_sched_entity

__calc_delta

数学模型: $\text{delta_exec} * \text{weight} / \text{lw.weight}$

实施方法: $(\text{delta_exec} * (\text{weight} * \text{lw} \rightarrow \text{inv_weight})) \gg \text{WMULT_SHIFT}$

calc_delta_fair(sched_slice(cfs_rq, se), se)

__calc_delta(delta,
NICE_0_LOAD, &se->load)

数学模型: $\text{runtime} * \text{NICE_0_LOAD} / \text{weight}$

计算vruntime

sched_vslice

以Task A为例:

$\text{sched_slice} = 6\text{ms} * (1586 * 801000) \gg 32$

$(\text{delta_exec} * (\text{weight} * \text{lw} \rightarrow \text{inv_weight})) \gg \text{WMULT_SHIFT}$

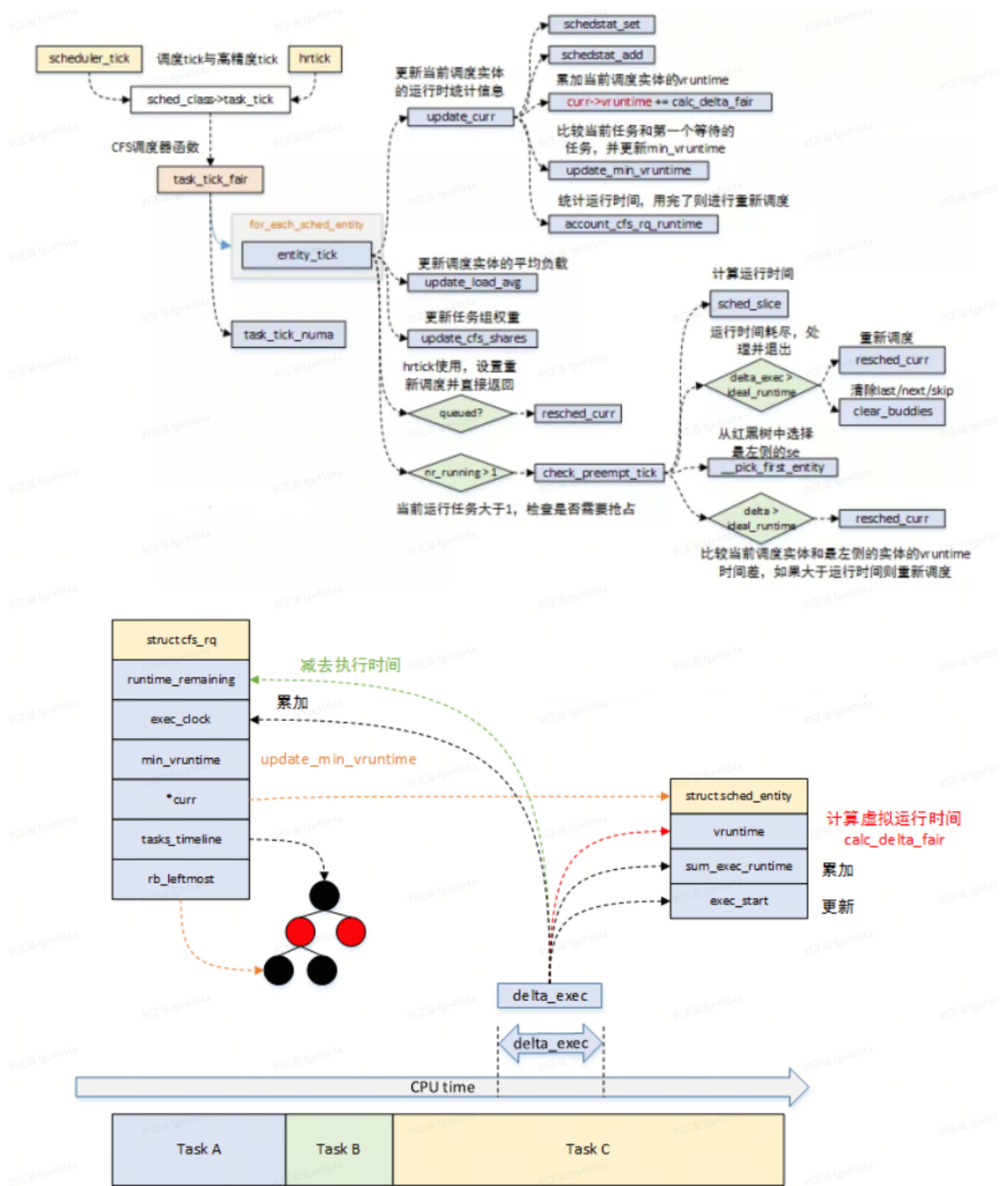
Task A nice = -2 weight=1586 inv_weight=2708050	Task B nice = -1 weight=1277 inv_weight=3363326	Task C nice = 0 weight=1024 inv_weight=4194304	Task D nice = 1 weight=820 inv_weight=523376	Task E nice = 2 weight=655 inv_weight=6557202
--	--	---	---	--

累加各个Task的weight得到 lw.weight=5362
 $\text{lw.inv_weight} = 0\text{xFFFFFFF} / \text{lw.weight} = 801000$

sysctl_sched_latency = 6ms

CFS调度tick

- task_tick_fair->entity->tick->update_curr->_update_curr



```

1  /*
2   * scheduler tick hitting a task of our scheduling class:
3   */
4  static void task_tick_fair(struct rq *rq, struct task_struct *curr, int
    queued)

```

```

5  {
6      struct cfs_rq *cfs_rq;
7      struct sched_entity *se = &curr->se;
8
9      for_each_sched_entity(se) {
10         cfs_rq = cfs_rq_of(se);
11         entity_tick(cfs_rq, se, queued);
12     }
13 }
14
15 static void
16 entity_tick(struct cfs_rq *cfs_rq, struct sched_entity *curr, int queued)
17 {
18     /*
19      * Update run-time statistics of the 'current'.
20      */
21     update_curr(cfs_rq);
22
23 #ifdef CONFIG_SCHED_HRTICK
24     /*
25      * queued ticks are scheduled to match the slice, so don't bother
26      * validating it and just reschedule.
27      */
28     if (queued) {
29         resched_task(rq_of(cfs_rq)->curr);
30         return;
31     }
32     /*
33      * don't let the period tick interfere with the hrtick preemption
34      */
35     if (!sched_feat(DOUBLE_TICK) &&
36         hrtimer_active(&rq_of(cfs_rq)->hrtick_timer))
37         return;
38 #endif
39
40     if (cfs_rq->nr_running > 1 || !sched_feat(WAKEUP_PREEMPT))
41         check_preempt_tick(cfs_rq, curr);
42 }
43
44 /*
45  * Update the current task's runtime statistics. Skip current tasks that
46  * are not in our scheduling class.
47  */
48 static inline void
49 __update_curr(struct cfs_rq *cfs_rq, struct sched_entity *curr,
50              unsigned long delta_exec)
51 {

```

```

52     unsigned long delta_exec_weighted;
53
54     schedstat_set(curr->exec_max, max((u64)delta_exec, curr->exec_max));
55
56     curr->sum_exec_runtime += delta_exec;
57     schedstat_add(cfs_rq, exec_clock, delta_exec);
58     delta_exec_weighted = calc_delta_fair(delta_exec, curr);
59
60     curr->vruntime += delta_exec_weighted;
61     update_min_vruntime(cfs_rq);
62 }
63
64
65 static void update_curr(struct cfs_rq *cfs_rq)
66 {
67     struct sched_entity *curr = cfs_rq->curr;
68     u64 now = rq_of(cfs_rq)->clock;
69     unsigned long delta_exec;
70
71     if (unlikely(!curr))
72         return;
73
74     /*
75      * Get the amount of time the current task was running
76      * since the last time we changed load (this cannot
77      * overflow on 32 bits):
78      */
79     delta_exec = (unsigned long)(now - curr->exec_start);
80     if (!delta_exec)
81         return;
82
83     __update_curr(cfs_rq, curr, delta_exec);
84     curr->exec_start = now;
85
86     if (entity_is_task(curr)) {
87         struct task_struct *curtask = task_of(curr);
88
89         trace_sched_stat_runtime(curtask, delta_exec, curr->vruntime);
90         cpuacct_charge(curtask, delta_exec);
91         account_group_exec_runtime(curtask, delta_exec);
92     }
93 }

```

任务出队入队

1. 当任务进入可运行状态时，需要将调度实体放入到红黑树中，完成入队操作；

2. 当任务退出可运行状态时，需要将调度实体从红黑树中移除，完成出队操作；
3. CFS 调度器，使用 `enqueue_entity` 函数将任务入队到 CFS 队列，使用 `dequeue_entity` 函数将任务从 CFS 队列中出队操作

```
1  static void
2  enqueue_entity(struct cfs_rq *cfs_rq, struct sched_entity *se, int flags)
3  {
4      /*
5       * Update the normalized vruntime before updating min_vruntime
6       * through callig update_curr().
7       */
8      if (!(flags & ENQUEUE_WAKEUP) || (flags & ENQUEUE_MIGRATE))
9          se->vruntime += cfs_rq->min_vruntime;
10
11     /*
12      * Update run-time statistics of the 'current'.
13      */
14     update_curr(cfs_rq);
15     account_entity_enqueue(cfs_rq, se);
16
17     if (flags & ENQUEUE_WAKEUP) {
18         place_entity(cfs_rq, se, 0);
19         enqueue_sleeper(cfs_rq, se);
20     }
21
22     update_stats_enqueue(cfs_rq, se);
23     check_spread(cfs_rq, se);
24     if (se != cfs_rq->curr)
25         __enqueue_entity(cfs_rq, se);
26 }
27
28 /*
29  * Enqueue an entity into the rb-tree:
30  */
31 static void __enqueue_entity(struct cfs_rq *cfs_rq, struct sched_entity *se)
32 {
33     struct rb_node **link = &cfs_rq->tasks_timeline.rb_node;
34     struct rb_node *parent = NULL;
35     struct sched_entity *entry;
36     s64 key = entity_key(cfs_rq, se);
37     int leftmost = 1;
38
39     /*
40      * Find the right place in the rbtree:
41      */
```

```

42 while (*link) {
43     parent = *link;
44     entry = rb_entry(parent, struct sched_entity, run_node);
45     /*
46      * We dont care about collisions. Nodes with
47      * the same key stay together.
48     */
49     if (key < entity_key(cfs_rq, entry)) {
50         link = &parent->rb_left;
51     } else {
52         link = &parent->rb_right;
53         leftmost = 0;
54     }
55 }
56
57 /*
58  * Maintain a cache of leftmost tree entries (it is frequently
59  * used):
60  */
61 if (leftmost)
62     cfs_rq->rb_leftmost = &se->run_node;
63
64 rb_link_node(&se->run_node, parent, link);
65 rb_insert_color(&se->run_node, &cfs_rq->tasks_timeline);
66 }

```

寻找下一个任务

- `__pick_next_entity()`：返回红黑树最左侧节点，该节点是vruntime最小的节点
- 使用`rb_leftmost`缓存最小vruntime节点

```

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

```