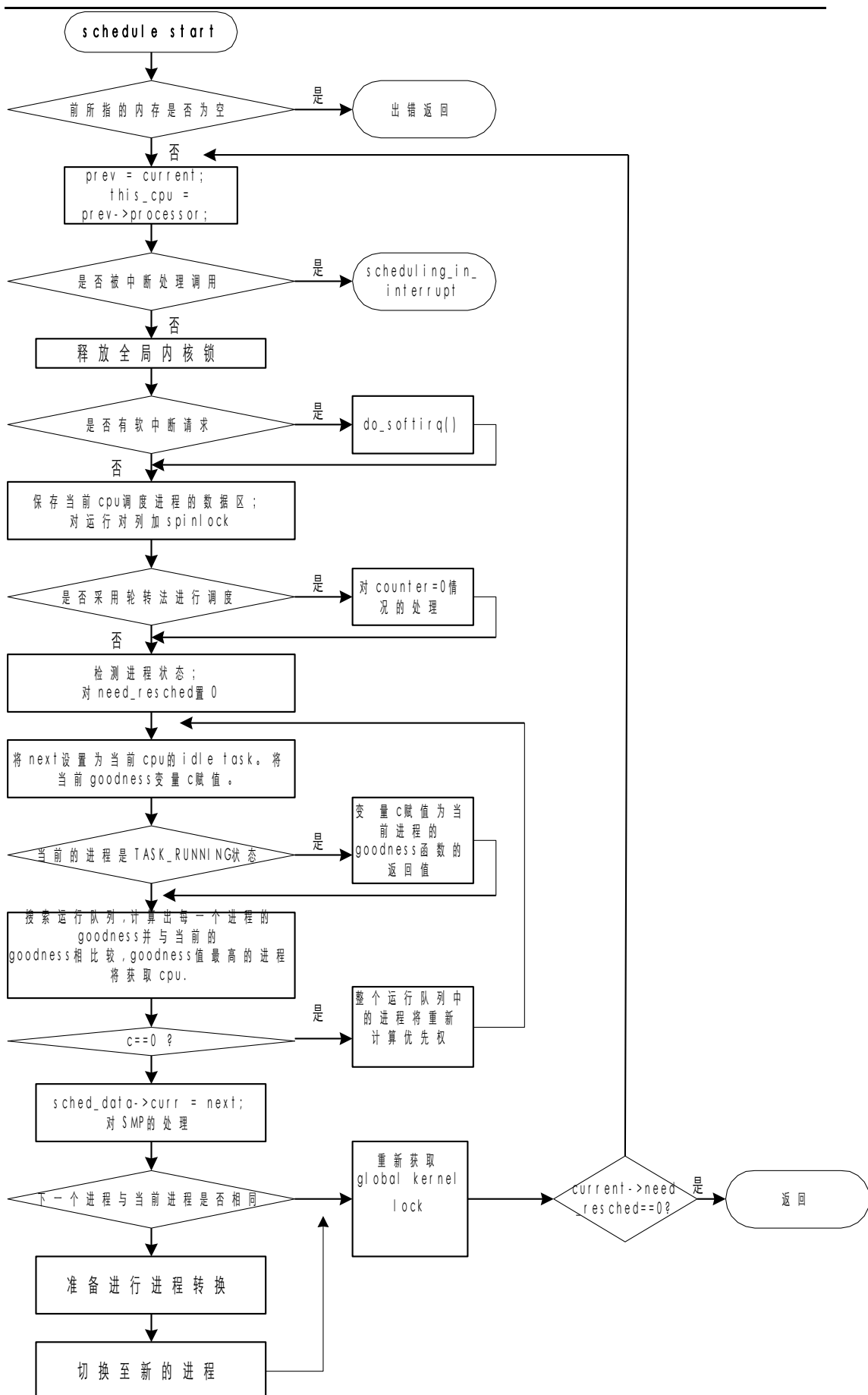


## Linux 2.4 版的 schedule() 函数的流程图



## kernel/sched.c

---

```
498 /*
499 * 'schedule()' is the scheduler function. It's a very simple and nice
500 * scheduler: it's not perfect, but certainly works for most things.
501 *
502 * The goto is "interesting".
503 *
504 * NOTE!! Task 0 is the 'idle' task, which gets called when no other
505 * tasks can run. It can not be killed, and it cannot sleep. The 'state'
506 * information in task[0] is never used.
507 */
508 asmlinkage void schedule(void)
509 {
510     struct schedule_data * sched_data;
511     struct task_struct *prev, *next, *p;
512     struct list_head *tmp;
513     int this_cpu, c;
514
515     if(!current->active_mm) BUG();
    如果 current->active_mm == NULL 那么进程必定有问题.任何进程都有一个确定的 active_mm.

516 need_resched_back:
517     prev = current;
518     this_cpu = prev->processor;
    准备交出 CPU 的进程, 它的 PCB 由 current 指示。将局部变量 prev 指向
    当前进程, this_cpu 指向当前进程拥有的 CPU (考虑对称多 CPU 的架
    构) 。

519
520     if(in_interrupt())
521         goto scheduling_in_interrupt;
    检查一下, 调用 schedule()函数时是否处于中断处理过程。若是, 则打印
    提示信息后, 直接返回。

522
523     release_kernel_lock(prev, this_cpu);
    释放全局内核锁.
```

```
524
525     /* Do "administrative" work here while we don't hold any locks */
526     if (softirq_active(this_cpu) & softirq_mask(this_cpu))
527         goto handle_softirq;
    如果有些工作需要由 softirq 机制完成的，那么立刻执行 do_softirq()
    (kernel/softirq.c)
```

```
528 handle_softirq_back:
```

```
529
530     /*
531      * 'sched_data' is protected by the fact that we can run
532      * only one process per CPU.
533      */
```

```
534     sched_data = & aligned_data[this_cpu].schedule_data;
    用局部变量 sched_data 保存当前 CPU 的调度进程的数据区。因为任一时刻，每个 CPU 只能运行一个 schedule()。schedule_data 结构包含一个指向当前进程 task_struct 结构的指针和最后一个调度进程的 TSC 值(Time Stamp Counter，CPU 附带一个 64 位的时间戳寄存器)。
```

```
535
536     spin_lock_irq(&runqueue_lock);
    对就绪队列加 spinlock。之所以用 spinlock，是因为 schedule()是允许中断的。这样在就绪队列解锁之后，我们可以直接使中断有效，而不必做一些保存标志、恢复标志的操作。
```

```
537
538     /* move an exhausted RR process to be last.. */
539     if (prev->policy == SCHED_RR)
540         goto move_rr_last;
    如果采用轮转法进行调度，则要重新检查一下 counter 是否为 0。若是 0，则将其挂到就绪队列的最后。
```

```
541 move_rr_back:
```

```
542
543     switch (prev->state) {
544         case TASK_INTERRUPTIBLE:
545             if (signal_pending(prev)) {
546                 prev->state = TASK_RUNNING;
547                 break;
548             }
```

```

549         default:
550             del_from_runqueue(prev);
551         case TASK_RUNNING:
552     }

```

检测进程状态：如果是 TASK\_RUNNING 则暂且放在一边不管；如果是 TASK\_INTERRUPTIBLE 状态并且能够唤醒它的信号已经迫近，则将其状态置为 TASK\_RUNNING。至于是别的状态，则将其从就绪队列中删去。

```

553     prev->need_resched = 0;

```

对 need\_resched 置 0，使下个时间片不会有新的 schedule()函数调用。

```

554
555     /*
556     * this is the scheduler proper:
557     */
558
559 repeat_schedule:
560     /*
561     * Default process to select..
562     */
563     next = idle_task(this_cpu);
564     c = -1000;

```

将 next 设置为当前 cpu 的 idle task。将当前 goodness(优先级)变量 c 赋值-1000。 -1000 是 goodness 函数返回值之中最小的权。

```

565     if(prev->state == TASK_RUNNING)
566         goto still_running;

```

如果当前的进程是 TASK\_RUNNING 状态，则将当前 goodness(优先级)变量 c 赋值为当前进程的 goodness 函数的返回值，这将优于 idle task。

```

567
568 still_running_back:
569     list_for_each(tmp, &runqueue_head) {
570         p = list_entry(tmp, struct task_struct, run_list);
571         if (can_schedule(p, this_cpu)) {
572             int weight = goodness(p, this_cpu, prev->active_mm);
573             if (weight > c)
574                 c = weight, next = p;
575         }
576     }

```

搜索就绪队列，计算出每一个进程的 goodness(优先级)并与当前的 goodness(优先级)相比较。goodness(优先级)值最高的进程将获取 CPU。

577

578     /\* Do we need to re-calculate counters? \*/

579     if(!c)

580         goto recalculate;

如果整个就绪队列中最高的 goodness(优先级)值是 0，那么整个运行队列中的进程将重新计算优先权。算法如下：

```
struct task_struct *p;
```

```
spin_unlock_irq(&runqueue_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_irq(&runqueue_lock);
```

581     /\*

582         \* from this point on nothing can prevent us from

583         \* switching to the next task, save this fact in

584         \* sched\_data.

585         \*/

586         sched\_data->curr = next;

从这一时刻起，要转换到下一个任务，将当前信息保存到 sched\_data 中。

587 #ifdef CONFIG\_SMP

588     next->has\_cpu = 1;

589     next->processor = this\_cpu;

590 #endif

next 进程拥有一个 CPU，它就是 this\_cpu

591     spin\_unlock\_irq(&runqueue\_lock);

592

593     if (prev == next)

594         goto same\_process;

处理一种特殊情况，可能选中的 next 进程就是当前进程。那么，我们只需重新获取 global kernel lock 并返回。

595

596 #ifdef CONFIG\_SMP

```

597  /*
598      * maintain the per-process 'last schedule' value.
599      * (this has to be recalculated even if we reschedule to
600      * the same process) Currently this is only used on SMP,
601      * and it's approximate, so we do not have to maintain
602      * it while holding the runqueue spinlock.
603      */
604  sched_data->last_schedule = get_cycles();
605
606  /*
607      * We drop the scheduler lock early (it's a global spinlock),
608      * thus we have to lock the previous process from getting
609      * rescheduled during switch_to().
610      */
611
612 #endif /* CONFIG_SMP */
613
614  kstat.context_switch++;
    累计上下文切换次数递增 1
615  /*
616      * there are 3 processes which are affected by a context switch:
617      *
618      * prev == .... ==> (last ==> next)
619      *
620      * It's the 'much more previous' 'prev' that is on next's stack,
621      * but prev is set to (the just run) 'last' process by switch_to().
622      * This might sound slightly confusing but makes tons of sense.
623      */
624  prepare_to_switch();
625  {
626      struct mm_struct *mm = next->mm;
627      struct mm_struct *oldmm = prev->active_mm;
628      if (!mm) {
629          if (next->active_mm) BUG();
630          next->active_mm = oldmm;
631          atomic_inc(&oldmm->mm_count);
632          enter_lazy_tlb(oldmm, next, this_cpu);
633      } else {

```

```

634         if (next->active_mm != mm) BUG();
635         switch_mm(oldmm, mm, next, this_cpu);
636     }

```

切换内存管理数据结构，从 prev 至 next

```

637
638     if (!prev->mm) {
639         prev->active_mm = NULL;
640         mmdrop(oldmm);
641     }

```

如果 prev 本来就共享了其它进程的 mm，那么，active\_mm 清 0，oldmm 的共享计数递减 1。

```

642     }
643
644     /*
645      * This just switches the register state and the
646      * stack.
647      */
648     switch_to(prev, next, prev);

```

切换寄存器和堆栈。从而，进程切换全部完成！

```

649     __schedule_tail(prev);

```

这次切换后失去 CPU 的 prev，仍然留在就绪队列。它终究会重新获得 CPU。一旦重新获得，其 PC 寄存器指向现在这个位置，也即 \_schedule\_tail(prev)。

```

650

```

```

651 same_process:

```

选中的 next 进程就是当前进程

```

652     reacquire_kernel_lock(current);
653     if (current->need_resched)
654         goto need_resched_back;

```

如果又有新一轮的调度需求，那么转 need\_resched\_back 去调度。

```

655

```

```

656     return;

```

否则，schedule()函数返回了 !!!

```

657

```

```

658 recalculate:

```

重新计算就绪队列中所有进程的优先权。前文已有说明

```
659     {
660         struct task_struct *p;
661         spin_unlock_irq(&runqueue_lock);
662         read_lock(&tasklist_lock);
663         for_each_task(p)
664             p->counter = (p->counter >> 1) + NICE_TO_TICKS(p->nice);
665         read_unlock(&tasklist_lock);
666         spin_lock_irq(&runqueue_lock);
667     }
668     goto repeat_schedule;
```

重新计算就绪队列中所有进程的优先权。前文已有说明

```
669
670 still_running:
671     c = goodness(prev, this_cpu, prev->active_mm);
672     next = prev;
673     goto still_running_back;
```

如果 prev 本身也是就绪进程，那么，先计算 prev 的优先权，先把 prev 设置成 next。后续会逐个计算机就绪队列里面之进程的优先权，如果有更高优先权的进程，会重置 c 和 next。

```
674
675 handle_softirq:
676     do_softirq();
677     goto handle_softirq_back;
678
679 move_rr_last:
680     if (!prev->counter) {
681         prev->counter = NICE_TO_TICKS(prev->nice);
682         move_last_runqueue(prev);
683     }
684     goto move_rr_back;
```

更新 prev 的优先权，并且把它放置就绪队列的末尾。

```
685
686 scheduling_in_interrupt:
687     printk("Scheduling in interrupt\n");
688     BUG();
689     return;
690 }
```

---





## Linux 2.4 版的 goodness() 函数

进程的 goodness 值通过 goodness() 函数计算。goodness 返回下面两类中的一个值：

- 1000 和 1000 以上的值只能赋给实时进程，
- 从 0-999 的值只能赋给普通进程。

有关这两类 goodness 结果的重要的一点是：该值在实时系统的范围肯定会比非实时系统的范围要高。POSIX 标准规定内核要确保在实时进程和非实时进程同时竞争 CPU 时，实时进程要优先于非实时进程。由于调度进程总是选择具有最大 goodness 值的进程，又由于任何尚未释放 CPU 的实时进程的 goodness 值总是比非实时进程的 goodness 大，Linux 是遵守 POSIX 标准的。

goodness() 函数从不会返回 -1000 的，也不会返回其它的负值。

而 idle 进程的 counter 值为负。

kernel/sched.c

---

```
123 /*
124  * This is the function that decides how desirable a process is..
125  * You can weigh different processes against each other depending
126  * on what CPU they've run on lately etc to try to handle cache
127  * and TLB miss penalties.
128  *
129  * Return values:
130  *   -1000: never select this
131  *     0: out of time, recalculate counters (but it might still be
132  *        selected)
133  *   +ve: "goodness" value (the larger, the better)
134  *   +1000: realtime process, select this.
135  */
136
137 static inline int goodness(struct task_struct * p, int this_cpu, struct mm_struct
*this_mm)
138 {
139     int weight;
140
141     /*
142      * select the current process after every other
143      * runnable process, but before the idle thread.
```

```

144     * Also, dont trigger a counter recalculation.
145     */
146     weight = -1;
147     if (p->policy & SCHED_YIELD)
148         goto out;
    如果要放弃这次调度，那么直接返回。
149
150     /*
151     * Non-RT process - normal case first.
152     */
153     if (p->policy == SCHED_OTHER) {
    如果是非实时进程调度，执行下列控制流

154         /*
155         * Give the process a first-approximation goodness value
156         * according to the number of clock-ticks it has left.
157         *
158         * Don't do any other calculations if the time slice is
159         * over..
160         */
161         weight = p->counter;
    时间片剩余值作为优先权。

162         if (!weight)
163             goto out;
    若时间片用完，则返回 0。

164
165 #ifdef CONFIG_SMP
166     /* Give a largish advantage to the same processor... */
167     /* (this is equivalent to penalizing other processors) */
168     if (p->processor == this_cpu)
169         weight += PROC_CHANGE_PENALTY;
170 #endif
171
172     /* .. and a slight advantage to the current MM */
173     if (p->mm == this_mm || !p->mm)
174         weight += 1;

```

如果进程的存储区与当前的一致，或者进程是内核线程，则优先权增1。

```
175         weight += 20 - p->nice;
```

最后，根据 p 进程的 nice 值适当地调整优先权。

```
176         goto out;
```

```
177     }
```

```
178
```

```
179     /*
```

```
180     * Realtime process, select the first one on the
```

```
181     * runqueue (taking priorities within processes
```

```
182     * into account).
```

```
183     */
```

```
184     weight = 1000 + p->rt_priority;
```

如果是实时进程调度，则取得实时进程优先权，而不是 counter。加上 1000，使其优先权高于任何普通进程。

```
185 out:
```

```
186     return weight;
```

```
187 }
```

---