

《操作系统综合实验》

实验报告

06 Linux操作系统实例研究报告

上课时间：	2022年5月28日
授课教师：	冯鹏斌
姓 名：	王栗政
班 级：	2018039
学 号：	20179100018
日 期：	2022年6月2日

1 实验题目

阅读教材第18章（Linux案例），并在互联网上查阅相关资料，对照操作系统课程中所讲的原理（进程管理，存储管理，文件系统，设备管理），了解Linux操作系统实例形成一份专题报告

- 可以是全面综述性报告
- 可以是侧重某一方面的报告（进程调度，进程间通信，存储管理，文件系统，安全）

2 相关原理与知识

2.1 进程基础知识

2.1.1 进程的定义

进程是资源分配的最小单位，线程是CPU调度的最小单位

进程^[1]（英语：process），是指计算机中已执行的程序。进程曾经是分时系统的基本运作单位。在面向进程设计的系统（如早期的UNIX，Linux 2.4及更早的版本）中，进程是程序的基本执行实体；在面向线程设计的系统（如当代多数操作系统、Linux 2.6及更新的版本）中，进程本身不是基本执行单位，而是线程的容器。

程序本身只是指令、数据及其组织形式的描述，相当于一个名词，进程才是程序（那些指令和数据）的真正执行实例，可以想像说是现在进行式。若干进程有可能与同一个程序相关系，且每个进程皆可以同步（循序）或异步（平行）的方式独立执行。现代计算机系统可在同一段时间内以进程的形式将多个程序加载到内存中，并借由时间共享（或称分时复用），以在一个处理器上表现出同时（平行性）执行的感觉。同样的，使用多线程技术的操作系统或计算机体系结构，同样程序的**平行线程**，可在多CPU主机或网络上真正**同时**执行（在不同的CPU上）。

进程不仅包括可执行程序的代码段，还包括一系列的**资源**，比如：打开的文件、内存、CPU时间、信号量、多个执行线程流等等。而线程可以共享进程内的资源空间。

2.1.2 进程的属性

Linux内核通过**进程描述符** `task_struct` 结构体来管理进程，该结构体包含了一个进程所需的所有信息，被定义在 `include/linux/sched.h` 文件中，接下来我将对 `task_struct` 结构体中与调度有关的内容进行简要的介绍^[2]。

```

1 struct task_struct {
2     /* ... */
3
4     /* 进程状态 */
5     volatile long    state;
6
7     /* 调度优先级相关，策略相关 */
8     int              prio;
9     int              static_prio;
10    int              normal_prio;

```

```

11     unsigned int      rt_priority;
12     unsigned int      policy;
13
14     /* 调度类, 调度实体相关, 任务组相关等 */
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     /* 进程之间的关系相关 */
24     /* Real parent process: */
25     struct task_struct __rcu *real_parent;
26
27     /* Recipient of SIGCHLD, wait4() reports: */
28     struct task_struct __rcu *parent;
29
30     /*
31      * Children/sibling form the list of natural children:
32      */
33     struct list_head children;
34     struct list_head sibling;
35     struct task_struct *group_leader;
36
37     /* ... */
38 }

```

2.1.2.1 进程标识符 (PID)

PID: Process Identifier, also known as process ID or PID, is a unique number to identify each process running in an operating system such as Linux, Windows, and Unix.

```

1  pid_t pid;
2  pid_t tgid;

```

Unix系统通过 `pid` 来标识进程, `linux` 把不同的 `pid` 与系统中每个进程或轻量级线程关联, 而 `unix` 程序员希望同一组线程具有共同的 `pid`, 遵照这个标准 `linux` 引入线程组的概念。一个线程组所有线程与领头线程具有相同的 `pid`, 存入 `tgid` 字段, `getpid()` 返回当前进程的 `tgid` 值而不是 `pid` 的值。

在 `CONFIG_BASE_SMALL` 配置为0的情况下, PID的取值范围是0到32767, 即系统中的进程数最大为32768个。

```

1  #define PID_MAX_DEFAULT (CONFIG_BASE_SMALL ? 0x1000 : 0x8000)

```

2.1.2.2 进程状态

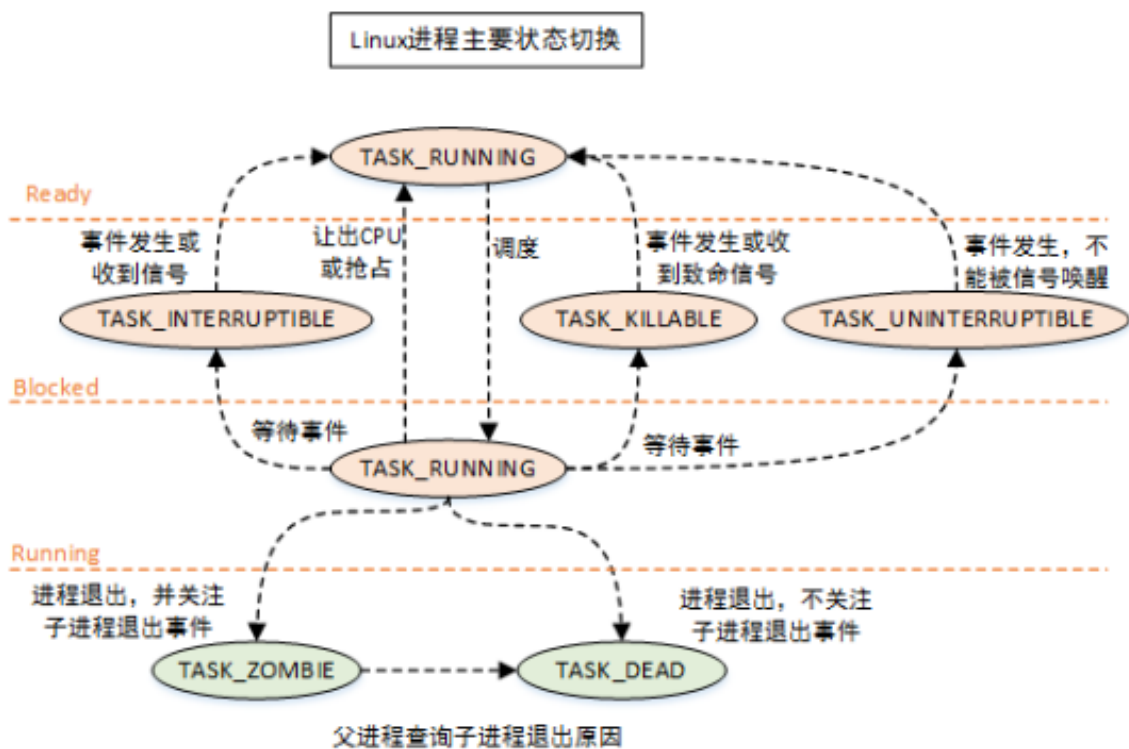
```
1 | volatile long state;    /* -1 unrunnable, 0 runnable, >0 stopped */
```

内核中主要定义的状态字段如下^[3]:

```
1 | /*
2 |  * Task state bitmask. NOTE! These bits are also
3 |  * encoded in fs/proc/array.c: get_task_state().
4 |  *
5 |  * We have two separate sets of flags: task->state
6 |  * is about runnability, while task->exit_state are
7 |  * about the task exiting. Confusing, but this way
8 |  * modifying one set can't modify the other one by
9 |  * mistake.
10 | */
11 | #define TASK_RUNNING          0
12 | #define TASK_INTERRUPTIBLE    1
13 | #define TASK_UNINTERRUPTIBLE  2
14 | #define __TASK_STOPPED        4
15 | #define __TASK_TRACED         8
16 |
17 | /* in tsk->exit_state */
18 | #define EXIT_DEAD             16
19 | #define EXIT_ZOMBIE           32
20 | #define EXIT_TRACE             (EXIT_ZOMBIE | EXIT_DEAD)
21 |
22 | /* in tsk->state again */
23 | #define TASK_DEAD             64
24 | #define TASK_WAKEKILL         128    /** wake on signals that are deadly
25 | **/
26 | #define TASK_WAKING           256
27 | #define TASK_PARKED           512
28 | #define TASK_NOLOAD           1024
29 | #define TASK_STATE_MAX        2048
30 |
31 | /* Convenience macros for the sake of set_task_state */
32 | #define TASK_KILLABLE          (TASK_WAKEKILL | TASK_UNINTERRUPTIBLE)
33 | #define TASK_STOPPED           (TASK_WAKEKILL | __TASK_STOPPED)
34 | #define TASK_TRACED            (TASK_WAKEKILL | __TASK_TRACED)
```

2.1.2.2.1 五个互斥状态

state域能够取以下5个互为排斥的值。系统中的每个进程都必然处于以上所列进程状态中的一种。



状态	描述
TASK_RUNNING	表示进程要么正在执行，要么正要准备执行（已经就绪），正在等待cpu时间片的调度
TASK_INTERRUPTIBLE	进程因为等待一些条件而被挂起（阻塞）而所处的状态。这些条件主要包括：硬中断、资源、一些信号……，一旦等待的条件成立，进程就会从该状态（阻塞）迅速转化成为就绪状态TASK_RUNNING
TASK_UNINTERRUPTIBLE	意义与TASK_INTERRUPTIBLE类似，除了不能通过接受一个信号来唤醒以外，对于处于TASK_UNINTERRUPTIBLE状态的进程，哪怕我们传递一个信号或者有一个外部中断都不能唤醒他们。只有它所等待的资源可用的时候，他才会被唤醒。这个标志很少用，但是并不代表没有任何用处，其实他的作用非常大，特别是对于驱动刺探相关的硬件过程很重要，这个刺探过程不能被一些其他的東西给中断，否则就会让进程进入不可预测的状态
TASK_STOPPED	进程被停止执行，当进程接收到SIGSTOP、SIGTTIN、SIGTSTP或者SIGTTOU信号之后就会进入该状态
TASK_TRACED	表示进程被debugger等进程监视，进程执行被调试程序所停止，当一个进程被另外的进程所监视，每一个信号都会让进程进入该状态

2.1.2.2.2 终止状态

还有两个附加的进程状态既可以被添加到state域中，又可以被添加到exit_state域中。只有当进程终止的时候，才会达到这两种状态。

```
1  /* task state */
2  int exit_state;
3  int exit_code, exit_signal;
```

状态	描述
EXIT_ZOMBIE	进程的执行被终止，但是其父进程还没有使用wait()等系统调用来获知它的终止信息，此时进程成为僵尸进程
EXIT_DEAD	进程的最终状态

2.1.2.2.3 睡眠状态

进程状态 TASK_UNINTERRUPTIBLE 和 TASK_INTERRUPTIBLE 都是睡眠状态，Linux 内核提供了两种方法将进程置为睡眠状态。

将进程置为睡眠状态的普通方法是 将进程状态设置为 TASK_INTERRUPTIBLE 或 TASK_UNINTERRUPTIBLE 并调用调度程序的 schedule() 函数。这样会将进程从 CPU 运行队列中移除。

- 如果进程处于可中断模式的睡眠状态（通过将其状态设置为 TASK_INTERRUPTIBLE），那么可以通过显式的唤醒呼叫（wakeup_process()）或需要处理的信号来唤醒它。
- 但是，如果进程处于非可中断模式的睡眠状态（通过将其状态设置为 TASK_UNINTERRUPTIBLE），那么只能通过显式的唤醒呼叫将其唤醒。除非万不得已，否则我们建议您将进程置为可中断睡眠模式，而不是不可中断睡眠模式（比如说在设备 I/O 期间，处理信号非常困难时）。

当处于可中断睡眠模式的任务接收到信号时，它需要处理该信号（除非它已被屏蔽），离开之前正在处理的任务（此处需要清除代码），并将 -EINTR 返回给用户空间。再一次，检查这些返回代码和采取适当操作的工作将由程序员完成。

因此，懒惰的程序员可能比较喜欢将进程置为不可中断模式的睡眠状态，因为信号不会唤醒这类任务。

但需要注意的一种情况是，对不可中断睡眠模式的进程的唤醒呼叫可能会由于某些原因不会发生，这会使进程无法被终止，从而最终引发问题，因为惟一的解决方法就是重启系统。一方面，您需要考虑一些细节，因为不这样做会在内核端和用户端引入 bug。另一方面，您可能会生成永远不会停止的进程（被阻塞且无法终止的进程）。

Linux Kernel 2.6.25 引入了一种新的进程睡眠状态。

状态	描述
TASK_KILLABLE	当进程处于这种可以终止的新睡眠状态中，它的运行原理类似于 TASK_UNINTERRUPTIBLE，只不过可以响应致命信号

它定义如下：

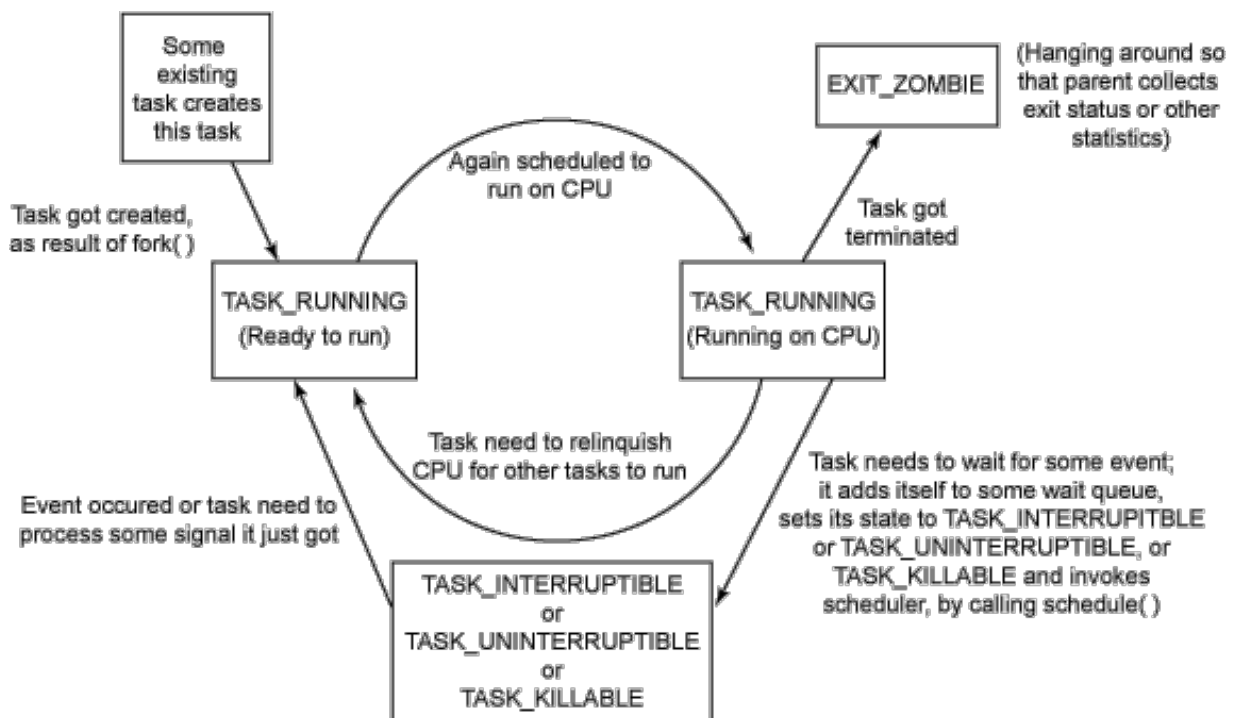
```

1 | #define TASK_WAKEKILL                128 /** wake on signals that are deadly **/
2 |
3 | /* Convenience macros for the sake of set_task_state */
4 | #define TASK_KILLABLE                (TASK_WAKEKILL | TASK_UNINTERRUPTIBLE)
5 | #define TASK_STOPPED                (TASK_WAKEKILL | __TASK_STOPPED)
6 | #define TASK_TRACED                 (TASK_WAKEKILL | __TASK_TRACED)

```

换句话说， $TASK_UNINTERRUPTIBLE + TASK_WAKEKILL = TASK_KILLABLE$ ，而 $TASK_WAKEKILL$ 用于在接收到致命信号时唤醒进程，新的睡眠状态允许 $TASK_UNINTERRUPTIBLE$ 响应致命信号。

进程状态的切换过程和原因大致如下图：



2.1.2.3 调度优先级

```

1 | int prio, static_prio, normal_prio;
2 | unsigned int rt_priority;

```

字段	描述
prio	用于保存动态优先级
static_prio	用于保存静态优先级，可以通过nice系统调用来进行修改
rt_priority	用于保存实时优先级
normal_prio	值取决于静态优先级和调度策略

实时优先级范围是0到 MAX_RT_PRIO-1 （即99），而普通进程的静态优先级范围是从 MAX_RT_PRIO 到 MAX_PRIO-1 （即100到139），值越大静态优先级越低，接下来将会对这些内容进行详细的介绍。

2.1.2.4 调度策略相关字段

```

1 unsigned int policy;
2
3 const struct sched_class *sched_class;
4 struct sched_entity se;
5 struct sched_rt_entity rt;
6
7 cpumask_t cpus_allowed;

```

字段	描述
policy	调度策略
sched_class	调度类
se	普通进程的调用实体，每个进程都有其中之一的实体
rt	实时进程的调用实体，每个进程都有其中之一的实体
cpus_allowed	用于控制进程可以在哪里处理器上运行

2.1.3 进程的分类

- I/O密集型(IO-bound)

在进程占用CPU期间频繁有IO操作，出现IO阻塞等待情况，比如负责监听socket的进程，真正使用CPU进行计算的时间并不多。

- CPU密集型(CPU-bound)

在进程占用CPU期间基本都在进行计算，很少进行IO操作，期间对CPU的真实使用率很高。

2.2 调度相关基础知识

2.2.1 scheduler 调度器

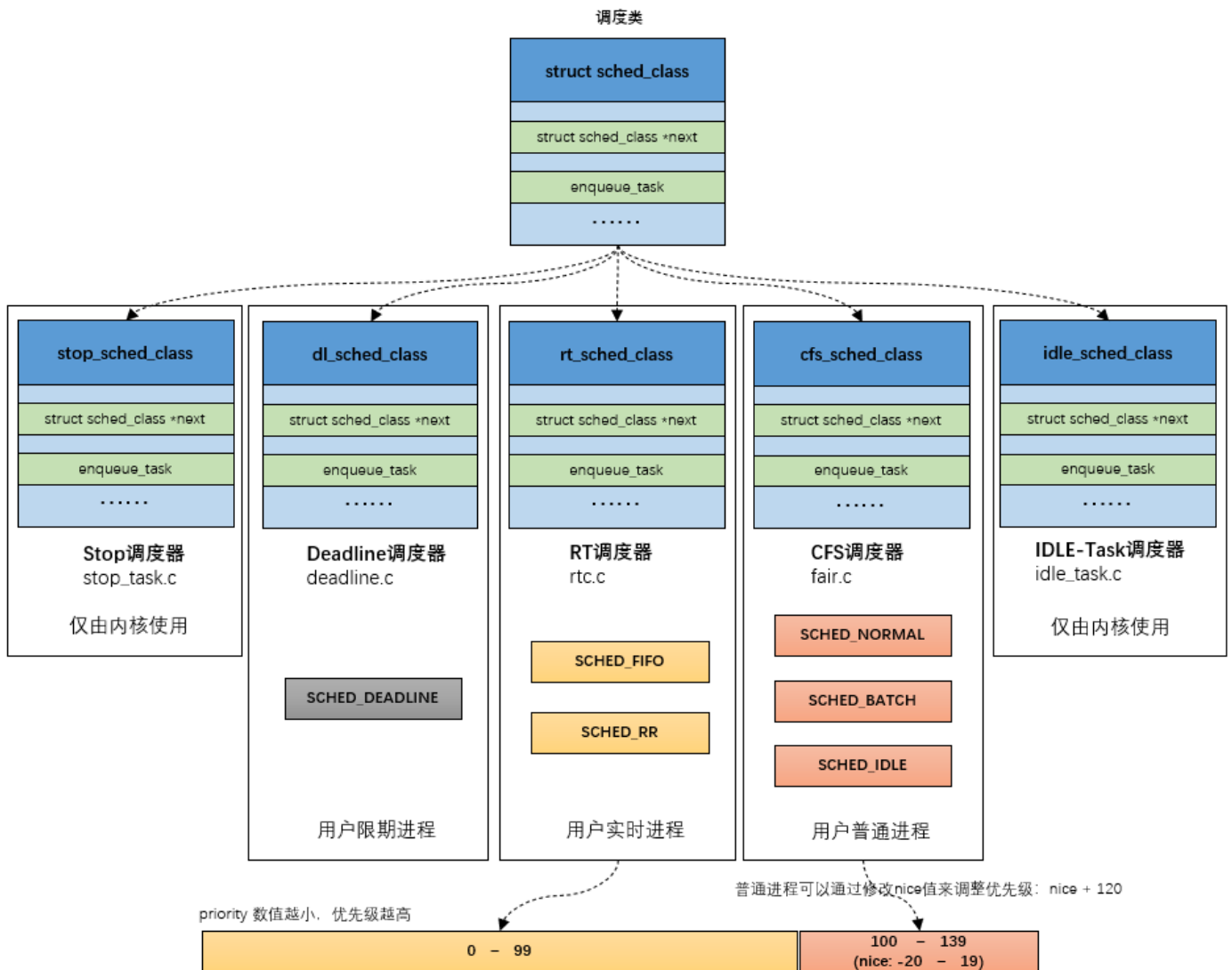
调度，就是按照某种调度的算法，从进程的就绪队列中选取进程分配CPU，主要是协调对CPU等的资源使用。进程调度的目标是最大限度利用CPU时间。

内核默认提供5类调度器，Linux 内核使用 `struct sched_class` 对调度器进行抽象，一共分为以下5类：

- **Stop调度器**, `stop_sched_class`：优先级最高的调度类，可以抢占其他所有进程，不能被其他进程抢占；
- **Deadline调度器**, `dl_sched_class`：使用红黑树，把进程按照绝对截止时间进行排序，选择最小进程进行调度运行；
- **RT调度器**, `rt_sched_class`：实时调度器，为每个优先级维护一个队列；
- **CFS调度器**, `cfs_sched_class`：完全公平调度器，采用完全公平调度算法，引入虚拟运行时间概念；
- **IDLE-Task调度器**, `idle_sched_class`：空闲调度器，每个CPU都会有一个idle线程，当没有其他进程可以调度时，调度运行idle线程；

Linux内核提供了一些调度策略供用户程序来选择调度器，其中 **Stop调度器** 和 **IDLE-Task调度器**，仅由内核使用，用户无法进行选择：

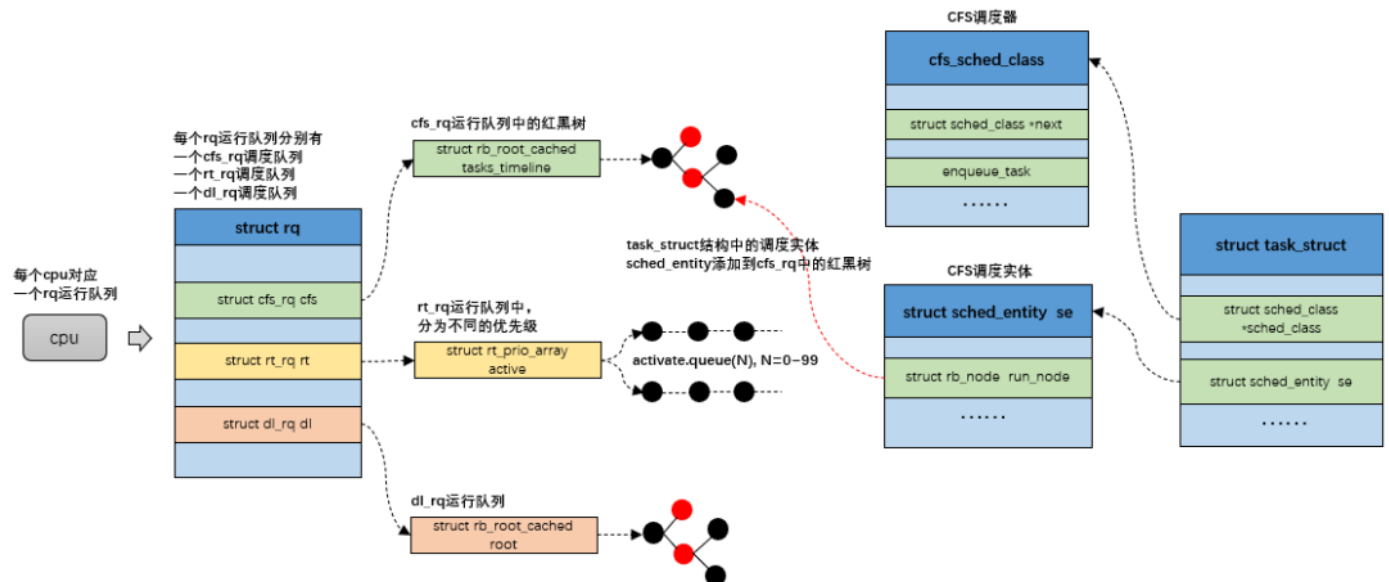
- `SCHED_DEADLINE`：限期进程调度策略，使task选择 `Deadline调度器` 来调度运行；
- `SCHED_RR`：实时进程调度策略，时间片轮转，进程用完时间片后加入优先级对应运行队列的尾部，把CPU让给同优先级的其他进程；
- `SCHED_FIFO`：实时进程调度策略，先进先出调度没有时间片，没有更高优先级的情况下，只能等待主动让出CPU；
- `SCHED_NORMAL`：普通进程调度策略，使task选择 `CFS调度器` 来调度运行；
- `SCHED_BATCH`：普通进程调度策略，批量处理，使task选择 `CFS调度器` 来调度运行；
- `SCHED_IDLE`：普通进程调度策略，使task以最低优先级选择 `CFS调度器` 来调度运行；



2.2.2 runqueue 运行队列

runqueue 运行队列是本 CPU 上所有可运行进程的队列集合。每个 CPU 都有一个运行队列，每个运行队列中有三个调度队列，task 作为调度实体加入到各自的调度队列中。

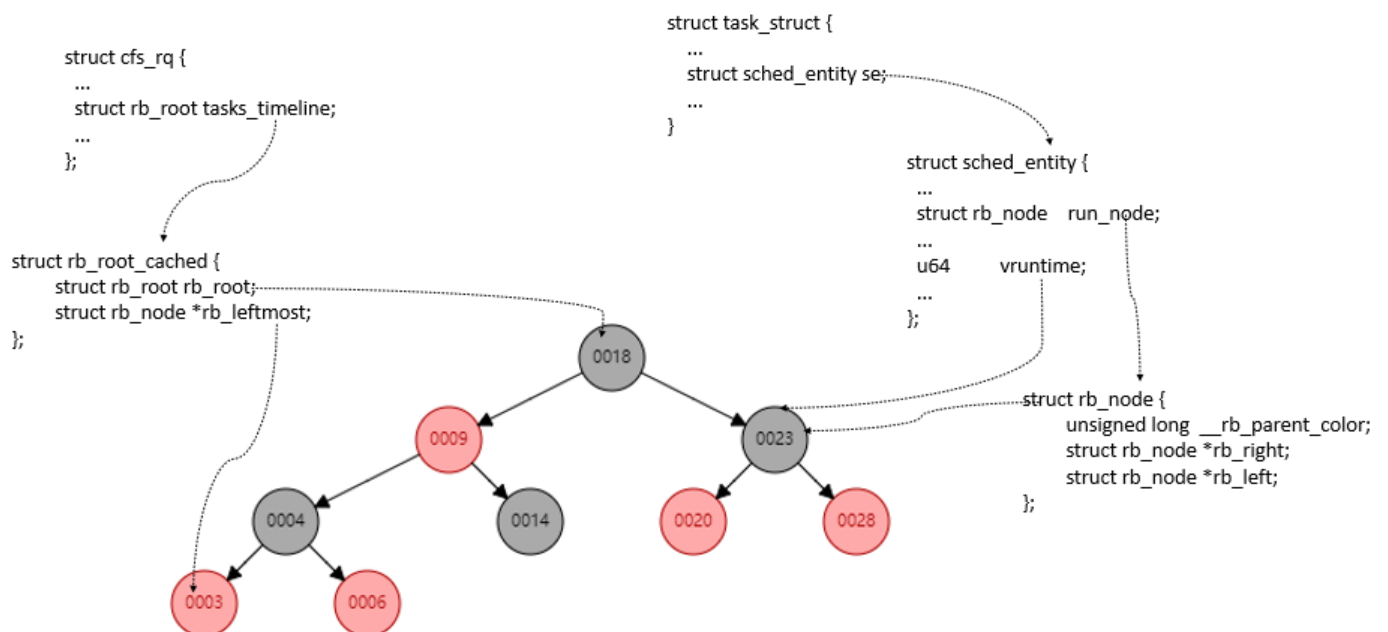
- `struct cfs_rq cfs`: CFS调度队列
- `struct rt_rq rt`: RT调度队列
- `struct dl_rq dl`: DL调度队列



cfs_rq: 跟踪就绪队列信息以及管理就绪态调度实体，并维护一棵按照虚拟时间排序的红黑树。tasks_timeline->rb_root是红黑树的根，tasks_timeline->rb_leftmost指向红黑树中最左边的调度实体，即虚拟时间最小的调度实体。

sched_entity: 可被内核调度的实体。每个就绪态的调度实体sched_entity包含插入红黑树中使用的节点rb_node，同时vruntime成员记录已经运行的虚拟时间。

这些数据结构的关系如下图所示：



2.2.3 调度算法

2.2.3.1 先来先服务FCFS

先来先服务First Come First Service可以说是最早最简单的调度算法，哪个进程先来就先让它占用CPU，释放CPU之后第二个进程再使用，依次类推。

- 场景一

假如有ABC三个进程依次进入等候使用CPU资源的队列FIFO，A进程占用CPU运行5ms，B进程10ms，C进程25ms，根据FCFS算法的规则，可知：

- - 周转时间 A-5ms B-15ms C-40ms 平均 $(5+15+40)/3=20\text{ms}$
 - 响应时间 A-0ms B-5ms C-15ms 平均 $(0+5+15)/3=6.67\text{ms}$

- 场景二

假如有ABC三个进程依次进入等候使用CPU资源的队列FIFO，A进程占用CPU运行100ms，B进程10ms，C进程25ms，根据FCFS算法的规则，可知：

- - 周转时间 A-100ms B-110ms C-135ms 平均 $(100+110+135)/3=115\text{ms}$
 - 响应时间 A-0ms B-100ms C-110ms 平均 $(0+100+110)/3=70\text{ms}$ 综上，在场景二中A进程的运行时间长达100ms，这样拉升了B和C的周转时间5倍多。

在FCFS中优先被调度的进程如果耗时很长，后续进程都必须等待这个大CPU消耗的进程，最终导致周转时间直线拉升，也就是护航效应。

2.2.3.2 最短任务优先SJF

最短任务优先Shortest Job First的思想是当多个进程同时出现时，优先安排执行时间最短的任务，然后是执行时间次短，以此类推。

SJF可以解决FCFS中同时到达进程执行时间不一致带来的护航效应问题：

- 场景一

假如有ABC三个进程同时进入等候使用CPU资源的队列FIFO，A进程占用CPU运行100ms，B进程10ms，C进程25ms，根据SJF算法的规则，可知：

- 进程执行顺序是B->C->A
- 周转时间 A-135ms B-10ms C-35ms 平均 $(135+10+35)/3=60\text{ms}$
- 响应时间 A-35ms B-0ms C-10ms 平均 $(35+0+10)/3=15\text{ms}$

相比于FCFS可能的执行顺序是A->C->B来说，周转时间和响应时间都得到很大的改善。

SJF的算法思想有些理想化，但并非一无是处，升级改进下也有用武之地：

- 大部分场景下进程都是有先后顺序进行等待队列的，同时出现的概率并不高
- 进程执行时间的长短并不能预测

2.2.3.3 抢占式最短任务优先PSJF

SJF算法最具启发的地方在于优先调度执行时间短的任务来解决护航效应，但是该算法属于非抢占式调度，对于先后到达且执行时间差别较大的场景也束手无策。

当向SJF算法增加抢占调度时能力就大大增强了，这就是PSJF(Preemptive Shortest Job First)算法。

- 场景一

假如有ABC三个进程间隔20ms(BC同时且晚于A)依次进入等候使用CPU资源的队列FIFO，A进程占用CPU运行100ms，B进程10ms，C进程25ms，根据PSJF算法的规则，可知：

- A先被执行20ms，再执行B10ms，再执行C25ms，再执行A80ms
- 周转时间 A-135ms B-10ms C-35ms 平均 $(135+10+35)/3=60\text{ms}$

- 响应时间 A-0ms B-0ms C-10ms 平均 $(0+0+10)/3=3.3\text{ms}$

抢占机制有效削弱了护航效应，周转时间和响应时间都降低了许多，但是还远不够完美。

PSJF算法对于进程A来说却不友好，进程A在被抢占之后只能等待B和C运行完成后，此时如果再来短任务DEF都会抢占A，就会产生饥饿效应。

PSJF算法是基于对任务运行时间长短来进行调度的，但在调度前任务的运行时间是未知的，因此首要问题是通过历史数据来预测运行时间。

2.2.3.4 时间片轮转RR

时间片轮转RR(Round-Robin)将CPU执行时间按照时钟脉冲进行切割称为时间切片Time-Slice，一个时间切片是时钟周期的数倍，时钟周期和CPU的主频呈倒数关系。

比如CPU的主频是1000hz，则时钟周期TimeClick=1ms，Time Slice = $n \times \text{Time Click}$ ，时间切片可以是2ms/4ms/10ms等。

在一个时间片内CPU分配给一个进程，时间片耗尽则调度选择下一个进程，如此循环。

进程往往需要多个循环获取多次时间片才能完成任务，因此需要操作系统记录上一次的数据，等进程下一次被分配时间片时再拿出来，这就是传说中的上下文Context。

进程上下文被存储和拿出的过程被称为上下文切换Context Switch，上下文切换是比较消耗资源的，因此时间片的划分就显得很重要：

- 时间片太短，上下文频繁切换，有效执行时间变少
- 时间片太大，无法保证多个进程可以雨露均沾，响应时间较大

RR算法在保障了多个进程都能占用CPU，属于公平策略的一种，但是RR算法将原本可以一次运行完的任务切分成多个部分来完成，从而拉升了周转时间。

2.2.3.5 多级反馈队列MLFQ

多级反馈队列(Multi-Level Feedback Queue)尝试去同时解决响应时间和周转时间两个问题，具体做法：

- 设置了多个任务队列，每个队列对应的优先级不同，队列内部的优先级相同
- 优先分配CPU给高优先级的任务，同优先级队列中的任务采用RR轮询机制
- 通过对任务运行状态的追踪来调整优先级，也就是所谓的Feedback反馈机制
 - 任务在运行期间有较多IO请求和等待，预测为交互进程，优先级保持或提升
 - 任务在运行期间一直进行CPU计算，预测为非交互进程，优先级保持或下降
- 最初将所以任务都设置为高优先级队列，随着后续的运行状态再进行调整
- - 运行期间有IO操作则保持优先级
 - 运行期间无IO操作则把任务放到低一级的队列中
- 不同队列分配不同的时间片
 - 高优先级队列往往是IO型任务，配置较小的时间片来保障响应时间
 - 低优先级队列往往是CPU型任务，配置较长时间片来保障任务一直运行

在实际应用中仍然会有一些问题：

- 饥饿问题
 - - CPU密集型的任务随着时间推移优先级会越来越低，在IO型进程多的场景下很容易出现饥饿问题，一直无法得到调度
 - - 任务是CPU密集型还是IO密集型可能是动态变化的，低优先级队列中的IO型任务的响应时间被拉升，调度频率下降
- 作弊问题
 - - 基于MLFQ对IO型任务的偏爱，用户可能为CPU密集型任务编写无用的IO空操作，从而人为提升任务优先级，相当于作弊

针对上述问题MLFQ还需增加几个补丁：

- 周期性提升所有任务的优先级到最高，避免饥饿问题
- 调度程序记录任务在某个层级队列中消耗的时间片，如果达到某个比例，无论期间是否有IO操作都降低任务的优先级，通过计时来确定真正的IO型任务

MLFQ的算法思想在1962年被提出，其作者也获得了图灵奖，可谓是影响深远。

在朴素MLFQ算法基础上出现一些变种，通过工程实现和经验配置最终被使用到操作系统中，成为真正的工业级进程调度器。

3 实验过程

对Linux的进程调度器进行研究：

Linux的进程调度器是不断演进的，先后出现过三种里程碑式的调度器：

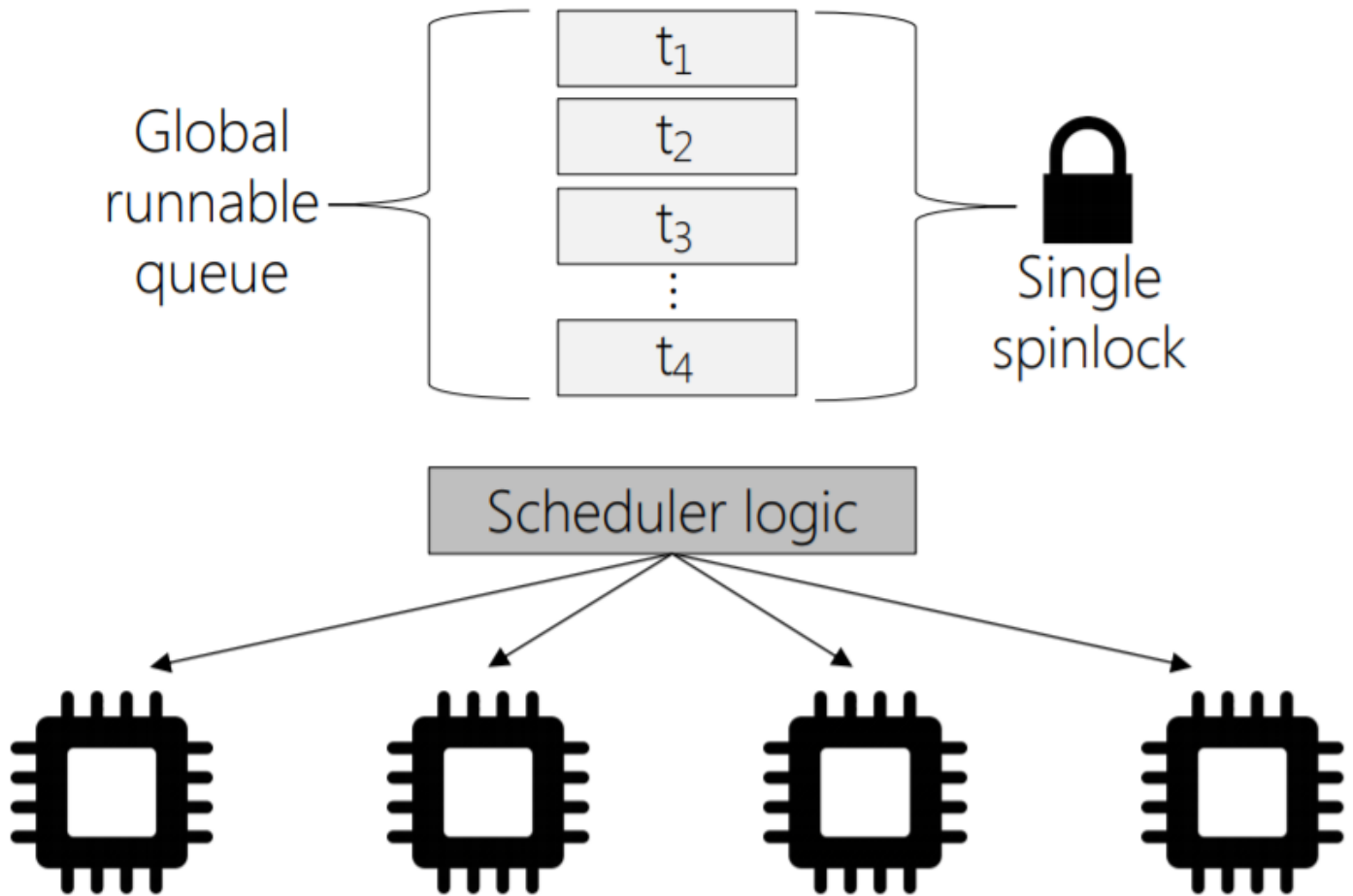
调度器种类	内核版本
O(n)调度器	linux0.11 - 2.4
O(1)调度器	linux2.6
CFS调度器	linux2.6至今

O(n)属于早期版本在pick next过程是需要遍历进程任务队列来实现，O(1)版本性能上有较大提升可以实现O(1)复杂度的pick next过程。

CFS调度器可以说是一种O(logn)调度器，但是其算法思想相比前两种有些不同，并且设计实现上也更加轻量，一直被Linux内核沿用至今。

3.1 O(n)调度器

Linux O(n) Scheduler



O(n)调度器采用一个全局队列runqueue作为核心数据结构，具备以下特点：

- 多个cpu共享全局队列，并非每个cpu有单独的队列
- 实时进程和普通进程混合且无序存放，寻找最合适进程需要遍历
- 就绪进程将被添加到队列，运行结束被删除
- 全局队列增删进程任务时需要加锁
- 进程被挂到不同CPU运行，缓存利用率低

在Linux中进程使用task_struct结构体来表示，其中有个counter表示进程剩余的CPU时间片：

```

1 struct task_struct {
2     long counter;
3     long nice;
4     unsigned long policy;
5     int processor;
6     unsigned long cpus_runnable, cpus_allowed;
7 }

```

counter值在调度周期开始时被赋值，随着调度的进行递减，直至counter=0表示无可用CPU时间，等待下一个调度周期。

O(n)调度器中调度权重是在goodness函数中完成计算的：

```

1  static inline int goodness(struct task_struct * p, int this_cpu, struct
mm_struct *this_mm)
2  {
3      int weight;
4      weight = -1;
5      /*进程可以设置调度策略为SCHED_YIELD即“礼让”策略，这时候它的权值为-1，权值相当
低*/
6      if (p->policy & SCHED_YIELD)
7          goto out;
8
9      /*
10     * Non-RT process - normal case first.
11     */
12
13     /*对于调度策略为SCHED_OTHER的进程，没有实时性要求，它的权值仅仅取决于
14     * 时间片的剩余和它的nice值，数值上nice越小，则优先级越高，总的权值 = 时间片剩
余 + (20 - nice)
15     * */
16     if (p->policy == SCHED_OTHER) {
17         /*
18         * Give the process a first-approximation goodness value
19         * according to the number of clock-ticks it has left.
20         *
21         * Don't do any other calculations if the time slice is
22         * over..
23         */
24         weight = p->counter;
25         if (!weight)
26             goto out;
27
28 #ifdef CONFIG_SMP
29         /* Give a largish advantage to the same processor... */
30         /* (this is equivalent to penalizing other processors) */
31         if (p->processor == this_cpu)
32             weight += PROC_CHANGE_PENALTY;
33 #endif
34
35         /* .. and a slight advantage to the current MM */
36         if (p->mm == this_mm || !p->mm)
37             weight += 1;
38         weight += 20 - p->nice;
39         goto out;

```

```

40         }
41
42         /*
43         *对于实时进程，也就是SCHED_FIFO或者SCHED_RR调度策略，
44         *具有一个实时优先级，总的权值仅仅取决于该实时优先级，
45         *总的权值 = 1000 + 实时优先级。
46         * */
47         weight = 1000 + p->rt_priority;
48     out:
49         return weight;
50     }

```

从代码可以看到：

- 当进程的剩余时间片Counter为0时，无论静态优先级是多少都不会被选中
- 普通进程的优先级=剩余时间片Counter+20-nice
- 实时进程的优先级=1000+进程静态优先级
- 实时进程的动态优先级远大于普通进程，更容易被选中
- 剩余时间片counter越多说明进程IO较多，分配给它的没用完，被调度的优先级需要高一些

```

1  #if HZ < 200
2  #define TICK_SCALE(x)      ((x) >> 2)
3  #elif HZ < 400
4  #define TICK_SCALE(x)      ((x) >> 1)
5  #elif HZ < 800
6  #define TICK_SCALE(x)      (x)
7  #elif HZ < 1600
8  #define TICK_SCALE(x)      ((x) << 1)
9  #else
10 #define TICK_SCALE(x)      ((x) << 2)
11 #endif
12
13 #define NICE_TO_TICKS(nice)    (TICK_SCALE(20-(nice))+1)

```

NICE_TO_TICKS是个宏函数，根据不同的调度频率HZ有对应的TICK_SCALE宏定义，这样就解决了不同优先级的进程的时间片分配问题。

O(n)调度器对实时进程和普通进程采用不同的调度策略：

- 实时进程采用的是SCHED_RR或者SCHED_FIFO，高级优先&同级轮转或者顺序执行
- 普通进程采用的是SCHED_OTHER
- 进程采用的策略在task_struct中policy体现

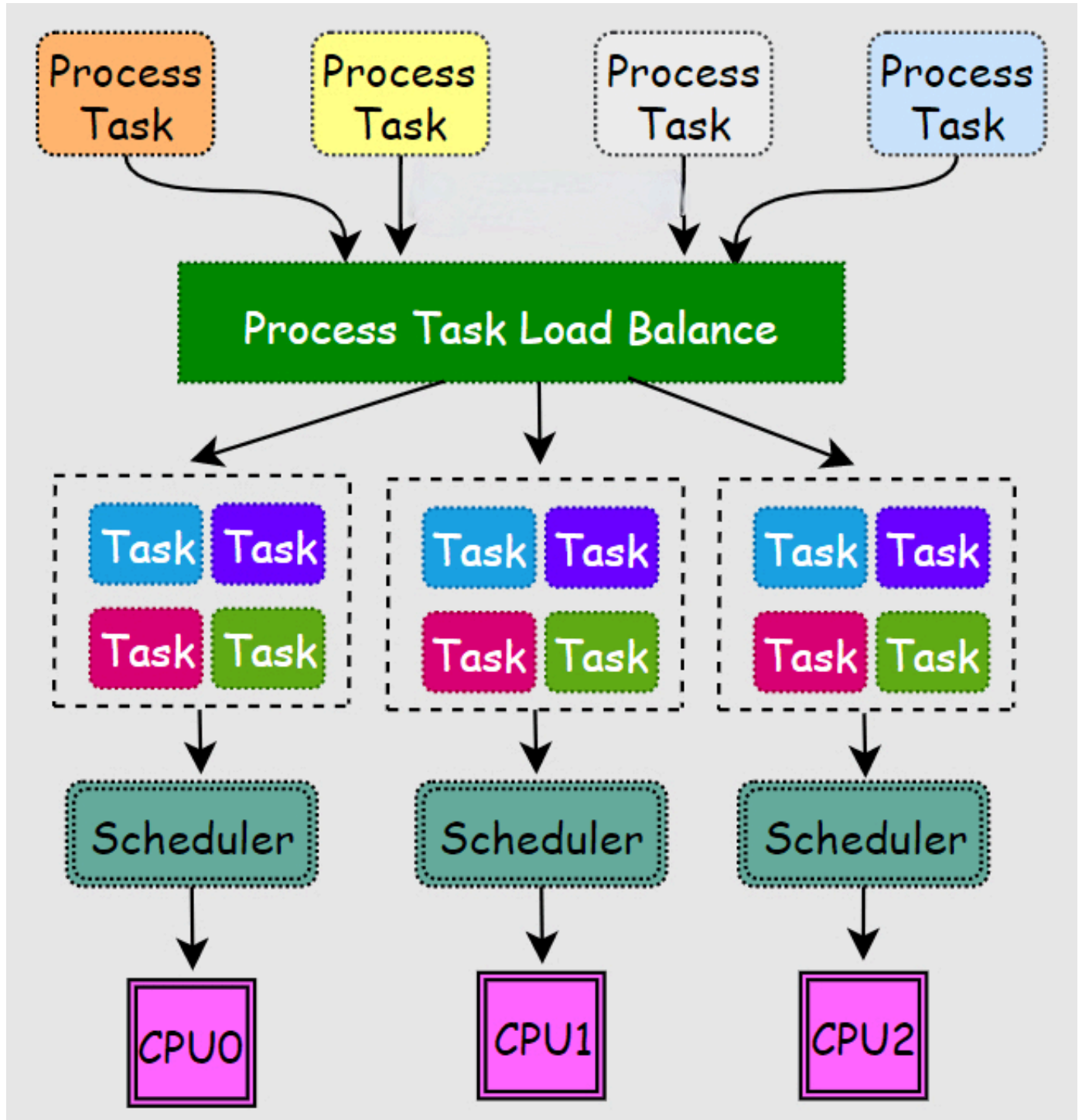
在runqueue中搜索下一个合适的进程是基于动态优先级来实现的，动态优先级最高的就是下一个被执行的进程。

O(n)调度器设计和实现上存在一些问题，但是其中的很多思想为后续调度器设计指明了方向，意义深远。

3.2 O(1)调度器

Linux 2.6.0采纳了Red Hat公司Ingo Molnar设计的O(1)调度算法，该调度算法的核心思想基于Corbato等人提出的多级反馈队列算法。

O(1)调度器引入了多个队列，并且增加了负载均衡机制，对新出现的进行任务分配到合适的cpu-runqueue中：



为了实现O(1)复杂度的pick-next算法，内核实现代码量增加了一倍多，其有以下几个特点：

- 实现了per-cpu-runqueue，每个CPU都有一个就绪进程任务队列
- 引入活跃数组active和过期数组expire，分别存储就绪进程和结束进程
- 采用全局优先级：实时进程0-99，普通进程100-139，数值越低优先级越高，更容易被调度
- 每个优先级对应一个链表，引入bitmap数组来记录140个链表中的活跃任务情况

任务队列的数据结构:

```

1 struct runqueue {
2     spinlock_t lock;
3     unsigned long nr_running;
4     unsigned long long nr_switches;
5     unsigned long expired_timestamp, nr_uninterruptible;
6     unsigned long long timestamp_last_tick;
7     task_t *curr, *idle;
8     struct mm_struct *prev_mm;
9     prio_array_t *active, *expired, arrays[2];
10    int best_expired_prio;
11    atomic_t nr_iowait;
12    .....
13 };

```

active和expired是指向prio_array_t的结构体指针

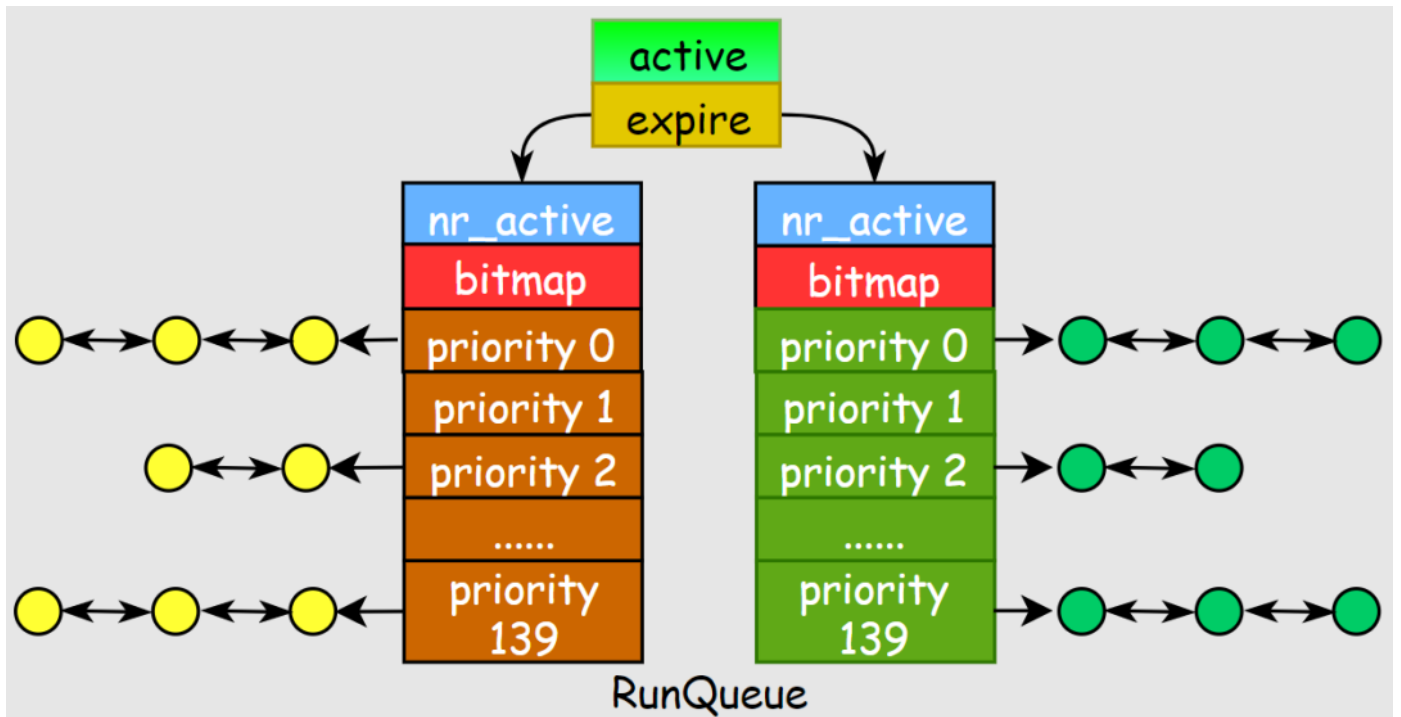
arrays是元素为prio_array_t的结构体数组

prio_array_t结构体的定义:

```

1 #define BITMAP_SIZE (((MAX_PRIO+1+7)/8)+sizeof(long)-1)/sizeof(long)
2 typedef struct prio_array prio_array_t;
3 struct prio_array {
4     unsigned int nr_active;
5     unsigned long bitmap[BITMAP_SIZE];
6     struct list_head queue[MAX_PRIO];
7 };

```



O(1)调度器对pick-next的实现:

- 在runqueue结构中有active和expire两个数组指针，active指向就绪进程的结构，从active-bitmap中寻找优先级最高且非空的数组元素，这个数组是元素是进程链表，找该链表中第1个进程即可。

```

1 | idx = sched_find_first_bit(array->bitmap);
2 | queue = array->queue + idx;
3 | next = list_entry(queue->next, task_t, run_list);

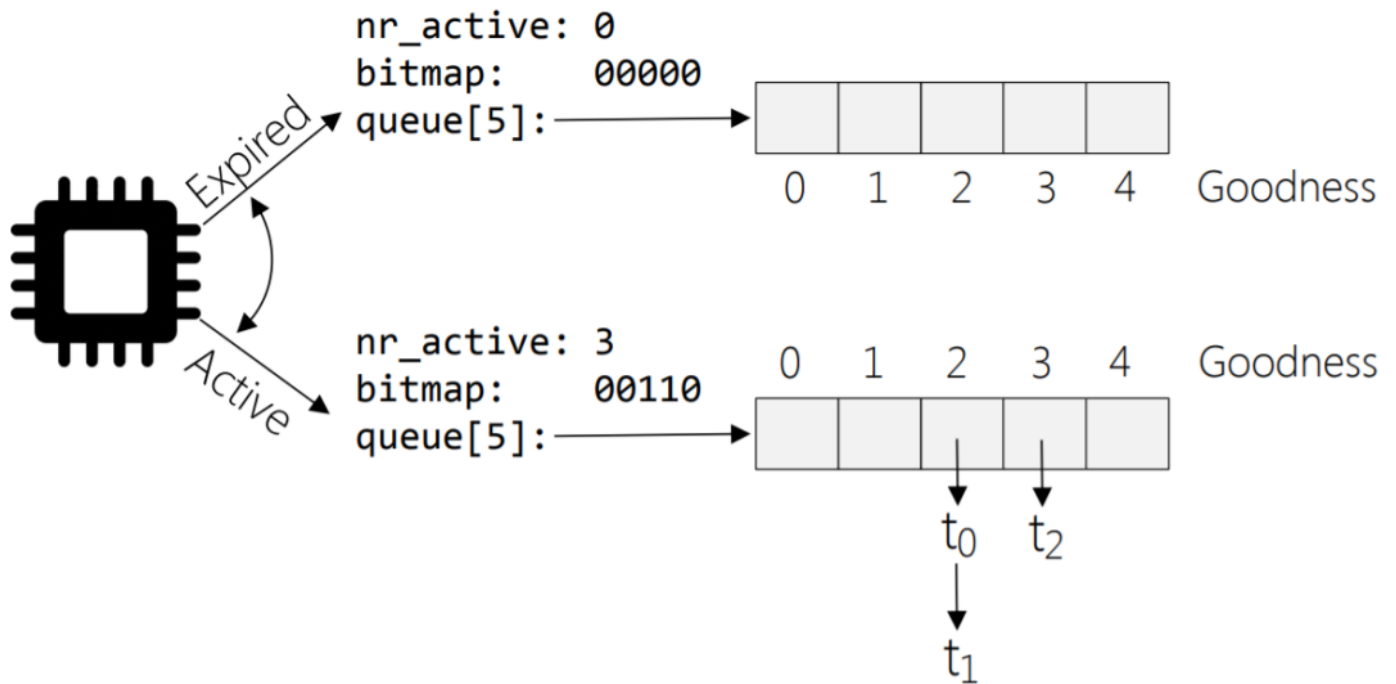
```

- 当active的nr_active=0时表示没有活跃任务，此时进行active和expire双指针互换，速度很快。

```

1 | array = rq->active;
2 | if (unlikely(!array->nr_active)) {
3 |     /*
4 |      * Switch the active and expired arrays.
5 |      */
6 |     rq->active = rq->expired;
7 |     rq->expired = array;
8 |     array = rq->active;
9 |     rq->expired_timestamp = 0;
10 |    rq->best_expired_prio = MAX_PRIO;
11 | }

```



O(1)和O(n)调度器确定进程优先级的方法不一样，O(1)借助了sleep_avg变量记录进程的睡眠时间，来识别IO密集型进程，计算bonus值来调整优先级：

```

1  #define NICE_TO_PRIO(nice)          (MAX_RT_PRIO + (nice) + 20)
2  #define NS_TO_JIFFIES(TIME)         ((TIME) / (1000000000 / HZ))
3  #define CURRENT_BONUS(p) \
4      (NS_TO_JIFFIES((p)->sleep_avg) * MAX_BONUS / \
5       MAX_SLEEP_AVG)
6
7  static int effective_prio(task_t *p)
8  {
9      int bonus, prio;
10
11     if (rt_task(p))
12         return p->prio;
13
14     bonus = CURRENT_BONUS(p) - MAX_BONUS / 2;
15
16     prio = p->static_prio - bonus;
17     if (prio < MAX_RT_PRIO)
18         prio = MAX_RT_PRIO;
19     if (prio > MAX_PRIO-1)
20         prio = MAX_PRIO-1;
21     return prio;
22 }

```

O(1)调度器为了实现复杂场景IO密集型任务的识别，做了大量的工作仍然无法到达100%的准确，但不可否认O(1)调度器是一款非常优秀的产品。

3.3 CFS调度器

O(1)调度器本质上是MLFQ算法的思想，随着时间的推移也暴露除了很多问题，主要集中在O(1)调度器对进程交互性的判断上积重难返。

无论是O(n)还是O(1)都需要去根据进程的运行状况判断它属于IO密集型还是CPU密集型，再做优先级奖励和惩罚，这种推测本身就会有误差，而且场景越复杂判断难度越大。

在2.6.23内核中引入scheduling class的概念，将调度器模块化，系统中可以有多种调度器，使用不同策略调度不同类型的进程：

- DL Scheduler 采用sched_deadline策略
- RT Scheduler 采用sched_rr和sched_fifo策略
- CFS Scheduler 采用sched_normal和sched_batch策略
- IDEL Scheduler 采用sched_idle策略

这样一来，CFS调度器就不关心实时进程了，专注于普通进程就可以了。

CFS(Completely Fair Scheduler)完全公平调度器，从实现思想上和之前的O(1)/O(n)很不一样。

3.3.1 优先级和权重

O(1)和O(n)都将CPU资源划分为时间片，采用了固定额度分配机制，在每个调度周期进程可使用的时间片是确定的，调度周期结束被重新分配。

CFS摒弃了固定时间片分配，采用动态时间片分配，本次调度中进程可占用的时间与进程总数、总CPU时间、进程权重等均有关关系，每个调度周期的值都可能会不一样。

CFS调度器从进程优先级出发，它建立了优先级prio和权重weight之间的映射关系，把优先级转换为权重来参与后续的计算：

```

1  const int sched_prio_to_weight[40] = {
2  /* -20 */      88761,      71755,      56483,      46273,      36291,
3  /* -15 */      29154,      23254,      18705,      14949,      11916,
4  /* -10 */      9548,       7620,       6100,       4904,       3906,
5  /*  -5 */      3121,       2501,       1991,       1586,       1277,
6  /*   0 */      1024,        820,        655,        526,        423,
7  /*   5 */       335,        272,        215,        172,        137,
8  /*  10 */       110,         87,         70,         56,         45,
9  /*  15 */        36,         29,         23,         18,         15,
10 };
```

普通进程的优先级范围是[100,139]，prio整体减小120就和代码左边的注释对上了，也就是nice值的范围[-20,19]，因此sched_prio=0相当于static_prio=120。

比如现有进程A sched_prio=0，进程B sched_prio=-5，通过sched_prio_to_weight的映射：

- 进程A weight=1024, 进程B weight = 3121
- 进程A的CPU占比 = $1024/(1024+3121) = 24.7\%$
- 进程B的CPU占比 = $3121/(1024+3121) = 75.3\%$
- 假如CPU总时间是10ms, 那么根据A占用2.47ms, B占用7.53ms

在CFS中引入sysctl_sched_latency(调度延迟)作为一个调度周期, 真实的CPU时间表示为:

$$\text{time_slice}_k = \frac{\text{weight}_k}{\sum_{n=0}^{n-1} \text{weight}_i} \cdot \text{sched_latency}$$

显然这样根据权重计算后的各个进程的运行时间是不等的, 也就违背了"完全公平"思想, 于是CFS引入了虚拟运行时间(virtual runtime)。

3.3.2 虚拟运行时间

每个进程的物理运行时间肯定不能一样的, CFS调度器只要保证的就是进程的虚拟运行时间相等即可。

那虚拟运行时间该如何计算呢?

$$\text{virtual_time} = \text{wall_time} * \text{nice_0_weight} / \text{sched_prio_to_weigh}$$

比如现有进程A sched_prio=0, 进程B sched_prio=-5:

- 调度延迟=10ms, A的运行时间=2.47ms B的运行时间=7.53ms, 也就是wall_time
- nice_0_weight表示sched_prio=0的权重为1024
- 进程A的虚拟时间: $2.47 * 1024 / 1024 = 2.47\text{ms}$
- 进程B的虚拟时间: $7.53 * 1024 / 3121 = 2.47\text{ms}$

经过这样映射, A和B的虚拟时间就相等了。

上述公式涉及了除法和浮点数运算, 因此需要转换为乘法来保证数据准确性, 再给出虚拟时间计算的变形等价公式:

$$\text{virtual_time} = (\text{wall_time} * \text{nice_0_weight} * 2^{32} / \text{sched_prio_to_weigh}) \gg 32$$

$$\text{令 } \text{inv_weight} = 2^{32} / \text{sched_prio_to_weigh}$$

$$\text{则 } \text{virtual_time} = (\text{wall_time} * 1024 * \text{inv_weight}) \gg 32$$

由于sched_prio_to_weigh的值存储在数组中, inv_weight同样可以:

```

1  const u32 sched_prio_to_wmult[40] = {
2      /* -20 */      48388,      59856,      76040,      92818,      118348,
3      /* -15 */      147320,      184698,      229616,      287308,      360437,
4      /* -10 */      449829,      563644,      704093,      875809,      1099582,
5      /*  -5 */      1376151,      1717300,      2157191,      2708050,      3363326,
6      /*   0 */      4194304,      5237765,      6557202,      8165337,      10153587,
7      /*   5 */      12820798,      15790321,      19976592,      24970740,      31350126,
8      /*  10 */      39045157,      49367440,      61356676,      76695844,      95443717,
9      /*  15 */      119304647,      148102320,      186737708,      238609294,      286331153,
10 };
```

经过一番计算，各个进程的虚拟运行时间一致了，似乎我们理解了"完全公平"的思想。

虚拟运行时间与优先级的衰减因子有关，也就是`inv_weight`随着`nice`值增大而增大，同时其作为分母也加速了低优先级进程的衰减。

- `nice=0` 虚拟运行时间 = 物理运行时间
- `nice>0` 虚拟运行时间 > 物理运行时间
- `nice<0` 虚拟运行时间 < 物理运行时间

简言之：**CFS**将物理运行时间在不同优先级进程中发生了不同的通胀。

摒弃了固定时间片机制也是CFS的亮点，系统负载高时大家都少用一些CPU，系统负载低时大家都多用一些CPU，让调度器有了一定的自适应能力。

3.3.3 pick-next和红黑树

那么这些进程应该采用哪种数据结构来实现pick-next算法呢？

CFS调度器采用了红黑树来保存活跃进程任务，红黑树的增删查复杂度是 $O(\log n)$ ，但是CFS引入了一些额外的数据结构，可以免去遍历获得下一个最合适的进程。

红黑树的key是进程已经使用的虚拟运行时间，并且把虚拟时间数值最小的放到最左的叶子节点，这个节点就是下一个被pick的进程了。

前面已经论证了，每个进程的虚拟运行时间是一样的，数值越小表示被调度的越少，因此需要更偏爱一些，当虚拟运行时间耗尽则从红黑树中删除，下个调度周期开始后再添加到红黑树上。

4 实验总结

1. $O(n)$ 调度器采用全局runqueue，导致多cpu加锁问题和cache利用率低的问题
2. $O(1)$ 调度器为每个cpu设计了一个runqueue，并且采用MLFQ算法思想设置140个优先级链表和active/expire两个双指针结构
3. CFS调度器采用红黑树来实现 $O(\log n)$ 复杂度的pick-next算法，摒弃固定时间片机制，采用调度周期内的动态时间机制
4. $O(1)$ 和 $O(n)$ 都在交互进程的识别算法上下了功夫，但是无法做的100%准确
5. CFS另辟蹊径采用完全公平思想以及虚拟运行时间来实现进行的调度，但并非完美，在某些方面可能不如 $O(1)$

5 参考文献

- [1] 进程.<https://zh.wikipedia.org/wiki/进程>
- [2] Linux进程描述符task_struct结构体详解--Linux进程的管理与调度 <https://blog.csdn.net/gatieme/article/details/51383272>
- [3] Linux进程管理之task_struct结构体 <https://blog.csdn.net/npylp/article/details/7292563>
- [4] Linux进程: task_struct结构体成员 <https://www.cnblogs.com/jiujuan/p/11715853.html>
- [5] Linux进程描述符task_struct结构体详解--Linux进程的管理与调度 <https://cloud.tencent.com/developer/article/1339556>
- [6] Linux进程调度器-基础 <https://www.cnblogs.com/LoyenWang/p/12249106.html>
- [7] 吐血整理 | 肝翻 Linux 进程调度所有知识点 <https://www.eefocus.com/embedded/498554>
- [8] The Linux Completely Fair Scheduler (CFS) https://cs334.cs.vassar.edu/slides/04_Linux_CFS.pdf
- [9] Scheduling: Case Studies <http://www.eecs.harvard.edu/~cs161/notes/scheduling-case-studies.pdf>
- [10] O(n)、O(1)和CFS调度器 http://www.wowotech.net/process_management/scheduler-history.html