# Linux schedule Cgoup

原创 伟林 Linux阅码场

作者简介

伟林，中年码农，从事过电信、手机、安全、芯片等行业，目前依旧从事Linux方向开发工作，个人爱好Linux相关知识分享，个人微博CSDN pwl999，欢迎大家关注！

> **Q** 学员问：我最近在看k8s对cgroup的管理部分，对于cfs对cgroup的调度有些疑惑。想搞明白cgroup里面的 period、quota是如何影响cfs的调度的
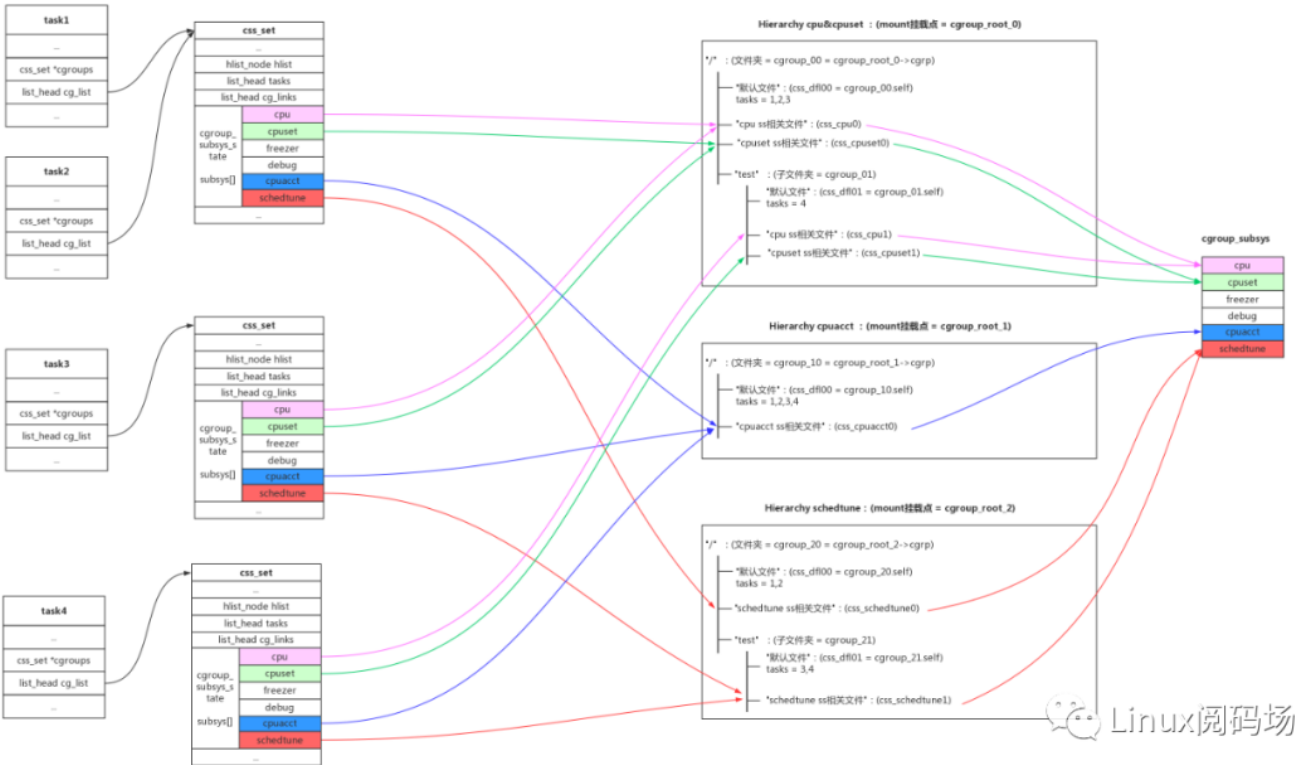
> **A** 伟林老师给出如下文章进行解答

# 1.Cgoup

## 1.1、cgroup概念

cgroup最基本的操作时我们可以使用以下命令创建一个cgroup文件夹：

```
1  mount -t cgroup -o cpu,cpuset cpu&cpuset /dev/cpu_cpuset_test
```

那么/dev/cpu_cpuset_test文件夹下就有一系列的cpu、cpuset cgroup相关的控制节点，tasks文件中默认加入了所有进程到这个cgroup中。可以继续创建子文件夹，子文件夹继承了父文件夹的结构形式，我们可以给子文件夹配置不同的参数，把一部分进程加入到子文件夹中的tasks文件当中，久可以实现分开的cgroup控制了。

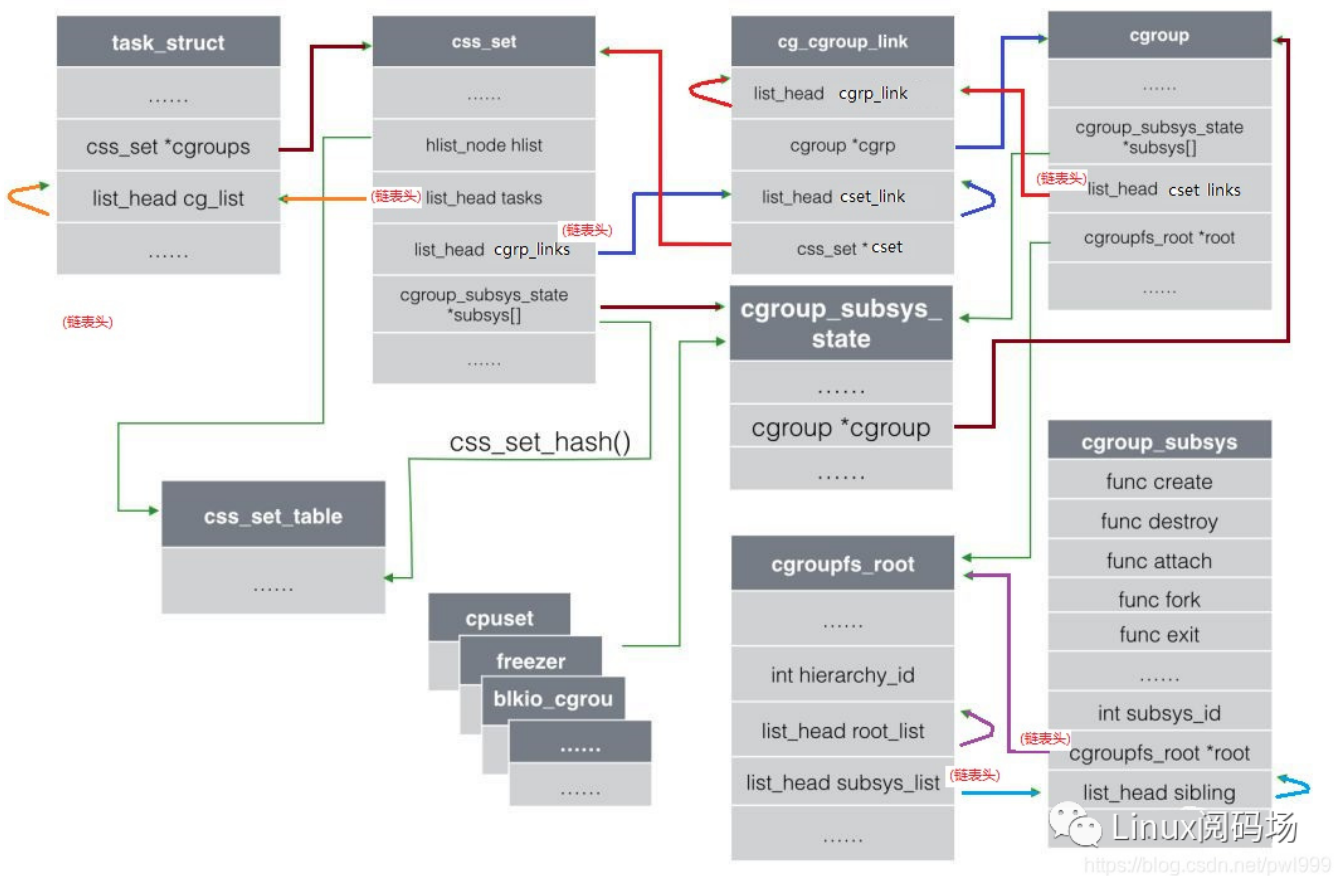

关于cgroup的结构有以下规则和规律：

- 1、cgroup有很多subsys，我们平时接触到的cpu、cpuset、cpuacct、memory、blkio都是cgroup_subsys；
- 2、一个cgroup hierarchy，就是使用mount命令挂载的一个cgroup文件系统，hierarchy对应mount的根cgroup_root；
- 3、一个hierarchy可以制定一个subsys，也可以制定多个subsys。可以是一个subsys，也可以是一个subsys组合；
- 4、一个subsys只能被一个hierarchy引用一次，如果subsys已经被hierarchy引用，新hierarchy创建时不能引用这个subsys；唯一例外的是，我们可以创建和旧的hierarchy相同的subsys组合，这其实没有创建新的hierarchy，只是简单的符号链接；

- 5、hierarchy对应一个文件系统，cgroup对应这个文件系统中的文件夹；subsys是基类，而css(cgroup_subsys_state)是cgroup引用subsys的实例；比如父目录和子目录分别是两个cgroup，他们都要引用相同的subsys，但是他们需要不同的配置，所以会创建不同的css供cgroup->subsys[]来引用；
- 6、一个任务对系统中不同的subsys一定会有引用，但是会引用到不同的hierarchy不同的cgroup即不同css当中；所以系统使用css_set结构来管理任务对css的引。如果任务引用的css组合相同，那他们开源使用相同的css_set；
- 7、还有cgroup到task的反向引用，系统引入了cg_group_link结构。这部分可以参考Docker背后的内核知识——cgroups资源限制一文的描述，如下图的结构关系：

## cgroup数据结构之间的关系



1、subsys是一组基类(cpu、blkio)，css(cgroup_subsys_state)是基类的实例化。

2、cgroup的一组css的集合。

3、hierarchy是多个cgoup的组合，它决定cgroup中能创建哪些subsys的css。hierarchy可以任意引用几种subsys，但是一个subsys只能被一个hierarchy引用。如果一个hierarchy已经引用某个subsys，那么其他hierarchy就不能再引用这个subsys了。hierarchy对应cgroupfs_root数据结构。
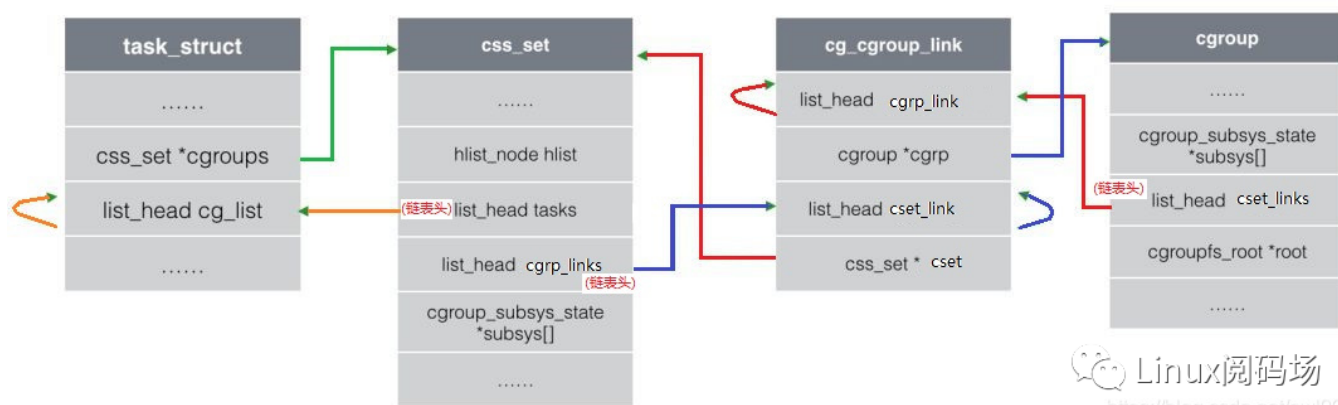
4、一旦hierarchy确定了subsys，那么它下面的cgroup只能创建对应的css实例。一个subsys只能存在于某个hierarchy中，hierarchy下的多个cgroup可以创建这个subsys对应的多个css。

5、hierarchy、cgroup、css三者还使用文件系统来表示层次关系：hierarchy是文件系统挂载点，cgroup是文件夹，css是文件夹中的文件。css的值，以及兄弟和父子关系，表示了subsys资源配额的关系。

6、cgoup是为了划分资源配额，配置的主体是进程task。每个task在每一类别的subsys上都有配额，所以每个task在每个类别的subsys上有一个唯一的css与之关联。

7、进程和css是一对多(1 x N)的关系。而系统中的多个进程和多个css，是多对多(M x N)的关系。为了收敛这种多对多的关系，系统把所有css属性都相同的一组进程放在一个css_set当中，把多个css放在一个cgroup当中，这样还是多对多但是已经收敛(M/a x N/b)。css_set根据属性组合，存入css_set_table当中。

8、css_set代表a个css属性相同的进程，cgroup代表引用的b个subsys。多对多的关系从task vs css的(M x N)，收敛到css_set vs cgroup的(M/a x N/b)。为了进一步简化css_set和cgroup之间多对多关系的双向查找，引入了cg_group_link数据结构：
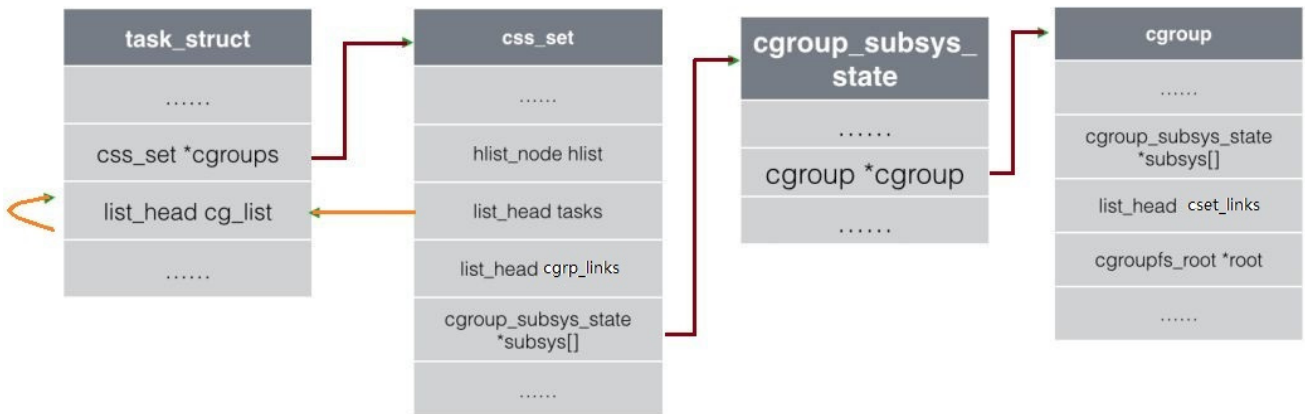


task_struct通过->cgroup成员找到css_set结构，css_set利用->tasks链表把所有css属性相同的进程链接到一起。

| dir | descript |
| --- | --- |
| css_set → cgroup | css_set的->cgrp_links链表上挂载了这组css相关cgroup对应的cg_cgroup_link，通过cg_cgro |

| | |
|---|---|
| | up_link->cgrp找到cgroup，再通过cgroup->subsys[]找到css。 |
| cgroup → css_set | cgroup的->cset_links链表上挂载了所有指向本cgoup的task对应的cg_cgroup_link，通过cg_cgroup_link->cset找到css_set，再通过css_set->tasks找到所有的task_struct。 |

9、还有一条task_struct → cgroup 的通路：



路径： task_struct->cgroup → css_set->subsys[] → cgroup_subsys_state->cgroup → cgroup

## 1.2、代码分析

1、"/proc/cgroups"

subsys的链表：for_each_subsys(ss, i)
一个susbsys对应一个hierarchy：ss->root
一个hierarchy有多少个cgroup：ss->root->nr_cgrps

```
1  # ount -t cgroup -o freezer,debug bbb freezer_test/
2
3  # cat /proc/cgroups
4  #subsys_name    hierarchy      num_cgroups     enabled
5  cpuset    4       6       1
6  cpu       3       2       1
```

```
 7  cpuacct 1       147     1
 8  schedtune       2       3       1
 9  freezer 6       1       1
10  debug   6       1       1
11
```

```c
 1  static int proc_cgroupstats_show(struct seq_file *m, void *v)
 2  {
 3    struct cgroup_subsys *ss;
 4    int i;
 5
 6    seq_puts(m, "#subsys_name\thierarchy\tnum_cgroups\tenabled\n");
 7    /*
 8     * ideally we don't want subsystems moving around while we do this.
 9     * cgroup_mutex is also necessary to guarantee an atomic snapshot of
10     * subsys/hierarchy state.
11     */
12    mutex_lock(&cgroup_mutex);
13
14    for_each_subsys(ss, i)
15      seq_printf(m, "%s\t%d\t%d\t%d\n",
16            ss->legacy_name, ss->root->hierarchy_id,
17            atomic_read(&ss->root->nr_cgrps),
18            cgroup_ssid_enabled(i));
19
20    mutex_unlock(&cgroup_mutex);
21    return 0;
22  }
```

## 2、"/proc/pid/cgroup"

每种subsys组合组成一个新的hierarchy，每个hierarchy在for_each_root(root)中创建一个root树；
每个hierarchy顶层目录和子目录都是一个cgroup，一个hierarchy可以有多个cgroup，对应的subsys组合一样，但是参数不一样
cgroup_root自带一个cgroup即root->cgrp，作为hierarchy的顶级目录

一个 cgroup 对应多个 subsys, 使用 cgroup_subsys_state 类型 (css) 的 cgroup->subsys[CGROUP_SUBSYS_COUNT]数组去和多个subsys链接;

一个cgroup自带一个cgroup_subsys_state即cgrp->self, 这个css的作用是css->parent指针, 建立起cgroup之间的父子关系;

css一个公用结构, 每个subsys使用自己的函数ss->css_alloc()分配自己的css结构, 这个结构包含公用css + subsys私有数据;

每个subsys只能存在于一个组合(hierarchy)当中, 如果一个subsys已经被一个组合引用, 其他组合不能再引用这个subsys。唯一例外的是, 我们可以重复mount相同的组合, 但是这样并没有创建新组合, 只是创建了一个链接指向旧组合;

进程对应每一种hierarchy, 一定有一个cgroup对应。

```
1   # cat /proc/832/cgroup
2   6:freezer,debug:/
3   4:cpuset:/
4   3:cpu:/
5   2:schedtune:/
6   1:cpuacct:/
```

```
1    int proc_cgroup_show(struct seq_file *m, struct pid_namespace *ns,
2            struct pid *pid, struct task_struct *tsk)
3    {
4      char *buf, *path;
5      int retval;
6      struct cgroup_root *root;
7
8      retval = -ENOMEM;
9      buf = kmalloc(PATH_MAX, GFP_KERNEL);
10     if (!buf)
11       goto out;
12
13     mutex_lock(&cgroup_mutex);
14     spin_lock_bh(&css_set_lock);
15
16     for_each_root(root) {
17       struct cgroup_subsys *ss;
18       struct cgroup *cgrp;
19       int ssid, count = 0;
```

```
20
21        if (root == &cgrp_dfl_root && !cgrp_dfl_root_visible)
22          continue;
23
24        seq_printf(m, "%d:", root->hierarchy_id);
25        if (root != &cgrp_dfl_root)
26          for_each_subsys(ss, ssid)
27            if (root->subsys_mask & (1 << ssid))
28              seq_printf(m, "%s%s", count++ ? "," : "",
29                    ss->legacy_name);
30        if (strlen(root->name))
31          seq_printf(m, "%sname=%s", count ? "," : "",
32                root->name);
33        seq_putc(m, ':');
34
35        cgrp = task_cgroup_from_root(tsk, root);
36
37        /*
38         * On traditional hierarchies, all zombie tasks show up as
39         * belonging to the root cgroup.  On the default hierarchy,
40         * while a zombie doesn't show up in "cgroup.procs" and
41         * thus can't be migrated, its /proc/PID/cgroup keeps
42         * reporting the cgroup it belonged to before exiting.  If
43         * the cgroup is removed before the zombie is reaped,
44         * " (deleted)" is appended to the cgroup path.
45         */
46        if (cgroup_on_dfl(cgrp) || !(tsk->flags & PF_EXITING)) {
47          path = cgroup_path(cgrp, buf, PATH_MAX);
48          if (!path) {
49            retval = -ENAMETOOLONG;
50            goto out_unlock;
51          }
52        } else {
53          path = "/";
54        }
55
56        seq_puts(m, path);
57
58        if (cgroup_on_dfl(cgrp) && cgroup_is_dead(cgrp))
59          seq_puts(m, " (deleted)\n");
```

```
60        else
61          seq_putc(m, '\n');
62      }
63
64      retval = 0;
65  out_unlock:
66      spin_unlock_bh(&css_set_lock);
67      mutex_unlock(&cgroup_mutex);
68      kfree(buf);
69  out:
70      return retval;
71  }
```

## 3、初始化

```
1   int __init cgroup_init_early(void)
2   {
3       static struct cgroup_sb_opts __initdata opts;
4       struct cgroup_subsys *ss;
5       int i;
6
7       /* (1) 初始化默认root cgrp_dfl_root，选项opts为空，初始了
8           root->cgrp          // cgrp->root = root;
9           root->cgrp.self     // cgrp->self.cgroup = cgrp; cgrp->self.flags |=
10      */
11      init_cgroup_root(&cgrp_dfl_root, &opts);
12      cgrp_dfl_root.cgrp.self.flags |= CSS_NO_REF;
13
14      RCU_INIT_POINTER(init_task.cgroups, &init_css_set);
15
16      /* (2) 轮询subsys进行初始化 */
17      for_each_subsys(ss, i) {
18          WARN(!ss->css_alloc || !ss->css_free || ss->name || ss->id,
19              "invalid cgroup_subsys %d:%s css_alloc=%p css_free=%p name:id=%d:%s
20              i, cgroup_subsys_name[i], ss->css_alloc, ss->css_free,
21              ss->id, ss->name);
22          WARN(strlen(cgroup_subsys_name[i]) > MAX_CGROUP_TYPE_NAMELEN,
23              "cgroup_subsys_name %s too long\n", cgroup_subsys_name[i]);
```

```
24
25          /* (3) 初始化ss->id、ss->name */
26      ss->id = i;
27      ss->name = cgroup_subsys_name[i];
28      if (!ss->legacy_name)
29        ss->legacy_name = cgroup_subsys_name[i];
30
31          /* (4) ss链接到默认root(cgrp_dfl_root)
32              默认css_set(init_css_set)指向ss
33          */
34      if (ss->early_init)
35        cgroup_init_subsys(ss, true);
36    }
37    return 0;
38  }
39
40  |→
41
42  static void __init cgroup_init_subsys(struct cgroup_subsys *ss, bool early)
43  {
44    struct cgroup_subsys_state *css;
45
46    printk(KERN_INFO "Initializing cgroup subsys %s\n", ss->name);
47
48    mutex_lock(&cgroup_mutex);
49
50    idr_init(&ss->css_idr);
51    INIT_LIST_HEAD(&ss->cfts);
52
53    /* Create the root cgroup state for this subsystem */
54    ss->root = &cgrp_dfl_root;
55
56    /* (4.1) subsys分配一个新的相关的cgroup_subsys_state */
57    css = ss->css_alloc(cgroup_css(&cgrp_dfl_root.cgrp, ss));
58    /* We don't handle early failures gracefully */
59    BUG_ON(IS_ERR(css));
60
61    /* (4.2) 初始化css的成员指向cgroup
62        cgroup为默认值cgrp_dfl_root.cgrp:
63        css->cgroup = cgrp;
```

```
64        css->ss = ss;
65        INIT_LIST_HEAD(&css->sibling);
66        INIT_LIST_HEAD(&css->children);
67     */
68    init_and_link_css(css, ss, &cgrp_dfl_root.cgrp);
69
70    /*
71     * Root csses are never destroyed and we can't initialize
72     * percpu_ref during early init.  Disable refcnting.
73     */
74    css->flags |= CSS_NO_REF;
75
76    if (early) {
77      /* allocation can't be done safely during early init */
78      css->id = 1;
79    } else {
80      css->id = cgroup_idr_alloc(&ss->css_idr, css, 1, 2, GFP_KERNEL);
81      BUG_ON(css->id < 0);
82    }
83
84    /* Update the init_css_set to contain a subsys
85     * pointer to this state - since the subsystem is
86     * newly registered, all tasks and hence the
87     * init_css_set is in the subsystem's root cgroup. */
88    /* (4.3) css_set指向新的css */
89    init_css_set.subsys[ss->id] = css;
90
91    have_fork_callback |= (bool)ss->fork << ss->id;
92    have_exit_callback |= (bool)ss->exit << ss->id;
93    have_free_callback |= (bool)ss->free << ss->id;
94    have_canfork_callback |= (bool)ss->can_fork << ss->id;
95
96    /* At system boot, before all subsystems have been
97     * registered, no tasks have been forked, so we don't
98     * need to invoke fork callbacks here. */
99    BUG_ON(!list_empty(&init_task.tasks));
100
101    /* (4.4) cgroup测指向css：
102        执行ss->css_online(css);
103        css->cgroup->subsys[ss->id] = css;
```

```
104        */
105     BUG_ON(online_css(css));
106
107     mutex_unlock(&cgroup_mutex);
108  }
109
110
111  int __init cgroup_init(void)
112  {
113     struct cgroup_subsys *ss;
114     int ssid;
115
116     BUG_ON(percpu_init_rwsem(&cgroup_threadgroup_rwsem));
117     BUG_ON(cgroup_init_cftypes(NULL, cgroup_dfl_base_files));
118     BUG_ON(cgroup_init_cftypes(NULL, cgroup_legacy_base_files));
119
120     /*
121      * The latency of the synchronize_sched() is too high for cgroups,
122      * avoid it at the cost of forcing all readers into the slow path.
123      */
124     rcu_sync_enter_start(&cgroup_threadgroup_rwsem.rss);
125
126     mutex_lock(&cgroup_mutex);
127
128     /*
129      * Add init_css_set to the hash table so that dfl_root can link to
130      * it during init.
131      */
132     hash_add(css_set_table, &init_css_set.hlist,
133        css_set_hash(init_css_set.subsys));
134
135     BUG_ON(cgroup_setup_root(&cgrp_dfl_root, 0));
136
137     mutex_unlock(&cgroup_mutex);
138
139     for_each_subsys(ss, ssid) {
140       if (ss->early_init) {
141         struct cgroup_subsys_state *css =
142           init_css_set.subsys[ss->id];
143
```

```c
            css->id = cgroup_idr_alloc(&ss->css_idr, css, 1, 2,
                        GFP_KERNEL);
            BUG_ON(css->id < 0);
        } else {
            cgroup_init_subsys(ss, false);
        }

        list_add_tail(&init_css_set.e_cset_node[ssid],
                &cgrp_dfl_root.cgrp.e_csets[ssid]);

        /*
         * Setting dfl_root subsys_mask needs to consider the
         * disabled flag and cftype registration needs kmalloc,
         * both of which aren't available during early_init.
         */
        if (cgroup_disable_mask & (1 << ssid)) {
            static_branch_disable(cgroup_subsys_enabled_key[ssid]);
            printk(KERN_INFO "Disabling %s control group subsystem\n",
                    ss->name);
            continue;
        }

            /* (1) 默认root(cgrp_dfl_root)，支持所有ss */
        cgrp_dfl_root.subsys_mask |= 1 << ss->id;

        if (!ss->dfl_cftypes)
            cgrp_dfl_root_inhibit_ss_mask |= 1 << ss->id;

            /* (2) 将cftypes(ss->legacy_cftypes/ss->legacy_cftypes)加入到ss->cfts
        if (ss->dfl_cftypes == ss->legacy_cftypes) {
            WARN_ON(cgroup_add_cftypes(ss, ss->dfl_cftypes));
        } else {
            WARN_ON(cgroup_add_dfl_cftypes(ss, ss->dfl_cftypes));
            WARN_ON(cgroup_add_legacy_cftypes(ss, ss->legacy_cftypes));
        }

        if (ss->bind)
            ss->bind(init_css_set.subsys[ssid]);
    }

```

```
184    /* init_css_set.subsys[] has been updated, re-hash */
185    hash_del(&init_css_set.hlist);
186    hash_add(css_set_table, &init_css_set.hlist,
187        css_set_hash(init_css_set.subsys));
188
189    WARN_ON(sysfs_create_mount_point(fs_kobj, "cgroup"));
190    WARN_ON(register_filesystem(&cgroup_fs_type));
191    WARN_ON(!proc_create("cgroups", 0, NULL, &proc_cgroupstats_operations));
192
193    return 0;
194 }
```

4、mount操作

创建新的root，因为ss默认都和默认root(cgrp_dfl_root)建立了关系，所以ss需要先解除旧的root链接，再和新root建立起链接。

```
1  static struct dentry *cgroup_mount(struct file_system_type *fs_type,
2          int flags, const char *unused_dev_name,
3          void *data)
4  {
5    struct super_block *pinned_sb = NULL;
6    struct cgroup_subsys *ss;
7    struct cgroup_root *root;
8    struct cgroup_sb_opts opts;
9    struct dentry *dentry;
10   int ret;
11   int i;
12   bool new_sb;
13
14   /*
15    * The first time anyone tries to mount a cgroup, enable the list
16    * linking each css_set to its tasks and fix up all existing tasks.
17    */
18   if (!use_task_css_set_links)
19     cgroup_enable_task_cg_lists();
20
```

```
21    mutex_lock(&cgroup_mutex);
22
23    /* First find the desired set of subsystems */
24    /* (1) 解析mount选项到opts */
25    ret = parse_cgroupfs_options(data, &opts);
26    if (ret)
27      goto out_unlock;
28
29    /* look for a matching existing root */
30    if (opts.flags & CGRP_ROOT_SANE_BEHAVIOR) {
31      cgrp_dfl_root_visible = true;
32      root = &cgrp_dfl_root;
33      cgroup_get(&root->cgrp);
34      ret = 0;
35      goto out_unlock;
36    }
37
38    /*
39     * Destruction of cgroup root is asynchronous, so subsystems may
40     * still be dying after the previous unmount.  Let's drain the
41     * dying subsystems.  We just need to ensure that the ones
42     * unmounted previously finish dying and don't care about new ones
43     * starting.  Testing ref liveliness is good enough.
44     */
45    /* (2) */
46    for_each_subsys(ss, i) {
47      if (!(opts.subsys_mask & (1 << i)) ||
48          ss->root == &cgrp_dfl_root)
49        continue;
50
51      if (!percpu_ref_tryget_live(&ss->root->cgrp.self.refcnt)) {
52        mutex_unlock(&cgroup_mutex);
53        msleep(10);
54        ret = restart_syscall();
55        goto out_free;
56      }
57      cgroup_put(&ss->root->cgrp);
58    }
59
60      /* (3) */
```

```
61    for_each_root(root) {
62      bool name_match = false;
63
64      if (root == &cgrp_dfl_root)
65        continue;
66
67      /*
68       * If we asked for a name then it must match.  Also, if
69       * name matches but sybsys_mask doesn't, we should fail.
70       * Remember whether name matched.
71       */
72      if (opts.name) {
73        if (strcmp(opts.name, root->name))
74          continue;
75        name_match = true;
76      }
77
78      /*
79       * If we asked for subsystems (or explicitly for no
80       * subsystems) then they must match.
81       */
82      if ((opts.subsys_mask || opts.none) &&
83          (opts.subsys_mask != root->subsys_mask)) {
84        if (!name_match)
85          continue;
86        ret = -EBUSY;
87        goto out_unlock;
88      }
89
90      if (root->flags ^ opts.flags)
91        pr_warn("new mount options do not match the existing superblock, will
92
93      /*
94       * We want to reuse @root whose lifetime is governed by its
95       * ->cgrp.  Let's check whether @root is alive and keep it
96       * that way.  As cgroup_kill_sb() can happen anytime, we
97       * want to block it by pinning the sb so that @root doesn't
98       * get killed before mount is complete.
99       *
100      * With the sb pinned, tryget_live can reliably indicate
```

```c
     * whether @root can be reused.  If it's being killed,
     * drain it.  We can use wait_queue for the wait but this
     * path is super cold.  Let's just sleep a bit and retry.
     */
    pinned_sb = kernfs_pin_sb(root->kf_root, NULL);
    if (IS_ERR(pinned_sb) ||
        !percpu_ref_tryget_live(&root->cgrp.self.refcnt)) {
      mutex_unlock(&cgroup_mutex);
      if (!IS_ERR_OR_NULL(pinned_sb))
        deactivate_super(pinned_sb);
      msleep(10);
      ret = restart_syscall();
      goto out_free;
    }

    ret = 0;
    goto out_unlock;
  }

  /*
   * No such thing, create a new one.  name= matching without subsys
   * specification is allowed for already existing hierarchies but we
   * can't create new one without subsys specification.
   */
  if (!opts.subsys_mask && !opts.none) {
    ret = -EINVAL;
    goto out_unlock;
  }

  /* (4) 分配新的root */
  root = kzalloc(sizeof(*root), GFP_KERNEL);
  if (!root) {
    ret = -ENOMEM;
    goto out_unlock;
  }

  /* (5) 初始化新的root，初始了
     root->cgrp           // cgrp->root = root;
     root->cgrp.self      // cgrp->self.cgroup = cgrp; cgrp->self.flags |=
     root->name = opts->name
```

```
141        */
142      init_cgroup_root(root, &opts);
143
144        /* (6) 将新的root和opts.subsys_mask指向的多个ss进行链接 */
145      ret = cgroup_setup_root(root, opts.subsys_mask);
146      if (ret)
147        cgroup_free_root(root);
148
149  out_unlock:
150      mutex_unlock(&cgroup_mutex);
151  out_free:
152      kfree(opts.release_agent);
153      kfree(opts.name);
154
155      if (ret)
156        return ERR_PTR(ret);
157
158        /* (7) mount新root对应的根目录 */
159      dentry = kernfs_mount(fs_type, flags, root->kf_root,
160            CGROUP_SUPER_MAGIC, &new_sb);
161      if (IS_ERR(dentry) || !new_sb)
162        cgroup_put(&root->cgrp);
163
164      /*
165       * If @pinned_sb, we're reusing an existing root and holding an
166       * extra ref on its sb.  Mount is complete.  Put the extra ref.
167       */
168      if (pinned_sb) {
169        WARN_ON(new_sb);
170        deactivate_super(pinned_sb);
171      }
172
173      return dentry;
174  }
175
176  |→
177
178  static int cgroup_setup_root(struct cgroup_root *root, unsigned long ss_mask
179  {
180      LIST_HEAD(tmp_links);
```

```c
    struct cgroup *root_cgrp = &root->cgrp;
    struct css_set *cset;
    int i, ret;

    lockdep_assert_held(&cgroup_mutex);

    ret = cgroup_idr_alloc(&root->cgroup_idr, root_cgrp, 1, 2, GFP_KERNEL);
    if (ret < 0)
      goto out;
    root_cgrp->id = ret;

    ret = percpu_ref_init(&root_cgrp->self.refcnt, css_release, 0,
              GFP_KERNEL);
    if (ret)
      goto out;

    /*
     * We're accessing css_set_count without locking css_set_lock here,
     * but that's OK - it can only be increased by someone holding
     * cgroup_lock, and that's us. The worst that can happen is that we
     * have some link structures left over
     */
    ret = allocate_cgrp_cset_links(css_set_count, &tmp_links);
    if (ret)
      goto cancel_ref;

    ret = cgroup_init_root_id(root);
    if (ret)
      goto cancel_ref;

      /* (6.1) 创建root对应的顶层root文件夹 */
    root->kf_root = kernfs_create_root(&cgroup_kf_syscall_ops,
              KERNFS_ROOT_CREATE_DEACTIVATED,
              root_cgrp);
    if (IS_ERR(root->kf_root)) {
      ret = PTR_ERR(root->kf_root);
      goto exit_root_id;
    }
    root_cgrp->kn = root->kf_root->kn;

```

```c
    /* (6.2) 创建cgroup自己对应的一些file，cgroup自己的file由cgroup自己的css(cgr
         后面cgroup会依次创建每个subsys的file，subsys的file由每个ss对应的css(cgrp
     */
    ret = css_populate_dir(&root_cgrp->self, NULL);
    if (ret)
      goto destroy_root;

    /* (6.3) 将新root需要的subsys和原默认root(cgrp_dfl_root)解除关系，
         并且把这些ss重新和新root建立关系
     */
    ret = rebind_subsystems(root, ss_mask);
    if (ret)
      goto destroy_root;

    /*
     * There must be no failure case after here, since rebinding takes
     * care of subsystems' refcounts, which are explicitly dropped in
     * the failure exit path.
     */
    list_add(&root->root_list, &cgroup_roots);
    cgroup_root_count++;

    /*
     * Link the root cgroup in this hierarchy into all the css_set
     * objects.
     */
    spin_lock_bh(&css_set_lock);
    hash_for_each(css_set_table, i, cset, hlist) {
      link_css_set(&tmp_links, cset, root_cgrp);
      if (css_set_populated(cset))
        cgroup_update_populated(root_cgrp, true);
    }
    spin_unlock_bh(&css_set_lock);

    BUG_ON(!list_empty(&root_cgrp->self.children));
    BUG_ON(atomic_read(&root->nr_cgrps) != 1);

    kernfs_activate(root_cgrp->kn);
    ret = 0;
    goto out;
```

```
261
262  destroy_root:
263      kernfs_destroy_root(root->kf_root);
264      root->kf_root = NULL;
265  exit_root_id:
266      cgroup_exit_root_id(root);
267  cancel_ref:
268      percpu_ref_exit(&root_cgrp->self.refcnt);
269  out:
270      free_cgrp_cset_links(&tmp_links);
271      return ret;
272  }
273
274  ||→
275
276  static int rebind_subsystems(struct cgroup_root *dst_root,
277              unsigned long ss_mask)
278  {
279      struct cgroup *dcgrp = &dst_root->cgrp;
280      struct cgroup_subsys *ss;
281      unsigned long tmp_ss_mask;
282      int ssid, i, ret;
283
284      lockdep_assert_held(&cgroup_mutex);
285
286      for_each_subsys_which(ss, ssid, &ss_mask) {
287          /* if @ss has non-root csses attached to it, can't move */
288          if (css_next_child(NULL, cgroup_css(&ss->root->cgrp, ss)))
289              return -EBUSY;
290
291          /* can't move between two non-dummy roots either */
292          if (ss->root != &cgrp_dfl_root && dst_root != &cgrp_dfl_root)
293              return -EBUSY;
294      }
295
296      /* skip creating root files on dfl_root for inhibited subsystems */
297      tmp_ss_mask = ss_mask;
298      if (dst_root == &cgrp_dfl_root)
299          tmp_ss_mask &= ~cgrp_dfl_root_inhibit_ss_mask;
300
```

```
301    for_each_subsys_which(ss, ssid, &tmp_ss_mask) {
302      struct cgroup *scgrp = &ss->root->cgrp;
303      int tssid;
304
305          /* (6.3.1) 在新root的根cgroup(dst_root->cgrp)下，
306              根据subsys的file链表(css->ss->cfts)创建subsys对应的file
307          */
308      ret = css_populate_dir(cgroup_css(scgrp, ss), dcgrp);
309      if (!ret)
310        continue;
311
312      /*
313       * Rebinding back to the default root is not allowed to
314       * fail.  Using both default and non-default roots should
315       * be rare.  Moving subsystems back and forth even more so.
316       * Just warn about it and continue.
317       */
318      if (dst_root == &cgrp_dfl_root) {
319        if (cgrp_dfl_root_visible) {
320          pr_warn("failed to create files (%d) while rebinding 0x%lx to defaul
321            ret, ss_mask);
322          pr_warn("you may retry by moving them to a different hierarchy and u
323        }
324        continue;
325      }
326
327      for_each_subsys_which(ss, tssid, &tmp_ss_mask) {
328        if (tssid == ssid)
329          break;
330        css_clear_dir(cgroup_css(scgrp, ss), dcgrp);
331      }
332      return ret;
333    }
334
335    /*
336     * Nothing can fail from this point on.  Remove files for the
337     * removed subsystems and rebind each subsystem.
338     */
339    for_each_subsys_which(ss, ssid, &ss_mask) {
340      struct cgroup_root *src_root = ss->root;
```

```c
    struct cgroup *scgrp = &src_root->cgrp;
    struct cgroup_subsys_state *css = cgroup_css(scgrp, ss);
    struct css_set *cset;

    WARN_ON(!css || cgroup_css(dcgrp, ss));

    css_clear_dir(css, NULL);

        /* (6.3.2) 取消原root cgroup对subsys的css的引用 */
    RCU_INIT_POINTER(scgrp->subsys[ssid], NULL);

    /* (6.3.3) 链接新root cgroup和subsys的css的引用 */
    rcu_assign_pointer(dcgrp->subsys[ssid], css);
    ss->root = dst_root;
    css->cgroup = dcgrp;

    spin_lock_bh(&css_set_lock);
    hash_for_each(css_set_table, i, cset, hlist)
      list_move_tail(&cset->e_cset_node[ss->id],
                &dcgrp->e_csets[ss->id]);
    spin_unlock_bh(&css_set_lock);

    src_root->subsys_mask &= ~(1 << ssid);
    scgrp->subtree_control &= ~(1 << ssid);
    cgroup_refresh_child_subsys_mask(scgrp);

    /* default hierarchy doesn't enable controllers by default */
    dst_root->subsys_mask |= 1 << ssid;
    if (dst_root == &cgrp_dfl_root) {
      static_branch_enable(cgroup_subsys_on_dfl_key[ssid]);
    } else {
      dcgrp->subtree_control |= 1 << ssid;
      cgroup_refresh_child_subsys_mask(dcgrp);
      static_branch_disable(cgroup_subsys_on_dfl_key[ssid]);
    }

    if (ss->bind)
      ss->bind(css);
}
```

```
381    kernfs_activate(dcgrp->kn);
382    return 0;
383 }
```

## 5、文件操作

创建一个新文件夹，相当于创建一个新的cgroup。我们重点来看看新建文件夹的操作：

```
1  static struct kernfs_syscall_ops cgroup_kf_syscall_ops = {
2    .remount_fs    = cgroup_remount,
3    .show_options    = cgroup_show_options,
4    .mkdir       = cgroup_mkdir,
5    .rmdir       = cgroup_rmdir,
6    .rename       = cgroup_rename,
7  };
8
9  static int cgroup_mkdir(struct kernfs_node *parent_kn, const char *name,
10        umode_t mode)
11 {
12    struct cgroup *parent, *cgrp;
13    struct cgroup_root *root;
14    struct cgroup_subsys *ss;
15    struct kernfs_node *kn;
16    int ssid, ret;
17
18    /* Do not accept '\n' to prevent making /proc/<pid>/cgroup unparsable.
19     */
20    if (strchr(name, '\n'))
21      return -EINVAL;
22
23    parent = cgroup_kn_lock_live(parent_kn);
24    if (!parent)
25      return -ENODEV;
26    root = parent->root;
27
28    /* allocate the cgroup and its ID, 0 is reserved for the root */
29    /* (1) 分配新的cgroup */
```

```c
30    cgrp = kzalloc(sizeof(*cgrp), GFP_KERNEL);
31    if (!cgrp) {
32      ret = -ENOMEM;
33      goto out_unlock;
34    }
35
36    ret = percpu_ref_init(&cgrp->self.refcnt, css_release, 0, GFP_KERNEL);
37    if (ret)
38      goto out_free_cgrp;
39
40    /*
41     * Temporarily set the pointer to NULL, so idr_find() won't return
42     * a half-baked cgroup.
43     */
44    cgrp->id = cgroup_idr_alloc(&root->cgroup_idr, NULL, 2, 0, GFP_KERNEL);
45    if (cgrp->id < 0) {
46      ret = -ENOMEM;
47      goto out_cancel_ref;
48    }
49
50    /* (2) 初始化cgroup */
51    init_cgroup_housekeeping(cgrp);
52
53    /* (3) 和父cgroup之间建立起关系 */
54    cgrp->self.parent = &parent->self;
55    cgrp->root = root;
56
57    if (notify_on_release(parent))
58      set_bit(CGRP_NOTIFY_ON_RELEASE, &cgrp->flags);
59
60    if (test_bit(CGRP_CPUSET_CLONE_CHILDREN, &parent->flags))
61      set_bit(CGRP_CPUSET_CLONE_CHILDREN, &cgrp->flags);
62
63    /* create the directory */
64    /* (3) 创建新的cgroup对应的文件夹 */
65    kn = kernfs_create_dir(parent->kn, name, mode, cgrp);
66    if (IS_ERR(kn)) {
67      ret = PTR_ERR(kn);
68      goto out_free_id;
69    }
```

```c
70    cgrp->kn = kn;
71
72    /*
73     * This extra ref will be put in cgroup_free_fn() and guarantees
74     * that @cgrp->kn is always accessible.
75     */
76    kernfs_get(kn);
77
78    cgrp->self.serial_nr = css_serial_nr_next++;
79
80    /* allocation complete, commit to creation */
81    list_add_tail_rcu(&cgrp->self.sibling, &cgroup_parent(cgrp)->self.children
82    atomic_inc(&root->nr_cgrps);
83    cgroup_get(parent);
84
85    /*
86     * @cgrp is now fully operational.  If something fails after this
87     * point, it'll be released via the normal destruction path.
88     */
89    cgroup_idr_replace(&root->cgroup_idr, cgrp, cgrp->id);
90
91    ret = cgroup_kn_set_ugid(kn);
92    if (ret)
93      goto out_destroy;
94
95     /* (4) 新cgroup文件夹下创建cgroup自己css对应的默认file */
96    ret = css_populate_dir(&cgrp->self, NULL);
97    if (ret)
98      goto out_destroy;
99
100   /* let's create and online css's */
101   /* (5) 针对root对应的各个susbsys，每个subsys创建新的css
102       并且在cgroup文件夹下创建css对应的file
103   */
104   for_each_subsys(ss, ssid) {
105     if (parent->child_subsys_mask & (1 << ssid)) {
106       ret = create_css(cgrp, ss,
107             parent->subtree_control & (1 << ssid));
108       if (ret)
109         goto out_destroy;
```

```c
      }
   }

   /*
    * On the default hierarchy, a child doesn't automatically inherit
    * subtree_control from the parent.  Each is configured manually.
    */
   if (!cgroup_on_dfl(cgrp)) {
     cgrp->subtree_control = parent->subtree_control;
     cgroup_refresh_child_subsys_mask(cgrp);
   }

   kernfs_activate(kn);

   ret = 0;
   goto out_unlock;

out_free_id:
   cgroup_idr_remove(&root->cgroup_idr, cgrp->id);
out_cancel_ref:
   percpu_ref_exit(&cgrp->self.refcnt);
out_free_cgrp:
   kfree(cgrp);
out_unlock:
   cgroup_kn_unlock(parent_kn);
   return ret;

out_destroy:
   cgroup_destroy_locked(cgrp);
   goto out_unlock;
}
```

cgroup默认文件，有一些重要的文件比如"tasks"，我们来看看具体的操作。

```c
static struct cftype cgroup_legacy_base_files[] = {
   {
     .name = "cgroup.procs",
```

```c
        .seq_start = cgroup_pidlist_start,
        .seq_next = cgroup_pidlist_next,
        .seq_stop = cgroup_pidlist_stop,
        .seq_show = cgroup_pidlist_show,
        .private = CGROUP_FILE_PROCS,
        .write = cgroup_procs_write,
    },
    {
        .name = "cgroup.clone_children",
        .read_u64 = cgroup_clone_children_read,
        .write_u64 = cgroup_clone_children_write,
    },
    {
        .name = "cgroup.sane_behavior",
        .flags = CFTYPE_ONLY_ON_ROOT,
        .seq_show = cgroup_sane_behavior_show,
    },
    {
        .name = "tasks",
        .seq_start = cgroup_pidlist_start,
        .seq_next = cgroup_pidlist_next,
        .seq_stop = cgroup_pidlist_stop,
        .seq_show = cgroup_pidlist_show,
        .private = CGROUP_FILE_TASKS,
        .write = cgroup_tasks_write,
    },
    {
        .name = "notify_on_release",
        .read_u64 = cgroup_read_notify_on_release,
        .write_u64 = cgroup_write_notify_on_release,
    },
    {
        .name = "release_agent",
        .flags = CFTYPE_ONLY_ON_ROOT,
        .seq_show = cgroup_release_agent_show,
        .write = cgroup_release_agent_write,
        .max_write_len = PATH_MAX - 1,
    },
    { }  /* terminate */
}
```

```
44
45  static ssize_t cgroup_tasks_write(struct kernfs_open_file *of,
46              char *buf, size_t nbytes, loff_t off)
47  {
48    return __cgroup_procs_write(of, buf, nbytes, off, false);
49  }
50
51  |→
52
53  static ssize_t __cgroup_procs_write(struct kernfs_open_file *of, char *buf,
54              size_t nbytes, loff_t off, bool threadgroup)
55  {
56    struct task_struct *tsk;
57    struct cgroup_subsys *ss;
58    struct cgroup *cgrp;
59    pid_t pid;
60    int ssid, ret;
61
62    if (kstrtoint(strstrip(buf), 0, &pid) || pid < 0)
63      return -EINVAL;
64
65    cgrp = cgroup_kn_lock_live(of->kn);
66    if (!cgrp)
67      return -ENODEV;
68
69    percpu_down_write(&cgroup_threadgroup_rwsem);
70    rcu_read_lock();
71    if (pid) {
72      tsk = find_task_by_vpid(pid);
73      if (!tsk) {
74        ret = -ESRCH;
75        goto out_unlock_rcu;
76      }
77    } else {
78      tsk = current;
79    }
80
81    if (threadgroup)
82      tsk = tsk->group_leader;
83
```

```
 84    /*
 85     * Workqueue threads may acquire PF_NO_SETAFFINITY and become
 86     * trapped in a cpuset, or RT worker may be born in a cgroup
 87     * with no rt_runtime allocated.  Just say no.
 88     */
 89    if (tsk == kthreadd_task || (tsk->flags & PF_NO_SETAFFINITY)) {
 90      ret = -EINVAL;
 91      goto out_unlock_rcu;
 92    }
 93
 94    get_task_struct(tsk);
 95    rcu_read_unlock();
 96
 97    ret = cgroup_procs_write_permission(tsk, cgrp, of);
 98    if (!ret) {
 99        /* (1) attach task到cgroup */
100      ret = cgroup_attach_task(cgrp, tsk, threadgroup);
101 #if defined(CONFIG_CPUSETS) && !defined(CONFIG_MTK_ACAO)
102      if (cgrp->id != SS_TOP_GROUP_ID && cgrp->child_subsys_mask == CSS_CPUSET
103      && excl_task_count > 0) {
104        remove_set_exclusive_task(tsk->pid, 0);
105      }
106 #endif
107    }
108    put_task_struct(tsk);
109    goto out_unlock_threadgroup;
110
111 out_unlock_rcu:
112    rcu_read_unlock();
113 out_unlock_threadgroup:
114    percpu_up_write(&cgroup_threadgroup_rwsem);
115    for_each_subsys(ss, ssid)
116      if (ss->post_attach)
117        ss->post_attach();
118    cgroup_kn_unlock(of->kn);
119    return ret ?: nbytes;
120 }
121
122 ||→
123
```

```c
124  static int cgroup_attach_task(struct cgroup *dst_cgrp,
125                  struct task_struct *leader, bool threadgroup)
126  {
127    LIST_HEAD(preloaded_csets);
128    struct task_struct *task;
129    int ret;
130
131    /* look up all src csets */
132    spin_lock_bh(&css_set_lock);
133    rcu_read_lock();
134    task = leader;
135
136    /* (1.1) 遍历task所在线程组，把需要迁移的进程的css_set加入到preloaded_csets链表
137    do {
138      cgroup_migrate_add_src(task_css_set(task), dst_cgrp,
139                  &preloaded_csets);
140      if (!threadgroup)
141        break;
142    } while_each_thread(leader, task);
143    rcu_read_unlock();
144    spin_unlock_bh(&css_set_lock);
145
146    /* (1.2) 去掉旧的css_set对css的应用，
147        分配新的css_set承担新的css组合的应用，并且给进程使用
148     */
149    /* prepare dst csets and commit */
150    ret = cgroup_migrate_prepare_dst(dst_cgrp, &preloaded_csets);
151    if (!ret)
152      ret = cgroup_migrate(leader, threadgroup, dst_cgrp);
153
154    cgroup_migrate_finish(&preloaded_csets);
155    return ret;
156  }
```

# 1.3、cgroup subsystem

我们关注cgroup子系统具体能提供的功能。

## 1.3.1、cpu

kernel/sched/core.c。会创建新的task_group，可以对cgroup对应的task_group进行cfs/rt类型的带宽控制。

```c
static struct cftype cpu_files[] = {
#ifdef CONFIG_FAIR_GROUP_SCHED
  {
    .name = "shares",
    .read_u64 = cpu_shares_read_u64,
    .write_u64 = cpu_shares_write_u64,
  },
#endif
#ifdef CONFIG_CFS_BANDWIDTH    // cfs 带宽控制
  {
    .name = "cfs_quota_us",
    .read_s64 = cpu_cfs_quota_read_s64,
    .write_s64 = cpu_cfs_quota_write_s64,
  },
  {
    .name = "cfs_period_us",
    .read_u64 = cpu_cfs_period_read_u64,
    .write_u64 = cpu_cfs_period_write_u64,
  },
  {
    .name = "stat",
    .seq_show = cpu_stats_show,
  },
#endif
#ifdef CONFIG_RT_GROUP_SCHED    // rt 带宽控制
  {
    .name = "rt_runtime_us",
    .read_s64 = cpu_rt_runtime_read,
    .write_s64 = cpu_rt_runtime_write,
  },
  {
    .name = "rt_period_us",
    .read_u64 = cpu_rt_period_read_uint,
```

```
34    .write_u64 = cpu_rt_period_write_uint,
35   },
36 #endif
37   { }  /* terminate */
38 };
39
40 struct cgroup_subsys cpu_cgrp_subsys = {
41   .css_alloc  = cpu_cgroup_css_alloc,         // 分配新的task_group
42   .css_released  = cpu_cgroup_css_released,
43   .css_free  = cpu_cgroup_css_free,
44   .fork     = cpu_cgroup_fork,
45   .can_attach  = cpu_cgroup_can_attach,
46   .attach    = cpu_cgroup_attach,
47   .legacy_cftypes  = cpu_files,
48   .early_init  = 1,
49 };
```

## 1.3.2、cpuset

kernel/cpusec.c。给cgroup分配不同的cpu和mem node节点，还可以配置一些flag。

```
1  static struct cftype files[] = {
2    {
3      .name = "cpus",
4      .seq_show = cpuset_common_seq_show,
5      .write = cpuset_write_resmask,
6      .max_write_len = (100U + 6 * NR_CPUS),
7      .private = FILE_CPULIST,
8    },
9
10   {
11     .name = "mems",
12     .seq_show = cpuset_common_seq_show,
13     .write = cpuset_write_resmask,
14     .max_write_len = (100U + 6 * MAX_NUMNODES),
15     .private = FILE_MEMLIST,
16   },
```

```c
	{
		.name = "effective_cpus",
		.seq_show = cpuset_common_seq_show,
		.private = FILE_EFFECTIVE_CPULIST,
	},

	{
		.name = "effective_mems",
		.seq_show = cpuset_common_seq_show,
		.private = FILE_EFFECTIVE_MEMLIST,
	},

	{
		.name = "cpu_exclusive",
		.read_u64 = cpuset_read_u64,
		.write_u64 = cpuset_write_u64,
		.private = FILE_CPU_EXCLUSIVE,
	},

	{
		.name = "mem_exclusive",
		.read_u64 = cpuset_read_u64,
		.write_u64 = cpuset_write_u64,
		.private = FILE_MEM_EXCLUSIVE,
	},

	{
		.name = "mem_hardwall",
		.read_u64 = cpuset_read_u64,
		.write_u64 = cpuset_write_u64,
		.private = FILE_MEM_HARDWALL,
	},

	{
		.name = "sched_load_balance",
		.read_u64 = cpuset_read_u64,
		.write_u64 = cpuset_write_u64,
		.private = FILE_SCHED_LOAD_BALANCE,
	},
```

```c
	{
		.name = "sched_relax_domain_level",
		.read_s64 = cpuset_read_s64,
		.write_s64 = cpuset_write_s64,
		.private = FILE_SCHED_RELAX_DOMAIN_LEVEL,
	},

	{
		.name = "memory_migrate",
		.read_u64 = cpuset_read_u64,
		.write_u64 = cpuset_write_u64,
		.private = FILE_MEMORY_MIGRATE,
	},

	{
		.name = "memory_pressure",
		.read_u64 = cpuset_read_u64,
	},

	{
		.name = "memory_spread_page",
		.read_u64 = cpuset_read_u64,
		.write_u64 = cpuset_write_u64,
		.private = FILE_SPREAD_PAGE,
	},

	{
		.name = "memory_spread_slab",
		.read_u64 = cpuset_read_u64,
		.write_u64 = cpuset_write_u64,
		.private = FILE_SPREAD_SLAB,
	},

	{
		.name = "memory_pressure_enabled",
		.flags = CFTYPE_ONLY_ON_ROOT,
		.read_u64 = cpuset_read_u64,
		.write_u64 = cpuset_write_u64,
		.private = FILE_MEMORY_PRESSURE_ENABLED,
```

```
 97      },
 98
 99      { }  /* terminate */
100  }
101
102  struct cgroup_subsys cpuset_cgrp_subsys = {
103      .css_alloc  = cpuset_css_alloc,
104      .css_online  = cpuset_css_online,
105      .css_offline  = cpuset_css_offline,
106      .css_free  = cpuset_css_free,
107      .can_attach  = cpuset_can_attach,
108      .cancel_attach  = cpuset_cancel_attach,
109      .attach     = cpuset_attach,
110      .post_attach  = cpuset_post_attach,
111      .bind     = cpuset_bind,
112      .fork     = cpuset_fork,
113      .legacy_cftypes  = files,
114      .early_init  = 1,
115  };
```

### 1.3.3、schedtune

kernel/sched/tune.c，可以进行schedle boost操作。

```
 1  static struct cftype files[] = {
 2    {
 3      .name = "boost",
 4      .read_u64 = boost_read,
 5      .write_u64 = boost_write,
 6    },
 7    {
 8      .name = "prefer_idle",
 9      .read_u64 = prefer_idle_read,
10      .write_u64 = prefer_idle_write,
11    },
12    { }  /* terminate */
13  };
```

```
14
15  struct cgroup_subsys schedtune_cgrp_subsys = {
16    .css_alloc  = schedtune_css_alloc,
17    .css_free = schedtune_css_free,
18    .legacy_cftypes  = files,
19    .early_init  = 1,
20  };
```

### 1.3.4、cpuacct

kernel/sched/cpuacct.c，可以按照cgroup的分组来统计cpu占用率。

```
1   static struct cftype files[] = {
2     {
3       .name = "usage",
4       .read_u64 = cpuusage_read,
5       .write_u64 = cpuusage_write,
6     },
7     {
8       .name = "usage_percpu",
9       .seq_show = cpuacct_percpu_seq_show,
10    },
11    {
12      .name = "stat",
13      .seq_show = cpuacct_stats_show,
14    },
15    { }  /* terminate */
16  };
17
18  struct cgroup_subsys cpuacct_cgrp_subsys = {
19    .css_alloc  = cpuacct_css_alloc,
20    .css_free  = cpuacct_css_free,
21    .legacy_cftypes  = files,
22    .early_init  = 1,
23  };
```

参考资料

1、linux 2.6 O(1)调度算法
2、linux cfs调度器_理论模型
3、linux cfs调度框图
4、linux cfs之特殊时刻vruntime的计算
5、entity级负载的计算
6、cpu级负载的计算update_cpu_load
7、系统级负载的计算:Linux Load Averages: Solving the Mystery
8、系统级负载的计算:UNIX Load Average
9、Linux Scheduling Domains
10、[MTK文档：CPU Utilization-scheduler(V1.1)]
11、Docker背后的内核知识——cgroups资源限制
12、Linux资源管理之cgroups简介

**往期课程可扫以下二维码试听与购买**



微信 Linux阅码场