

阅读新闻

背景: □□□□□□□□

Linux内核部件分析

设备驱动模型之device

[日期: 2011-10-06]来源: Linux社区 作者: qb_2008[字体: 大 中 小]

linux的设备驱动模型，是建立在sysfs和kobject之上的，由总线、设备、驱动、类所组成的关系结构。从本节开始，我们将对linux这一设备驱动模型进行深入分析。

头文件是include/linux/device.h，实现在drivers/base目录中。本节要分析的，是其中的设备，主要在core.c中。

1. struct device {

2. struct device *parent;

3.

4. struct device_private *p;

5.

6. struct kobject kobj;

7. const char *init_name; /* initial name of the device */

8. struct device_type *type;

9.

10. struct semaphore sem; /* semaphore to synchronize calls to

11. * its driver.

12. */

13.

14. struct bus_type *bus; /* type of bus device is on */

15. struct device_driver *driver; /* which driver has allocated this

16. device */

17. void *platform_data; /* Platform specific data, device

18. core doesn't touch it */

19. struct dev_pm_info power;

20.

21. #ifdef CONFIG_NUMA

22. int numa_node; /* NUMA node this device is close to */

23. #endif

24. u64 *dma_mask; /* dma mask (if dma'able device) */

25. u64 coherent_dma_mask; /* Like dma_mask, but for

26. alloc_coherent mappings as

27. not all hardware supports

28. 64 bit addresses for consistent

29. allocations such descriptors. */

30.

31. struct device_dma_parameters *dma_parms;

32.

33. struct list_head dma_pools; /* dma pools (if dma'ble) */

34.

35. struct dma_coherent_mem *dma_mem; /* internal for coherent mem

36. override */

37. /* arch specific additions */

38. struct dev_archdata archdata;

39.

40. dev_t devt; /* dev_t, creates the sysfs "dev" */

41.

42. spinlock_t devres_lock;

43. struct list_head devres_head;

44.

45. struct klist_node knode_class;

46. struct class *class;

47. const struct attribute_group **groups; /* optional groups */

48.

49. void (*release)(struct device *dev);

50. };

先来分析下struct device的结构变量。首先是指向父节点的指针parent，kobj是内嵌在device中的kobject，用于把它联系到sysfs中。bus是对设备所在总线的指针，driver是对设备所用驱动的指针。还有DMA需要的数据，表示设备号的devt，表示设备资源的devres_head和保护它的devres_lock。指向类的指针class，knode_e_class是被连入class链表时所用的klist节点。group是设备的属性集合。release应该是设备释放时调用的函数。

```

1. struct device_private {
2.     struct klist klist_children;
3.     struct klist_node knode_parent;
4.     struct klist_node knode_driver;
5.     struct klist_node knode_bus;
6.     void *driver_data;
7.     struct device *device;
8. };
9. #define to_device_private_parent(obj) \
10.     container_of(obj, struct device_private, knode_parent)
11. #define to_device_private_driver(obj) \
12.     container_of(obj, struct device_private, knode_driver)
13. #define to_device_private_bus(obj) \
14.     container_of(obj, struct device_private, knode_bus)

```

struct device中有一部分不愿意让外界看到，所以做出struct device_private结构，包括了设备驱动模型内部的链接。klist_children是子设备的链表，knode_parent是连入父设备的klist_children时所用的节点，knode_driver是连入驱动的设备链表所用的节点，knode_bus是连入总线的设备链表时所用的节点。driver_data用于在设备结构中存放相关的驱动信息，也许是驱动专门为设备建立的结构实例。device则是指向struct device_private所属的device。

下面还有一些宏，to_device_private_parent()是从父设备的klist_children上节点，获得相应的device_private。to_device_private_driver()是从驱动的设备链表上节点，获得对应的device_private。to_device_private_bus()是从总线的设备链表上节点，获得对应的device_private。

或许会奇怪，为什么knode_class没有被移入struct device_private，或许有外部模块需要用到它。

```

1. /*
2.  * The type of device, "struct device" is embedded in. A class
3.  * or bus can contain devices of different types
4.  * like "partitions" and "disks", "mouse" and "event".
5.  * This identifies the device type and carries type-specific
6.  * information, equivalent to the kobj_type of a kobject.
7.  * If "name" is specified, the uevent will contain it in
8.  * the DEVTTYPE variable.
9.  */
10. struct device_type {
11.     const char *name;
12.     const struct attribute_group **groups;
13.     int (*uevent)(struct device *dev, struct kobj_uevent_env *env);
14.     char *(*devnode)(struct device *dev, mode_t *mode);
15.     void (*release)(struct device *dev);
16.
17.     const struct dev_pm_ops *pm;
18. };

```

device竟然有device_type，类似于与kobject相对的kobj_type，之后我们再看它怎么用。

```

1. /* interface for exporting device attributes */
2. struct device_attribute {
3.     struct attribute attr;
4.     ssize_t (*show)(struct device *dev, struct device_attribute *attr,
5.                     char *buf);
6.     ssize_t (*store)(struct device *dev, struct device_attribute *attr,
7.                     const char *buf, size_t count);
8. };
9.
10. #define DEVICE_ATTR(_name, _mode, _show, _store) \
11.     struct device_attribute dev_attr_##_name = __ATTR(_name, _mode, _show, _store)

```

这个device_attribute显然就是device对struct attribute的封装，新加的show()、store()函数都是以与设备相关的结构调用的。

至于device中其它的archdata、dma、devres，都是作为设备特有的，我们现在主要关心设备驱动模型的建立，这些会尽量忽略。

下面来看看device的实现，这主要在core.c中。

```

1. int __init devices_init(void)
2. {
3.     devices_kset = kset_create_and_add("devices", &device_uevent_ops, NULL);
4.     if (!devices_kset)
5.         return -ENOMEM;
6.     dev_kobj = kobject_create_and_add("dev", NULL);
7.     if (!dev_kobj)

```

```

8.     goto dev_kobj_err;
9.     sysfs_dev_block_kobj = kobject_create_and_add("block", dev_kobj);
10.    if (!sysfs_dev_block_kobj)
11.        goto block_kobj_err;
12.    sysfs_dev_char_kobj = kobject_create_and_add("char", dev_kobj);
13.    if (!sysfs_dev_char_kobj)
14.        goto char_kobj_err;
15.
16.    return 0;
17.
18. char_kobj_err:
19.     kobject_put(sysfs_dev_block_kobj);
20. block_kobj_err:
21.     kobject_put(dev_kobj);
22. dev_kobj_err:
23.     kset_unregister(devices_kset);
24.     return -ENOMEM;
25. }

```

这是在设备驱动模型初始化时调用的device部分初始的函数devices_init()。它干的事情我们都很熟悉，就是建立sysfs中的devices目录，和dev目录。还在dev目录下又建立了block和char两个子目录。因为dev目录只打算存放辅助的设备号，所以没必要使用kset。

```

1. static ssize_t dev_attr_show(struct kobject *kobj, struct attribute *attr,
2.     char *buf)
3. {
4.     struct device_attribute *dev_attr = to_dev_attr(attr);
5.     struct device *dev = to_dev(kobj);
6.     ssize_t ret = -EIO;
7.
8.     if (dev_attr->show)
9.         ret = dev_attr->show(dev, dev_attr, buf);
10.    if (ret >= (ssize_t)PAGE_SIZE) {
11.        print_symbol("dev_attr_show: %s returned bad count\n",
12.            (unsigned long)dev_attr->show);
13.    }
14.    return ret;
15. }
16.
17. static ssize_t dev_attr_store(struct kobject *kobj, struct attribute *attr,
18.     const char *buf, size_t count)
19. {
20.     struct device_attribute *dev_attr = to_dev_attr(attr);
21.     struct device *dev = to_dev(kobj);
22.     ssize_t ret = -EIO;
23.
24.     if (dev_attr->store)
25.         ret = dev_attr->store(dev, dev_attr, buf, count);
26.     return ret;
27. }
28.
29. static struct sysfs_ops dev_sysfs_ops = {
30.     .show = dev_attr_show,
31.     .store = dev_attr_store,
32. };

```

看到这里是不是很熟悉，dev_sysfs_ops就是device准备注册到sysfs中的操作函数。dev_attr_show()和dev_attr_store()都会再调用与属性相关的函数。

```

1. static void device_release(struct kobject *kobj)
2. {
3.     struct device *dev = to_dev(kobj);
4.     struct device_private *p = dev->p;
5.
6.     if (dev->release)
7.         dev->release(dev);
8.     else if (dev->type && dev->type->release)
9.         dev->type->release(dev);
10.    else if (dev->class && dev->class->dev_release)
11.        dev->class->dev_release(dev);
12.    else
13.        WARN(1, KERN_ERR "Device '%s' does not have a release() "
14.            "function, it is broken and must be fixed.\n",

```

```

15.     dev_name(dev));
16.     kfree(p);
17. }
18.
19. static struct kobj_type device_ktype = {
20.     .release    = device_release,
21.     .sysfs_ops  = &dev_sysfs_ops,
22. };

```

使用的release函数是device_release。在释放device时，会依次调用device结构中定义的release函数，device_type中定义的release函数，device所属的class中所定义的release函数，最后会吧device_private结构释放掉。

```

1. static int dev_uevent_filter(struct kset *kset, struct kobject *kobj)
2. {
3.     struct kobj_type *ktype = get_ktype(kobj);
4.
5.     if (ktype == &device_ktype) {
6.         struct device *dev = to_dev(kobj);
7.         if (dev->bus)
8.             return 1;
9.         if (dev->class)
10.            return 1;
11.    }
12.    return 0;
13. }
14.
15. static const char *dev_uevent_name(struct kset *kset, struct kobject *kobj)
16. {
17.     struct device *dev = to_dev(kobj);
18.
19.     if (dev->bus)
20.         return dev->bus->name;
21.     if (dev->class)
22.         return dev->class->name;
23.     return NULL;
24. }
25.
26. static int dev_uevent(struct kset *kset, struct kobject *kobj,
27.     struct kobj_uevent_env *env)
28. {
29.     struct device *dev = to_dev(kobj);
30.     int retval = 0;
31.
32.     /* add device node properties if present */
33.     if (MAJOR(dev->devt)) {
34.         const char *tmp;
35.         const char *name;
36.         mode_t mode = 0;
37.
38.         add_uevent_var(env, "MAJOR=%u", MAJOR(dev->devt));
39.         add_uevent_var(env, "MINOR=%u", MINOR(dev->devt));
40.         name = device_get_devnode(dev, &mode, &tmp);
41.         if (name) {
42.             add_uevent_var(env, "DEVNAME=%s", name);
43.             kfree(tmp);
44.             if (mode)
45.                 add_uevent_var(env, "DEVMODE=%#o", mode & 0777);
46.         }
47.     }
48.
49.     if (dev->type && dev->type->name)
50.         add_uevent_var(env, "DEVTYPE=%s", dev->type->name);
51.
52.     if (dev->driver)
53.         add_uevent_var(env, "DRIVER=%s", dev->driver->name);
54.
55. #ifdef CONFIG_SYSFS_DEPRECATED
56.     if (dev->class) {
57.         struct device *parent = dev->parent;
58.

```

```

59.  /* find first bus device in parent chain */
60.  while (parent && !parent->bus)
61.      parent = parent->parent;
62.  if (parent && parent->bus) {
63.      const char *path;
64.
65.      path = kobject_get_path(&parent->kobj, GFP_KERNEL);
66.      if (path) {
67.          add_uevent_var(env, "PHYSDEVPATH=%s", path);
68.          kfree(path);
69.      }
70.
71.      add_uevent_var(env, "PHYSDEVBUS=%s", parent->bus->name);
72.
73.      if (parent->driver)
74.          add_uevent_var(env, "PHYSDEVDRIVER=%s",
75.                          parent->driver->name);
76.  }
77.  } else if (dev->bus) {
78.      add_uevent_var(env, "PHYSDEVBUS=%s", dev->bus->name);
79.
80.      if (dev->driver)
81.          add_uevent_var(env, "PHYSDEVDRIVER=%s",
82.                          dev->driver->name);
83.  }
84. #endif
85.
86.  /* have the bus specific function add its stuff */
87.  if (dev->bus && dev->bus->uevent) {
88.      retval = dev->bus->uevent(dev, env);
89.      if (retval)
90.          pr_debug("device: '%s': %s: bus uevent() returned %d\n",
91.                  dev_name(dev), __func__, retval);
92.  }
93.
94.  /* have the class specific function add its stuff */
95.  if (dev->class && dev->class->dev_uevent) {
96.      retval = dev->class->dev_uevent(dev, env);
97.      if (retval)
98.          pr_debug("device: '%s': %s: class uevent() "
99.                  "returned %d\n", dev_name(dev),
100.                     __func__, retval);
101.  }
102.
103.  /* have the device type specific fuction add its stuff */
104.  if (dev->type && dev->type->uevent) {
105.      retval = dev->type->uevent(dev, env);
106.      if (retval)
107.          pr_debug("device: '%s': %s: dev_type uevent() "
108.                  "returned %d\n", dev_name(dev),
109.                     __func__, retval);
110.  }
111.
112.  return retval;
113. }
114.
115. static struct kset_uevent_ops device_uevent_ops = {
116.     .filter = dev_uevent_filter,
117.     .name = dev_uevent_name,
118.     .uevent = dev_uevent,
119. };

```

前面在讲到kset时，我们并未关注其中的kset_event_ops结构变量。但这里device既然用到了，我们就对其中的三个函数做简单介绍。kset_uevent_ops中的函数是用于管理kset内部kobject的uevent操作。其中filter函数用于阻止一个kobject向用户空间发送uevent，返回值为0表示阻止。这里dev_uevent_filter()检查device所属的bus或者class是否存在，如果都不存在，也就没有发送uevent的必要了。name函数是用于覆盖kset发送给用户空间的名称。这里dev_uevent_name()选择使用bus或者class的名称。uevent()函数是在uevent将被发送到用户空间之前调用的，用于向uevent中增加新的环境变量。dev_uevent()的实现很热闹，向uevent中添加了各种环境变量。

1. static ssize_t show_uevent(struct device *dev, struct device_attribute *attr,
2. char *buf)

```

3. {
4.     struct kobject *top_kobj;
5.     struct kset *kset;
6.     struct kobj_uevent_env *env = NULL;
7.     int i;
8.     size_t count = 0;
9.     int retval;
10.
11.     /* search the kset, the device belongs to */
12.     top_kobj = &dev->kobj;
13.     while (!top_kobj->kset && top_kobj->parent)
14.         top_kobj = top_kobj->parent;
15.     if (!top_kobj->kset)
16.         goto out;
17.
18.     kset = top_kobj->kset;
19.     if (!kset->uevent_ops || !kset->uevent_ops->uevent)
20.         goto out;
21.
22.     /* respect filter */
23.     if (kset->uevent_ops && kset->uevent_ops->filter)
24.         if (!kset->uevent_ops->filter(kset, &dev->kobj))
25.             goto out;
26.
27.     env = kzalloc(sizeof(struct kobj_uevent_env), GFP_KERNEL);
28.     if (!env)
29.         return -ENOMEM;
30.
31.     /* let the kset specific function add its keys */
32.     retval = kset->uevent_ops->uevent(kset, &dev->kobj, env);
33.     if (retval)
34.         goto out;
35.
36.     /* copy keys to file */
37.     for (i = 0; i < env->envp_idx; i++)
38.         count += sprintf(&buf[count], "%s\n", env->envp[i]);
39. out:
40.     kfree(env);
41.     return count;
42. }
43.
44. static ssize_t store_uevent(struct device *dev, struct device_attribute *attr,
45.                             const char *buf, size_t count)
46. {
47.     enum kobject_action action;
48.
49.     if (kobject_action_type(buf, count, &action) == 0) {
50.         kobject_uevent(&dev->kobj, action);
51.         goto out;
52.     }
53.
54.     dev_err(dev, "uevent: unsupported action-string; this will "
55.             "be ignored in a future kernel version\n");
56.     kobject_uevent(&dev->kobj, KOBJ_ADD);
57. out:
58.     return count;
59. }
60.
61. static struct device_attribute uevent_attr =
62.     __ATTR(uevent, S_IRUGO | S_IWUSR, show_uevent, store_uevent);

```

device不仅在对kset中添加了对uevent的管理，而且还把uevent信息做成设备的一个属性uevent。其中show_event()是显示uevent中环境变量的，store_uevent()是发送uevent的。

```

1. static int device_add_attributes(struct device *dev,
2.                                 struct device_attribute *attrs)
3. {
4.     int error = 0;
5.     int i;
6.

```

```

7.  if (attrs) {
8.      for (i = 0; attr_name(attrs[i]); i++) {
9.          error = device_create_file(dev, &attrs[i]);
10.         if (error)
11.             break;
12.     }
13.     if (error)
14.         while (--i >= 0)
15.             device_remove_file(dev, &attrs[i]);
16. }
17. return error;
18. }
19.
20. static void device_remove_attributes(struct device *dev,
21.                                     struct device_attribute *attrs)
22. {
23.     int i;
24.
25.     if (attrs)
26.         for (i = 0; attr_name(attrs[i]); i++)
27.             device_remove_file(dev, &attrs[i]);
28. }
29.
30. static int device_add_groups(struct device *dev,
31.                              const struct attribute_group **groups)
32. {
33.     int error = 0;
34.     int i;
35.
36.     if (groups) {
37.         for (i = 0; groups[i]; i++) {
38.             error = sysfs_create_group(&dev->kobj, groups[i]);
39.             if (error) {
40.                 while (--i >= 0)
41.                     sysfs_remove_group(&dev->kobj,
42.                                         groups[i]);
43.                 break;
44.             }
45.         }
46.     }
47.     return error;
48. }
49.
50. static void device_remove_groups(struct device *dev,
51.                                  const struct attribute_group **groups)
52. {
53.     int i;
54.
55.     if (groups)
56.         for (i = 0; groups[i]; i++)
57.             sysfs_remove_group(&dev->kobj, groups[i]);
58. }

```

以上四个内部函数是用来向device中添加或删除属性与属性集合的。

device_add_attributes、device_remove_attributes、device_add_groups、device_remove_groups，都是直接通过sysfs提供的API实现。

```

1. static int device_add_attrs(struct device *dev)
2. {
3.     struct class *class = dev->class;
4.     struct device_type *type = dev->type;
5.     int error;
6.
7.     if (class) {
8.         error = device_add_attributes(dev, class->dev_attrs);
9.         if (error)
10.            return error;
11.     }
12.
13.     if (type) {
14.         error = device_add_groups(dev, type->groups);

```

```

15.     if (error)
16.         goto err_remove_class_attrs;
17. }
18.
19. error = device_add_groups(dev, dev->groups);
20. if (error)
21.     goto err_remove_type_groups;
22.
23. return 0;
24.
25. err_remove_type_groups:
26. if (type)
27.     device_remove_groups(dev, type->groups);
28. err_remove_class_attrs:
29. if (class)
30.     device_remove_attributes(dev, class->dev_attrs);
31.
32. return error;
33. }
34.
35. static void device_remove_attrs(struct device *dev)
36. {
37.     struct class *class = dev->class;
38.     struct device_type *type = dev->type;
39.
40.     device_remove_groups(dev, dev->groups);
41.
42.     if (type)
43.         device_remove_groups(dev, type->groups);
44.
45.     if (class)
46.         device_remove_attributes(dev, class->dev_attrs);
47. }

```

device_add_attrs()实际负责device中的属性添加。也是几个部分的集合，包括class中的dev_attrs，device_type中的groups，还有device本身的groups。

device_remove_attrs()则负责对应的device属性删除工作。

```

1. #define print_dev_t(buffer, dev)          \
2.     sprintf((buffer), "%u:%u\n", MAJOR(dev), MINOR(dev))
3.
4. static ssize_t show_dev(struct device *dev, struct device_attribute *attr,
5.     char *buf)
6. {
7.     return print_dev_t(buf, dev->devt);
8. }
9.
10. static struct device_attribute devt_attr =
11.     __ATTR(dev, S_IRUGO, show_dev, NULL);

```

这里又定义了一个名为dev的属性，就是显示设备的设备号。

```

1. /**
2.  * device_create_file - create sysfs attribute file for device.
3.  * @dev: device.
4.  * @attr: device attribute descriptor.
5.  */
6. int device_create_file(struct device *dev, struct device_attribute *attr)
7. {
8.     int error = 0;
9.     if (dev)
10.         error = sysfs_create_file(&dev->kobj, &attr->attr);
11.     return error;
12. }
13.
14. /**
15.  * device_remove_file - remove sysfs attribute file.
16.  * @dev: device.
17.  * @attr: device attribute descriptor.
18.  */
19. void device_remove_file(struct device *dev, struct device_attribute *attr)

```



```

20. {
21.     if (dev)
22.         sysfs_remove_file(&dev->kobj, &attr->attr);
23. }
24.
25. /**
26.  * device_create_bin_file - create sysfs binary attribute file for device.
27.  * @dev: device.
28.  * @attr: device binary attribute descriptor.
29.  */
30. int device_create_bin_file(struct device *dev, struct bin_attribute *attr)
31. {
32.     int error = -EINVAL;
33.     if (dev)
34.         error = sysfs_create_bin_file(&dev->kobj, attr);
35.     return error;
36. }
37.
38. /**
39.  * device_remove_bin_file - remove sysfs binary attribute file
40.  * @dev: device.
41.  * @attr: device binary attribute descriptor.
42.  */
43. void device_remove_bin_file(struct device *dev, struct bin_attribute *attr)
44. {
45.     if (dev)
46.         sysfs_remove_bin_file(&dev->kobj, attr);
47. }
48.
49. int device_schedule_callback_owner(struct device *dev,
50.     void (*func)(struct device *), struct module *owner)
51. {
52.     return sysfs_schedule_callback(&dev->kobj,
53.         (void (*)(void *)) func, dev, owner);
54. }

```

这里的五个函数，也是对sysfs提供的API的简单封装。

device_create_file()和device_remove_file()提供直接的属性文件管理方法。

device_create_bin_file()和device_remove_bin_file()则是提供设备管理二进制文件的方法。

device_schedule_callback_owner()也是简单地将func加入工作队列。

```

1. static void klist_children_get(struct klist_node *n)
2. {
3.     struct device_private *p = to_device_private_parent(n);
4.     struct device *dev = p->device;
5.
6.     get_device(dev);
7. }
8.
9. static void klist_children_put(struct klist_node *n)
10. {
11.     struct device_private *p = to_device_private_parent(n);
12.     struct device *dev = p->device;
13.
14.     put_device(dev);
15. }

```

如果之前认真看过klist的实现，应该知道，klist_children_get()和klist_children_put()就是在设备挂入和删除父设备的klist_children链表时调用的函数。在父设备klist_children链表上的指针，相当于对device的一个引用计数。

```

1. struct device *get_device(struct device *dev)
2. {
3.     return dev ? to_dev(kobject_get(&dev->kobj)) : NULL;
4. }
5.
6. /**
7.  * put_device - decrement reference count.
8.  * @dev: device in question.
9.  */

```

```

10. void put_device(struct device *dev)
11. {
12.     /* might_sleep(); */
13.     if (dev)
14.         kobject_put(&dev->kobj);
15. }

```

device中的引用计数，完全交给内嵌的kobject来做。如果引用计数降为零，自然是调用之前说到的包含甚广的device_release函数。

```

1. void device_initialize(struct device *dev)
2. {
3.     dev->kobj.kset = devices_kset;
4.     kobject_init(&dev->kobj, &device_ktype);
5.     INIT_LIST_HEAD(&dev->dma_pools);
6.     init_MUTEX(&dev->sem);
7.     spin_lock_init(&dev->devres_lock);
8.     INIT_LIST_HEAD(&dev->devres_head);
9.     device_init_wakeup(dev, 0);
10.    device_pm_init(dev);
11.    set_dev_node(dev, -1);
12. }

```

device_initialize()就是device结构的初始化函数，它把device中能初始化的部分全初始化。它的界限在其中kobj的位置与device在设备驱动模型中的位置，这些必须由外部设置。可以看到，调用kobject_init()时，object的kobj_type选择了device_ktype，其中主要是sysops的两个函数，还有device_release函数。

```

1. static struct kobject *virtual_device_parent(struct device *dev)
2. {
3.     static struct kobject *virtual_dir = NULL;
4.
5.     if (!virtual_dir)
6.         virtual_dir = kobject_create_and_add("virtual",
7.             &devices_kset->kobj);
8.
9.     return virtual_dir;
10. }
11.
12. static struct kobject *get_device_parent(struct device *dev,
13.     struct device *parent)
14. {
15.     int retval;
16.
17.     if (dev->class) {
18.         struct kobject *kobj = NULL;
19.         struct kobject *parent_kobj;
20.         struct kobject *k;
21.
22.         /*
23.          * If we have no parent, we live in "virtual".
24.          * Class-devices with a non class-device as parent, live
25.          * in a "glue" directory to prevent namespace collisions.
26.          */
27.         if (parent == NULL)
28.             parent_kobj = virtual_device_parent(dev);
29.         else if (parent->class)
30.             return &parent->kobj;
31.         else
32.             parent_kobj = &parent->kobj;
33.
34.         /* find our class-directory at the parent and reference it */
35.         spin_lock(&dev->class->p->class_dirs.list_lock);
36.         list_for_each_entry(k, &dev->class->p->class_dirs.list, entry)
37.             if (k->parent == parent_kobj) {
38.                 kobj = kobject_get(k);
39.                 break;
40.             }
41.         spin_unlock(&dev->class->p->class_dirs.list_lock);
42.         if (kobj)
43.             return kobj;
44.
45.         /* or create a new class-directory at the parent device */
46.         k = kobject_create();

```

```

47.     if (!k)
48.         return NULL;
49.     k->kset = &dev->class->p->class_dirs;
50.     retval = kobject_add(k, parent_kobj, "%s", dev->class->name);
51.     if (retval < 0) {
52.         kobject_put(k);
53.         return NULL;
54.     }
55.     /* do not emit an uevent for this simple "glue" directory */
56.     return k;
57. }
58.
59. if (parent)
60.     return &parent->kobj;
61. return NULL;
62. }

```

这里的get_device_parent()就是获取父节点的kobject，但也并非就如此简单。get_device_parent()的返回值直接决定了device将被挂在哪个目录下。到底该挂在哪儿，是由dev->class、dev->parent、dev->parent->class等因素综合决定的。我们看get_device_parent()中是如何判断的。如果dev->class为空，表示一切随父设备，有parent则返回parent->kobj，没有则返回NULL。如果有dev->class呢，情况就比较复杂了，也许device有着与parent不同的class，也许device还没有一个parent，等等。我们看具体的情况。如果parent不为空，而且存在parent->class，则还放在parent目录下。不然，要么parent不存在，要么parent没有class，很难直接将有class的device放在parent下面。目前的解决方法很简单，在parent与device之间，再加一层表示class的目录。如果parent都没有，那就把/sys/devices/virtual当parent。class->p->class_dirs就是专门存放这种中间kobject的kset。思路理清后，再结合实际的sysfs，代码就很容易看懂了。

```

1. static void cleanup_glue_dir(struct device *dev, struct kobject *glue_dir)
2. {
3.     /* see if we live in a "glue" directory */
4.     if (!glue_dir || !dev->class ||
5.         glue_dir->kset != &dev->class->p->class_dirs)
6.         return;
7.
8.     kobject_put(glue_dir);
9. }
10.
11. static void cleanup_device_parent(struct device *dev)
12. {
13.     cleanup_glue_dir(dev, dev->kobj.parent);
14. }

```

cleanup_device_parent()是取消对parent引用时调用的函数，看起来只针对这种glue形式的目录起作用。

```

1. static void setup_parent(struct device *dev, struct device *parent)
2. {
3.     struct kobject *kobj;
4.     kobj = get_device_parent(dev, parent);
5.     if (kobj)
6.         dev->kobj.parent = kobj;
7. }

```

setup_parent()就是调用get_device_parent()获得应该存放的父目录kobj，并把dev->kobj.parent设为它。

```

1. static int device_add_class_symlinks(struct device *dev)
2. {
3.     int error;
4.
5.     if (!dev->class)
6.         return 0;
7.
8.     error = sysfs_create_link(&dev->kobj,
9.                             &dev->class->p->class_subsys.kobj,
10.                            "subsystem");
11.     if (error)
12.         goto out;
13.     /* link in the class directory pointing to the device */
14.     error = sysfs_create_link(&dev->class->p->class_subsys.kobj,
15.                             &dev->kobj, dev_name(dev));
16.     if (error)
17.         goto out_subsys;
18.
19.     if (dev->parent && device_is_not_partition(dev)) {
20.         error = sysfs_create_link(&dev->kobj, &dev->parent->kobj,
21.                                "device");

```

```

22.     if (error)
23.         goto out_busid;
24. }
25. return 0;
26.
27. out_busid:
28.     sysfs_remove_link(&dev->class->p->class_subsys.kobj, dev_name(dev));
29. out_subsys:
30.     sysfs_remove_link(&dev->kobj, "subsystem");
31. out:
32.     return error;
33. }

```

device_add_class_symlinks()在device和class直接添加一些软链接。在device目录下创建指向class的subsystem文件，在class目录下创建指向device的同名文件。如果device有父设备，而且device不是块设备分区，则在device目录下建立一个指向父设备的device链接文件。这一点在usb设备和usb接口间很常见。

```

1. static void device_remove_class_symlinks(struct device *dev)
2. {
3.     if (!dev->class)
4.         return;
5.
6. #ifdef CONFIG_SYSFS_DEPRECATED
7.     if (dev->parent && device_is_not_partition(dev)) {
8.         char *class_name;
9.
10.        class_name = make_class_name(dev->class->name, &dev->kobj);
11.        if (class_name) {
12.            sysfs_remove_link(&dev->parent->kobj, class_name);
13.            kfree(class_name);
14.        }
15.        sysfs_remove_link(&dev->kobj, "device");
16.    }
17.
18.    if (dev->kobj.parent != &dev->class->p->class_subsys.kobj &&
19.        device_is_not_partition(dev))
20.        sysfs_remove_link(&dev->class->p->class_subsys.kobj,
21.            dev_name(dev));
22. #else
23.     if (dev->parent && device_is_not_partition(dev))
24.         sysfs_remove_link(&dev->kobj, "device");
25.
26.     sysfs_remove_link(&dev->class->p->class_subsys.kobj, dev_name(dev));
27. #endif
28.
29.     sysfs_remove_link(&dev->kobj, "subsystem");
30. }

```

device_remove_class_symlinks()删除device和class之间的软链接。

```

1. static inline const char *dev_name(const struct device *dev)
2. {
3.     return kobject_name(&dev->kobj);
4. }
5.
6. int dev_set_name(struct device *dev, const char *fmt, ...)
7. {
8.     va_list vargs;
9.     int err;
10.
11.     va_start(vargs, fmt);
12.     err = kobject_set_name_vargs(&dev->kobj, fmt, vargs);
13.     va_end(vargs);
14.     return err;
15. }

```

dev_name()获得设备名称，dev_set_name()设置设备名称。但这里的dev_set_name()只能在设备未注册前使用。device的名称其实是完全靠dev->kobj管理的。

```

1. static struct kobject *device_to_dev_kobj(struct device *dev)
2. {
3.     struct kobject *kobj;
4.
5.     if (dev->class)

```

```

6.     kobj = dev->class->dev_kobj;
7.     else
8.         kobj = sysfs_dev_char_kobj;
9.
10.    return kobj;
11. }

```

device_to_dev_kobj()为dev选择合适的/sys/dev下的kobject，或者是块设备，或者是字符设备，或者没有。

```

1. #define format_dev_t(buffer, dev)          \
2.     ({                                     \
3.         sprintf(buffer, "%u:%u", MAJOR(dev), MINOR(dev)); \
4.         buffer;                          \
5.     })
6.
7. static int device_create_sys_dev_entry(struct device *dev)
8. {
9.     struct kobject *kobj = device_to_dev_kobj(dev);
10.    int error = 0;
11.    char devt_str[15];
12.
13.    if (kobj) {
14.        format_dev_t(devt_str, dev->devt);
15.        error = sysfs_create_link(kobj, &dev->kobj, devt_str);
16.    }
17.
18.    return error;
19. }
20.
21. static void device_remove_sys_dev_entry(struct device *dev)
22. {
23.     struct kobject *kobj = device_to_dev_kobj(dev);
24.     char devt_str[15];
25.
26.     if (kobj) {
27.         format_dev_t(devt_str, dev->devt);
28.         sysfs_remove_link(kobj, devt_str);
29.     }
30. }

```

device_create_sys_dev_entry()是在/sys/dev相应的目录下建立对设备的软链接。先是通过device_to_dev_kobj()获得父节点的kobj，然后调用sysfs_create_link()建立软链接。

device_remove_sys_dev_entry()与其操作正相反，删除在/sys/dev下建立的软链接。

```

1. int device_private_init(struct device *dev)
2. {
3.     dev->p = kzalloc(sizeof(*dev->p), GFP_KERNEL);
4.     if (!dev->p)
5.         return -ENOMEM;
6.     dev->p->device = dev;
7.     klist_init(&dev->p->klist_children, klist_children_get,
8.               klist_children_put);
9.     return 0;
10. }

```

device_private_init()分配并初始化dev->p。至于空间的释放，是等到释放设备时调用的device_release()中。

之前的函数比较散乱，或许找不出一个整体的印象。但下面马上就要看到重要的部分了，因为代码终于攒到了爆发的程度！

```

1. /**
2.  * device_register - register a device with the system.
3.  * @dev: pointer to the device structure
4.  *
5.  * This happens in two clean steps - initialize the device
6.  * and add it to the system. The two steps can be called
7.  * separately, but this is the easiest and most common.
8.  * I.e. you should only call the two helpers separately if
9.  * have a clearly defined need to use and refcount the device
10. * before it is added to the hierarchy.
11. *
12. * NOTE: _Never_ directly free @dev after calling this function, even
13. * if it returned an error! Always use put_device() to give up the

```

```

14. * reference initialized in this function instead.
15. */
16. int device_register(struct device *dev)
17. {
18.     device_initialize(dev);
19.     return device_add(dev);
20. }

```

device_register()是提供给外界注册设备的接口。它先是调用device_initialize()初始化dev结构，然后调用device_add()将其加入系统中。但要注意，在调用device_register()注册dev之前，有一些dev结构变量是需要自行设置的。这其中有指明设备位置的struct device *parent, struct bus_type *bus, struct class *class，有指明设备属性的 const char *init_name, struct device_type *type, const struct attribute_group **groups, void (*release)(struct device *dev), dev_t devt，等等。不同设备的使用方法不同，我们留待之后再具体分析。device_initialize()我们已经看过，下面重点看看device_add()是如何实现的。

```

1. int device_add(struct device *dev)
2. {
3.     struct device *parent = NULL;
4.     struct class_interface *class_intf;
5.     int error = -EINVAL;
6.
7.     dev = get_device(dev);
8.     if (!dev)
9.         goto done;
10.
11.     if (!dev->p) {
12.         error = device_private_init(dev);
13.         if (error)
14.             goto done;
15.     }
16.
17.     /*
18.      * for statically allocated devices, which should all be converted
19.      * some day, we need to initialize the name. We prevent reading back
20.      * the name, and force the use of dev_name()
21.      */
22.     if (dev->init_name) {
23.         dev_set_name(dev, "%s", dev->init_name);
24.         dev->init_name = NULL;
25.     }
26.
27.     if (!dev_name(dev))
28.         goto name_error;
29.
30.     pr_debug("device: '%s': %s\n", dev_name(dev), __func__);
31.
32.     parent = get_device(dev->parent);
33.     setup_parent(dev, parent);
34.
35.     /* use parent numa_node */
36.     if (parent)
37.         set_dev_node(dev, dev_to_node(parent));
38.
39.     /* first, register with generic layer. */
40.     /* we require the name to be set before, and pass NULL */
41.     error = kobject_add(&dev->kobj, dev->kobj.parent, NULL);
42.     if (error)
43.         goto Error;
44.
45.     /* notify platform of device entry */
46.     if (platform_notify)
47.         platform_notify(dev);
48.
49.     error = device_create_file(dev, &uevent_attr);
50.     if (error)
51.         goto attrError;
52.
53.     if (MAJOR(dev->devt)) {
54.         error = device_create_file(dev, &devt_attr);
55.         if (error)
56.             goto ueventattrError;
57.

```

```

58.     error = device_create_sys_dev_entry(dev);
59.     if (error)
60.         goto devtattrError;
61.
62.     devtmpfs_create_node(dev);
63. }
64.
65. error = device_add_class_symlinks(dev);
66. if (error)
67.     goto SymlinkError;
68. error = device_add_attrs(dev);
69. if (error)
70.     goto AttrsError;
71. error = bus_add_device(dev);
72. if (error)
73.     goto BusError;
74. error = dpm_sysfs_add(dev);
75. if (error)
76.     goto DPMErrors;
77. device_pm_add(dev);
78.
79. /* Notify clients of device addition. This call must come
80.  * after dpm_sysfs_add() and before kobject_uevent().
81.  */
82. if (dev->bus)
83.     blocking_notifier_call_chain(&dev->bus->p->bus_notifier,
84.                                  BUS_NOTIFY_ADD_DEVICE, dev);
85.
86. kobject_uevent(&dev->kobj, KOBJ_ADD);
87. bus_probe_device(dev);
88. if (parent)
89.     klist_add_tail(&dev->p->knode_parent,
90.                   &parent->p->klist_children);
91.
92. if (dev->class) {
93.     mutex_lock(&dev->class->p->class_mutex);
94.     /* tie the class to the device */
95.     klist_add_tail(&dev->knode_class,
96.                   &dev->class->p->class_devices);
97.
98.     /* notify any interfaces that the device is here */
99.     list_for_each_entry(class_intf,
100.                        &dev->class->p->class_interfaces, node)
101.         if (class_intf->add_dev)
102.             class_intf->add_dev(dev, class_intf);
103.     mutex_unlock(&dev->class->p->class_mutex);
104. }
105. done:
106.     put_device(dev);
107.     return error;
108. DPMErrors:
109.     bus_remove_device(dev);
110. BusError:
111.     device_remove_attrs(dev);
112. AttrsError:
113.     device_remove_class_symlinks(dev);
114. SymlinkError:
115.     if (MAJOR(dev->devt))
116.         device_remove_sys_dev_entry(dev);
117. devtattrError:
118.     if (MAJOR(dev->devt))
119.         device_remove_file(dev, &devt_attr);
120. ueventattrError:
121.     device_remove_file(dev, &uevent_attr);
122. attrError:
123.     kobject_uevent(&dev->kobj, KOBJ_REMOVE);
124.     kobject_del(&dev->kobj);
125. Error:
126.     cleanup_device_parent(dev);
127.     if (parent)

```

```

128.     put_device(parent);
129. name_error:
130.     kfree(dev->p);
131.     dev->p = NULL;
132.     goto done;
133. }

```

device_add()将dev加入设备驱动模型。它先是调用get_device(dev)增加dev的引用计数，然后调用device_private_init()分配和初始化dev->p，调用dev_set_name()设置dev名字。然后是准备将dev加入sysfs，先是用get_device(parent)增加对parent的引用计数(无论是直接挂在parent下还是通过一个类层挂在parent下都要增加parent的引用计数)，然后调用setup_parent()找到实际要加入的父kobject，通过kobject_add()加入其下。然后是添加属性和属性集合的操作，调用device_create_file()添加uevent属性，调用device_add_attrs()添加device/type/class预定义的属性与属性集合。如果dev有被分配设备号，再用device_create_file()添加dev属性，并用device_create_sys_dev_entry()在/sys/dev下添加相应的软链接，最后调用devtmpfs_create_node()在/dev下创建相应的设备文件。然后调用device_add_class_symlinks()添加dev与class间的软链接，调用bus_add_device()添加dev与bus间的软链接，并将dev挂入bus的设备链表。调用dpm_sysfs_add()增加dev下的power属性集合，调用device_pm_add()将dev加入dpm_list链表。

调用kobject_uevent()发布KOBJ_ADD消息，调用bus_probe_device()为dev寻找合适的驱动。如果有parent节点，把dev->p->knode_parent挂入parent->p->klist_children链表。如果dev有所属的class，将dev->knode_class挂在class->p->class_devices上，并调用可能的类设备接口的add_dev()方法。可能对于直接在bus上的设备来说，自然可以调用bus_probe_device()查找驱动，而不与总线直接接触的设备，则要靠class来发现驱动，这里的class_interface中的add_dev()方法，就是一个绝好的机会。最后会调用put_device(dev)释放在函数开头增加的引用计数。

device_add()要做的事很多，但想想每件事都在情理之中。device是设备驱动模型的基本元素，在class、bus、dev、devices中都有它的身影。device_add()要适应各种类型的设备注册，自然会越来越复杂。可以说文件开头定义的内部函数，差不多都是为了这里服务的。

```

1. void device_unregister(struct device *dev)
2. {
3.     pr_debug("device: '%s': %s\n", dev_name(dev), __func__);
4.     device_del(dev);
5.     put_device(dev);
6. }

```

有注册自然又注销。device_unregister()就是用于将dev从系统中注销，并释放创建时产生的引用计数。

```

1. void device_del(struct device *dev)
2. {
3.     struct device *parent = dev->parent;
4.     struct class_interface *class_intf;
5.
6.     /* Notify clients of device removal. This call must come
7.      * before dpm_sysfs_remove().
8.      */
9.     if (dev->bus)
10.         blocking_notifier_call_chain(&dev->bus->p->bus_notifier,
11.                                     BUS_NOTIFY_DEL_DEVICE, dev);
12.     device_pm_remove(dev);
13.     dpm_sysfs_remove(dev);
14.     if (parent)
15.         klist_del(&dev->p->knode_parent);
16.     if (MAJOR(dev->devt)) {
17.         devtmpfs_delete_node(dev);
18.         device_remove_sys_dev_entry(dev);
19.         device_remove_file(dev, &devt_attr);
20.     }
21.     if (dev->class) {
22.         device_remove_class_symlinks(dev);
23.
24.         mutex_lock(&dev->class->p->class_mutex);
25.         /* notify any interfaces that the device is now gone */
26.         list_for_each_entry(class_intf,
27.                             &dev->class->p->class_interfaces, node)
28.             if (class_intf->remove_dev)
29.                 class_intf->remove_dev(dev, class_intf);
30.         /* remove the device from the class list */
31.         klist_del(&dev->knode_class);
32.         mutex_unlock(&dev->class->p->class_mutex);
33.     }
34.     device_remove_file(dev, &uevent_attr);
35.     device_remove_attrs(dev);
36.     bus_remove_device(dev);
37.
38.     /*
39.      * Some platform devices are driven without driver attached
40.      * and managed resources may have been acquired. Make sure

```



```

41.  * all resources are released.
42.  */
43. devres_release_all(dev);
44.
45. /* Notify the platform of the removal, in case they
46.  * need to do anything...
47.  */
48. if (platform_notify_remove)
49.     platform_notify_remove(dev);
50. kobject_uevent(&dev->kobj, KOBJ_REMOVE);
51. cleanup_device_parent(dev);
52. kobject_del(&dev->kobj);
53. put_device(parent);
54. }

```

device_del()是与device_add()相对的函数，进行实际的将dev从系统中脱离的工作。这其中既有将dev从设备驱动模型各种链表中脱离的工作，又有将dev从sysfs的各个角落删除的工作。大致流程与dev_add()相对，就不一一介绍。

爆发结束，下面来看一些比较轻松的函数。

```

1. /**
2.  * device_get_devnode - path of device node file
3.  * @dev: device
4.  * @mode: returned file access mode
5.  * @tmp: possibly allocated string
6.  *
7.  * Return the relative path of a possible device node.
8.  * Non-default names may need to allocate a memory to compose
9.  * a name. This memory is returned in tmp and needs to be
10. * freed by the caller.
11. */
12. const char *device_get_devnode(struct device *dev,
13.                                mode_t *mode, const char **tmp)
14. {
15.     char *s;
16.
17.     *tmp = NULL;
18.
19.     /* the device type may provide a specific name */
20.     if (dev->type && dev->type->devnode)
21.         *tmp = dev->type->devnode(dev, mode);
22.     if (*tmp)
23.         return *tmp;
24.
25.     /* the class may provide a specific name */
26.     if (dev->class && dev->class->devnode)
27.         *tmp = dev->class->devnode(dev, mode);
28.     if (*tmp)
29.         return *tmp;
30.
31.     /* return name without allocation, tmp == NULL */
32.     if (strchr(dev_name(dev), '!') == NULL)
33.         return dev_name(dev);
34.
35.     /* replace '!' in the name with '/' */
36.     *tmp = kstrdup(dev_name(dev), GFP_KERNEL);
37.     if (!*tmp)
38.         return NULL;
39.     while ((s = strchr(*tmp, '!')))
40.         s[0] = '/';
41.     return *tmp;
42. }

```

device_get_devnode()返回设备的路径名。不过似乎可以由device_type或者class定义一些独特的返回名称。

```

1. static struct device *next_device(struct klist_iter *i)
2. {
3.     struct klist_node *n = klist_next(i);
4.     struct device *dev = NULL;
5.     struct device_private *p;
6.

```

```

7.  if (n) {
8.      p = to_device_private_parent(n);
9.      dev = p->device;
10. }
11. return dev;
12. }
13.
14. int device_for_each_child(struct device *parent, void *data,
15.     int (*fn)(struct device *dev, void *data))
16. {
17.     struct klist_iter i;
18.     struct device *child;
19.     int error = 0;
20.
21.     if (!parent->p)
22.         return 0;
23.
24.     klist_iter_init(&parent->p->klist_children, &i);
25.     while ((child = next_device(&i)) && !error)
26.         error = fn(child, data);
27.     klist_iter_exit(&i);
28.     return error;
29. }
30.
31. struct device *device_find_child(struct device *parent, void *data,
32.     int (*match)(struct device *dev, void *data))
33. {
34.     struct klist_iter i;
35.     struct device *child;
36.
37.     if (!parent)
38.         return NULL;
39.
40.     klist_iter_init(&parent->p->klist_children, &i);
41.     while ((child = next_device(&i)))
42.         if (match(child, data) && get_device(child))
43.             break;
44.     klist_iter_exit(&i);
45.     return child;
46. }

```

device_for_each_child()对dev下的每个子device，都调用一遍特定的处理函数。

device_find_child()则是查找dev下特点的子device，查找使用特定的match函数。

这两个遍历过程都使用了klist特有的遍历函数，支持遍历过程中的节点删除等功能。next_device()则是为了遍历方便封装的一个内部函数。

下面本该是root_device注册相关的代码。但经过检查，linux内核中使用到的root_device很少见，而且在sysfs中也未能找到一个实际的例子。所以root_device即使还未被弃用，也并非主流，我们将其跳过。

与kobject和kset类似，device也为我们提供了快速device创建方法，下面就看看吧。

```

1. static void device_create_release(struct device *dev)
2. {
3.     pr_debug("device: '%s': %s\n", dev_name(dev), __func__);
4.     kfree(dev);
5. }
6.
7. struct device *device_create_vars(struct class *class, struct device *parent,
8.     dev_t devt, void *drvdata, const char *fmt,
9.     va_list args)
10. {
11.     struct device *dev = NULL;
12.     int retval = -ENODEV;
13.
14.     if (class == NULL || IS_ERR(class))
15.         goto error;
16.
17.     dev = kzalloc(sizeof(*dev), GFP_KERNEL);
18.     if (!dev) {
19.         retval = -ENOMEM;
20.         goto error;

```

```

21. }
22.
23. dev->devt = devt;
24. dev->class = class;
25. dev->parent = parent;
26. dev->release = device_create_release;
27. dev_set_drvdata(dev, drvdata);
28.
29. retval = kobject_set_name_vargs(&dev->kobj, fmt, args);
30. if (retval)
31.     goto error;
32.
33. retval = device_register(dev);
34. if (retval)
35.     goto error;
36.
37. return dev;
38.
39. error:
40. put_device(dev);
41. return ERR_PTR(retval);
42. }
43.
44. struct device *device_create(struct class *class, struct device *parent,
45.                             dev_t devt, void *drvdata, const char *fmt, ...)
46. {
47.     va_list vars;
48.     struct device *dev;
49.
50.     va_start(vars, fmt);
51.     dev = device_create_vargs(class, parent, devt, drvdata, fmt, vars);
52.     va_end(vars);
53.     return dev;
54. }

```

这里的device_create()提供了一个快速的dev创建注册方法。只是中间没有提供设置device_type的方法，或许是这样的device已经够特立独行了，不需要搞出一类来。

```

1. static int __match_devt(struct device *dev, void *data)
2. {
3.     dev_t *devt = data;
4.
5.     return dev->devt == *devt;
6. }
7.
8. void device_destroy(struct class *class, dev_t devt)
9. {
10.    struct device *dev;
11.
12.    dev = class_find_device(class, NULL, &devt, __match_devt);
13.    if (dev) {
14.        put_device(dev);
15.        device_unregister(dev);
16.    }
17. }

```

device_destroy()就是与device_create()相对的注销函数。至于这里为什么会多一个put_device(dev)，也很简单，因为在class_find_device()找到dev时，调用了get_device()。

```

1. struct device *class_find_device(struct class *class, struct device *start,
2.                                 void *data,
3.                                 int (*match)(struct device *, void *))
4. {
5.     struct class_dev_iter iter;
6.     struct device *dev;
7.
8.     if (!class)
9.         return NULL;
10.    if (!class->p) {
11.        WARN(1, "%s called for class '%s' before it was initialized",
12.             __func__, class->name);

```

```

13.     return NULL;
14. }
15.
16. class_dev_iter_init(&iter, class, start, NULL);
17. while ((dev = class_dev_iter_next(&iter))) {
18.     if (match(dev, data)) {
19.         get_device(dev);
20.         break;
21.     }
22. }
23. class_dev_iter_exit(&iter);
24.
25. return dev;
26. }

```

class_find_device()本来是class.c中的内容，其实现也于之前将的遍历dev->p->klist_children类似，无非是在klist提供的遍历方法上加以封装。但我们这里列出class_find_device()的实现与使用它的device_destroy()，却是为了更好地分析这个调用流程中dev是如何被保护的。它实际上是经历了三个保护手段：首先在class_dev_iter_next()->klist_next()中，是受到struct klist中 spinlock_t k_lock保护的。在找到下一点并解锁之前，就增加了struct klist_node中的struct kref n_ref引用计数。在当前的next()调用完，到下一个next()调用之前，都是受这个增加的引用计数保护的。再看class_find_device()中，使用get_device(dev)增加了dev本身的引用计数保护(当然也要追溯到kobj->kref中)，这是第三种保护。知道device_destroy()中主动调用put_device(dev)才去除了这种保护。

本来对dev的保护，应该完全是由dev中的引用计数完成的。但实际上这种保护很多时候是间接完成的。例如这里的klist中的自旋锁，klist_node中的引用计数，都不过是为了保持class的设备链表中对dev的引用计数不消失，这是一种间接保护的手段，保证了这中间即使外界主动释放class设备链表对dev的引用计数，dev仍然不会被实际注销。这种曲折的联系，才真正发挥了引用计数的作用，构成设备驱动模型独特的魅力。

```

1. int device_rename(struct device *dev, char *new_name)
2. {
3.     char *old_device_name = NULL;
4.     int error;
5.
6.     dev = get_device(dev);
7.     if (!dev)
8.         return -EINVAL;
9.
10.    pr_debug("device: '%s': %s: renaming to '%s'\n", dev_name(dev),
11.        __func__, new_name);
12.    old_device_name = kstrdup(dev_name(dev), GFP_KERNEL);
13.    if (!old_device_name) {
14.        error = -ENOMEM;
15.        goto out;
16.    }
17.
18.    error = kobject_rename(&dev->kobj, new_name);
19.    if (error)
20.        goto out;
21.    if (dev->class) {
22.        error = sysfs_create_link_nowarn(&dev->class->p->class_subsys.kobj,
23.            &dev->kobj, dev_name(dev));
24.        if (error)
25.            goto out;
26.        sysfs_remove_link(&dev->class->p->class_subsys.kobj,
27.            old_device_name);
28.    }
29. out:
30.    put_device(dev);
31.
32.    kfree(old_device_name);
33.
34.    return error;
35. }

```

device_rename()是供设备注册后改变名称用的，除了改变/sys/devices下地名称，还改变了/sys/class下地软链接名称。前者很自然，但后者却很难想到。即使简单的地方，经过重重调试，我们也会惊讶于linux的心细如发。

```

1. static int device_move_class_links(struct device *dev,
2.     struct device *old_parent,
3.     struct device *new_parent)
4. {
5.     int error = 0;
6.     if (old_parent)
7.         sysfs_remove_link(&dev->kobj, "device");

```

```

8.  if (new_parent)
9.      error = sysfs_create_link(&dev->kobj, &new_parent->kobj,
10.          "device");
11.  return error;
12. #endif
13. }

```

device_move_class_links()只是一个内部函数，后面还有操纵它的那只手。这里的device_move_class_links显得很名不副实，并没用操作class中软链接的举动。这很正常，因为在sysfs中软链接是针对kobject来说的，所以即使位置变掉了，软链接还是很很准确地定位。

```

1.  /**
2.   * device_move - moves a device to a new parent
3.   * @dev: the pointer to the struct device to be moved
4.   * @new_parent: the new parent of the device (can be NULL)
5.   * @dpm_order: how to reorder the dpm_list
6.   */
7.  int device_move(struct device *dev, struct device *new_parent,
8.      enum dpm_order dpm_order)
9.  {
10.     int error;
11.     struct device *old_parent;
12.     struct kobject *new_parent_kobj;
13.
14.     dev = get_device(dev);
15.     if (!dev)
16.         return -EINVAL;
17.
18.     device_pm_lock();
19.     new_parent = get_device(new_parent);
20.     new_parent_kobj = get_device_parent(dev, new_parent);
21.
22.     pr_debug("device: '%s': %s: moving to '%s'\n", dev_name(dev),
23.         __func__, new_parent ? dev_name(new_parent) : "<NULL>");
24.     error = kobject_move(&dev->kobj, new_parent_kobj);
25.     if (error) {
26.         cleanup_glue_dir(dev, new_parent_kobj);
27.         put_device(new_parent);
28.         goto out;
29.     }
30.     old_parent = dev->parent;
31.     dev->parent = new_parent;
32.     if (old_parent)
33.         klist_remove(&dev->p->knode_parent);
34.     if (new_parent) {
35.         klist_add_tail(&dev->p->knode_parent,
36.             &new_parent->p->klist_children);
37.         set_dev_node(dev, dev_to_node(new_parent));
38.     }
39.
40.     if (!dev->class)
41.         goto out_put;
42.     error = device_move_class_links(dev, old_parent, new_parent);
43.     if (error) {
44.         /* We ignore errors on cleanup since we're hosed anyway... */
45.         device_move_class_links(dev, new_parent, old_parent);
46.         if (!kobject_move(&dev->kobj, &old_parent->kobj)) {
47.             if (new_parent)
48.                 klist_remove(&dev->p->knode_parent);
49.             dev->parent = old_parent;
50.             if (old_parent) {
51.                 klist_add_tail(&dev->p->knode_parent,
52.                     &old_parent->p->klist_children);
53.                 set_dev_node(dev, dev_to_node(old_parent));
54.             }
55.         }
56.         cleanup_glue_dir(dev, new_parent_kobj);
57.         put_device(new_parent);
58.         goto out;
59.     }
60.     switch (dpm_order) {

```

```
61. case DPM_ORDER_NONE:
62.     break;
63. case DPM_ORDER_DEV_AFTER_PARENT:
64.     device_pm_move_after(dev, new_parent);
65.     break;
66. case DPM_ORDER_PARENT_BEFORE_DEV:
67.     device_pm_move_before(new_parent, dev);
68.     break;
69. case DPM_ORDER_DEV_LAST:
70.     device_pm_move_last(dev);
71.     break;
72. }
73. out_put:
74.     put_device(old_parent);
75. out:
76.     device_pm_unlock();
77.     put_device(dev);
78.     return error;
79. }
```

device_move()就是将dev移到一个新的parent下。但也有可能这个parent是空的。大部分操作围绕在引用计数上，get_device()，put_device()。而且换了新的parent，到底要加到sysfs中哪个目录下，还要再调用get_device_parent()研究一下。主要的操作就是kobject_move()和device_move_class_links()。因为在sysfs中软链接是针对kobject来说的，所以即使位置变掉了，软链接还是很很准确地定位，所以在/sys/dev、/sys/bus、/sys/class中的软链接都不用变，这实在是sysfs的一大优势。除此之外，device_move()还涉及到电源管理的问题，device移动影响到dev在dpm_list上的位置，我们对此不了解，先忽略之。

```
1. void device_shutdown(void)
2. {
3.     struct device *dev, *devn;
4.
5.     list_for_each_entry_safe_reverse(dev, devn, &devices_kset->list,
6.         kobj.entry) {
7.         if (dev->bus && dev->bus->shutdown) {
8.             dev_dbg(dev, "shutdown\n");
9.             dev->bus->shutdown(dev);
10.        } else if (dev->driver && dev->driver->shutdown) {
11.            dev_dbg(dev, "shutdown\n");
12.            dev->driver->shutdown(dev);
13.        }
14.    }
15.    kobject_put(sysfs_dev_char_kobj);
16.    kobject_put(sysfs_dev_block_kobj);
17.    kobject_put(dev_kobj);
18.    async_synchronize_full();
19. }
```

这个device_shutdown()是在系统关闭时才调用的。它动用了很少使用的devices_kset，从而可以遍历到每个注册到sysfs上的设备，调用相应的总线或驱动定义的shutdown()函数。提起这个，还是在device_initialize()中将dev->kobj->kset统一设为devices_kset的。原来设备虽然有不同的parent，但kset还是一样的。这样我们就能理解/sys/devices下的顶层设备目录是怎么来的，因为没用parent，就在调用kobject_add()时将kset->kobj当成了parent，所以会直接挂在顶层目录下。这样的目录大致有pci0000:00、virtual等等。

看完了core.c，我有种明白机器人也是由零件组成的的感觉。linux设备驱动模型的大门已经打开了四分之一。随着分析的深入，我们大概也会越来越明白linux的良苦用心。



关注Linux公社（LinuxIDC.com）官方微信与QQ群，随机发放邀请码

上一页

1

2

3

4

5

6

7

8

9

下一页

7

【内容导航】

- 第1页：连通世界的list
- 第2页：原子性操作atomic_t
- 第3页：记录生命周期的kref
- 第4页：更强的链表klist
- 第5页：设备驱动模型的基石kobject
- 第6页：设备驱动模型之device
- 第7页：设备驱动模型之driver
- 第8页：设备驱动模型之bus
- 第9页：设备驱动模型之device-driver