

Learning Linux Device Drivers

Jeremiah Mahler (jmmahler@gmail.com)

August 11, 2013

Contents

1	Introduction	2
2	Hello, World	2
2.1	hello	2
2.2	param	4
3	Read/Write Data	5
3.1	data_chr	5
3.2	data_rw	8
3.3	data_sk	11
3.4	data_ioctl	12
3.5	/dev/null, /dev/zero	15
4	Sysfs	16
4.1	sysx_file	16
4.2	sysx_file2	16
4.3	sysx_group	16
4.4	sysx_ktype	16
4.5	sysx_ktype2	16
5	Concurrency	16
5.1	fifo_rw	16
5.2	fifo_sysfs	16
5.3	fifo_xxx	16
5.4	fifo_fix	16

1 Introduction

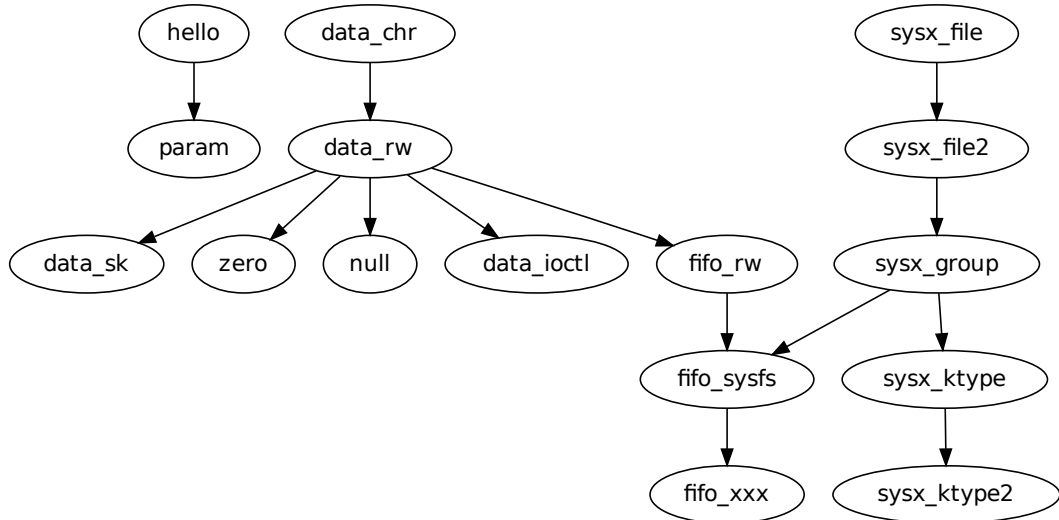


Figure 1: Hierarchy of kernel module examples. Simplest at the top downward to the more complex.

There are many excellent books about Linux device drivers^{123.4} However, in this authors experience, they were difficult to learn from. It certainly was not from their lack of detail. Clearly each of the authors have a profound understanding of the Linux kernel and their books reflect this. If anything it is due to this lack of simplicity.

The drivers described in this document aim to be simple and concise. Each one introduces as few concepts as possible. And each driver is a fully working example⁵. Many of the drivers are built in stages. Each stage introduces a new concept. And the changes are concisely described showing the differences (diff). Figure 1 shows the hierarchy of driver examples.

2 Hello, World

2.1 hello

The hello module (Listing 1) simply prints message when it is loaded and unload.

```
hello$ make
(should compile without error, resulting in hello.ko)
hello$ sudo insmod hello.ko
Hello, World
hello$ sudo rmmod hello
Goodbye, cruel world
```

¹J. Corbet, A. Rubini, and Greg. Kroah-Hartman. *Linux Device Drivers*. O'Reilly Media, 2009. ISBN: 9780596555382.

²S. Venkateswaran. *Essential Linux Device Drivers*. Pearson Education, 2008. ISBN: 9780132715812.

³R. Love. *Linux Kernel Development*. Developer's Library. Pearson Education, 2010. ISBN: 9780768696790.

⁴Robert. Love. *Linux System Programming: Talking Directly to the Kernel and C Library*. O'Reilly Media, 2013. ISBN: 9781449341541.

⁵Be creative with the examples. Try changing something and see what happens. Actively exploring in this way is a great way to solidify your understanding.

```
1 #include <linux/init.h>
2 #include <linux/module.h>
3
4 static int __init hello_init(void)
5 {
6     printk(KERN_ALERT "Hello , World\n");
7     return 0;
8 }
9
10 static void __exit hello_exit(void)
11 {
12     printk(KERN_ALERT "Goodbye, cruel world\n");
13 }
14
15 MODULE_AUTHOR("Jeremiah Mahler <jmmahler@gmail.com>");
16 MODULE_LICENSE("GPL");
17
18 module_init( hello_init );
19 module_exit( hello_exit );
```

Listing 1: Hello, World module in hello/hello.c

The `module_init` (line 18) and `module_exit` (line 19) tell the kernel which functions to call when this module is loaded (`insmod`) and unloaded (`rmmod`).

The `__init` (line 4) and `__exit` (line 10) are optional hints for the compiler. For example in the case of `__init`, this tells the kernel that it may discard the code after initialization has been completed.

Both the init function (line 4) and the exit function (line 10) are declared `static`. Since these functions are not meant to be used outside the scope of this file, declaring them `static` enforces this constraint.⁶

The `printk` statements are the `printf` of the kernel domain. There are various levels, in this case `KERN_ALERT` is used which will cause the messages to appear on the console. Notice that there is no comma between the level and the message.

The `MODULE_AUTHOR` and `MODULE_LICENSE` on lines 15 and 16 are optional but recommended.⁷ There are various other `MODULE_*` options as well (`linux/module.h`).

⁶Corbet, Rubini, and Kroah-Hartman, see n. 1, Pg. 52.

⁷Ibid., Pg. 51.

2.2 param

The `param` module expands upon the `hello` module to take a parameter specifying how many times to print the message.

```
param$ sudo insmod hello.ko howmany=2
Hello, World
Hello, World
param$ sudo rmmod hello
Goodbye, cruel world
Goodbye, cruel world
```

Listing 2 shows the differences between this parameterized hello world module and the previous `hello` module.

```
1  --- ../hello/hello.c      2013-08-09 12:23:58.222416131 -0700
2  +++ hello.c 2013-08-09 12:53:38.082434726 -0700
3  @@ -1,15 +1,28 @@
4  #include <linux/init.h>
5  #include <linux/module.h>
6  +#include <linux/moduleparam.h>
7  +
8  +static int howmany = 1;
9  +module_param(howmany, int, S_IRUGO);
10
11 static int __init hello_init(void)
12 {
13 - printk(KERN_ALERT "Hello , World\n");
14 + int i;
15 +
16 + for (i = 0; i < howmany; i++) {
17 +     printk(KERN_ALERT "Hello , World\n");
18 + }
19 +
20     return 0;
21 }
22
23 static void __exit hello_exit(void)
24 {
25 - printk(KERN_ALERT "Goodbye, cruel world\n");
26 + int i;
27 +
28 + for (i = 0; i < howmany; i++) {
29 +     printk(KERN_ALERT "Goodbye, cruel world\n");
30 + }
31 }
32
33 MODULE_AUTHOR("Jeremiah Mahler <jmmahler@gmail.com>");
```

Listing 2: `param$ diff -u hello.c ../hello/hello.c`

To use a parameter a global variable has been created named `howmany` on line 8. And on line 9 the `module_param` function is used to tell the kernel about this parameter ⁸.

On lines 13-19 and 25-30 it can be seen that the same message is printed `howmany` times.

⁸The `module_param` function create a `sysfs` entry in `/sys/module/parameters/howmany`. `sysfs` will be discussed in detail in later modules.

3 Read/Write Data

The `data` module allocates some memory which can then be read from and written to. This is accomplished as a character device and supports all the usual file operations.

3.1 `data_chr`

The first step is to construct the basic infrastructure for a character driver as shown in Listing 3. It doesn't do anything useful but it will simplify the description of upcoming drivers.

The `DEVICE_NAME` (line 8) is just a shortcut for the name which is used in several places.

Lines 10-17 are the global variables that will be used. The `struct data_dev` is the per device structure. Notice that a character device is placed inside.

The `file_operations` (line 19-21) in this case only defines the `.owner`. Upcoming modules will add references to the `open`, `close`, `read`, `write`, and `seek` functions to this structure.

`alloc_chrdev_region` (line 26) allocates a major and minor number for the character device.⁹ In this case only one major and minor pair is needed.

Functions such as `alloc_chrdev_region` may fail and when they do anything that has been created up to that point must be undone to ensure the kernel is left in a consistent state. A common way this is done is using `goto` statements which branch to different steps in the exit sequence.¹⁰ It can be seen that if `alloc_chrdev_region` fails its `goto` (line 29) will branch to line 66. Since nothing was created up to that point nothing has to be undone.

`class_create` (line 32) establishes a “class” for this module which is also represented in sysfs under `/sys/class/data`. This object will be used later as an argument to `device_create`.

Since the per device structure is just a pointer it must be allocated before it is used (line 32).

To establish the character device it must be initialized and added (lines 41-43). And finally the device is created (line 49). This device will now appear under `/dev/data0`.

⁹Corbet, Rubini, and Kroah-Hartman, see n. 1, Pg. 66.

¹⁰Ibid., Pg. 53.

```

1  #include <linux/cdev.h>
2  #include <linux/device.h>
3  #include <linux/fs.h>
4  #include <linux/module.h>
5  #include <linux/slab.h>
6  #include <linux/uaccess.h>
7
8  #define DEVICE_NAME "data"
9
10 static dev_t data_major;
11 struct class *data_class;
12 struct device *data_device;
13
14 struct data_dev {
15     struct cdev cdev;
16 } *data_devp;
17
18 struct file_operations data_fops = {
19     .owner = THIS_MODULE,
20 };
21
22 static int __init data_init(void)
23 {
24     int err = 0;
25
26     err = alloc_chrdev_region(&data_major, 0, 1, DEVICE_NAME);
27     if (err < 0) {
28         printk(KERN_WARNING "Unable to register device\n");
29         goto err_chrdev_region;
30     }
31
32     data_class = class_create(THIS_MODULE, DEVICE_NAME);
33
34     data_devp = kmalloc(sizeof(struct data_dev), GFP_KERNEL);
35     if (!data_devp) {
36         printk(KERN_WARNING "Unable to kmalloc data_devp\n");
37         err = -ENOMEM;
38         goto err_malloc_data_devp;
39     }
40
41     cdev_init(&data_devp->cdev, &data_fops);
42     data_devp->cdev.owner = THIS_MODULE;
43     err = cdev_add(&data_devp->cdev, data_major, 1);
44     if (err) {
45         printk(KERN_WARNING "cdev_add failed\n");
46         goto err_cdev_add;
47     }
48
49     data_device = device_create(data_class, NULL,
50                               MKDEV(MAJOR(data_major), 0), NULL, "data%d", 0);
51     if (IS_ERR(data_device)) {
52         printk(KERN_WARNING "device_create failed\n");
53         err = PTR_ERR(data_device);
54         goto err_device_create;
55     }

```

```

56
57     return 0;  /* success */
58
59 err_device_create:
60     cdev_del(&data_devp->cdev);
61 err_cdev_add:
62     kfree(data_devp);
63 err_malloc_data_devp:
64     class_destroy(data_class);
65     unregister_chrdev_region(data_major, 1);
66 err_chrdev_region:
67
68     return err;
69 }
70
71 static void __exit data_exit(void)
72 {
73     device_destroy(data_class, data_major);
74
75     cdev_del(&data_devp->cdev);
76
77     kfree(data_devp);
78
79     class_destroy(data_class);
80
81     unregister_chrdev_region(data_major, 1);
82 }
83
84 MODULE_AUTHOR("Jeremiah Mahler <jmmahler@gmail.com>");
85 MODULE_LICENSE("GPL");
86
87 module_init(data_init);
88 module_exit(data_exit);

```

Listing 3: data_chr driver.

3.2 data_rw

With the addition of read/write operations the device can be operated upon just like any other file. As an example the driver source code can be copied in to the device and then read back out. The result should be the same up to the maximum amount which in this case was 128 bytes. This maximum size is a `#define` inside the driver.

```
$ sudo dd if=data.c of=/dev/data0 bs=128 count=1
```

```
$ sudo dd if=/dev/data0 of=output bs=128 count=1
```

Listing 4 shows the differences compared to the previous `data_chr` driver. An array of bytes has been added to the per device structure along with the current offset (lines 7-17).

File operations for open, read, write and release have been added (lines 82-88). The release operation is called when a process closes the device file.

When the file is opened the open function (lines 19-29) is called. The `container_of` function (line 23) is used to obtain a parent structure from a child structure^{11,12}. Recall that the per device structure, `data_devp` contains a `cdev` (line 14). `container_of` makes it possible to obtain the `data_devp` from the `cdev`.

The open functions sets the offset to zero (line 24) when is the usual behavior when opening a file.

The open function also stores the device structure under `private_data` (line 26) so it is easy to access in the read/write functions.

When data is read from the device file the read function (lines 31-52) is called. Since the amount of data that can be read is limited by `MAX_DATA` the amount requested will be reduced if it is too large (lines 39-43). Then the `copy_to_user` function is used to attempt to transfer the data in to user space (lines 45-47). The `copy_to_user` also checks to make sure that destination it is transferring to is valid for the given process. If the transfer was a success the new offset is stored (line 49) and then the number of bytes that were successfully transferred are returned (line 51).

The write function (lines 54-70) has the same operation as read except in the opposite directory. Notice that the `copy_from_user` (line 68) function is used in this case.

And since nothing needs to be done when the device is closed, the release function (lines 77-80) simply returns success.

¹¹Corbet, Rubini, and Kroah-Hartman, see n. 1, Pg. 79.

¹²Greg. Kroah-Hartman. *container_of()*. [Online; accessed 10-August-2013]. 2005. URL: http://www.kroah.com/log/linux/container_of.html.

```

1  ——— ../data_chr/data.c  2013-08-10 10:17:01.359016284 -0700
2  +++ data.c  2013-08-11 00:28:57.774964943 -0700
3  @@ -6,6 +6,7 @@
4  #include <linux/uaccess.h>
5
6  #define DEVICENAME "data"
7  +#define MAXDATA 128
8
9  static dev_t data_major;
10 struct class *data_class;
11 @@ -13,10 +14,79 @@
12
13 struct data_dev {
14     struct cdev cdev;
15 +   char data[MAXDATA];
16 +   loff_t cur_ofs; // current offset
17 } *data_devp;
18
19 +static int data_open(struct inode* inode, struct file* filp)
20 +{
21 +   struct data_dev *data_devp;
22 +
23 +   data_devp = container_of(inode->i_cdev, struct data_dev, cdev);
24 +   data_devp->cur_ofs = 0;
25 +
26 +   filp->private_data = data_devp;
27 +
28 +   return 0;
29 +}
30 +
31 +static ssize_t data_read(struct file *filp, char __user *buf, size_t count,
32 +                        loff_t *f_pos)
33 +{
34 +   struct data_dev *data_devp = filp->private_data;
35 +   loff_t cur_ofs;
36 +   char *datp;
37 +   size_t left;
38 +
39 +   cur_ofs = data_devp->cur_ofs;
40 +   datp = data_devp->data;
41 +   left = MAXDATA - cur_ofs;
42 +
43 +   count = (count > left) ? left : count;
44 +
45 +   if (copy_to_user(buf, (void *) (datp + cur_ofs), count) != 0) {
46 +       return -EIO;
47 +   }
48 +
49 +   data_devp->cur_ofs = cur_ofs + count;
50 +
51 +   return count;
52 +}
53 +
54 +static ssize_t data_write(struct file *filp, const char __user *buf, size_t count,
55 +                        loff_t *f_pos)

```

```

56 +{
57 +   struct data_dev *data_devp = filp->private_data;
58 +   loff_t cur_ofs;
59 +   char *datp;
60 +   size_t left;
61 +
62 +   cur_ofs = data_devp->cur_ofs;
63 +   datp = data_devp->data;
64 +   left = MAXDATA - cur_ofs;
65 +
66 +   count = (count > left) ? left : count;
67 +
68 +   if (copy_from_user((void *) (datp + cur_ofs), buf, count) != 0) {
69 +       return -EIO;
70 +   }
71 +
72 +   data_devp->cur_ofs = cur_ofs + count;
73 +
74 +   return count;
75 +}
76 +
77 +static int data_release(struct inode *inode, struct file *filp)
78 +{
79 +   return 0;
80 +}
81 +
82 +struct file_operations data_fops = {
83 +   .owner = THIS_MODULE,
84 +   .open = data_open,
85 +   .read = data_read,
86 +   .write = data_write,
87 +   .release = data_release,
88 +};
89
90 static int __init data_init(void)

```

Listing 4: data_rw\$ diff -u ../data_chr/data.c data.c

3.3 data_sk

To add support for the seek operation requires the addition of one more function along with its corresponding entry in the file operations. Listing 5 shows the differences.

```
1  --- ../data_rw/data.c    2013-08-11 00:28:57.774964943 -0700
2  +++ data.c    2013-08-11 00:28:18.202964529 -0700
3  @@ -76,6 +76,35 @@
4      return count;
5  }
6
7  +static loff_t data_llseek(struct file *filp, loff_t offset, int orig)
8  +{
9  +    struct data_dev *data_devp = filp->private_data;
10 +    loff_t cur_ofs;
11 +
12 +    cur_ofs = data_devp->cur_ofs;
13 +
14 +    switch (orig) {
15 +        case SEEK_SET:
16 +            cur_ofs = offset;
17 +            break;
18 +        case SEEK_CUR:
19 +            cur_ofs += offset;
20 +            break;
21 +        case SEEK_END:
22 +            cur_ofs = MAXDATA + offset;
23 +            break;
24 +        default:
25 +            return -EINVAL;
26 +    }
27 +
28 +    if (cur_ofs < 0 || cur_ofs >= MAXDATA)
29 +        return -EINVAL;
30 +
31 +    data_devp->cur_ofs = cur_ofs;
32 +
33 +    return cur_ofs;
34 +}
35 +
36  static int data_release(struct inode *inode, struct file *filp)
37  {
38      return 0;
39  @@ -86,6 +115,7 @@
40      .open = data_open,
41      .read = data_read,
42      .write = data_write,
43  +    .llseek = data_llseek,
44      .release = data_release,
45  };
```

Listing 5: data_sk\$ diff -u ../data_rw/data.c data.c

3.4 data_ioctl

Using `ioctl()` for new designs is not recommended^{13,14}. Instead `sysfs` is preferred. While `/proc` is another option, it is also becoming obsolete in favor of `sysfs`. Nonetheless it is still used so it is worth knowing how it works.

Include with this driver is a test program (`data_ioctl/test/ioctlx.c`). The driver allows the reset, read and write of a single global variable (`x`) in the driver. The test program allows these operations to be performed.

```
data_ioctl$ cd test/
test$ sudo ./ioctlx 10      # set value
test$ sudo ./ioctlx        # read current value
10
test$ sudo ./ioctlx 0      # reset
test$ sudo ./ioctlx
0
test$
```

The `data_ioctl` driver has the fewest number of differences compared to the `data_chr` driver (Section 3.1) as shown in Listing 6. The open (lines 10-19) and release (lines 75-78) are the same as those for the `data_rw` driver (Section 3.2). The only code of interest is for the `data_ioctl` function (lines 29-73).

The first thing to notice is the use of "magic" (lines 23-25, 37). Magic is used to make the `ioctl` calls unique across the entire system which helps prevent inadvertent configuration if the wrong device is opened.¹⁵ The user program must also contain the corresponding magic values as is done in the test program (Listing 7).

The three defines (lines 23-24) describe the supported `ioctl` operations. Any number of additional operations can be added. And it can be seen that there is an operation that has no data (`_IO`), reads data (`_IOR`) and writes data (`_IOW`). The `DATA_IOC_MAXNR` (line 27) is used as a sanity check later (line 41).

And the switch statement (lines 53-70) process the `ioctl` commands. In this case all the operations involve the global variable `x`. It is either reset, read or written.

¹³Corbet, Rubini, and Kroah-Hartman, see n. 1, Pg. 156.

¹⁴Love, *Linux Kernel Development*, see n. 3.

¹⁵Corbet, Rubini, and Kroah-Hartman, see n. 1, Pg. 158.

```

1  --- ../data_chr/data.c  2013-08-10 10:17:01.359016284 -0700
2  +++ data.c  2013-08-11 11:08:06.287053188 -0700
3  @@ -15,14 +15,92 @@
4      struct cdev cdev;
5  } *data_devp;
6
7  +
8  +int x;
9  +
10 +static int data_open(struct inode *inode, struct file *filp)
11 +{
12 +    struct data_dev *data_devp;
13 +
14 +    data_devp = container_of(inode->i_cdev, struct data_dev, cdev);
15 +
16 +    filp->private_data = data_devp;
17 +
18 +    return 0;
19 +}
20 +
21 +#define DATA_IOC_MAGIC 'm'
22 +
23 +#define DATA_IOCRESET _IO(DATA_IOC_MAGIC, 1)
24 +#define DATA_IOCWX _IOW(DATA_IOC_MAGIC, 2, int)
25 +#define DATA_IOCRIX _IOR(DATA_IOC_MAGIC, 3, int)
26 +
27 +#define DATA_IOC_MAXNR 3
28 +
29 +static long data_ioctl(struct file *filp, unsigned int cmd,
30 +    unsigned long arg)
31 +{
32 +    int err = 0;
33 +    int retval = 0;
34 +
35 +    //struct data_dev *data_dev = filp->private_data;
36 +
37 +    if (_IOC_TYPE(cmd) != DATA_IOC_MAGIC) {
38 +        printk(KERN_ALERT "invalid ioctl magic\n");
39 +        return -ENOTTY;
40 +    }
41 +    if (_IOC_NR(cmd) > DATA_IOC_MAXNR) {
42 +        printk(KERN_ALERT "ioctl beyond maximum\n");
43 +        return -ENOTTY;
44 +    }
45 +
46 +    if (_IOC_DIR(cmd) & _IOC_READ)
47 +        err = !access_ok(VERIFY_WRITE, (void __user *) arg, _IOC_SIZE(cmd));
48 +    else if (_IOC_DIR(cmd) & _IOC_WRITE)
49 +        err = !access_ok(VERIFY_READ, (void __user *) arg, _IOC_SIZE(cmd));
50 +    if (err)
51 +        return -EFAULT;
52 +
53 +    switch (cmd) {
54 +        case DATA_IOCRESET:
55 +            /* takes no argument, sets values to default */

```

```

56 +         x = 0;
57 +         break;
58 +     case DATAIOCRX:
59 +         /* read integer */
60 +         retval = __put_user(x, (int __user *) arg);
61 +         break;
62 +     case DATAIOCWX:
63 +         /* write integer */
64 +         retval = __get_user(x, (int __user *) arg);
65 +         break;
66 +     default:
67 +         return -ENOTTY; /* POSIX standard */
68 +         //return -EINVAL; /* common */
69 +         /* Pg. 161 Linux Device Drivers (2005) */
70 +     }
71 +
72 +     return retval;
73 + }
74 +
75 + int data_release(struct inode *inode, struct file *filp)
76 + {
77 +     return 0;
78 + }
79 +
80 + struct file_operations data_fops = {
81 +     .owner = THIS_MODULE,
82 +     .open = data_open,
83 +     .unlocked_ioctl = data_ioctl,
84 +     .release = data_release,
85 + };
86 +
87 + static int __init data_init(void)
88 + {
89 +     int err = 0;
90 +
91 +     x = 0;
92 +
93 +     err = alloc_chrdev_region(&data_major, 0, 1, DEVICENAME);
94 +     if (err < 0) {
95 +         printk(KERN_WARNING "Unable to register device\n");

```

Listing 6: data_ioctl\$ diff -u ../data_chr/data.c data.c

```

9 #define DEVFILE "/dev/data0"
10
11 #define DATAIOC_MAGIC 'm'
12
13 #define DATAIOCRESET _IO(DATAIOC_MAGIC, 1)
14 #define DATAIOCWX _IOW(DATAIOC_MAGIC, 2, int)
15 #define DATAIOCRX _IOR(DATAIOC_MAGIC, 3, int)
16
17 int main(int argc, char* argv[])
18 {
19     int devfd;

```

Listing 7: Corresponding "magic" in user program.

3.5 /dev/null, /dev/zero

From what has been described so far it is easy construct a driver for the well known /dev/null and /dev/zero devices.

The zero device is even simpler than the data_rw example (Section 3.2). The only real difference, other than names (data -> null), is the read and write operations as shown in Listing 8.

```
30 static ssize_t null_read(struct file *filp, char __user *buf,
31                          size_t count, loff_t *f_pos)
32 {
33     return 0;
34 }
35
36 static ssize_t null_write(struct file *filp, const char __user *buf,
37                           size_t count, loff_t *f_pos)
38 {
39     return count;
40 }
```

Listing 8: /dev/null read and write functions.

The read and write functions for the zero driver are also quite simple (Listing 9). The one new addition is the `clear_user` function. It behaves like the `copy_to_user` function except that it simply zeros out the users buffer.

```
30 static ssize_t zero_read(struct file *filp, char __user *buf,
31                          size_t count, loff_t *f_pos)
32 {
33     if (clear_user((void __user *) buf, count) > 0) {
34         return -EFAULT;
35     }
36
37     return count;
38 }
39
40 ssize_t zero_write(struct file *filp, const char __user *buf,
41                   size_t count, loff_t *f_pos)
42 {
43     return count;
44 }
```

Listing 9: /dev/zero read and write functions.

4 Sysfs

4.1 sysx_file

4.2 sysx_file2

4.3 sysx_group

4.4 sysx_ktype

4.5 sysx_ktype2

5 Concurrency

5.1 fifo_rw

5.2 fifo_sysfs

5.3 fifo_xxx

5.4 fifo_fix

References

Corbet, J., A. Rubini, and Greg. Kroah-Hartman. *Linux Device Drivers*. O'Reilly Media, 2009. ISBN: 9780596555382.

Kroah-Hartman, Greg. *container_of()*. [Online; accessed 10-August-2013]. 2005. URL: http://www.kroah.com/log/linux/container_of.html.

Love, R. *Linux Kernel Development*. Developer's Library. Pearson Education, 2010. ISBN: 9780768696790.

Love, Robert. *Linux System Programming: Talking Directly to the Kernel and C Library*. O'Reilly Media, 2013. ISBN: 9781449341541.

Venkateswaran, S. *Essential Linux Device Drivers*. Pearson Education, 2008. ISBN: 9780132715812.