

20.03.2016

seq-file

/proc Filesystem » ADMIN Magazine

april 2014

Eva-Katharina Kunst

before dppa 3.10

proc creatadata()

# ADMIN

Network & Security

Subscribe To Our Newsletter:

Email

Search

Search

News

Articles

Tech Tools

Subscribe

Archive

Whitepapers

Digisub

Write for Us!

Newsletter

Shop

Cloud Computing

Virtualization

HPC

Linux

Windows

Security

Monitoring

Databases

all Topics...

Home » Archive » 2014 » Issue 23: 10 Ti... » Kernel and driv...

→ → → Login

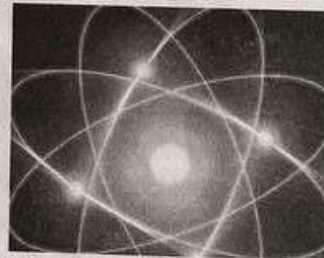
## Kernel and driver development for the Linux kernel

### Core Technology

Article from ADMIN 23/2014

By Jürgen Quade and Eva-Katharina Kunst

The /proc filesystem facilitates the exchange of current data between the system and user. To access the data, you simply read and write to a file. This mechanism is the first step for understanding kernel programming. ü



Lead image © psdesign1, Fotolia.com

In keeping with the central Unix philosophy that everything is a file, Linux systems publish system information through the virtual filesystems /proc and /sys (see the "Proc Filesystem" box). This brilliant mechanism gives the user read and write access to system internals with the read() and write() functions.

### Proc Filesystem

The /proc virtual filesystem folders and files are not stored on a hard disk; rather, the kernel creates them dynamically on access. The term "proc" derives from "processes"; thus, it is clear that the /proc filesystem primarily provides information about computing processes (Figure 1).

For each computational process, the kernel creates a new directory with the process identification number as its name. Below the directory, you find extensive information about the job, including the call parameters (cmdline), shared file descriptors (fd) and environment variables (environ), and process statistics.

In addition to information about computing processes, the kernel uses the virtual filesystem to share



information with the user system and receive configurations. All data for the current CPU is in /proc/cpuinfo, interrupt sources and the frequency of their occurrence is under /proc/interrupts, and info for activated device drivers with their device numbers is under /proc/devices. Writing 1 to /proc/sys/net/ipv4/ip\_forward enables routing in the Linux kernel and configures the watchdog feature when written to /proc/sys/kernel/watchdog. KL: Navigating the /proc filesystem is both useful and instructive. A proc file name reveals how important a file in the /proc filesystem is. Hardware-related information is located – with the exception of CPU information – in the /sys filesystem. This is useful for driver programmers as a platform for exchanging information.

*but it's not a good measure of CPU usage*

*/proc \$ ls*

quedeget-mobil:/proc \$ ls									
1	1213	147	12	2281	2473	3	42	5828	9
10	12347	1470	1000	2204	2487	30	4213	5823	903
10495	12368	148	2002	2204	23	3014	41	628	9024
10519	12365	146	2003	2204	2300	3038	440	6108	904
10540	12376	15	2000	2226	2309	3039	4412	69	9136
10561	12381	150	2019	2226	2331	31	45	7	932
10582	12404	151	2023	2255	2394	3112	44	73	9330
10584	12525	152	2019	2300	34	3132	47	71	943
10586	1234	16	2582	2392	2402	2306	479	7183	95
10547	12543	1519	2019	23	3418	2381	48	72	97
10588	11594	17	254	2318	2412	35	49	7206	9445
1064	1240	1740	1640	2326	2437	335	5	73	9347
10737	12617	18	2047	2329	2443	34	59	74	9451
10736	1278	1834	2049	2335	2471	345	5012	743	9444
10940	13	1889	205	2336	2738	35	5104	7450	9450
13	1304	18	2050	2352	2751	36	525	75	945
1282	1257	1844	2054	2313	2752	37	534	76	946
1278	12109	1944	2064	2314	28	2787	5472	75	9472
1278	1222	1943	2068	2342	2837	38	55	8	948
12792	1053	1947	21	24	2845	39	56	871	949
1285	1552	1970	2177	249	2881	294	568	872	950
1283	1299	1964	2100	2481	2904	397	57	813	951
1293	34	1990	2188	2411	29	40	1720	837	952
13	1437	1999	2190	2417	2903	4072	3751	843	953
1224	144	2	2191	2468	2918	41	3752	891	954
quedeget-mobil:/proc \$									

Figure 1: Directories and files in the /proc directory.

The shell, in turn, maps these access functions to the commands `cat` (read) and `echo` (write). To read the number of interrupts that have been triggered since the last boot, then, you can simply enter the `cat /proc/interrupts` command in a terminal. To pass network traffic through, on the other hand, the superuser only needs to write 1 to the `/proc/sys/net/ipv4/ip_forward` file:

```
sudo echo 1 > /proc/sys/net/ipv4/ip_forward
```

Such functions are very easy to use in scripts.

The /proc filesystem is also useful for your first steps in kernel programming: With fewer than 50 lines of code, you can make the compiler generate a module that outputs the famous "Hello World" string when accessing a proc file (Listing 1).

### Listing 1

#### Simple Proc File

```
01 #include <linux/module.h>
02 #include <linux/proc_fs.h>
03 #include <linux/seq_file.h>
04
05 #define PROC_FILE_NAME "Hello_World"
06 static struct proc_dir_entry *proc_file;
```



```

07 static char *output_string;
08
09 static int prochello_show( struct seq_file *m, void *v )
10 {
11     int error = 0;
12     error = seq_printf( m, "%s\n", output_string);
13     return error;
14 }
15
16 static int prochello_open(struct inode *inode, struct file *file)
17 {
18     return single_open(file, prochello_show, NULL);
19 }
20
21
22 static const struct file_operations prochello_fops = {
23     .owner = THIS_MODULE,
24     .open = prochello_open,
25     .release = single_release,
26     .read = seq_read,
27 };
28
29 static int __init prochello_init(void)
30 {
31     output_string = "Hello World";
32     proc_file = proc_create_data( PROC_FILE_NAME, S_IRUGO, NULL, &prochello_fops, NULL );
33     if (!proc_file)
34         return -ENOMEM;
35     return 0;
36 }
37
38 static void __exit prochello_exit(void)
39 {
40     if( proc_file )
41         remove_proc_entry( PROC_FILE_NAME, NULL );
42 }
43
44 module_init( prochello_init );
45 module_exit( prochello_exit );
46 MODULE_LICENSE("GPL");

```

*биднавогара / seq-file.h*

*файл - последовательность*

*структура данных  
содержит  
ссылки на  
функции*

*1) 2) 3) 4) 5)*

*права доступа*

*определяет  
уровень доступа  
к файлу*

*создает файл*

*если NULL  
ниже /proc*

*удаляет созданный файл*

Modules are known extensions of the Linux kernel and are always constructed in a similar way (Figure 2).



*Чтобы управлять модулем, вам нужно 2 вещи: подпрограммы, которые реализуют функциональность модуля, и структуру данных, которая принимает адреса идентификаторов подпрограмм.*

To manage a module you need two things: routines that implement the functionality of the module and a data structure that receives the addresses of the main routines.

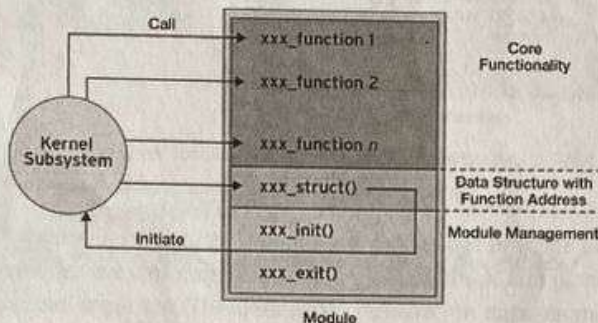


Figure 2: Kernel modules are divided into three categories.

*Для управленческих функций — произвольное имя, когда команда insmod загружает модуль в ядро, и когда модуль выгружается из ядра с помощью rmmod.*

The two management functions are `<xxx>_init()` and `<xxx>_exit()`, where `<xxx>` is the module or an arbitrary name. The `<xxx>_init()` function runs when the `insmod` command loads the module in the kernel, and `<xxx>_exit()` runs when the module is unloaded using `rmmod`.

*Значение*

The main task of the `<xxx>_init()` function is to pass the data structure containing the addresses of the relevant module routines to a kernel subsystem. The module logs in to the desired subsystem(s) and itself. In this way, the kernel knows the relevant module routines and can activate them. Similarly `<xxx>_exit()` is for signing off.

*Ка 1. три области типичных модулей: административные функции, модульные функции и структура данных для регистрации.*

In Listing 1, the three areas of a typical Linux module are clearly visible: the administrative functions `prochello_init()` and `prochello_exit()` (lines 29 and 38), the module routines `prochello_show()` and `prochello_open()` (lines 9 and 17), and the data structure used for logging in to the kernel `prochello_fops` (line 22).

## Transferring Access Rights

*Она передает имя файла в proc и разрешения для этого файла. Потому что права доступа закодированы в битовых шаблонах, вы используете здесь (S\_IRUGO). В примере, владелец (user), группа, и все остальные получают права на чтение. Третий параметр определяет, где в дереве каталогов файловой системы /proc файл должен быть создан. NULL создает файл ниже /proc. Наконец, четвертый параметр передает структуру данных, которая содержит функции доступа к новому файлу. Всего требуется три функции для доступа: open(), read(), и release().*

Within the `prochello_init()` init function, the module logs in to the Linux kernel's proc subsystem by using the `proc_create_data()` function. It passes in the name of the proc file and the permissions for this file. Because the access rights are encoded as bit patterns, you use defines here (`S_IRUGO`). In the example, the owner (user), the group, and all others are given read permissions. The third parameter specifies where in the directory tree of the `/proc` filesystem the file is to be generated. `NULL` creates the file below `/proc`. Finally, the fourth parameter transfers the data structure that contains the access functions to the new proc file. A total of three functions are needed for read access: `open()`, `read()`, and `release()`.

*Использовать*

However, you only need to implement the open routine yourself (`prochello_open()`). If you want, you can also implement the other routines (`read()` and `release()`), but this is complicated and only useful in exceptional cases. The kernel developers have built a fault-resistant intermediate layer for this, which reduces the complexity of the data exchange to something like `printf`. The basic idea of the intermediate layer is that the module programmer writes all the data to a sufficiently large memory area (Figure 3).

*использовать proc\_create\_data(), чтобы избежать сложной "гойки"*

use `proc_create_data()`, to avoid the complex "goiky"



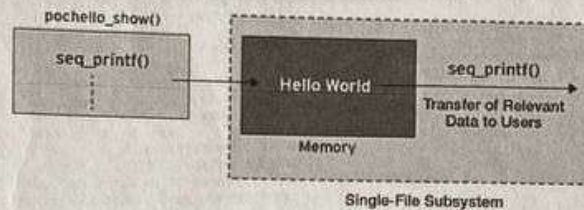


Figure 3: The show function uses seq\_printf() to write the data to RAM; the single-file subsystem copies the relevant to the user.

Further processing of the data – in particular, transferring data to the proc file users – is handled by the intermediate layer. This transfer is more complicated than it appears at first glance because the user might not read all the data, just a subset. The user might retrieve the data via multiple read calls.

The intermediate layer is called the "single-file subsystem" and is a special form of sequence file. Above all, it implements the read and the release functions (equivalent to the close() system call), so they can be used without changes when accessing the proc file. To create the single-file instance of the module, programmers call the single\_open() function, which passes the address to a routine, typically called show(). This show routine is given the address of a memory page (in this example, the data to be output by seq\_printf()); seq\_printf() can be called a more or less arbitrary number of times within the show function.

### Limited Single Files

Single-file instances are not intended for writing large amounts of data. Instead, the space is limited to 64KB, but that's enough for most tasks. Also, seq\_printf() monitors overflows on every call. If that were to happen, the subsystem would automatically grow the memory buffer (e.g., by creating a new one and copying) and would then write the data. If scaling is impossible, seq\_printf() returns a negative error code. Professional programmers evaluate the return value, of course. The single\_open() function, which receives the address of the show function, is called in the open function for the proc file (Listing 1, line 19).

Figure 4 shows how you can compile the source code in Listing 1 using the Makefile (Listing 2) followed by insmod, which loads the generated module into the kernel, and finally cat, which tests the results. If you don't want to create the proc file directly below /proc but in a subfolder, you can use proc\_mkdir to create a directory. The function returns a pointer to a data structure of the type proc\_dir\_entry, which represents the newly created directory.



```

quade@ezs-nobil:~/tmp/procfs$ ls
Makefile  prochello.c
quade@ezs-nobil:~/tmp/procfs$ make
make -C /lib/modules/3.15.0-rc5/build M=/tmp/procfs modules
make[1]: Betrete Verzeichnis '/usr/src/mainline/linux'
CC [M] /tmp/procfs/prochello.o
Building modules, stage 2.
MODPOST 1 modules
CC /tmp/procfs/prochello.mod.o
LD [M] /tmp/procfs/prochello.ko
make[1]: Verlasse Verzeichnis '/usr/src/mainline/linux'
quade@ezs-nobil:~/tmp/procfs$ sudo su
[sudo] password for quade:
root@ezs-nobil:~/tmp/procfs# insmod prochello.ko
root@ezs-nobil:~/tmp/procfs# ls -l /proc/Hello_World
-r--r--r-- 1 root root 0 Jul 30 18:44 /proc/Hello_World
root@ezs-nobil:~/tmp/procfs# cat /proc/Hello_World
Hello World
root@ezs-nobil:~/tmp/procfs#

```

Figure 4: Generating and using the kernel module.

## Listing 2

### Make

```

01 ifneq ($(KERNELRELEASE),)
02 obj-m    := prochello.o
03
04 else
05 KDIR     := /lib/modules/$(shell uname -r)/build
06 PWD      := $(shell pwd)
07
08 default:
09     $(MAKE) -C $(KDIR) M=$(PWD) modules
10 endif

```

The `proc_mkdir()` function expects two parameters. The first contains the directory name as a string, and the second is an identifier that says where to create the new directory. If the second parameter is equal to `NULL`, the new folder is created below the `/proc` directory.

A cardinal error is to forget to delete the `proc` directories and files you created if you no longer need them or to remove the entire module. This is what the `remove_proc_entry()` function does. In addition to the name of the file to remove, it also supplies an identifier for the subdirectory that contains the `proc` file.

## Hands On

The subsystem does not support the ability to transfer configuration data to the kernel by writing to a `proc` file. This is where the module programmer needs to lend a hand. To start, you need memory to cache the data to be written to the kernel for evaluation.

If enough memory is available, you copy the data to be written from userspace to the memory buffer. You can then analyze the data and perform the action requested by the user. It is important not to transfer more data than you can buffer between userspace and kernel-space. The `min()` function ensures that







Finally, four lines need to be added to the program header:

```
#include <asm/uaccess.h>
static char kernel_buffer[256];
#define TEXT_GERMAN "Hallo Welt"
#define TEXT_ENGLISH "Hello World"
```

After making these changes and compiling, then unloading the old module version and reloading the driver, write access should be possible.

The code for implementing proc files is a good template for your own development. Essentially, you need to change the name of the proc file and the show function.

If you have extensive output that changes frequently, you will probably want to look to sequence files instead. Unfortunately, earlier writings on this topic are not really that useful as a guide (see the box "Changes in Kernel 3.10").

### Changes in Kernel 3.10

Kernel 3.10 no longer supports the `create_proc_entry()` function:

```
proc_file = create_proc_entry("example_file", S_IRUGO, proc_dir );
if (proc_file) {
    proc_file->read_proc = proc_read;
    proc_file->data = NULL;
}
```

Instead, programmers use the function presented in this article, `proc_create_data()`, which uses different parameterization. Whereas individual elements of the `proc_dir_entry` data structure had to be initialized in some of the earlier versions of the kernel, developers working with a more recent kernel version reserve a structure that is already familiar from driver development, `struct file_operations`, and assign the access methods (`open()`, `read()`, `release()`):

```
static struct file_operations example_proc_fops = {
    .owner = THIS_MODULE,
    .open = example_proc_open,
    .read = example_proc_read,
    .release = example_proc_release,
}
[...]
static int __init example_proc_init(void)
{
    proc_file = proc_create_data( "example_file", S_IRUGO, proc_dir,
                                &example_proc_fops, NULL );
    [...]
}
```






метод организации от предыдущих версий.

The read() and write() access methods differ compared with previous versions. The former peof parameter, which used to indicate that the system had read all the data, no longer exists. In contrast, the current version writes data to the memory address passed in as a parameter and then returns the number of characters written (if the single-file subsystem is not used, unlike the description in this article).

### The Author

Eva-Katharina Kunst has been an open source fan since the early days of Linux. Jürgen Quade is a professor at the Niederrhein University; he published his third Linux book *Learning Embedded Linux with the Raspberry Pi* in April 2014.

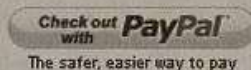
**Buy this article as PDF**

Share / Save   

Express-Checkout as PDF

Price \$2.95

(incl. VAT)



**Buy ADMIN Magazine**

SINGLE ISSUES

Print Issues

Digital Issues

SUBSCRIPTIONS

Print Subs