

The Linux Kernel Module Programming Guide

Peter Jay Salzman, Michael Burian, Ori Pomerantz, Bob Mottram, Jim Huang

Introduction

The Linux Kernel Module Programming Guide is a free book; you may reproduce and/or modify it under the terms of the Open Software License, version 3.0.

This book is distributed in the hope it will be useful, but without any warranty, without even the implied warranty of merchantability or fitness for a particular purpose.

The author encourages wide distribution of this book for personal or commercial use, provided the above copyright notice remains intact and the method adheres to the provisions of the Open Software License. In summary, you may copy and distribute this book free of charge or for a profit. No explicit permission is required from the author for reproduction of this book in any medium, physical or electronic.

Derivative works and translations of this document must be placed under the Open Software License, and the original copyright notice must remain intact. If you have contributed new material to this book, you must make the material and source code available for your revisions. Please make revisions and updates available directly to the document maintainer, Peter Jay Salzman <p@dirac.org>. This will allow for the merging of updates and provide consistent revisions to the Linux community.

If you publish or distribute this book commercially, donations, royalties, and/or printed copies are greatly appreciated by the author and the Linux Documentation Project (LDP). Contributing in this way shows your support for free software and the LDP. If you have questions or comments, please contact the address above.

Authorship

The Linux Kernel Module Programming Guide was originally written for the 2.2 kernels by Ori Pomerantz. Eventually, Ori no longer had time to maintain the document. After all, the Linux kernel is a fast moving target. Peter Jay Salzman took over maintenance and updated it for the 2.4 kernels. Eventually, Peter no longer had time to follow developments with the 2.6 kernel, so Michael

Burian became a co-maintainer to update the document for the 2.6 kernels. Bob Mottram updated the examples for 3.8+ kernels. Jim Huang upgraded to recent kernel versions (v5.x) and revise LaTeX scripts.

Acknowledgements

The following people have contributed corrections or good suggestions: Ignacio Martin, David Porter, Daniele Paolo Scarpazza, Dimo Velez, Francois Audeon, Horst Schirmeier, and Roman Lakeev.

What Is A Kernel Module?

So, you want to write a kernel module. You know C, you have written a few normal programs to run as processes, and now you want to get to where the real action is, to where a single wild pointer can wipe out your file system and a core dump means a reboot.

What exactly is a kernel module? Modules are pieces of code that can be loaded and unloaded into the kernel upon demand. They extend the functionality of the kernel without the need to reboot the system. For example, one type of module is the device driver, which allows the kernel to access hardware connected to the system. Without modules, we would have to build monolithic kernels and add new functionality directly into the kernel image. Besides having larger kernels, this has the disadvantage of requiring us to rebuild and reboot the kernel every time we want new functionality.

Kernel module package

Linux distributions provide the commands *modprobe*, *insmod* and *depmod* within a package.

On Ubuntu/Debian:

```
sudo apt-get install build-essential kmod
```

On Arch Linux:

```
sudo pacman -S gcc kmod
```

What Modules are in my Kernel?

To discover what modules are already loaded within your current kernel use the command **lsmod**.

```
sudo lsmod
```

Modules are stored within the file `/proc/modules`, so you can also see them with:

```
sudo cat /proc/modules
```

This can be a long list, and you might prefer to search for something particular. To search for the *fat* module:

```
sudo lsmod | grep fat
```

Do I need to download and compile the kernel?

For the purposes of following this guide you don't necessarily need to do that. However, it would be wise to run the examples within a test distribution running on a virtual machine in order to avoid any possibility of messing up your system.

Before We Begin

Before we delve into code, there are a few issues we need to cover. Everyone's system is different and everyone has their own groove. Getting your first "hello world" program to compile and load correctly can sometimes be a trick. Rest assured, after you get over the initial hurdle of doing it for the first time, it will be smooth sailing thereafter.

1. Modversioning. A module compiled for one kernel will not load if you boot a different kernel unless you enable `CONFIG_MODVERSIONS` in the kernel. We will not go into module versioning until later in this guide. Until we cover modversions, the examples in the guide may not work if you are running a kernel with modversioning turned on. However, most stock Linux distribution kernels come with it turned on. If you are having trouble loading the modules because of versioning errors, compile a kernel with modversioning turned off.
2. Using X Window System. It is highly recommended that you extract, compile and load all the examples this guide discusses. It is also highly recommended you do this from a console. You should not be working on this stuff in X Window System.

Modules can not print to the screen like `printf()` can, but they can log information and warnings, which ends up being printed on your screen, but only on a console. If you `insmod` a module from an `xterm`, the information and warnings will be logged, but only to your `systemd` journal. You will not see it unless you look through your `journalctl`. See 4 for details. To have immediate access to this information, do all your work from the console.

Headers

Before you can build anything you'll need to install the header files for your kernel.

On Ubuntu/Debian:

```
sudo apt-get update apt-cache search linux-headers-‘uname -r‘
```

On Arch Linux:

```
sudo pacman -S linux-libre-headers
```

This will tell you what kernel header files are available. Then for example:

```
sudo apt-get install kmod linux-headers-5.4.0-80-generic
```

Examples

All the examples from this document are available within the *examples* subdirectory.

If there are any compile errors then you might have a more recent kernel version or need to install the corresponding kernel header files.

Hello World

The Simplest Module

Most people learning programming start out with some sort of "*hello world*" example. I don't know what happens to people who break with this tradition, but I think it is safer not to find out. We will start with a series of hello world programs that demonstrate the different aspects of the basics of writing a kernel module.

Here is the simplest module possible.

Make a test directory:

```
mkdir -p /develop/kernel/hello-1 cd /develop/kernel/hello-1
```

Paste this into your favorite editor and save it as **hello-1.c**:

Now you will need a Makefile. If you copy and paste this, change the indentation to use *tabs*, not spaces.

```
obj-m += hello-1.o
```

```
all: make -C /lib/modules/$(shelluname -r)/buildM=$(PWD) modules
```

```
clean: make -C /lib/modules/$(shelluname -r)/buildM=$(PWD) clean
```

And finally just:

```
make
```

If all goes smoothly you should then find that you have a compiled **hello-1.ko** module. You can find info on it with the command:

```
sudo modinfo hello-1.ko
```

At this point the command:

```
sudo lsmod | grep hello
```

should return nothing. You can try loading your shiny new module with:

```
sudo insmod hello-1.ko
```

The dash character will get converted to an underscore, so when you again try:

```
sudo lsmod | grep hello
```

you should now see your loaded module. It can be removed again with:

```
sudo rmmod hello_1
```

Notice that the dash was replaced by an underscore. To see what just happened in the logs:

```
journalctl -since "1 hour ago" | grep kernel
```

You now know the basics of creating, compiling, installing and removing modules. Now for more of a description of how this module works.

Kernel modules must have at least two functions: a "start" (initialization) function called **init_module()** which is called when the module is insmoded into the kernel, and an "end" (cleanup) function called **cleanup_module()** which is called just before it is removed from the kernel. Actually, things have changed starting with kernel 2.3.13. You can now use whatever name you like for the start and end functions of a module, and you will learn how to do this in Section 2.3. In fact, the new method is the preferred method. However, many people still use **init_module()** and **cleanup_module()** for their start and end functions.

Typically, **init_module()** either registers a handler for something with the kernel, or it replaces one of the kernel functions with its own code (usually code to do something and then call the original function). The **cleanup_module()** function is supposed to undo whatever **init_module()** did, so the module can be unloaded safely.

Lastly, every kernel module needs to include **linux/module.h**. We needed to include **linux/kernel.h** only for the macro expansion for the **pr_alert()** log level, which you'll learn about in Section 2.1.1.

1. A point about coding style. Another thing which may not be immediately obvious to anyone getting started with kernel programming is that indentation within your code should be using **tabs** and **not spaces**. It is one of the coding conventions of the kernel. You may not like it, but you'll need to get used to it if you ever submit a patch upstream.
2. Introducing print macros. In the beginning there was **printk**, usually followed by a priority such as **KERN_INFO** or **KERN_DEBUG**. More recently this can also be expressed in abbreviated form using a set of print macros,

such as **pr_info** and **pr_debug**. This just saves some mindless keyboard bashing and looks a bit neater. They can be found within **linux/printk.h**. Take time to read through the available priority macros.

3. About Compiling. Kernel modules need to be compiled a bit differently from regular userspace apps. Former kernel versions required us to care much about these settings, which are usually stored in Makefiles. Although hierarchically organized, many redundant settings accumulated in sublevel Makefiles and made them large and rather difficult to maintain. Fortunately, there is a new way of doing these things, called kbuild, and the build process for external loadable modules is now fully integrated into the standard kernel build mechanism. To learn more on how to compile modules which are not part of the official kernel (such as all the examples you will find in this guide), see file **Documentation/kbuild/modules.rst**.

Additional details about Makefiles for kernel modules are available in **Documentation/kbuild/makefiles.rst**. Be sure to read this and the related files before starting to hack Makefiles. It will probably save you lots of work.

Here is another exercise for the reader. See that comment above the return statement in **init_module()**? Change the return value to something negative, recompile and load the module again. What happens?

Hello and Goodbye

In early kernel versions you had to use the **init_module** and **cleanup_module** functions, as in the first hello world example, but these days you can name those anything you want by using the **module_init** and **module_exit** macros. These macros are defined in **linux/init.h**. The only requirement is that your init and cleanup functions must be defined before calling the those macros, otherwise you'll get compilation errors. Here is an example of this technique:

So now we have two real kernel modules under our belt. Adding another module is as simple as this:

```
obj-m += hello-1.o obj-m += hello-2.o
```

```
all: make -C /lib/modules/$(shelluname -r)/buildM=$(PWD) modules
```

```
clean: make -C /lib/modules/$(shelluname -r)/buildM=$(PWD) clean
```

Now have a look at **linux/drivers/char/Makefile** for a real world example. As you can see, some things get hardwired into the kernel (**obj-y**) but where are all those **obj-m** gone? Those familiar with shell scripts will easily be able to spot them. For those not, the **obj-\$(CONFIG_FOO)** entries you see everywhere expand into **obj-y** or **obj-m**, depending on whether the **CONFIG_FOO** variable has been set to y or m. While we are at it, those were exactly the kind of variables

that you have set in the `linux/.config` file, the last time when you said make `menuconfig` or something like that.

The `__init` and `__exit` Macros

This demonstrates a feature of kernel 2.2 and later. Notice the change in the definitions of the `init` and `cleanup` functions. The `__init` macro causes the `init` function to be discarded and its memory freed once the `init` function finishes for built-in drivers, but not loadable modules. If you think about when the `init` function is invoked, this makes perfect sense.

There is also an `__initdata` which works similarly to `__init` but for `init` variables rather than functions.

The `__exit` macro causes the omission of the function when the module is built into the kernel, and like `__init`, has no effect for loadable modules. Again, if you consider when the `cleanup` function runs, this makes complete sense; built-in drivers do not need a `cleanup` function, while loadable modules do.

These macros are defined in `linux/init.h` and serve to free up kernel memory. When you boot your kernel and see something like `Freeing unused kernel memory: 236k freed`, this is precisely what the kernel is freeing.

Licensing and Module Documentation

Honestly, who loads or even cares about proprietary modules? If you do then you might have seen something like this:

```
$ sudo insmod xxxxxx.ko
loading out-of-tree module taints kernel.
module license 'unspecified' taints kernel.
```

You can use a few macros to indicate the license for your module. Some examples are "GPL", "GPL v2", "GPL and additional rights", "Dual BSD/GPL", "Dual MIT/GPL", "Dual MPL/GPL" and "Proprietary". They are defined within `linux/module.h`.

To reference what license you're using a macro is available called **MODULE_LICENSE**. This and a few other macros describing the module are illustrated in the below example.

Passing Command Line Arguments to a Module

Modules can take command line arguments, but not with the `argc/argv` you might be used to.

To allow arguments to be passed to your module, declare the variables that will take the values of the command line arguments as global and then use the `module_param()` macro, (defined in `linux/moduleparam.h`) to set the mechanism up.

At runtime, insmod will fill the variables with any command line arguments that are given, like `./insmod mymodule.ko myvariable=5`. The variable declarations and macros should be placed at the beginning of the module for clarity. The example code should clear up my admittedly lousy explanation.

The `module_param()` macro takes 3 arguments: the name of the variable, its type and permissions for the corresponding file in sysfs. Integer types can be signed as usual or unsigned. If you'd like to use arrays of integers or strings see `module_param_array()` and `module_param_string()`.

```
int myint = 3; module_param(myint, int, 0);
```

Arrays are supported too, but things are a bit different now than they were in the olden days. To keep track of the number of parameters you need to pass a pointer to a count variable as third parameter. At your option, you could also ignore the count and pass NULL instead. We show both possibilities here:

```
int myintarray[2]; module_param_array(myintarray, int, NULL, 0); /* not
interested in count */
```

```
short myshortarray[4]; int count; module_param_array(myshortarray, short,
&count, 0); /* put count into "count" variable */
```

A good use for this is to have the module variable's default values set, like an port or IO address. If the variables contain the default values, then perform autodetection (explained elsewhere). Otherwise, keep the current value. This will be made clear later on.

Lastly, there's a macro function, **MODULE_PARM_DESC()**, that is used to document arguments that the module can take. It takes two parameters: a variable name and a free form string describing that variable.

I would recommend playing around with this code:

Modules Spanning Multiple Files

Sometimes it makes sense to divide a kernel module between several source files.

Here is an example of such a kernel module.

The next file:

And finally, the makefile:

```
obj-m += hello-1.o obj-m += hello-2.o obj-m += hello-3.o obj-m += hello-4.o
obj-m += hello-5.o obj-m += startstop.o startstop-objs := start.o stop.o
```

```
all: make -C /lib/modules/$(shelluname -r)/buildM=$(PWD) modules
```

```
clean: make -C /lib/modules/$(shelluname -r)/buildM=$(PWD) clean
```

This is the complete makefile for all the examples we have seen so far. The first five lines are nothing special, but for the last example we will need two lines.

First we invent an object name for our combined module, second we tell make what object files are part of that module.

Building modules for a precompiled kernel

Obviously, we strongly suggest you to recompile your kernel, so that you can enable a number of useful debugging features, such as forced module unloading (**MODULE_FORCE_UNLOAD**): when this option is enabled, you can force the kernel to unload a module even when it believes it is unsafe, via a **sudo rmmod -f module** command. This option can save you a lot of time and a number of reboots during the development of a module. If you do not want to recompile your kernel then you should consider running the examples within a test distribution on a virtual machine. If you mess anything up then you can easily reboot or restore the virtual machine (VM).

There are a number of cases in which you may want to load your module into a precompiled running kernel, such as the ones shipped with common Linux distributions, or a kernel you have compiled in the past. In certain circumstances you could require to compile and insert a module into a running kernel which you are not allowed to recompile, or on a machine that you prefer not to reboot. If you can't think of a case that will force you to use modules for a precompiled kernel you might want to skip this and treat the rest of this chapter as a big footnote.

Now, if you just install a kernel source tree, use it to compile your kernel module and you try to insert your module into the kernel, in most cases you would obtain an error as follows:

```
insmod: error inserting 'poet_atkm.ko': -1 Invalid module format
```

Less cryptical information are logged to the systemd journal:

```
Jun  4 22:07:54 localhost kernel: poet_atkm: version magic '2.6.5-1.358custom 686  
REGPARM 4KSTACKS gcc-3.3' should be '2.6.5-1.358 686 REGPARM 4KSTACKS gcc-3.3'
```

In other words, your kernel refuses to accept your module because version strings (more precisely, version magics) do not match. Incidentally, version magics are stored in the module object in the form of a static string, starting with `vermagic:`. Version data are inserted in your module when it is linked against the **init/vermagic.o** file. To inspect version magics and other strings stored in a given module, issue the `modinfo module.ko` command:

```
$ modinfo hello-4.ko
description:    A sample driver
author:         LKMPG
license:        GPL
srcversion:     B2AA7FBFCC2C39AED665382
depends:
retpoline:     Y
```

```
name:          hello_4
vermagic:      5.4.0-70-generic SMP mod_unload modversions
```

To overcome this problem we could resort to the **-force-vermagic** option, but this solution is potentially unsafe, and unquestionably unacceptable in production modules. Consequently, we want to compile our module in an environment which was identical to the one in which our precompiled kernel was built. How to do this, is the subject of the remainder of this chapter.

First of all, make sure that a kernel source tree is available, having exactly the same version as your current kernel. Then, find the configuration file which was used to compile your precompiled kernel. Usually, this is available in your current *boot* directory, under a name like *config-2.6.x*. You may just want to copy it to your kernel source tree: `cp /boot/config-`uname -r` .config`.

Let's focus again on the previous error message: a closer look at the version magic strings suggests that, even with two configuration files which are exactly the same, a slight difference in the version magic could be possible, and it is sufficient to prevent insertion of the module into the kernel. That slight difference, namely the custom string which appears in the module's version magic and not in the kernel's one, is due to a modification with respect to the original, in the makefile that some distribution include. Then, examine your `/usr/src/linux/Makefile`, and make sure that the specified version information matches exactly the one used for your current kernel. For example, you makefile could start as follows:

```
VERSION = 5
PATCHLEVEL = 14
SUBLEVEL = 0
EXTRAVERSION = -rc2
```

In this case, you need to restore the value of symbol **EXTRAVERSION** to **-rc2**. We suggest to keep a backup copy of the makefile used to compile your kernel available in `/lib/modules/5.14.0-rc2/build`. A simple `cp /lib/modules/`uname -r`/build/Makefile /usr/src/linux-`uname -r`` should suffice. Additionally, if you already started a kernel build with the previous (wrong) Makefile, you should also rerun make, or directly modify symbol `UTS_RELEASE` in file `/usr/src/linux-5.14.0/include/linux/version.h` according to contents of file `/lib/modules/5.14.0/build/include/linux/version.h`, or overwrite the latter with the first.

Now, please run make to update configuration and version headers and objects:

```
$ make
CHK      include/linux/version.h
UPD      include/linux/version.h
SYMLINK  include/asm -> include/asm-i386
SPLIT    include/linux/autoconf.h -> include/config/*
HOSTCC   scripts/basic/fixdep
```

```
HOSTCC  scripts/basic/split-include
HOSTCC  scripts/basic/docproc
HOSTCC  scripts/conmakehash
HOSTCC  scripts/kallsyms
CC       scripts/empty.o
```

If you do not desire to actually compile the kernel, you can interrupt the build process (CTRL-C) just after the SPLIT line, because at that time, the files you need will be ready. Now you can turn back to the directory of your module and compile it: It will be built exactly according your current kernel settings, and it will load into it without any errors.

Preliminaries

How modules begin and end

A program usually begins with a **main()** function, executes a bunch of instructions and terminates upon completion of those instructions. Kernel modules work a bit differently. A module always begin with either the `init_module` or the function you specify with `module_init` call. This is the entry function for modules; it tells the kernel what functionality the module provides and sets up the kernel to run the module's functions when they're needed. Once it does this, entry function returns and the module does nothing until the kernel wants to do something with the code that the module provides.

All modules end by calling either **cleanup_module** or the function you specify with the **module_exit** call. This is the exit function for modules; it undoes whatever entry function did. It unregisters the functionality that the entry function registered.

Every module must have an entry function and an exit function. Since there's more than one way to specify entry and exit functions, I'll try my best to use the terms 'entry function' and 'exit function', but if I slip and simply refer to them as `init_module` and `cleanup_module`, I think you'll know what I mean.

Functions available to modules

Programmers use functions they don't define all the time. A prime example of this is **printf()**. You use these library functions which are provided by the standard C library, `libc`. The definitions for these functions don't actually enter your program until the linking stage, which insures that the code (for `printf()` for example) is available, and fixes the call instruction to point to that code.

Kernel modules are different here, too. In the hello world example, you might have noticed that we used a function, **pr_info()** but didn't include a standard I/O library. That's because modules are object files whose symbols get resolved upon `insmod`'ing. The definition for the symbols comes from the kernel itself;

the only external functions you can use are the ones provided by the kernel. If you're curious about what symbols have been exported by your kernel, take a look at `/proc/kallsyms`.

One point to keep in mind is the difference between library functions and system calls. Library functions are higher level, run completely in user space and provide a more convenient interface for the programmer to the functions that do the real work — system calls. System calls run in kernel mode on the user's behalf and are provided by the kernel itself. The library function `printf()` may look like a very general printing function, but all it really does is format the data into strings and write the string data using the low-level system call `write()`, which then sends the data to standard output.

Would you like to see what system calls are made by `printf()`? It's easy! Compile the following program:

```
#include <stdio.h>

int main(void) { printf("hello"); return 0; }
```

with **`gcc -Wall -o hello hello.c`**. Run the executable with **`strace ./hello`**. Are you impressed? Every line you see corresponds to a system call. `strace` is a handy program that gives you details about what system calls a program is making, including which call is made, what its arguments are and what it returns. It's an invaluable tool for figuring out things like what files a program is trying to access. Towards the end, you'll see a line which looks like `write(1, "hello", 5hello)`. There it is. The face behind the `printf()` mask. You may not be familiar with `write`, since most people use library functions for file I/O (like `fopen`, `fputs`, `fclose`). If that's the case, try looking at `man 2 write`. The 2nd man section is devoted to system calls (like `kill()` and `read()`). The 3rd man section is devoted to library calls, which you would probably be more familiar with (like `cosh()` and `random()`).

You can even write modules to replace the kernel's system calls, which we'll do shortly. Crackers often make use of this sort of thing for backdoors or trojans, but you can write your own modules to do more benign things, like have the kernel write Tee hee, that tickles! everytime someone tries to delete a file on your system.

User Space vs Kernel Space

A kernel is all about access to resources, whether the resource in question happens to be a video card, a hard drive or even memory. Programs often compete for the same resource. As I just saved this document, `updatedb` started updating the locate database. My vim session and `updatedb` are both using the hard drive concurrently. The kernel needs to keep things orderly, and not give users access to resources whenever they feel like it. To this end, a CPU can run in different modes. Each mode gives a different level of freedom to do what you want on the system. The Intel 80386 architecture had 4 of these modes, which

were called rings. Unix uses only two rings; the highest ring (ring 0, also known as ‘supervisor mode’ where everything is allowed to happen) and the lowest ring, which is called ‘user mode’.

Recall the discussion about library functions vs system calls. Typically, you use a library function in user mode. The library function calls one or more system calls, and these system calls execute on the library function’s behalf, but do so in supervisor mode since they are part of the kernel itself. Once the system call completes its task, it returns and execution gets transferred back to user mode.

Name Space

When you write a small C program, you use variables which are convenient and make sense to the reader. If, on the other hand, you’re writing routines which will be part of a bigger problem, any global variables you have are part of a community of other peoples’ global variables; some of the variable names can clash. When a program has lots of global variables which aren’t meaningful enough to be distinguished, you get namespace pollution. In large projects, effort must be made to remember reserved names, and to find ways to develop a scheme for naming unique variable names and symbols.

When writing kernel code, even the smallest module will be linked against the entire kernel, so this is definitely an issue. The best way to deal with this is to declare all your variables as static and to use a well-defined prefix for your symbols. By convention, all kernel prefixes are lowercase. If you don’t want to declare everything as static, another option is to declare a symbol table and register it with a kernel. We’ll get to this later.

The file `/proc/kallsyms` holds all the symbols that the kernel knows about and which are therefore accessible to your modules since they share the kernel’s codespace.

Code space

Memory management is a very complicated subject and the majority of O’Reilly’s *"Understanding The Linux Kernel"* exclusively covers memory management! We’re not setting out to be experts on memory managements, but we do need to know a couple of facts to even begin worrying about writing real modules.

If you haven’t thought about what a segfault really means, you may be surprised to hear that pointers don’t actually point to memory locations. Not real ones, anyway. When a process is created, the kernel sets aside a portion of real physical memory and hands it to the process to use for its executing code, variables, stack, heap and other things which a computer scientist would know about. This memory begins with 0x00000000 and extends up to whatever it needs to be. Since the memory space for any two processes don’t overlap, every process that can access a memory address, say 0xbffff978, would be accessing a different location in real physical memory! The processes would be accessing an index

named 0xbffff978 which points to some kind of offset into the region of memory set aside for that particular process. For the most part, a process like our Hello, World program can't access the space of another process, although there are ways which we'll talk about later.

The kernel has its own space of memory as well. Since a module is code which can be dynamically inserted and removed in the kernel (as opposed to a semi-autonomous object), it shares the kernel's codespace rather than having its own. Therefore, if your module segfaults, the kernel segfaults. And if you start writing over data because of an off-by-one error, then you're trampling on kernel data (or code). This is even worse than it sounds, so try your best to be careful.

By the way, I would like to point out that the above discussion is true for any operating system which uses a monolithic kernel. This isn't quite the same thing as *"building all your modules into the kernel"*, although the idea is the same. There are things called microkernels which have modules which get their own codespace. The GNU Hurd and the Magenta kernel of Google Fuchsia are two examples of a microkernel.

Device Drivers

One class of module is the device driver, which provides functionality for hardware like a serial port. On Unix, each piece of hardware is represented by a file located in /dev named a device file which provides the means to communicate with the hardware. The device driver provides the communication on behalf of a user program. So the es1370.o sound card device driver might connect the /dev/sound device file to the Ensoniq IS1370 sound card. A userspace program like mp3blaster can use /dev/sound without ever knowing what kind of sound card is installed.

1. Major and Minor Numbers [sec:org63bb013]

Let's look at some device files. Here are device files which represent the first three partitions on the primary master IDE hard drive:

```
$ ls -l /dev/hda[1-3]
brw-rw---- 1 root disk 3, 1 Jul 5 2000 /dev/hda1
brw-rw---- 1 root disk 3, 2 Jul 5 2000 /dev/hda2
brw-rw---- 1 root disk 3, 3 Jul 5 2000 /dev/hda3
```

Notice the column of numbers separated by a comma. The first number is called the device's major number. The second number is the minor number. The major number tells you which driver is used to access the hardware. Each driver is assigned a unique major number; all device files with the same major number are controlled by the same driver. All the above major numbers are 3, because they're all controlled by the same driver.

The minor number is used by the driver to distinguish between the various hardware it controls. Returning to the example above, although all three

devices are handled by the same driver they have unique minor numbers because the driver sees them as being different pieces of hardware.

Devices are divided into two types: character devices and block devices. The difference is that block devices have a buffer for requests, so they can choose the best order in which to respond to the requests. This is important in the case of storage devices, where it's faster to read or write sectors which are close to each other, rather than those which are further apart. Another difference is that block devices can only accept input and return output in blocks (whose size can vary according to the device), whereas character devices are allowed to use as many or as few bytes as they like. Most devices in the world are character, because they don't need this type of buffering, and they don't operate with a fixed block size. You can tell whether a device file is for a block device or a character device by looking at the first character in the output of `ls -l`. If it's 'b' then it's a block device, and if it's 'c' then it's a character device. The devices you see above are block devices. Here are some character devices (the serial ports):

```
crw-rw---- 1 root dial 4, 64 Feb 18 23:34 /dev/ttyS0
crw-r----- 1 root dial 4, 65 Nov 17 10:26 /dev/ttyS1
crw-rw---- 1 root dial 4, 66 Jul 5 2000 /dev/ttyS2
crw-rw---- 1 root dial 4, 67 Jul 5 2000 /dev/ttyS3
```

If you want to see which major numbers have been assigned, you can look at `/usr/src/linux/Documentation/devices.txt`.

When the system was installed, all of those device files were created by the `mknod` command. To create a new char device named 'coffee' with major/minor number 12 and 2, simply do `mknod /dev/coffee c 12 2`. You don't have to put your device files into `/dev`, but it's done by convention. Linus put his device files in `/dev`, and so should you. However, when creating a device file for testing purposes, it's probably OK to place it in your working directory where you compile the kernel module. Just be sure to put it in the right place when you're done writing the device driver.

I would like to make a few last points which are implicit from the above discussion, but I'd like to make them explicit just in case. When a device file is accessed, the kernel uses the major number of the file to determine which driver should be used to handle the access. This means that the kernel doesn't really need to use or even know about the minor number. The driver itself is the only thing that cares about the minor number. It uses the minor number to distinguish between different pieces of hardware.

By the way, when I say "*hardware*", I mean something a bit more abstract than a PCI card that you can hold in your hand. Look at these two device files:

```
$ ls -l /dev/sda /dev/sdb
```

```
brw-rw---- 1 root disk 8,  0 Jan  3 09:02 /dev/sda
brw-rw---- 1 root disk 8, 16 Jan  3 09:02 /dev/sdb
```

By now you can look at these two device files and know instantly that they are block devices and are handled by same driver (block major 8). Sometimes two device files with the same major but different minor number can actually represent the same piece of physical hardware. So just be aware that the word "hardware" in our discussion can mean something very abstract.

Character Device drivers

The `proc_ops` Structure

The `proc_ops` structure is defined in `/usr/include/linux/fs.h`, and holds pointers to functions defined by the driver that perform various operations on the device. Each field of the structure corresponds to the address of some function defined by the driver to handle a requested operation.

For example, every character driver needs to define a function that reads from the device. The `proc_ops` structure holds the address of the module's function that performs that operation. Here is what the definition looks like for kernel 3.0:

```
struct proc_ops {
    struct module *owner;
    loff_t (*llseek) (struct file *, loff_t,
int);
    ssize_t (*read) (struct file *, char __user *, size_t, loff_t *);
    ssize_t (*write) (struct file *, const char __user *, size_t, loff_t *);
    ssize_t (*aio_read) (struct kiocb *, const struct iovec *, unsigned long, loff_t);
    ssize_t (*aio_write) (struct kiocb *, const struct iovec *, unsigned long, loff_t);
    int (*iterate) (struct file *, struct dir_context *);
    unsigned int (*poll) (struct file *, struct poll_table_struct *);
    long (*unlocked_ioctl) (struct file *, unsigned int, unsigned long);
    long (*compat_ioctl) (struct file *, unsigned int, unsigned long);
    int (*mmap) (struct file *, struct vm_area_struct *);
    int (*open) (struct inode *, struct file *);
    int (*flush) (struct file *, fl_owner_t id);
    int (*release) (struct inode *, struct file *);
    int (*fsync) (struct file *, loff_t, loff_t, int datasync);
    int (*aio_fsync) (struct kiocb *, int datasync);
    int (*fasync) (int, struct file *, int);
    int (*lock) (struct file *, int, struct file_lock *);
    ssize_t (*sendpage) (struct file *, struct page *, int, size_t, loff_t *, int);
    unsigned long (*get_unmapped_area) (struct file *, unsigned long, unsigned long, unsigned long, unsigned long);
    int (*check_flags) (int);
    int (*flock) (struct file *, int, struct file_lock *);
    ssize_t (*splice_write) (struct pipe_inode_info *, struct file *, loff_t *, size_t, unsigned int);
    ssize_t (*splice_read) (struct file *, loff_t *, struct pipe_inode_info *, size_t, unsigned int);
    int (*setlease) (struct file *, long, struct file_lock **);
    long (*fallocate) (struct file *, int mode, loff_t offset, loff_t len);
    int (*show_fdinfo) (struct seq_file *, struct file *);
};
```

Some operations are not implemented by a driver. For example, a driver that

handles a video card won't need to read from a directory structure. The corresponding entries in the `proc_ops` structure should be set to `NULL`.

There is a gcc extension that makes assigning to this structure more convenient. You'll see it in modern drivers, and may catch you by surprise. This is what the new way of assigning to the structure looks like:

```
struct proc_ops fops = proc_read: device_read, proc_write: device_write,  
proc_open: device_open, proc_release: device_release ;
```

However, there's also a C99 way of assigning to elements of a structure, and this is definitely preferred over using the GNU extension. The version of gcc the author used when writing this, 2.95, supports the new C99 syntax. You should use this syntax in case someone wants to port your driver. It will help with compatibility:

```
struct proc_ops fops = .proc_read = device_read, .proc_write = device_write,  
.proc_open = device_open, .proc_release = device_release ;
```

The meaning is clear, and you should be aware that any member of the structure which you don't explicitly assign will be initialized to `NULL` by gcc.

An instance of `struct proc_ops` containing pointers to functions that are used to implement read, write, open, ... syscalls is commonly named `fops`.

The file structure

Each device is represented in the kernel by a file structure, which is defined in **linux/fs.h**. Be aware that a file is a kernel level structure and never appears in a user space program. It's not the same thing as a **FILE**, which is defined by glibc and would never appear in a kernel space function. Also, its name is a bit misleading; it represents an abstract open 'file', not a file on a disk, which is represented by a structure named `inode`.

An instance of `struct file` is commonly named `filp`. You'll also see it referred to as `struct file file`. Resist the temptation.

Go ahead and look at the definition of `file`. Most of the entries you see, like `struct dentry` aren't used by device drivers, and you can ignore them. This is because drivers don't fill `file` directly; they only use structures contained in `file` which are created elsewhere.

Registering A Device

As discussed earlier, char devices are accessed through device files, usually located in `/dev`. This is by convention. When writing a driver, it's OK to put the device file in your current directory. Just make sure you place it in `/dev` for a production driver. The major number tells you which driver handles which device file. The minor number is used only by the driver itself to differentiate which device it's operating on, just in case the driver handles more than one device.

Adding a driver to your system means registering it with the kernel. This is synonymous with assigning it a major number during the module's initialization. You do this by using the `register_chrdev` function, defined by `linux/fs.h`.

```
int register_chrdev(unsigned int major, const char *name, struct proc_ops
*fops);
```

where unsigned int `major` is the major number you want to request, *const char *name* is the name of the device as it'll appear in **/proc/devices** and *struct proc_ops *fops* is a pointer to the `proc_ops` table for your driver. A negative return value means the registration failed. Note that we didn't pass the minor number to `register_chrdev`. That's because the kernel doesn't care about the minor number; only our driver uses it.

Now the question is, how do you get a major number without hijacking one that's already in use? The easiest way would be to look through `Documentation/devices.txt` and pick an unused one. That's a bad way of doing things because you'll never be sure if the number you picked will be assigned later. The answer is that you can ask the kernel to assign you a dynamic major number.

If you pass a major number of 0 to `register_chrdev`, the return value will be the dynamically allocated major number. The downside is that you can't make a device file in advance, since you don't know what the major number will be. There are a couple of ways to do this. First, the driver itself can print the newly assigned number and we can make the device file by hand. Second, the newly registered device will have an entry in **/proc/devices**, and we can either make the device file by hand or write a shell script to read the file in and make the device file. The third method is we can have our driver make the device file using the **device_create** function after a successful registration and **device_destroy** during the call to `cleanup_module`.

Unregistering A Device

We can't allow the kernel module to be `rmmod`'ed whenever root feels like it. If the device file is opened by a process and then we remove the kernel module, using the file would cause a call to the memory location where the appropriate function (read/write) used to be. If we're lucky, no other code was loaded there, and we'll get an ugly error message. If we're unlucky, another kernel module was loaded into the same location, which means a jump into the middle of another function within the kernel. The results of this would be impossible to predict, but they can't be very positive.

Normally, when you don't want to allow something, you return an error code (a negative number) from the function which is supposed to do it. With `cleanup_module` that's impossible because it's a void function. However, there's a counter which keeps track of how many processes are using your module. You can see what its value is by looking at the 3rd field of **/proc/modules**. If this number isn't zero, `rmmod` will fail. Note that you don't have to check the counter

from within `cleanup_module` because the check will be performed for you by the system call `sys_delete_module`, defined in **linux/module.c**. You shouldn't use this counter directly, but there are functions defined in **linux/module.h** which let you increase, decrease and display this counter:

- `try_module_get(THIS_MODULE)`: Increment the use count.
- `module_put(THIS_MODULE)`: Decrement the use count.

It's important to keep the counter accurate; if you ever do lose track of the correct usage count, you'll never be able to unload the module; it's now reboot time, boys and girls. This is bound to happen to you sooner or later during a module's development.

chardev.c

The next code sample creates a char driver named `chardev`. You can cat its device file.

```
cat /proc/devices
```

(or open the file with a program) and the driver will put the number of times the device file has been read from into the file. We don't support writing to the file (like `echo "hi" > /dev/hello`), but catch these attempts and tell the user that the operation isn't supported. Don't worry if you don't see what we do with the data we read into the buffer; we don't do much with it. We simply read in the data and print a message acknowledging that we received it.

Writing Modules for Multiple Kernel Versions

The system calls, which are the major interface the kernel shows to the processes, generally stay the same across versions. A new system call may be added, but usually the old ones will behave exactly like they used to. This is necessary for backward compatibility – a new kernel version is not supposed to break regular processes. In most cases, the device files will also remain the same. On the other hand, the internal interfaces within the kernel can and do change between versions.

The Linux kernel versions are divided between the stable versions (n.\$<evennumber>\$.m) and the development versions (n.\$<oddnumber>\$.m). The development versions include all the cool new ideas, including those which will be considered a mistake, or reimplemented, in the next version. As a result, you can't trust the interface to remain the same in those versions (which is why I don't bother to support them in this book, it's too much work and it would become dated too quickly). In the stable versions, on the other hand, we can expect the interface to remain the same regardless of the bug fix version (the m number).

There are differences between different kernel versions, and if you want to

support multiple kernel versions, you'll find yourself having to code conditional compilation directives. The way to do this to compare the macro `LINUX_VERSION_CODE` to the macro `KERNEL_VERSION`. In version `a.b.c` of the kernel, the value of this macro would be $2^{16}a + 2^8b + c$.

While previous versions of this guide showed how you can write backward compatible code with such constructs in great detail, we decided to break with this tradition for the better. People interested in doing such might now use a LKMPG with a version matching to their kernel. We decided to version the LKMPG like the kernel, at least as far as major and minor number are concerned. We use the patchlevel for our own versioning so use LKMPG version `2.4.x` for kernels `2.4.x`, use LKMPG version `2.6.x` for kernels `2.6.x` and so on. Also make sure that you always use current, up to date versions of both, kernel and guide.

You might already have noticed that recent kernels look different. In case you haven't they look like `2.6.x.y` now. The meaning of the first three items basically stays the same, but a subpatchlevel has been added and will indicate security fixes till the next stable patchlevel is out. So people can choose between a stable tree with security updates and use the latest kernel as developer tree. Search the kernel mailing list archives if you're interested in the full story.

The `/proc` File System

In Linux, there is an additional mechanism for the kernel and kernel modules to send information to processes — the `/proc` file system. Originally designed to allow easy access to information about processes (hence the name), it is now used by every bit of the kernel which has something interesting to report, such as `/proc/modules` which provides the list of modules and `/proc/meminfo` which stats memory usage statistics.

The method to use the `proc` file system is very similar to the one used with device drivers — a structure is created with all the information needed for the `/proc` file, including pointers to any handler functions (in our case there is only one, the one called when somebody attempts to read from the `/proc` file). Then, `init_module` registers the structure with the kernel and `cleanup_module` unregisters it.

Normal file systems are located on a disk, rather than just in memory (which is where `/proc` is), and in that case the inode number is a pointer to a disk location where the file's index-node (inode for short) is located. The inode contains information about the file, for example the file's permissions, together with a pointer to the disk location or locations where the file's data can be found.

Because we don't get called when the file is opened or closed, there's nowhere for us to put `try_module_get` and `try_module_put` in this module, and if the file is opened and then the module is removed, there's no way to avoid the consequences.

Here a simple example showing how to use a `/proc` file. This is the HelloWorld for the `/proc` filesystem. There are three parts: create the file `proc helloworld` in the function `init_module`, return a value (and a buffer) when the file `/proc/helloworld` is read in the callback function `procfile_read`, and delete the file `/proc/helloworld` in the function `cleanup_module`.

The `/proc/helloworld` is created when the module is loaded with the function `proc_create`. The return value is a `struct proc_dir_entry`, and it will be used to configure the file `/proc/helloworld` (for example, the owner of this file). A null return value means that the creation has failed.

Each time, everytime the file `/proc/helloworld` is read, the function `procfile_read` is called. Two parameters of this function are very important: the buffer (the first parameter) and the offset (the third one). The content of the buffer will be returned to the application which read it (for example the cat command). The offset is the current position in the file. If the return value of the function isn't null, then this function is called again. So be careful with this function, if it never returns zero, the read function is called endlessly.

```
$ cat /proc/helloworld
HelloWorld!
```

Read and Write a /proc File

We have seen a very simple example for a `/proc` file where we only read the file `/proc/helloworld`. It's also possible to write in a `/proc` file. It works the same way as read, a function is called when the `/proc` file is written. But there is a little difference with read, data comes from user, so you have to import data from user space to kernel space (with `copy_from_user` or `get_user`)

The reason for `copy_from_user` or `get_user` is that Linux memory (on Intel architecture, it may be different under some other processors) is segmented. This means that a pointer, by itself, does not reference a unique location in memory, only a location in a memory segment, and you need to know which memory segment it is to be able to use it. There is one memory segment for the kernel, and one for each of the processes.

The only memory segment accessible to a process is its own, so when writing regular programs to run as processes, there's no need to worry about segments. When you write a kernel module, normally you want to access the kernel memory segment, which is handled automatically by the system. However, when the content of a memory buffer needs to be passed between the currently running process and the kernel, the kernel function receives a pointer to the memory buffer which is in the process segment. The `put_user` and `get_user` macros allow you to access that memory. These functions handle only one character, you can handle several characters with `copy_to_user` and `copy_from_user`. As the buffer (in read or write function) is in kernel space, for write function you need to import data because it comes from user space, but not for the read function

because data is already in kernel space.

Manage /proc file with standard filesystem

We have seen how to read and write a /proc file with the /proc interface. But it's also possible to manage /proc file with inodes. The main concern is to use advanced functions, like permissions.

In Linux, there is a standard mechanism for file system registration. Since every file system has to have its own functions to handle inode and file operations, there is a special structure to hold pointers to all those functions, struct **inode_operations**, which includes a pointer to struct **proc_ops**.

The difference between file and inode operations is that file operations deal with the file itself whereas inode operations deal with ways of referencing the file, such as creating links to it.

In /proc, whenever we register a new file, we're allowed to specify which struct **inode_operations** will be used to access to it. This is the mechanism we use, a struct **inode_operations** which includes a pointer to a struct **proc_ops** which includes pointers to our **procfs_read** and **procfs_write** functions.

Another interesting point here is the **module_permission** function. This function is called whenever a process tries to do something with the /proc file, and it can decide whether to allow access or not. Right now it is only based on the operation and the uid of the current user (as available in **current**, a pointer to a structure which includes information on the currently running process), but it could be based on anything we like, such as what other processes are doing with the same file, the time of day, or the last input we received.

It's important to note that the standard roles of read and write are reversed in the kernel. Read functions are used for output, whereas write functions are used for input. The reason for that is that read and write refer to the user's point of view — if a process reads something from the kernel, then the kernel needs to output it, and if a process writes something to the kernel, then the kernel receives it as input.

Still hungry for **procfs** examples? Well, first of all keep in mind, there are rumors around, claiming that **procfs** is on it's way out, consider using **sysfs** instead. Second, if you really can't get enough, there's a highly recommendable bonus level for **procfs** below `linux/Documentation/DocBook/`. Use `make help` in your toplevel kernel directory for instructions about how to convert it into your favourite format. Example: `make htmldocs`. Consider using this mechanism, in case you want to document something kernel related yourself.

Manage /proc file with seq_file

As we have seen, writing a /proc file may be quite "complex". So to help people writting /proc file, there is an API named **seq_file** that helps forming a /proc

file for output. It's based on sequence, which is composed of 3 functions: `start()`, `next()`, and `stop()`. The `seq_file` API starts a sequence when a user read the `/proc` file.

A sequence begins with the call of the function `start()`. If the return is a non NULL value, the function `next()` is called. This function is an iterator, the goal is to go through all the data. Each time `next()` is called, the function `show()` is also called. It writes data values in the buffer read by the user. The function `next()` is called until it returns NULL. The sequence ends when `next()` returns NULL, then the function `stop()` is called.

BE CAREFUL: when a sequence is finished, another one starts. That means that at the end of function `stop()`, the function `start()` is called again. This loop finishes when the function `start()` returns NULL. You can see a scheme of this in the figure "How `seq_file` works".

`Seq_file` provides basic functions for `proc_ops`, as `seq_read`, `seq_lseek`, and some others. But nothing to write in the `/proc` file. Of course, you can still use the same way as in the previous example.

If you want more information, you can read this web page:

- <http://lwn.net/Articles/22355/>
- http://www.kernelnewbies.org/documents/seq_file_howto.txt

You can also read the code of `fs/seq_file.c` in the linux kernel.

sysfs: Interacting with your module

sysfs allows you to interact with the running kernel from userspace by reading or setting variables inside of modules. This can be useful for debugging purposes, or just as an interface for applications or scripts. You can find *sysfs* directories and files under the *sys* directory on your system.

```
ls -l /sys
```

An example of a hello world module which includes the creation of a variable accessible via *sysfs* is given below.

Make and install the module:

```
make sudo insmod hello-sysfs.ko
```

Check that it exists:

```
sudo lsmod | grep hello_sysfs
```

What is the current value of *myvariable* ?

```
cat /sys/kernel/mymodule/myvariable
```

Set the value of *myvariable* and check that it changed.

```
echo "32" > /sys/kernel/mymodule/myvariable cat /sys/kernel/mymodule/myvariable
```

Finally, remove the test module:

```
sudo rmmod hello_sysfs
```

Talking To Device Files

Device files are supposed to represent physical devices. Most physical devices are used for output as well as input, so there has to be some mechanism for device drivers in the kernel to get the output to send to the device from processes. This is done by opening the device file for output and writing to it, just like writing to a file. In the following example, this is implemented by `device_write`.

This is not always enough. Imagine you had a serial port connected to a modem (even if you have an internal modem, it is still implemented from the CPU's perspective as a serial port connected to a modem, so you don't have to tax your imagination too hard). The natural thing to do would be to use the device file to write things to the modem (either modem commands or data to be sent through the phone line) and read things from the modem (either responses for commands or the data received through the phone line). However, this leaves open the question of what to do when you need to talk to the serial port itself, for example to send the rate at which data is sent and received.

The answer in Unix is to use a special function called **ioctl** (short for Input Output ConTroL). Every device can have its own ioctl commands, which can be read ioctls (to send information from a process to the kernel), write ioctls (to return information to a process), both or neither. Notice here the roles of read and write are reversed again, so in ioctl's read is to send information to the kernel and write is to receive information from the kernel.

The ioctl function is called with three parameters: the file descriptor of the appropriate device file, the ioctl number, and a parameter, which is of type long so you can use a cast to use it to pass anything. You won't be able to pass a structure this way, but you will be able to pass a pointer to the structure.

The ioctl number encodes the major device number, the type of the ioctl, the command, and the type of the parameter. This ioctl number is usually created by a macro call (`_IO`, `_IOR`, `_IOW` or `_IOWR` — depending on the type) in a header file. This header file should then be included both by the programs which will use ioctl (so they can generate the appropriate ioctls) and by the kernel module (so it can understand it). In the example below, the header file is `chardev.h` and the program which uses it is `ioctl.c`.

If you want to use ioctls in your own kernel modules, it is best to receive an official ioctl assignment, so if you accidentally get somebody else's ioctls, or if they get yours, you'll know something is wrong. For more information, consult the kernel source tree at `Documentation/ioctl-number.txt`.

System Calls

So far, the only thing we've done was to use well defined kernel mechanisms to register **/proc** files and device handlers. This is fine if you want to do something the kernel programmers thought you'd want, such as write a device driver. But what if you want to do something unusual, to change the behavior of the system in some way? Then, you're mostly on your own.

If you're not being sensible and using a virtual machine then this is where kernel programming can become hazardous. While writing the example below, I killed the **open()** system call. This meant I couldn't open any files, I couldn't run any programs, and I couldn't shutdown the system. I had to restart the virtual machine. No important files got annihilated, but if I was doing this on some live mission critical system then that could have been a possible outcome. To ensure you don't lose any files, even within a test environment, please run **sync** right before you do the **insmod** and the **rmmod**.

Forget about **/proc** files, forget about device files. They're just minor details. Minutiae in the vast expanse of the universe. The real process to kernel communication mechanism, the one used by all processes, is *system calls*. When a process requests a service from the kernel (such as opening a file, forking to a new process, or requesting more memory), this is the mechanism used. If you want to change the behaviour of the kernel in interesting ways, this is the place to do it. By the way, if you want to see which system calls a program uses, run **strace <arguments>**.

In general, a process is not supposed to be able to access the kernel. It can't access kernel memory and it can't call kernel functions. The hardware of the CPU enforces this (that's the reason why it's called 'protected mode' or 'page protection').

System calls are an exception to this general rule. What happens is that the process fills the registers with the appropriate values and then calls a special instruction which jumps to a previously defined location in the kernel (of course, that location is readable by user processes, it is not writable by them). Under Intel CPUs, this is done by means of interrupt 0x80. The hardware knows that once you jump to this location, you are no longer running in restricted user mode, but as the operating system kernel — and therefore you're allowed to do whatever you want.

The location in the kernel a process can jump to is called `system_call`. The procedure at that location checks the system call number, which tells the kernel what service the process requested. Then, it looks at the table of system calls (`sys_call_table`) to see the address of the kernel function to call. Then it calls the function, and after it returns, does a few system checks and then return back to the process (or to a different process, if the process time ran out). If you want to read this code, it's at the source file `arch/$<architecture>/kernel/entry.S`, after the line `ENTRY(system_call)`.

So, if we want to change the way a certain system call works, what we need to do is to write our own function to implement it (usually by adding a bit of our own code, and then calling the original function) and then change the pointer at `sys_call_table` to point to our function. Because we might be removed later and we don't want to leave the system in an unstable state, it's important for `cleanup_module` to restore the table to its original state.

The source code here is an example of such a kernel module. We want to "spy" on a certain user, and to `pr_info()` a message whenever that user opens a file. Towards this end, we replace the system call to open a file with our own function, called `our_sys_open`. This function checks the uid (user's id) of the current process, and if it's equal to the uid we spy on, it calls `pr_info()` to display the name of the file to be opened. Then, either way, it calls the original `open()` function with the same parameters, to actually open the file.

The `init_module` function replaces the appropriate location in `sys_call_table` and keeps the original pointer in a variable. The `cleanup_module` function uses that variable to restore everything back to normal. This approach is dangerous, because of the possibility of two kernel modules changing the same system call. Imagine we have two kernel modules, A and B. A's open system call will be `A_open` and B's will be `B_open`. Now, when A is inserted into the kernel, the system call is replaced with `A_open`, which will call the original `sys_open` when it's done. Next, B is inserted into the kernel, which replaces the system call with `B_open`, which will call what it thinks is the original system call, `A_open`, when it's done.

Now, if B is removed first, everything will be well — it will simply restore the system call to `A_open`, which calls the original. However, if A is removed and then B is removed, the system will crash. A's removal will restore the system call to the original, `sys_open`, cutting B out of the loop. Then, when B is removed, it will restore the system call to what it thinks is the original, `A_open`, which is no longer in memory. At first glance, it appears we could solve this particular problem by checking if the system call is equal to our open function and if so not changing it at all (so that B won't change the system call when it's removed), but that will cause an even worse problem. When A is removed, it sees that the system call was changed to `B_open` so that it is no longer pointing to `A_open`, so it won't restore it to `sys_open` before it is removed from memory. Unfortunately, `B_open` will still try to call `A_open` which is no longer there, so that even without removing B the system would crash.

Note that all the related problems make syscall stealing unfeasible for production use. In order to keep people from doing potential harmful things `sys_call_table` is no longer exported. This means, if you want to do something more than a mere dry run of this example, you will have to patch your current kernel in order to have `sys_call_table` exported. In the example directory you will find a README and the patch. As you can imagine, such modifications are not to be taken lightly. Do not try this on valuable systems (ie systems that you do not own - or cannot restore easily). You'll need to get the complete sourcecode of

this guide as a tarball in order to get the patch and the README. Depending on your kernel version, you might even need to hand apply the patch. Still here? Well, so is this chapter. If Wyle E. Coyote was a kernel hacker, this would be the first thing he'd try. ;)

Blocking Processes and threads

Sleep

What do you do when somebody asks you for something you can't do right away? If you're a human being and you're bothered by a human being, the only thing you can say is: "*Not right now, I'm busy. Go away!*". But if you're a kernel module and you're bothered by a process, you have another possibility. You can put the process to sleep until you can service it. After all, processes are being put to sleep by the kernel and woken up all the time (that's the way multiple processes appear to run on the same time on a single CPU).

This kernel module is an example of this. The file (called `/proc/sleep`) can only be opened by a single process at a time. If the file is already open, the kernel module calls `wait_event_interruptible`. The easiest way to keep a file open is to open it with:

```
tail -f
```

This function changes the status of the task (a task is the kernel data structure which holds information about a process and the system call it's in, if any) to **TASK_INTERRUPTIBLE**, which means that the task will not run until it is woken up somehow, and adds it to WaitQ, the queue of tasks waiting to access the file. Then, the function calls the scheduler to context switch to a different process, one which has some use for the CPU.

When a process is done with the file, it closes it, and `module_close` is called. That function wakes up all the processes in the queue (there's no mechanism to only wake up one of them). It then returns and the process which just closed the file can continue to run. In time, the scheduler decides that that process has had enough and gives control of the CPU to another process. Eventually, one of the processes which was in the queue will be given control of the CPU by the scheduler. It starts at the point right after the call to **module_interruptible_sleep_on**.

This means that the process is still in kernel mode - as far as the process is concerned, it issued the open system call and the system call hasn't returned yet. The process doesn't know somebody else used the CPU for most of the time between the moment it issued the call and the moment it returned.

It can then proceed to set a global variable to tell all the other processes that the file is still open and go on with its life. When the other processes get a piece of the CPU, they'll see that global variable and go back to sleep.

So we'll use `tail -f` to keep the file open in the background, while trying to access it with another process (again in the background, so that we need not switch to a different vt). As soon as the first background process is killed with `kill %1`, the second is woken up, is able to access the file and finally terminates.

To make our life more interesting, `module_close` doesn't have a monopoly on waking up the processes which wait to access the file. A signal, such as `Ctrl +c` (**SIGINT**) can also wake up a process. This is because we used `module_interruptible_sleep_on`. We could have used `module_sleep_on` instead, but that would have resulted in extremely angry users whose `Ctrl+c`'s are ignored.

In that case, we want to return with **-EINTR** immediately. This is important so users can, for example, kill the process before it receives the file.

There is one more point to remember. Some times processes don't want to sleep, they want either to get what they want immediately, or to be told it cannot be done. Such processes use the **O_NONBLOCK** flag when opening the file. The kernel is supposed to respond by returning with the error code **-EAGAIN** from operations which would otherwise block, such as opening the file in this example. The program `cat_noblock`, available in the source directory for this chapter, can be used to open a file with **O_NONBLOCK**.

```
$ sudo insmod sleep.ko
$ cat_noblock /proc/sleep
Last input:
$ tail -f /proc/sleep &
Last input:
Last input:
Last input:
Last input:
Last input:
Last input:
Last input:
tail: /proc/sleep: file truncated
[1] 6540
$ cat_noblock /proc/sleep
Open would block
$ kill %1
[1]+  Terminated                  tail -f /proc/sleep
$ cat_noblock /proc/sleep
Last input:
$
```

Completions

Sometimes one thing should happen before another within a module having multiple threads. Rather than using `/proc/sleep` commands the kernel has

another way to do this which allows timeouts or interrupts to also happen.

In the following example two threads are started, but one needs to start before another.

The *machine* structure stores the completion states for the two threads. At the exit point of each thread the respective completion state is updated, and *wait_for_completion* is used by the flywheel thread to ensure that it doesn't begin prematurely.

So even though *flywheel_thread* is started first you should notice if you load this module and run *dmesg* that turning the crank always happens first because the flywheel thread waits for it to complete.

There are other variations upon the *wait_for_completion* function, which include timeouts or being interrupted, but this basic mechanism is enough for many common situations without adding a lot of complexity.

Avoiding Collisions and Deadlocks

If processes running on different CPUs or in different threads try to access the same memory then it's possible that strange things can happen or your system can lock up. To avoid this various types of mutual exclusion kernel functions are available. These indicate if a section of code is "locked" or "unlocked" so that simultaneous attempts to run it can't happen.

Mutex

You can use kernel mutexes (mutual exclusions) in much the same manner that you might deploy them in userland. This may be all that's needed to avoid collisions in most cases.

Spinlocks

As the name suggests, spinlocks lock up the CPU that the code is running on, taking 100% of its resources. Because of this you should only use the spinlock mechanism around code which is likely to take no more than a few milliseconds to run and so won't noticeably slow anything down from the user's point of view.

The example here is *"irq safe"* in that if interrupts happen during the lock then they won't be forgotten and will activate when the unlock happens, using the *flags* variable to retain their state.

Read and write locks

Read and write locks are specialised kinds of spinlocks so that you can exclusively read from something or write to something. Like the earlier spinlocks example the one below shows an *"irq safe"* situation in which if other functions were triggered

from irqs which might also read and write to whatever you are concerned with then they wouldn't disrupt the logic. As before it's a good idea to keep anything done within the lock as short as possible so that it doesn't hang up the system and cause users to start revolting against the tyranny of your module.

Of course if you know for sure that there are no functions triggered by irqs which could possibly interfere with your logic then you can use the simpler `read_lock(&myrwlock)` and `read_unlock(&myrwlock)` or the corresponding write functions.

Atomic operations

If you're doing simple arithmetic: adding, subtracting or bitwise operations then there's another way in the multi-CPU and multi-hyperthreaded world to stop other parts of the system from messing with your mojo. By using atomic operations you can be confident that your addition, subtraction or bit flip did actually happen and wasn't overwritten by some other shenanigans. An example is shown below.

Replacing Print Macros

Replacement

In Section 1.2.1.2, I said that X and kernel module programming don't mix. That's true for developing kernel modules, but in actual use, you want to be able to send messages to whichever tty the command to load the module came from.

"tty" is an abbreviation of *teletype*: originally a combination keyboard-printer used to communicate with a Unix system, and today an abstraction for the text stream used for a Unix program, whether it's a physical terminal, an xterm on an X display, a network connection used with ssh, etc.

The way this is done is by using `current`, a pointer to the currently running task, to get the current task's tty structure. Then, we look inside that tty structure to find a pointer to a string write function, which we use to write a string to the tty.

Flashing keyboard LEDs

In certain conditions, you may desire a simpler and more direct way to communicate to the external world. Flashing keyboard LEDs can be such a solution: It is an immediate way to attract attention or to display a status condition. Keyboard LEDs are present on every hardware, they are always visible, they do not need any setup, and their use is rather simple and non-intrusive, compared to writing to a tty or a file.

The following source code illustrates a minimal kernel module which, when loaded, starts blinking the keyboard LEDs until it is unloaded.

If none of the examples in this chapter fit your debugging needs there might yet be some other tricks to try. Ever wondered what CONFIG_LL_DEBUG in make menuconfig is good for? If you activate that you get low level access to the serial port. While this might not sound very powerful by itself, you can patch **kernel/printk.c** or any other essential syscall to use printascii, thus making it possible to trace virtually everything what your code does over a serial line. If you find yourself porting the kernel to some new and former unsupported architecture this is usually amongst the first things that should be implemented. Logging over a netconsole might also be worth a try.

While you have seen lots of stuff that can be used to aid debugging here, there are some things to be aware of. Debugging is almost always intrusive. Adding debug code can change the situation enough to make the bug seem to disappear. Thus you should try to keep debug code to a minimum and make sure it does not show up in production code.

Scheduling Tasks

There are two main ways of running tasks: tasklets and work queues. Tasklets are a quick and easy way of scheduling a single function to be run, for example when triggered from an interrupt, whereas work queues are more complicated but also better suited to running multiple things in a sequence.

Tasklets

Here's an example tasklet module. The *tasklet_fn* function runs for a few seconds and in the mean time execution of the *example_tasklet_init* function continues to the exit point.

So with this example loaded *dmesg* should show:

```
tasklet example init
Example tasklet starts
Example tasklet init continues...
Example tasklet ends
```

Work queues

To add a task to the scheduler we can use a workqueue. The kernel then uses the Completely Fair Scheduler (CFS) to execute work within the queue.

Interrupt Handlers

Interrupt Handlers

Except for the last chapter, everything we did in the kernel so far we've done as a response to a process asking for it, either by dealing with a special file, sending

an `ioctl()`, or issuing a system call. But the job of the kernel isn't just to respond to process requests. Another job, which is every bit as important, is to speak to the hardware connected to the machine.

There are two types of interaction between the CPU and the rest of the computer's hardware. The first type is when the CPU gives orders to the hardware, the other is when the hardware needs to tell the CPU something. The second, called interrupts, is much harder to implement because it has to be dealt with when convenient for the hardware, not the CPU. Hardware devices typically have a very small amount of RAM, and if you don't read their information when available, it is lost.

Under Linux, hardware interrupts are called IRQ's (Interrupt ReQuests). There are two types of IRQ's, short and long. A short IRQ is one which is expected to take a very short period of time, during which the rest of the machine will be blocked and no other interrupts will be handled. A long IRQ is one which can take longer, and during which other interrupts may occur (but not interrupts from the same device). If at all possible, it's better to declare an interrupt handler to be long.

When the CPU receives an interrupt, it stops whatever it's doing (unless it's processing a more important interrupt, in which case it will deal with this one only when the more important one is done), saves certain parameters on the stack and calls the interrupt handler. This means that certain things are not allowed in the interrupt handler itself, because the system is in an unknown state. The solution to this problem is for the interrupt handler to do what needs to be done immediately, usually read something from the hardware or send something to the hardware, and then schedule the handling of the new information at a later time (this is called the "bottom half") and return. The kernel is then guaranteed to call the bottom half as soon as possible – and when it does, everything allowed in kernel modules will be allowed.

The way to implement this is to call `request_irq()` to get your interrupt handler called when the relevant IRQ is received.

In practice IRQ handling can be a bit more complex. Hardware is often designed in a way that chains two interrupt controllers, so that all the IRQs from interrupt controller B are cascaded to a certain IRQ from interrupt controller A. Of course that requires that the kernel finds out which IRQ it really was afterwards and that adds overhead. Other architectures offer some special, very low overhead, so called "fast IRQ" or FIQs. To take advantage of them requires handlers to be written in assembler, so they do not really fit into the kernel. They can be made to work similar to the others, but after that procedure, they're no longer any faster than "common" IRQs. SMP enabled kernels running on systems with more than one processor need to solve another truckload of problems. It's not enough to know if a certain IRQs has happend, it's also important for what CPU(s) it was for. People still interested in more details, might want to do a web search for "APIC" now ;)

This function receives the IRQ number, the name of the function, flags, a name for `/proc/interrupts` and a parameter to pass to the interrupt handler. Usually there is a certain number of IRQs available. How many IRQs there are is hardware-dependent. The flags can include `SA_SHIRQ` to indicate you're willing to share the IRQ with other interrupt handlers (usually because a number of hardware devices sit on the same IRQ) and `SA_INTERRUPT` to indicate this is a fast interrupt. This function will only succeed if there isn't already a handler on this IRQ, or if you're both willing to share.

Detecting button presses

Many popular single board computers, such as Raspberry Pis or Beagleboards, have a bunch of GPIO pins. Attaching buttons to those and then having a button press do something is a classic case in which you might need to use interrupts so that instead of having the CPU waste time and battery power polling for a change in input state it's better for the input to trigger the CPU to then run a particular handling function.

Here's an example where buttons are connected to GPIO numbers 17 and 18 and an LED is connected to GPIO 4. You can change those numbers to whatever is appropriate for your board.

Bottom Half

Suppose you want to do a bunch of stuff inside of an interrupt routine. A common way to do that without rendering the interrupt unavailable for a significant duration is to combine it with a tasklet. This pushes the bulk of the work off into the scheduler.

The example below modifies the previous example to also run an additional task when an interrupt is triggered.

Crypto

At the dawn of the internet everybody trusted everybody completely... but that didn't work out so well. When this guide was originally written it was a more innocent era in which almost nobody actually gave a damn about crypto - least of all kernel developers. That's certainly no longer the case now. To handle crypto stuff the kernel has its own API enabling common methods of encryption, decryption and your favourite hash functions.

Hash functions

Calculating and checking the hashes of things is a common operation. Here is a demonstration of how to calculate a sha256 hash within a kernel module.

Make and install the module:

```
make sudo insmod cryptosha256.ko dmesg
```

And you should see that the hash was calculated for the test string.

Finally, remove the test module:

```
sudo rmmod cryptosha256
```

Symmetric key encryption

Here is an example of symmetrically encrypting a string using the AES algorithm and a password.

Standardising the interfaces: The Device Model

Up to this point we've seen all kinds of modules doing all kinds of things, but there was no consistency in their interfaces with the rest of the kernel. To impose some consistency such that there is at minimum a standardised way to start, suspend and resume a device a device model was added. An example is show below, and you can use this as a template to add your own suspend, resume or other interface functions.

Optimizations

Likely and Unlikely conditions

Sometimes you might want your code to run as quickly as possible, especially if it's handling an interrupt or doing something which might cause noticeable latency. If your code contains boolean conditions and if you know that the conditions are almost always likely to evaluate as either *true* or *false*, then you can allow the compiler to optimise for this using the *likely* and *unlikely* macros.

For example, when allocating memory you're almost always expecting this to succeed.

```
bvl = bvec_alloc(gfp_mask, nr_iovecs, &idx); if (unlikely(!bvl)) mem-  
pool_free(bio, bio_pool); bio = NULL; goto out;
```

When the *unlikely* macro is used the compiler alters its machine instruction output so that it continues along the false branch and only jumps if the condition is true. That avoids flushing the processor pipeline. The opposite happens if you use the *likely* macro.

Common Pitfalls

Before I send you on your way to go out into the world and write kernel modules, there are a few things I need to warn you about. If I fail to warn you and

something bad happens, please report the problem to me for a full refund of the amount I was paid for your copy of the book.

Using standard libraries

You can't do that. In a kernel module you can only use kernel functions, which are the functions you can see in `/proc/kallsyms`.

Disabling interrupts

You might need to do this for a short time and that is OK, but if you don't enable them afterwards, your system will be stuck and you'll have to power it off.

Sticking your head inside a large carnivore

I probably don't have to warn you about this, but I figured I will anyway, just in case.

Where To Go From Here?

I could easily have squeezed a few more chapters into this book. I could have added a chapter about creating new file systems, or about adding new protocol stacks (as if there's a need for that – you'd have to dig underground to find a protocol stack not supported by Linux). I could have added explanations of the kernel mechanisms we haven't touched upon, such as bootstrapping or the disk interface.

However, I chose not to. My purpose in writing this book was to provide initiation into the mysteries of kernel module programming and to teach the common techniques for that purpose. For people seriously interested in kernel programming, I recommend kernelnewbies.org and the *Documentation* subdirectory within the kernel source code which isn't always easy to understand but can be a starting point for further investigation. Also, as Linus said, the best way to learn the kernel is to read the source code yourself.

If you're interested in more examples of short kernel modules then searching on sites such as Github and Gitlab is a good way to start, although there is a lot of duplication of older LKMPG examples which may not compile with newer kernel versions. You will also be able to find examples of the use of kernel modules to attack or compromise systems or exfiltrate data and those can be useful for thinking about how to defend systems and learning about existing security mechanisms within the kernel.

I hope I have helped you in your quest to become a better programmer, or at least to have fun through technology. And, if you do write useful kernel modules, I hope you publish them under the GPL, so I can use them too.

If you'd like to contribute to this guide, notice anything glaringly wrong, or just want to add extra sarcastic remarks perhaps involving monkeys or some other kind of animal then fire an email to bob@freedombone.net and you may be able to get commit access to <https://code.freedombone.net/bashrc/LKMPG>.

Happy hacking.