

Unreliable Guide To Hacking The Linux Kernel

Rusty Russell <rusty@rustcorp.com.au>

Unreliable Guide To Hacking The Linux Kernel

by Rusty Russell

This is the first release of this document as part of the kernel tarball.

Copyright © 2005 Rusty Russell

This documentation is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; either version 2 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with this program; if not, write to the Free Software Foundation, Inc., 59 Temple Place, Suite 330, Boston, MA 02111-1307 USA

For more details see the file COPYING in the source distribution of Linux.

Chapter 1. Introduction

Welcome, gentle reader, to Rusty's Remarkably Unreliable Guide to Linux Kernel Hacking. This document describes the common routines and general requirements for kernel code: its goal is to serve as a primer for Linux kernel development for experienced C programmers. I avoid implementation details: that's what the code is for, and I ignore whole tracts of useful routines.

Before you read this, please understand that I never wanted to write this document, being grossly under-qualified, but I always wanted to read it, and this was the only way. I hope it will grow into a compendium of best practice, common starting points and random information.

Chapter 2. The Players

At any time each of the CPUs in a system can be:

- not associated with any process, serving a hardware interrupt;
- not associated with any process, serving a softirq or tasklet;
- running in kernel space, associated with a process (user context);
- running a process in user space.

There is an ordering between these. The bottom two can preempt each other, but above that is a strict hierarchy: each can only be preempted by the ones above it. For example, while a softirq is running on a CPU, no other softirq will preempt it, but a hardware interrupt can. However, any other CPUs in the system execute independently.

We'll see a number of ways that the user context can block interrupts, to become truly non-preemptable.

User Context

User context is when you are coming in from a system call or other trap: like userspace, you can be preempted by more important tasks and by interrupts. You can sleep, by calling `schedule()`.

Note

You are always in user context on module load and unload, and on operations on the block device layer.

In user context, the current pointer (indicating the task we are currently executing) is valid, and `in_interrupt()` (`include/linux/interrupt.h`) is false.

Caution

Beware that if you have preemption or softirqs disabled (see below), `in_interrupt()` will return a false positive.

Hardware Interrupts (Hard IRQs)

Timer ticks, network cards and keyboard are examples of real hardware which produce interrupts at any time. The kernel runs interrupt handlers, which services the hardware. The kernel guarantees that this handler is never re-entered: if the same interrupt arrives, it is queued (or dropped). Because it disables interrupts, this handler has to be fast: frequently it simply acknowledges the interrupt, marks a 'software interrupt' for execution and exits.

You can tell you are in a hardware interrupt, because `in_irq()` returns true.

Caution

Beware that this will return a false positive if interrupts are disabled (see below).

Software Interrupt Context: Softirqs and Tasklets

Whenever a system call is about to return to userspace, or a hardware interrupt handler exits, any 'software interrupts' which are marked pending (usually by hardware interrupts) are run (`kernel/softirq.c`).

Much of the real interrupt handling work is done here. Early in the transition to SMP, there were only 'bottom halves' (BHs), which didn't take advantage of multiple CPUs. Shortly after we switched from wind-up computers made of match-sticks and snot, we abandoned this limitation and switched to 'softirqs'.

`include/linux/interrupt.h` lists the different softirqs. A very important softirq is the timer softirq (`include/linux/timer.h`): you can register to have it call functions for you in a given length of time.

Softirqs are often a pain to deal with, since the same softirq will run simultaneously on more than one CPU. For this reason, tasklets (`include/linux/interrupt.h`) are more often used: they are dynamically-registrable (meaning you can have as many as you want), and they also guarantee that any tasklet will only run on one CPU at any time, although different tasklets can run simultaneously.

Caution

The name 'tasklet' is misleading: they have nothing to do with 'tasks', and probably more to do with some bad vodka Alexey Kuznetsov had at the time.

You can tell you are in a softirq (or tasklet) using the `in_softirq()` macro (`include/linux/interrupt.h`).

Caution

Beware that this will return a false positive if a bh lock (see below) is held.

Chapter 3. Some Basic Rules

No memory protection

If you corrupt memory, whether in user context or interrupt context, the whole machine will crash. Are you sure you can't do what you want in userspace?

No floating point or MMX

The FPU context is not saved; even in user context the FPU state probably won't correspond with the current process: you would mess with some user process' FPU state. If you really want to do this, you would have to explicitly save/restore the full FPU state (and avoid context switches). It is generally a bad idea; use fixed point arithmetic first.

A rigid stack limit

Depending on configuration options the kernel stack is about 3K to 6K for most 32-bit architectures: it's about 14K on most 64-bit archs, and often shared with interrupts so you can't use it all. Avoid deep recursion and huge local arrays on the stack (allocate them dynamically instead).

The Linux kernel is portable

Let's keep it that way. Your code should be 64-bit clean, and endi-an-independent. You should also minimize CPU specific stuff, e.g. inline assembly should be cleanly encapsulated and minimized to ease porting. Generally it should be restricted to the architecture-de-pendent part of the kernel tree.

Chapter 4. ioctls: Not writing a new system call

A system call generally looks like this

```
asmlinkage long sys_mycall(int arg)
{
    return 0;
}
```

First, in most cases you don't want to create a new system call. You create a character device and implement an appropriate ioctl for it. This is much more flexible than system calls, doesn't have to be entered in every architecture's `include/asm/unistd.h` and `arch/kernel/entry.S` file, and is much more likely to be accepted by Linus.

If all your routine does is read or write some parameter, consider implementing a `sysfs` interface instead.

Inside the ioctl you're in user context to a process. When an error occurs you return a negated `errno` (see `include/linux/errno.h`), otherwise you return 0.

After you slept you should check if a signal occurred: the Unix/Linux way of handling signals is to temporarily exit the system call with the `-ERESTARTSYS` error. The system call entry code will switch back to user context, process the signal handler and then your system call will be restarted (unless the user disabled that). So you should be prepared to process the restart, e.g. if you're in the middle of manipulating some data structure.

```
if (signal_pending(current))
    return -ERESTARTSYS;
```

If you're doing longer computations: first think userspace. If you *really* want to do it in kernel you should regularly check if you need to give up the CPU (remember there is cooperative multitasking per CPU). Idiom:

```
cond_resched(); /* Will sleep */
```

A short note on interface design: the UNIX system call motto is "Provide mechanism not policy".

Chapter 5. Recipes for Deadlock

You cannot call any routines which may sleep, unless:

- You are in user context.
- You do not own any spinlocks.
- You have interrupts enabled (actually, Andi Kleen says that the scheduling code will enable them for you, but that's probably not what you wanted).

Note that some functions may sleep implicitly: common ones are the user space access functions (`*_user`) and memory allocation functions without `GFP_ATOMIC`.

You should always compile your kernel `CONFIG_DEBUG_ATOMIC_SLEEP` on, and it will warn you if you break these rules. If you *do* break the rules, you will eventually lock up your box.

Really.

`__initcall()/module_init() include/linux/init.h`

Many parts of the kernel are well served as a module (dynamically-loadable parts of the kernel). Using the `module_init()` and `module_exit()` macros it is easy to write code without `#ifdefs` which can operate both as a module or built into the kernel.

The `module_init()` macro defines which function is to be called at module insertion time (if the file is compiled as a module), or at boot time: if the file is not compiled as a module the `module_init()` macro becomes equivalent to `__initcall()`, which through linker magic ensures that the function is called on boot.

The function can return a negative error number to cause module loading to fail (unfortunately, this has no effect if the module is compiled into the kernel). This function is called in user context with interrupts enabled, so it can sleep.

`module_exit() include/linux/init.h`

This macro defines the function to be called at module removal time (or never, in the case of the file compiled into the kernel). It will only be called if the module usage count has reached zero. This function can also sleep, but cannot fail: everything must be cleaned up by the time it returns.

Note that this macro is optional: if it is not present, your module will not be removable (except for `'rmmod -f'`).

`try_module_get()/module_put() include/linux/module.h`

These manipulate the module usage count, to protect against removal (a module also can't be removed if another module uses one of its exported symbols: see below). Before calling into module code, you should call `try_module_get()` on that module: if it fails, then the module is being removed and you should act as if it wasn't there. Otherwise, you can safely enter the module, and call `module_put()` when you're finished.

Most registerable structures have an *owner* field, such as in the `file_operations` structure. Set this field to the macro `THIS_MODULE`.

Chapter 7. Wait Queues `include/linux/wait.h`

[SLEEPS]

A wait queue is used to wait for someone to wake you up when a certain condition is true. They must be used carefully to ensure there is no race condition. You declare a `wait_queue_head_t`, and then processes which want to wait for that condition declare a `wait_queue_t` referring to themselves, and place that in the queue.

Declaring

You declare a `wait_queue_head_t` using the `DECLARE_WAIT_QUEUE_HEAD()` macro, or using the `init_waitqueue_head()` routine in your initialization code.

Queuing

Placing yourself in the waitqueue is fairly complex, because you must put yourself in the queue before checking the condition. There is a macro to do this: `wait_event_interruptible()` in `include/linux/wait.h`. The first argument is the wait queue head, and the second is an expression which is evaluated; the macro returns 0 when this expression is true, or `-ERESTARTSYS` if a signal is received. The `wait_event()` version ignores signals.

Waking Up Queued Tasks

Call `wake_up()` in `include/linux/wait.h`, which will wake up every process in the queue. The exception is if one has `TASK_EXCLUSIVE` set, in which case the remainder of the queue will not be woken. There are other variants of this basic function available in the same header.

Chapter 8. Atomic Operations

Certain operations are guaranteed atomic on all platforms. The first class of operations work on `atomic_t` include/asm/atomic.h; this contains a signed integer (at least 32 bits long), and you must use these functions to manipulate or read `atomic_t` variables. `atomic_read()` and `atomic_set()` get and set the counter, `atomic_add()`, `atomic_sub()`, `atomic_inc()`, `atomic_dec()`, and `atomic_dec_and_test()` (returns true if it was decremented to zero).

Yes. It returns true (i.e. `!= 0`) if the atomic variable is zero.

Note that these functions are slower than normal arithmetic, and so should not be used unnecessarily.

The second class of atomic operations is atomic bit operations on an unsigned long, defined in include/linux/bitops.h. These operations generally take a pointer to the bit pattern, and a bit number: 0 is the least significant bit. `set_bit()`, `clear_bit()` and `change_bit()` set, clear, and flip the given bit. `test_and_set_bit()`, `test_and_clear_bit()` and `test_and_change_bit()` do the same thing, except return true if the bit was previously set; these are particularly useful for atomically setting flags.

It is possible to call these operations with bit indices greater than `BITS_PER_LONG`. The resulting behavior is strange on big-endian platforms though so it is a good idea not to do this.

Chapter 9. Symbols

Within the kernel proper, the normal linking rules apply (ie. unless a symbol is declared to be file scope with the `static` keyword, it can be used anywhere in the kernel). However, for modules, a special exported symbol table is kept which limits the entry points to the kernel proper. Modules can also export symbols.

`EXPORT_SYMBOL()` `include/linux/export.h`

This is the classic method of exporting a symbol: dynamically loaded modules will be able to use the symbol as normal.

`EXPORT_SYMBOL_GPL()` `include/linux/export.h`

Similar to `EXPORT_SYMBOL()` except that the symbols exported by `EXPORT_SYMBOL_GPL()` can only be seen by modules with a `MODULE_LICENSE()` that specifies a GPL compatible license. It implies that the function is considered an internal implementation issue, and not really an interface. Some maintainers and developers may however require `EXPORT_SYMBOL_GPL()` when adding any new APIs or functionality.

- Inline functions
- Statement expressions (ie. the `{ { and } }` constructs).
- Declaring attributes of a function / variable / type (`__attribute__`)
- `typeof`
- Zero length arrays
- Macro `varargs`
- Arithmetic on void pointers
- Non-Constant initializers
- Assembler Instructions (not outside `arch/` and `include/asm/`)
- Function names as strings (`__func__`).
- `__builtin_constant_p()`

Be wary when using `long long` in the kernel, the code gcc generates for it is horrible and worse: division and multiplication does not work on i386 because the GCC runtime functions for it are missing from the kernel environment.

C++

Using C++ in the kernel is usually a bad idea, because the kernel does not provide the necessary runtime environment and the include files are not tested for it. It is still possible, but not recommended. If you really want to do this, forget about exceptions at least.

#if

It is generally considered cleaner to use macros in header files (or at the top of `.c` files) to abstract away functions rather than using `#if` pre-processor statements throughout the source code.

Chapter 11. Putting Your Stuff in the Kernel

In order to get your stuff into shape for official inclusion, or even to make a neat patch, there's administrative work to be done:

- Figure out whose pond you've been pissing in. Look at the top of the source files, inside the MAINTAINERS file, and last of all in the CREDITS file. You should coordinate with this person to make sure you're not duplicating effort, or trying something that's already been rejected.

Make sure you put your name and EMail address at the top of any files you create or mangle significantly. This is the first place people will look when they find a bug, or when *they* want to make a change.

- Usually you want a configuration option for your kernel hack. Edit `Kconfig` in the appropriate directory. The Config language is simple to use by cut and paste, and there's complete documentation in `Documentation/kbuild/kconfig-language.txt`.

In your description of the option, make sure you address both the expert user and the user who knows nothing about your feature. Mention incompatibilities and issues here. *Definitely* end your description with “ if in doubt, say N ” (or, occasionally, `Y'); this is for people who have no idea what you are talking about.

- Edit the Makefile: the CONFIG variables are exported here so you can usually just add a "obj-\$(CONFIG_XXX) += xxx.o" line. The syntax is documented in `Documentation/kbuild/makefiles.txt`.
- Put yourself in CREDITS if you've done something noteworthy, usually beyond a single file (your name should be at the top of the source files anyway). MAINTAINERS means you want to be consulted when changes are made to a subsystem, and hear about bugs; it implies a more-than-passing commitment to some part of the code.
- Finally, don't forget to read `Documentation/SubmittingPatches` and possibly `Documentation/SubmittingDrivers`.

Chapter 12. Kernel Cantrips

Some favorites from browsing the source. Feel free to add to this list.

arch/x86/include/asm/delay.h:

```
#define ndelay(n) (__builtin_constant_p(n) ? \
    ((n) > 20000 ? __bad_ndelay() : __const_udelay((n) * 5ul)) : \
    __ndelay(n))
```

include/linux/fs.h:

```
/*
 * Kernel pointers have redundant information, so we can use a
 * scheme where we can return either an error code or a dentry
 * pointer with the same return value.
 *
 * This should be a per-architecture thing, to allow different
 * error and pointer decisions.
 */
#define ERR_PTR(err)      ((void *)((long)(err)))
#define PTR_ERR(ptr)      ((long)(ptr))
#define IS_ERR(ptr)      ((unsigned long)(ptr) > (unsigned long)(-1000))
```

arch/x86/include/asm/uaccess_32.h:

```
#define copy_to_user(to, from, n) \
    (__builtin_constant_p(n) ? \
    __constant_copy_to_user((to), (from), (n)) : \
    __generic_copy_to_user((to), (from), (n)))
```

arch/sparc/kernel/head.S:

```
/*
 * Sun people can't spell worth damn. "compatability" indeed.
 * At least we *know* we can't spell, and use a spell-checker.
 */

/* Uh, actually Linus it is I who cannot spell. Too much murky
 * Sparc assembly will do this to ya.
 */
C_LABEL(cputypvar):
    .asciz "compatibility"

/* Tested on SS-5, SS-10. Probably someone at Sun applied a spell-checker. */
    .align 4
C_LABEL(cputypvar_sun4m):
    .asciz "compatible"
```

arch/sparc/lib/checksum.S:

```
/* Sun, you just can't beat me, you just can't.  Stop trying,  
 * give up.  I'm serious, I am going to kick the living shit  
 * out of you, game over, lights out.  
 */
```

Chapter 13. Thanks

Thanks to Andi Kleen for the idea, answering my questions, fixing my mistakes, filling content, etc. Philipp Rumpf for more spelling and clarity fixes, and some excellent non-obvious points. Werner Almesberger for giving me a great summary of `disable_irq()`, and Jes Sorensen and Andrea Arcangeli added caveats. Michael Elizabeth Chastain for checking and adding to the Configure section. Telsa Gwynne for teaching me DocBook.