

Часть 1. Первые шаги



Олег Цилюрик

Опубликовано 01.01.2012

FacebookTwitterLinked InGoogle+Отправить эту страницу по электронной почте



6

Серия контента:

ЭТОТ КОНТЕНТ ЯВЛЯЕТСЯ ЧАСТЬЮ СЕРИИ: **Разработка модулей ядра Linux**

Эта статья является первой в протяжённом цикле статей, который в последствии должен превратиться в **пошаговое руководство** по написанию модулей ядра Linux. Но для того, чтобы такая амбициозная задача имела шансы на успех, а её результат был измеримым и конкретным, потребуется с самого начала разграничить области интересов и выделить зоны ответственности при разработке целевых модулей. В данном цикле не будут рассматриваться следующие вопросы.

- детали внутреннего устройства и функционирования ядра Linux, а также выявляющиеся в нём проблемы. Эти вопросы относятся к сфере компетенции команды разработчиков ядра, возглавляемой Линусом Торвальдсом.
- утилиты и библиотеки, поставляемые в составе Linux. Разработчику достаточно уметь использовать эти инструменты при построении модулей ядра, но более глубокие знания — это уже прерогатива сообществ GNU, FSF и разработчиков независимых проектов.
- вопросы интеграции создаваемого модуля в дерево исходных кодов Linux или любого его дистрибутива. Эти вопросы должны решаться системотехниками или внедренцами, берущими на себя ответственность за дальнейшую судьбу разрабатываемого проекта. Поэтому демонстрируемые примеры модулей будут создаваться не в дереве исходных кодов системы, а в отдельных каталогах целевых проектов.

Разработчика модулей ядра должны интересовать исключительно вопросы создания собственного модуля (драйвера) для использования в рамках некоего целевого прикладного программного проекта. Поэтому все вопросы, выходящие за рамки этой практической задачи, не будут рассматриваться в рамках данного цикла статей.

Эта короткое, но необходимое, введение поясняет, что можно ожидать от этой и последующих статей, а на какие вопросы они не смогут ответить.

Создание первого модуля ядра

Лучший способ научиться плавать — начать плавать. Поэтому, вместо скучных обсуждений терминологии, систематизации и архитектуры цикл начинается сразу с написания кода модулей ядра. Такой код будет интуитивно понятен любому программисту без особых пояснений. Хотя в дальнейшем, конечно, обязательно

придётся вернуться к обсуждению скучных вещей в объёме, необходимом для достижения поставленной цели. Обычно любую иллюстрацию из мира программирования начинают с примера "Hello world!". На самом деле, "Hello world!" модуль - это не самый лучший пример, так как все пишущие о ядре авторы используют именно его. Вместо этого лучше создать несколько простейших модулей, в конечном итоге имеющих суммарно такую же потребительскую ценность, как и "Hello world!", но при этом сразу иллюстрирующих гораздо больше понятий из мира модулей ядра и открывающих дорогу для дальнейших экспериментов. В листинге 1 приведен код вызываемого модуля (файл **md1.c**), а в листинге 2 — код вызывающего модуля (файл **md2.c**)

Листинг 1. Вызываемый модуль ядра Linux

```
1
2
3     #include <linux/init.h>
4     #include <linux/module.h>
5     #include "md.h"
6     MODULE_LICENSE( "GPL" );
7     MODULE_AUTHOR( "Oleg Tsiliuric <olej@front.ru>" );
8     char* mdl_data = "Привет мир!";
9     extern char* mdl_proc( void ) {
10         return mdl_data;
11     }
12     static char* mdl_local( void ) {
13         return mdl_data;
14     }
15     extern char* mdl_noexport( void ) {
16         return mdl_data;
17     }
18     EXPORT_SYMBOL( mdl_data );
19     EXPORT_SYMBOL( mdl_proc );
20     static int __init md_init( void ) {
21         printk( "+ module md1 start!\n" );
22         return 0;
23     }
24     static void __exit md_exit( void ) {
25         printk( "+ module md1 unloaded!\n" );
26     }
27     module_init( md_init );
28     module_exit( md_exit );
```

Листинг 2. Вызывающий модуль ядра Linux

```
1     #include <linux/init.h>
2     #include <linux/module.h>
3     #include "md.h"
4     MODULE_LICENSE( "GPL" );
5     MODULE_AUTHOR( "Oleg Tsiliuric <olej@front.ru>" );
6     static int __init md_init( void ) {
7         printk( "+ module md2 start!\n" );
8         printk( "+ data string exported from md1 : %s\n", mdl_data );
9         printk( "+ string returned mdl_proc() is : %s\n", mdl_proc() );
10        return 0;
11    }
12    static void __exit md_exit( void ) {
13        printk( "+ module md2 unloaded!\n" );
14    }
15    module_init( md_init );
16    module_exit( md_exit );
```

14
15
16

Также все модули, включают заголовочный файл (только обеспечивающий их синтаксическую связь), приведенный ниже:

```
1 extern char* mdl_data;  
2 extern char* mdl_proc( void );
```

Как и для сборки любого проекта, в данном случае понадобится файл сценария сборки с именем **Makefile**, приведенный в листинге 3. Опытные пользователи легко поймут, что в нем собираются 3 независимых модуля.

Листинг 3. Типовой сценарий сборки

```
1  
2 CURRENT = $(shell uname -r)  
3 KDIR = /lib/modules/$(CURRENT)/build  
4 PWD = $(shell pwd)  
5 TARGET1 = md1  
6 TARGET2 = md2  
7 TARGET3 = md3  
8 obj-m := $(TARGET1).o $(TARGET2).o $(TARGET3).o  
9 default:  
10 $(MAKE) -C $(KDIR) M=$(PWD) modules  
11 clean:  
12 @rm -f *.o *.cmd *.flags *.mod.c *.order  
13 @rm -f *.*.cmd *~ *.*~ TODO.*  
14 @rm -fR .tmp*  
15 @rm -rf .tmp_versions  
16 disclean: clean  
17 @rm *.ko *.symvers
```

Примечание. Согласно синтаксическим правилам утилиты `make`, каждая строка, имеющая отступ в листинге 3 и не начинающаяся с первого символа строки, должна начинаться с одного или нескольких символов ТАБ, но ни в коем случае не с пробелов.

Более детально разнообразные варианты содержимого **Makefile** для различных случаев сборки модулей будут рассмотрены позже, но и представленной информации уже достаточно, чтобы успешно выполнить команду сборки.

Листинг 4. Результат успешной сборки модуля

```
1 $ make  
2 make -C /lib/modules/2.6.32.9-70.fc12.i686.PAE/build M=/home/olej/TMP modules  
3 make[1]: Entering directory `/usr/src/kernels/2.6.32.9-70.fc12.i686.PAE'  
4 CC [M] /home/olej/TMP/md1.o  
5 /home/olej/TMP/md1.c:14: предупреждение: `mdl_local' определена, но нигде не исп  
6 CC [M] /home/olej/TMP/md2.o  
7 CC [M] /home/olej/TMP/md3.o  
8 Building modules, stage 2.  
9 MODPOST 3 modules  
10 CC /home/olej/TMP/md1.mod.o  
11 LD [M] /home/olej/TMP/md1.ko  
12 CC /home/olej/TMP/md2.mod.o  
13 LD [M] /home/olej/TMP/md2.ko  
14 CC /home/olej/TMP/md3.mod.o  
15 LD [M] /home/olej/TMP/md3.ko  
16 make[1]: Leaving directory `/usr/src/kernels/2.6.32.9-70.fc12.i686.PAE'
```

В листинге 4 приведен пример успешной сборки модуля (предупреждение касательно строки 14 файла **md1.c** не существенно), результаты сборок других модулей должны выглядеть аналогично. На этом вся работа по программированию и сборке завершена, и можно переходить к дальнейшим экспериментам и наблюдениям. Все ключевые понятия и термины, возникающие в ходе этих экспериментов, имеющие существенное значение для понимания природы и техники модулей, буду выделяться **таким шрифтом**. В результате сборки были созданы три файла модуля ядра:

```
1 $ ls *.ko
2 md1.ko md2.ko md3.ko
```

Это файлы объектного формата компилятора gcc, расширенные некоторыми дополнительными символами (в терминологии объектных модулей):

```
1 $ file md1.ko
2 md1.ko: ELF 32-bit LSB relocatable, Intel 80386, version 1 (SYSV), not stripped
```

Попробуем установить (загрузить) один из новых модулей в системе:

```
1 $ sudo insmod md2.ko
2 insmod: error inserting 'md2.ko': -1 Unknown symbol in module
```

Возникла ошибка, потому что модуль **содержит** ссылки на неизвестные ядру имена (хотя известно, что эти недостающие имена определены в другом модуле md1). В системном журнале об этой ошибке будет выведена следующая информация:

```
1 $ dmesg | tail -n30 | grep md
2 md2: Unknown symbol md1_data
3 md2: Unknown symbol md1_proc
```

Следует выполнить загрузку модулей в другом, на этот раз правильном, порядке:

```
1 $ sudo insmod md1.ko<br>$ sudo insmod md2.ko<br>$ lsmod | grep md
2 md2                                646  0
3 md1                                860  1 md2
```

Всё прошло нормально, но в разработанных модулях присутствовали вызовы функции `printk()`, до сих пор не рассматривавшейся, но, по аналогии с `printf()` она должна была выводить текстовые сообщения по ходу загрузки модулей. Всё это так, но в **пространстве пользователя**, при запуске приложения и вызове `printf()` вывод осуществляется на **управляющий терминал**, а таким терминалом является текстовая консоль или приложение графического терминала, если выполнение происходит в среде X Window System. Загрузка модулей в свою очередь выполнялась в пространстве ядра, где нет и не может быть никакого управляющего терминала, поэтому вывод `printk()` направляется демону системного журналирования, который помещает его, в частности, в **системный журнал (/var/log/messages)**. Этот вопрос еще будет обсуждаться дальше, а пока достаточно просто воспользоваться командой чтения системного журнала (`dmesg`), как показано ниже:

```
1 $ dmesg | tail -n60 | grep +
2 + module md1 start!
3 + module md2 start!
4 + data string exported from md1 : Привет мир!
5 + string returned md1_proc() is : Привет мир!
```

Другой способ найти выведенные сообщения, это изучить файл системного журнала, но в некоторых дистрибутивах для этого могут потребоваться права **root**:

```
1 $ sudo cat /var/log/messages | tail -n150 | grep +
```

```

2   Dec 17 20:08:03 notebook kernel: + module md1 start!
3   Dec 17 20:08:09 notebook kernel: + module md2 start!
4   Dec 17 20:08:09 notebook kernel: + data string exported from md1 : Привет мир!
5   Dec 17 20:08:09 notebook kernel: + string returned md1_proc() is : Привет мир!

```

После успешного создания и загрузки модулей можно их выгрузить с помощью команды `rmmmod`:

```

1   $ sudo rmmmod md1
2   ERROR: Module md1 is in use by md2

```

Однако на этом шаге снова возникает ошибка. Рассмотрим листинг выполнения команды `lsmod`, расположенный несколькими абзацами выше:

- на модуль **md1** ссылаются некоторые другие модули или объекты ядра: цифра **1** — это число таких ссылающихся модулей, которое называется **счётчиком ссылок**;
- дальше за счётчиком ссылок указывается **список** тех модулей, откуда исходят такие ссылки, в данном случае, это один модуль **md2**;
- до тех пор, пока число ссылок на любой модуль в системе **не станет нулевым**, модуль **не может быть выгружен**;
- другими словами, модуль может быть выгружен только после того, как будут выгружены все ссылающиеся на него другие модули (загруженные **после** него — здесь не может возникнуть цикличности или перекрёстности ссылок);
- бывают случаи (как будет показано дальше), когда модуль вообще не может быть выгружен, в частности, когда счётчик ссылок для модуля не может быть сделан нулевым в силу каких-то причин.

Следующая попытка выгрузить модули уже учитывает эти правила, сначала выгружая модуль **md2**, и только потом **md1**:

```

1   $ sudo rmmmod md2<br>$ sudo rmmmod md1<br>$ lsmod | grep md
2   $

```

К этому моменту модули были созданы, загружены и использованы (по крайней мере, удалось обнаружить их диагностический вывод), а затем выгружены, проделав тем самым полный жизненный цикл. В конечном итоге операционная система вернулась в **исходное состояние**. Теперь можно сделать некоторые комментарии к отдельным фрагментам исходного кода модулей, которые могут показаться не совсем очевидными:

1. Модуль **md1** **экспортирует** для использования другими модулями имя процедуры `md1_proc()` и, что далеко не так очевидно, имя структуры данных `md1_data`. Любой другой модуль (**md2**) может использовать в своём коде любые **экспортируемые** имена. Это могут быть имена, экспортируемые ранее загруженными модулями, но гораздо чаще это имена, **экспортируемые ядром**. Это множество экспортируемых имён ядра далее будет называться **API ядра**. Примером одного из вызовов из набора API ядра в показанных фрагментах кода является вызов `printk()`.
2. Модуль **md2**, использующий экспортируемое имя, связывается с этим именем по прямому **абсолютному адресу**. Как следствие этого, любые изменения (новая сборка), вносимые в ядро или экспортирующие модули, делают собранный модуль непригодным для использования. Именно поэтому бессмысленно предоставлять модуль в собранном виде — он должен собираться только на месте использования.
3. Модуль сможет использовать только те имена, которые явно экспортированы. В модуле **md1** специально показаны два других имени: `md1_local()` является локальным именем (модификатор `static`), непригодным для связывания, а

имя `mdl_noexport()` не объявлено как экспортируемое имя и также не может быть использовано вне модуля.

4. Почему в качестве строки, выводимой `mdl_proc()`, была выбрана строка "Привет мир!" в русскоязычном написании? Для того, чтобы сразу проверить прозрачность настроек самых разных подсистем Linux на работу с UNICODE представлением символьных данных в кодировке UTF-8 — в ранних версиях Linux всё было не так однозначно. Что в данном случае понимается под прозрачностью? Это единообразное и слаженное поведение на таких кодировках и подсистемы клавиатуры, и отображение в графических терминалах и текстовой консоли, и поведение системного журнала. Кроме того, для большей степени общности, интересно работать со строковыми представлениями, для которых значения `strlen()` значительно больше визуальной длины строки.
5. Зачем каждую выводимую строку предварять строкой "+"? Это **маркер**, отмечающий вывод из **собственных** модулей. В качестве него можно выбрать любой символ или вообще отказаться от него (что чаще всего и происходит). Но если настройки Linux таковы, что работают различные сервисы аудита или подобные службы, то они могут «засыпать» системный журнал достаточно плотным потоком своих сообщений, а сообщения собственных модулей будут сильно разрежены таким потоком. Так что их придётся потом разыскивать в этом потоке. Заблаговременно предварять сообщения собственных модулей фиксированными маркерами — это простейший способ позже осуществить их отбор и группировку, по крайней мере, в иллюстрационных целях, как и было показано. Отобрать собственные сообщения можно с помощью команд, подобных приведенной ниже:

```
1 $ dmesg | tail -n60 | grep +
```

Существует и другой способ отобрать интересующие сообщения из системного журнала. Если знать внешний вид сообщений аудита, можно напротив: **отбраковать** засоряющие листинг сообщения по соответствующему регулярному выражению. Достаточно часто журнал может быть засорён повторяющимися сообщениями вида:

```
1 $ dmesg | tail -n2
2 audit(:0): major=340 name_count=0: freeing multiple contexts (16)
3 audit: freed 16 contexts
```

В этом случае используется фильтр, который в дальнейшем иногда будет фигурировать в демонстрируемых примерах:

```
1 $ dmesg | tail -n50 | grep -v ^audit
```

В листинге 5 приведен исходный код последнего из модулей, представленных в архиве проекта (файл **md3.c**):

Листинг 5. Модифицированная версия модуля md2

```
1
2 #include <linux/init.h>
3 #include <linux/module.h>
4 #include "md.h"
5 MODULE_LICENSE( "GPL" );
6 MODULE_AUTHOR( "Oleg Tsiliuric <olej@front.ru>" );
7 static int __init md_init( void ) {
8     printk( "+ module md3 start!\n" );
9     printk( "+ data string exported from mdl : %s\n", mdl_data );
10    printk( "+ string returned mdl_proc() is : %s\n", mdl_proc() );
11    return -1;
12 }
```

Это ещё одна форма модуля, которая не часто упоминается в других статьях, посвященных модулям ядра, но очень удобная для отладочных действий с модулями и для других специальных действий:

- функция инициализации этого модуля, выполнив все предписанные ей действия, **преднамеренно** возвращает ненулевое значение, что означает ошибку инициализации модуля;
- тогда такой модуль **не будет подгружен** к ядру, но произойдёт это уже **после** выполнения кода инициализирующей функции модуля **в пространстве ядра**;
- а если такой модуль по замыслу не загружается, то он, в принципе, **может** не иметь функции выгрузки.

Такой модуль очень напоминает привычные приложения пространства пользователя (процессы), но только код этот выполняется в пространстве ядра со всеми **ядерными** привилегиями. Ниже приведен вывод после запуска этого модуля:

```
1 $ sudo insmod md3.ko
2 insmod: error inserting 'md3.ko': -1 Operation not permitted
3 $ dmesg | tail -n60 | grep +
4 + module md3 start!
5 + data string exported from mdl : Привет мир!
6 + string returned mdl_proc() is : Привет мир!
7 $ lsmod | grep md
8 mdl                  860    0
```

Заключение

В первой статье этого цикла был представлен краткий обзор техники модульного программирования, который можно подытожить следующими утверждениями:

- код модуля объявляет две функции: инициализации и финализации (завершения);
- имена функций, выполняющих эти задачи, объявляются в макросах `module_init()` и `module_exit()`;
- эти функции должны в точности соответствовать прототипу, показанному в листингах примеров;
- функция инициализации выполняется при загрузке модуля в ядро, если эта функция возвращает нулевое значение (успех), то код модуля остаётся резидентным и выполняет дальнейшую работу;
- функция завершения вызывается при выгрузке модуля командой `rmmmod`;
- функция завершения по своему прототипу **не имеет возвращаемого значения**, поэтому, начавшись, она уже не имеет механизмов сообщить о своём неудачном выполнении;
- большинство операций над модулями (`insmod`, `rmmmod`, `modprobe`) требуют прав **root**, но некоторые индикативные команды (`lsmod`) могут выполняться и от имени обычного пользователя.

Названные выше соглашения по объявлению функций инициализации и завершения являются общепринятыми. Но существует ещё один, не документированный, способ описания этих функций: воспользоваться их предопределёнными именами, а именно `init_module()` и `cleanup_module()`, как показано ниже:

```
1 int init_module( void ) {
2     ...
3 }
4 void cleanup_module( void ) {
5     ...
6 }
```

При такой записи необходимость в использовании макросов `module_init()` и `module_exit()` отпадает, а использовать модификатор `static` с этими функциями нельзя (именно эти имена и используются при загрузке и удалении модуля). Конечно, такая запись не способствует улучшению читаемости кода, но она может существенно сократить количество рутинного кода, особенно в коротких иллюстративных примерах, так что она будет использоваться в демонстрируемых примерах.

Ресурсы для скачивания

- [этот контент в PDF](#)
 - [исходный код примеров](#) (export-data-1.tgz | 6KB)
- Похожие темы
- [Инструкция по работе с примерами исходного кода в цикле статей "Разработка модулей ядра Linux"](#)
 - В [разделе Linux](#) сайта developerWorks находятся [сотни практических руководств и учебных пособий](#), материалы для скачивания, форумы для обсуждения и множество других ресурсов для разработчиков и администраторов Linux.
 - [Знакомьтесь с продуктами IBM](#) различными способами: загружайте ознакомительные версии, испытывайте продукты в онлайн-режиме или в облачной среде или проведите несколько часов в [SOA Sandbox](#), чтобы узнать, как эффективно создавать SOA-приложения.
 - Следите за [публикациями developerWorks в Твиттере](#) или подпишитесь на [канал твитов по Linux на developerWorks](#).