

LINUX DEVICE DRIVERS

LINUX

设备驱动程序



O'REILLY®
中国电力出版社

ALESSANDRO RUBINI 著
LISOLEG 译

LINUX

设备驱动程序

ALESSANDRO RUBINI

LISOLEG 译

O'REILLY®

Beijing • Cambridge • Farnham • Köln • Paris • Sebastopol • Taipei • Tokyo

中国电力出版社

图书在版编目 (CIP) 数据

Linux 设备驱动程序: / (美) 鲁宾尼 (Rubini, A.) 编著; 聊鸿斌等译. - 北京: 中国电力出版社, 2000. 4

(开源软件丛书)

书名原文: Linux Device Drivers

ISBN 7-5083-0221-4

I .L ... II .①鲁 ... ②聊 ... III .操作系统 (软件), Linux IV .TP316

中国版本图书馆 CIP 数据核字 (1999) 第 75034 号

北京市版权局著作权合同登记

图字: 01-1999-3753 号

© 1998 by O'Reilly & Associates, Inc.

Simplified Chinese Edition, co-published by O'Reilly & Associates, Inc. and Chinese Electric Power Press, 2000. Authorized translation of the English edition, 1998 O'Reilly & Associates, Inc., the owner of all rights to publish and sell the same.

All rights reserved including the rights of reproduction in whole or in part in any form.

简体中文版 中国电力出版社 2000。授权英文译文, 1998, 奥莱理有限公司。此译本的出版和销售得到出版权和销售权的所有者——奥莱理有限公司的许可。

版权所有, 未得书面许可, 本书的任何部分和全部不得以任何形式复制。

书 名 / Linux 设备驱动程序

书 号 / ISBN 7-5083-0221-4

责任编辑 / 胡顺增, 杨伟国

审 校 / 章远琳

封面设计 / Ellie Volckhausen, Hanna Dyer, 张健

出版发行 / 中国电力出版社

地 址 / 北京三里河路 6 号 (邮政编码 100044)

经 销 / 全国新华书店

印 刷 / 北京市地矿印刷厂

开 本 / 787 毫米 × 1092 毫米 16 开本 30.5 印张 550 千字

版 次 / 2000 年 4 月第一版 2000 年 4 月第一次印刷

印 数 / 0001-5000 册

定 价 / 59.00 元 (册)

目录

前言	1
第一章 Linux 内核简介	11
驱动程序作者的作用	12
划分内核	13
设备和模块的分类	16
安全问题	18
版本编号	19
许可证术语	21
全书概貌	22
第二章 编写和运行模块	25
模块与应用程序	26
编译和加载	31
版本相关性	33
内核符号表	35
初始化和终止	37
使用资源	42

自动和手动配置	48
在用户空间编写驱动程序	50
快速索引	52
第三章 字符设备驱动程序	57
scull 的设计	57
主设备号和次设备号	59
文件操作	66
file 结构	70
Open 和 Close	71
Scull 的内存使用	76
读和写	80
试试新设备	86
快速索引	86
第四章 调试技术	89
用打印信息调试	89
通过查询调试	94
通过监视进行调试	98
调试系统故障	100
使用调试器	111
第五章 字符设备驱动程序的扩展操作	117
ioctl	118
阻塞型 I/O	130
Select	140
异步触发	145
定位设备	148
设备文件的访问控制	150
快速索引	156

第六章 时间流	161
内核中的时间间隔	161
获取当前时间	162
延迟执行	164
任务队列	168
内核定时器	179
快速索引	183
第七章 获取内存	185
kmalloc 函数的内幕	185
get_free_page 和相关函数	188
vmalloc 和相关函数	191
“脏”的处理方法 (Playing Dirty)	194
快速索引	195
第八章 硬件管理	197
使用 I/O 端口	198
使用并口	203
访问设备卡上的内存	206
访问字符模式的视频缓冲区	211
快速索引	212
第九章 中断处理	215
准备并口	215
安装中断处理程序	216
实现中断处理程序	228
下半部	233
共享中断	239
中断驱动的 I/O	243

竞争条件	244
中断处理的版本相关性	253
快速索引	255
第十章 合理使用数据类型	259
使用标准 C 类型	260
分配确定的空间大小给数据项	261
接口特定的类型	262
其他与移植有关的问题	263
快速索引	266
第十一章 kerneld 和高级模块化	269
按需加载模块	269
模块中的版本控制	275
跨过卸载 / 装载的持久存储	280
快速索引	282
第十二章 加载块设备驱动程序	285
注册驱动程序	285
头文件 blk.h	290
处理请求	293
挂载 (Mounting) 是如何工作的	300
ioctl 方法	301
可拆卸的设备	305
可分区设备	308
中断驱动的块设备驱动程序	317
快速索引	320

第十三章 MMAP 和 DMA	323
Linux 中的内存管理	323
mmap 设备操作	332
直接内存访问 (DMA)	347
快速索引	358
第十四章 网络驱动程序	361
snull 如何设计	362
与内核相连	366
设备结构的细节	371
打开和关闭	380
包发送	382
包接收	384
中断驱动的操作	386
插座缓冲区	388
地址解析	391
加载时配置	394
运行时配置	395
自定义 ioctl 命令	397
统计信息	399
选播 (multicasting)	399
快速索引	403
第十五章 外设总线概览	407
PCI 接口	407
回顾: ISA	423
其他 PC 总线	425
Sbus	427
快速索引	428

第十六章 内核源代码的物理布局	431
引导内核	431
引导之前	433
Init 进程	438
kernel 目录	439
mm 目录	441
fs 目录	443
网络	448
IPC 和 lib 函数	449
Drivers	450
体系结构相关性	452
第十七章 最新进展	453
模块化	454
文件操作	459
访问用户空间	463
任务队列	469
中断管理	469
位操作	470
转换函数	470
vremap	471
虚拟内存	472
处理内核空间错误	472
其他变化	474

前言



作为一名电子工程师，而且是一个什么都想自己做的人，我一向愿意用计算机来控制外部设备。甚至还在我们父辈的Apple-2e时代，我就已经开始寻找别的平台，希望可以与我定制的环境相连，并且可以写我自己的驱动程序软件。不幸的是，20世纪80年代的个人电脑的功能并没有那么强大，不论在软件层还是硬件层上，个人电脑的内部设计都远不如Apple-2e，并且在很长时间内，可以获得的文档都不能令人满意。但是，接着发生的事是Linux出现了，我决定试一试，于是买了昂贵的386主板，并且没有跑任何商业软件。

那时候，我正在大学里使用Unix系统，并为这样一个聪明的操作系统所震动，特别是又补充了GNU工程提供给用户的更智能的工具后，我更是为Unix所深深陶醉。在我自己的个人电脑主板上跑Linux实在是一个有趣的经历，我甚至可以自己写设备驱动程序，而且又可以玩烙铁了。我不断告诉别人：“当我长大了，我一定会成为一个黑客”，Linux是实现这个梦想的绝佳平台。这就是说，用不着长大我就可以实现梦想了。

当Linux成熟起来以后，越来越多的人对编写特制电路的设备驱动程序和商业设备的驱动程序感兴趣。正如Linus Torvalds所说的那样：“我们退回到这样一个时代——人人都为自己编写设备驱动程序。”

当我不能再写出有创意的编码以后，我就开始为《Linux Journal》写一些技术性文章了，这也算是为Linux社团的一点贡献吧。后来，O'Reilly的Andy Oram有

意让我编写一本讲设备驱动程序的书，我接受了这个任务。尽管真正的黑客可以在正式发布的内核代码中找到所有必要的信息，但是写出来的文本对提高编程技巧还是有用的。你拿到的这个东西是我花几小时的时间对内核资源耐心grep的结果，我希望最后的结果对得起我的努力。如果本书能作为那些想成为内核黑客又不知道从何下手的人的一个起点的话，就足以实现我的愿望了。

本书读者

从技术角度来讲，本书应该为您提供唾手可得的途径去理解内核内幕以及 Linus 本人在开发时所做的设计抉择。尽管本书的主要目的是讲述如何写设备驱动程序，但它所包含的内容应该也对内核的运行给出一个完整的概貌。

对那些想玩电脑的人和那些涉及 Linux 机器内部的专业程序员来讲，本书都会是一个很有趣的资料来源。注意：“Linux 机器”这个概念比“跑 Linux 的 PC”这个概念还要广泛，因为我们的操作系统支持很多平台，而且内核编程也不局限于某一种平台上。

Linux 的狂热支持者会发现本书提供了大量的精神食粮。开始可以玩一玩编程，然后就可以加入到开发者小组中了，他们可是在孜孜不倦地工作以提供新的功能，提高系统性能。Linux 仍处在不断完善的过程当中，并且总能为新加入的程序员提供新的空间。

换句话说，如果你只是想给你自己的设备编一个设备驱动程序，而不想在内核内幕上下什么工夫，本书的模块化结构也绝对可以满足你的要求。如果你不想深入细节，你可以跳过大多数技术章节，而直接查找设备驱动程序的标准 API，它们与系统的其他部分是无缝连接的。

本书主要目的是为 Linux 2.0 写一个内核模块。模块是对象代码，可以在运行的内核中动态加载新功能。讨论也会涉及到 1.2 版本的内核。最后一章描述从 2.0 到 2.1.43（在对本书进行技术审校时的最新版本）的驱动程序接口的变化。

材料的组织

本书介绍的主题会越来越复杂，它们可以分成两个部分。第一部分（第一章到第十章）从内核模块的正确安装开始介绍，会涉及到写字符设备驱动程序需要面对的各方面问题。每章会讨论一个独立主题，并且在末尾包含一个“符号表”，在实际开发时，可以用作参考。当我写自己的设备驱动程序时，我发现我会回头去查自己写的章节，我希望你也可以充分利用这个符号表。

贯穿本书第一部分材料的组织方式大致是从面向软件到面向硬件。这意味着你可以在没有附加设备的情况下在你的机器上测试软件。每章都包含源代码，并指出可以使用的驱动程序的例子，这些例子在所有 Linux 机器上都可以实现。在第八章和第九章，我会要求你在并口上连一根金属线，用来测试一下硬件，但这种要求对所有人来讲都是很容易实现的。

本书第二部分描述了一些块设备驱动程序和网络接口，以及更进一步的深入讨论。这里讨论的大部分东西很有可能在你写实际驱动程序中用不到，但我希望第一部分能够引起你的足够兴趣来阅读第二部分。

事实上，我罗列的大部分材料都很有意思，而且与实际编写设备驱动程序时需要的材料没有什么关系。在我写这本书的时候，很多学生就他们的需求问了我一些关于 Linux 的问题。他们一定很高兴看到这些章节对他们的工作有帮助，即使他们的工作与编写驱动程序无关。

背景信息

要想阅读本书，你应该熟悉 C 语言编程，也要有点 Unix 的专业知识，例如我会经常提到 Unix 命令和管道。

在硬件层，不需要预先的专业知识，前提只要求总体概念清晰就可以了。本书不基于什么特殊的 PC 硬件设备，而当我用到什么特殊设备的时候，我会提供所有相关信息。

如果能够连接到互联网上，读者可以享受很多便利，因为从网络上可以获得很多新鲜有趣的文档和升级软件。当然，有上网条件并不是必须的，我自己的上网条件就很有限（主要得托 Italian 电话通讯公司快速网络速度的福）。

只要一涉及到软件，你就需要安装 Linux 系统了，这样就可以运行例子的驱动程序。注意，任何发布版本都可以使用（并且所有的硬件平台也都适用）。第一章完整地列出了所需软件包的清单，因为如果放到“序”中，很多读者会漏掉这些信息（我希望能有很多人可以跳过第一章，大多数读者不都是黑客吗？）。

深入信息来源

本书中涉及的大多数信息都是直接取材于内核的。只要你的系统上装了 Linux，就几乎不需要什么文档来做补充。在写设备驱动程序时，几乎没有什么富有趣味的书可供参考，而主要信息来源就是内核的源代码和你设备的技术文档。不用再说什么了，你应该很感激描述你机器平台的手册。

至于了解内核内部机制的工作，最好的信息来自互联网（仅次于源文件）。《Linux Journal》也有一些有趣的技术性文章。查看“内核之角”卷，但跳过我的文章——因为我可能重复自己。没有“内核之角”标识的文章实际上也挺有意思，但是一般技术性没有本书读者要求的那样高。

在互联网上，我建议查找下面网址：

<http://www.redhat.com:8080/>

小红帽上的超新闻（HYPERNEWS）服务器提供“内核黑客指南”，这是关于内核内幕很有趣的文档。其中的一些章节已经很老了，但是最近更新了其中一部分东西。依我之见，这些材料相当有趣。

<http://www.kernel.org/>

<ftp://ftp.kernel.org>

本站点是 Linux 内核开发的中心，可以获得最新发行版和相关信息。注意，这个 FTP 站点在全球都有镜像，所以你可以找最近的。

<ftp://sunsite.unc.edu/pub/Linux/docs/>

<ftp://tsx-11.mit.edu/pub/Linux/docs/>

“Linux 文档计划”中有大量称为“HOWTO”的有趣文档，一些与内核相关的主题极具技术性。Sunsite 和 tsx-11 还有大量在 Linux 上可以应用的程序。总的来讲，不仅仅是 docs/ 目录下的文档，他们都相当有趣。我敢肯定你已经知道这些文件了，但我觉得还是有必要提到他们。

<http://www.ssc.com/>

SSC，专业系统顾问，是《Linux Journal》的出版商，他们的站点有他们出版的大部分文章的 HTML 版。他们发表的有趣文章在出版不久后就转换成 HTML 文件，在 Web 上发行。

<http://www.conceta.it/linux/>

这是个意大利站点，Linux 的拥护者积累了大量信息，这些信息是关于所有正在运行的和 Linux 有关的项目。也许你已经知道一些有关 Linux 开发的 HTTP 连接的站点，如果你不知道，这个站点是个很好的起点。

相关书目

除了源代码和互联网资源，很多好书也涉及到本书讨论的一些主题。下面的列表是我个人在这个领域内选择的一些书籍。我列的这些书或者是 Unix 系统软件功能文档，或者描述了有趣的硬件主题。我没有列出任何关于 PC 结构的书，因为现在的书太多了。不幸的是我也没法建议任何关于 Sparc 结构的书，因为我找不到这样的书。如果你需要有关信息，我绝对相信通过 Web 你可以填补这个空缺。

[0] Bach, Maurice. *The Design of the Unix Operating System*, Prentice Hall. 1986.

本书尽管内容相当陈旧，但涵盖了所有运行 Unix 的主题。它可是 Linus 本人编写 Linux 第一版时主要的灵感来源。

[1] Beck, Michael. *Linux Kernel Internals*. Addison-Wesley. 1997.

本书重点在于 Linux 的内部数据结构和算法。如果你喜爱这些详细介绍，你

会很喜欢这本书的。第一版对应 Linux 1.2 版，我不知道最新版本有什么进展。2.0 版及其后继版本与 1.2 版有很大差别。

- [2] Stevens, Richard. *Advanced Programming in The Unix Environment*. Addison-Wesley. 1992.

这里介绍了所有 Unix 系统调用的详细资料。在使用设备高级功能的方法时，本书会是很好的参考。对 Unix 语义中任何可能的不明之处，参考本书都可以得到解决。

- [3] Stevens, Richard. *Unix Network Programming*. Prentice Hall. 1990.

如你所想，本书是网络主题的高效参考书。在主题选择范围和质量上和“高级编程”相匹配。本书包含有各种源代码可以测试用户网络编程空间的各个角落。

- [4] Comer, Douglas, and Stevens, David. *Internetworking with TCP/IP Vol I,II,III*. Prentice Hall. 1991.

本书搜集了所有关于 Internet 的网络信息。描述了 Internet 协议族和一些它们的实现。

- [5] Shanley, Tom, and Anderson, Don. *PCI System Architecture*. Addison-Wesley. 1995.

本书详述 PCI 总线和标准接口。在大多数硬件主题中都可以找到类似“系统结构”这样的标题，这些都是由一个作者写的。这些书都很有趣，尽管有那么点偏向于 PC。我最喜欢 PCI 那一卷。这些书中至少有一本我不喜欢，但如果仔细研究，可以看出这本书不错，就是所描述系统结构不怎么样。

- [6] Digital Semiconductor. *Alpha AXP Architecture Handbook*. Digital Semiconductor. 1994.

从 Digital Semiconductor 可以免费获得本书和“Alpha AXP Reference Manual”。它们介绍了 Alpha 处理器的机器语言，以及所涉及使用的设计主题。本书的订货号码是 EC-QD2KA-TE。

本书使用的约定

下面给出本书所使用的排版字体约定。

Italic

用于文件和目录名、程序和命令名、命令行选项、电子邮件地址以及路径名、URL 和突出表示新名词。

Boldface

用来表示按键（如 **Ctrl-N**）。

Constant Width

用来表示变量选项、关键字，或是用户用来替代实际值的文本。

Constant Bold

在例子用来表示应该由用户键入的命令或其他文本。

我们很愿意听您的反馈

我们已尽全力调整本书内容，但您仍可能发现有些内容不对（甚至我们可能出了错误！）。如果您的建议与我们以后版本有关，请告诉我们您找到的错误以及您的建议，写信到：

美国：

O'Reilly & Associates, Inc.
101 Morris Street
Sebastopol, CA 95472

中国：

100031 北京市西城区复兴门内大街 160 号 2411 室
奥莱理软件（北京）有限公司

询问技术问题或对本书的评论，请发电子邮件到：

info@mail.oreilly.com.cn

最后，您可以在 WWW 上找到我们：

http://www.oreilly.com

http://www.oreilly.com.cn

致谢

本书不仅仅是我个人努力的结果：许多人不仅在物质上给了我充分的帮助，精神上也给了我巨大的支持。我要感谢 Quant-X 的 Dreyer 先生，他借给我一部 Alpha 计算机，这样我可以测试本书例子中的代码的可移植性。Sun-Italia 对我也很好，他们借了我一部 Sparc 机器，这样我可以把他们机器的操作系统升级成我需要的。ImageNation 赠送给我一个 PCI 视频捕捉卡，我可以用来研究 PCI 和 DMA 特性。

如果没有 Andy Oram 和 Michael Johnson 的支持，没有 Federica——我的女朋友、我妻子在心理上对我的支持，本书是不可能完成的。Andy 是给我强有力支持的编辑，而正是 Michael 要求我给《Linux Journal》写东西，并且把我介绍给 Andy——如果有什么人对本书觉得内疚，那就是 Michael 了。我还要感谢 Georg van Zezschwitz，他介绍给我这个奇妙的内核模块世界，并且在给《Linux Journal》写文章时，给了我很大帮助。我想感谢 Silvana Ranzoli，我高中时代的英语老师，由于她无情地（有时简直感觉像是残酷）承诺可以得益于她的班级。我感激 Ellen Siever，她纠正了我在高中以后学到的所有不正规语法，由于我对重写从来不满意，所以每当我有那种黑客主义和极端细致倾向的时候，她对我总是格外耐心。

Alan Cox、Greg Hankins、Hans Lermen、Heiko Eissfeldt 和 Miguel de Icaza（按照首字母顺序）从技术角度评论了本书。他们的意见和建议对我的小错和不足很有用。我要感谢他们在我写书的过程中所花费的宝贵时间，这看起来跟他们这些大拿要做的事情毫不相干。