

- 全面指导Linux内核的设计与实现
- 填补Linux内核理论与实践之间的鸿沟
- 畅销技术图书新版

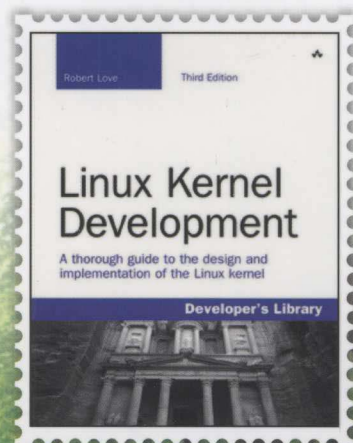
原书第3版

# Linux

# 内核设计与实现

## Linux Kernel Development Third Edition

(美) Robert Love 著  
陈莉君 康华 译



华章专业开发者书库

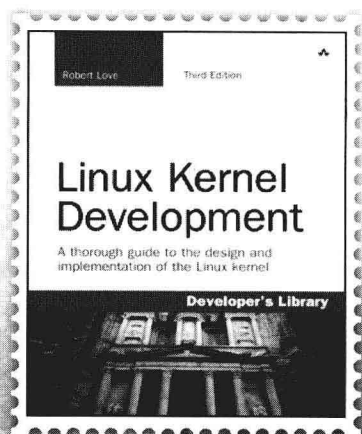
原书第3版

# Linux

## 内核设计与实现

Linux Kernel Development Third Edition

(美) Robert Love 著  
陈莉君 康华 译



机械工业出版社  
China Machine Press

本书基于 Linux 2.6.34 内核详细介绍了 Linux 内核系统，覆盖了从核心内核系统的应用到内核设计与实现等方面的内容。本书主要内容包括：进程管理、进程调度、时间管理和定时器、系统调用接口、内存寻址、内存管理和页缓存、VFS、内核同步以及调试技术等。同时本书也涵盖了 Linux 2.6 内核中颇具特色的内容，包括 CFS 调度程序、抢占式内核、块 I/O 层以及 I/O 调度程序等。本书采用理论与实践相结合的路线，能够带领读者快速走进 Linux 内核世界，真正开发内核代码。

本书适合作为高等院校操作系统课程的教材或参考书，也可供相关技术人员参考。

Simplified Chinese edition copyright © 2011 by Pearson Education Asia Limited and China Machine Press.  
Original English language title : *Linux Kernel Development Third Edition* (ISBN 978-0-672-32946-3) by Robert Love , Copyright © 2010.

All rights reserved.

Published by arrangement with the original publisher, Pearson Education, Inc., publishing as Addison-Wesley.

本书封面贴有 Pearson Education（培生教育出版集团）激光防伪标签，无标签者不得销售。

封底无防伪标均为盗版

版权所有，侵权必究

本书法律顾问 北京市展达律师事务所

本书版权登记号：图字：01-2010-4824

图书在版编目（CIP）数据

Linux 内核设计与实现（原书第3版）/（美）拉美\_（Love, R）著；陈莉君，康华译. —北京：机械工业出版社，2011.6

（华章专业开发者书库）

书名原文：Linux Kernel Development, Third Edition

ISBN 978-7-111-33829-1

I . L… II . ①拉… ②陈… ③康… III . Linux 操作系统—程序设计 IV . TP316.89

中国版本图书馆 CIP 数据核字（2011）第 047916 号

机械工业出版社（北京市西城区百万庄大街 22 号 邮政编码 100037）

责任编辑：秦 健

北京市荣盛彩色印刷有限公司印刷

2011 年 6 月第 1 版第 1 次印刷

186mm×240mm·22 印张

标准书号：ISBN 978-7-111-33829-1

定价：69.00 元

凡购本书，如有缺页、倒页、脱页，由本社发行部调换

客服热线：（010）88378991；88361066

购书热线：（010）68326294；88376949；68995259

投稿热线：（010）88379604

读者信箱：hzjsj@hzbook.com

# 译者序

不知不觉涉足 Linux 内核已经十多个年头了，与其他有志（兴趣）于此的朋友一样，我们也经历了学习—实用—追踪—再学习的过程。也就是说，我们也是从漫无边际到茫然无措，再到初窥门径，转而觉得心有戚戚焉这一路走下来的。其中甘苦，犹然在心。

Linux 最为人称道的莫过于它的自由精神，所有源代码唾手可得。侯捷先生云：“源码在前，了无秘密。”是的，但是我们在面对它的时候，为什么却总是因为这种规模和层面所造就的陡峭学习曲线陷入困顿呢？很多朋友就此倒下，纵然 Linux 世界繁花似锦，纵然内核天空无边广阔，但是，眼前的迷雾重重，心中的阴霾又怎能被阳光驱散呢？纵有雄心壮志，拔剑四顾心茫然，脚下路在何方？

Linux 内核入门是不容易，它之所以难学，在于庞大的规模和复杂的层面。规模一大，就不易现出本来面目，浑然一体，自然不容易找到着手之处；层面一多，就会让人眼花缭乱，盘根错节，怎能让人提纲挈领？

“如果有这样一本书，既能提纲挈领，为我理顺思绪，指引方向，同时又能照顾小节，阐述细微，帮助我们更好更快地理解 STL 源码，那该有多好。”孟岩先生如此说，虽然针对的是 C++，但道出的是研习源码的人们共同的心声。然而，Linux 源码研究的方法却不大相同。这还是由于规模和层面决定的。比如说，在语言学习中，我们可以采取小步快跑的方法，通过一个个小程序和小尝试，就可以取得渐进的成果，就能从新技术中有所收获。而掌握 Linux 呢？如果没有对整体的把握，即使你对某个局部的算法、技术或是代码再熟悉，也无法将其融入实用。其实，像内核这样的大规模的软件，正是编程技术施展身手的舞台（当然，目前的内核虽然包含了一些面向对象思想，但还不能让 C++ 一展身手）。

那么，我们能不能做点什么，让 Linux 的内核学习过程更符合程序员的习惯呢？

Robert Love 回答了这个问题。Robert Love 是一个狂热的内核爱好者，所以他的想法自然贴近程序员。是的，我们注定要在对所有核心的子系统有了全面认识之后，才能开始自己的实践，但却完全可以舍弃细枝末节，将行李压缩到最小，自然可以轻装快走，迅速进入动手阶段。

因此，相对于 Daniel P. Bovet 和 Marco Cesati 的内核巨著《Understanding the Linux Kernel》，它少了五分细节；相对于实践经典《Linux Device Drivers》，多了五分说理。可以说，本书填补了 Linux 内核理论和实践之间的鸿沟，真可谓“一桥飞架南北，天堑变通途”。

就我们的经验，内核初学者（不是编程初学者）可以从本书着手，对内核各个核心子系统有个整体把握，包括它们提供什么样的服务，为什么要提供这样的服务，又是怎样实现的。而且，本书还包含了 Linux 内核开发者在开发时需要用到的很多信息，包括调试技术、编程风格、注意事项等。在消化本书的基础上，如果你侧重于了解内核，可以进一步研究《Understanding the Linux Kernel》和源代码本身；如果你侧重于实际编程，可以研读《Linux Device Drivers》，

直接开始动手工作；如果你想有一个轻松的内核学习和实践环节，请访问我们的网站 [www.kerneltravel.net](http://www.kerneltravel.net)。

Linux 内核是一艘永不停息的轮船，它将驶向何方我们并不知晓，但在这些变化的背后，总有一些原理是恒定不变的，总有一些变化是我们想知晓的，比如调度程序的大幅度改进，内核性能的不断提升，本书第 3 版虽然针对的是较新的 2.6.34 Linux 内核版本，但在旧版本上积累的知识 and 经验依然有效，而新增内容将使读者在应对变化了的内核代码时更加从容。

感谢牛涛和武特，他们在第 2 版和第 3 版差异的校对中花费了大量精力。感谢素不相识的网友 Cheng Renquan，他主动承担了其中一章的修订。还要感谢苏锦绣、黄伟、王泽宇、赵格娟、刘周平、周永飞、曹江峰、陈白虎和孟阿龙，他们参与了后期的校对和查错补漏。

最后，特别感谢我的合作者康华，从十年前一块分析 Linux 内核代码到今天，他对技术孜孜不倦的追求不但在业界赢得声誉，也使我们在翻译过程中所遇到的技术难点和晦涩句子被一一迎刃而解。感谢合作者张波，他流畅有趣的文笔让本书少了份枯燥，多了份趣味。

陈莉君

2010 年 10 月

# 序 言

随着 Linux 内核和 Linux 应用程序越来越成熟，越来越多的系统软件工程师涉足 Linux 开发和维护领域。他们中有些人纯粹是出于个人爱好，有些人是为 Linux 公司工作，有些人是为硬件厂商做开发，还有一些是为内部项目工作的。

但是所有人都必须直面一个问题：内核的学习曲线变得越来越长，也越来越陡峭。系统规模不断扩大，复杂程度不断提高。虽然现在的内核开发者对内核的掌握越发炉火纯青，但新手却无法跟上内核发展的步伐，长此以往将出现青黄不接的断层。

我认为这种新老鸿沟已经成为内核质量的一个隐患，而且问题将继续恶化。所以那些真正关心内核的人已经开始致力于扩大内核开发群体。

解决上述问题的一个方法是尽量保证代码简洁：接口定义合理，代码风格一致，“一次做一件事，做到完美”等。这也就是 Linus Torvalds 倡导的解决办法。

我提倡的解决办法是对代码慷慨地加上注释，即能够让读者立刻了解代码开发者意图的文字（识别意图和实现之间差异的工作称为调试。如果意图不明确显然调试就难以进行）。

可是，即使有注解，也没办法清楚地展现内核的各个主要子系统的全景，说明它们到底要做什么。那么，这些开发者又该从何下手呢？

由文字材料来说明这些在起步阶段就该理解的材料，其实是最合适的。

Robert Love 的贡献就在于此，有经验的开发者可以通过本书全面了解内核子系统提供的服务，同时还可以了解这些服务是怎么实现的。对不少人来说，这些知识就已经足够了：那些好奇的人，那些应用程序开发者，那些想对内核的设计品头论足一番的人，都有足够的谈资了。

但是学习本书同样可以作为那些有抱负的内核开发者更上一层楼的契机，可以帮他们更改内核代码以达到预定的目标。我建议有抱负的开发者能够亲身实践：理解内核某部分的捷径就是对它做些修改，这样能为开发者揭示仅仅通过看内核代码无法看到的深层机理。

严肃认真的内核开发者应该加入开发邮件列表，不断和其他开发者交流。这是内核开发者相互切磋和并肩前进的最好方法。而 Robert 在书中对内核生活中至关重要的文化和技巧都做了精彩介绍。

请学习和欣赏 Robert 的书吧。想必你也希望能精益求精，继续探索，成为内核开发社区中的一员，那么首先你要清楚的是：社区欢迎你。我们评价和衡量一个人是根据他所作的贡献，当你投身于 Linux 时，你要明白：虽然你仅仅贡献了一小份力，但马上就会有数千万或上亿人受益。这是我们的欢乐之源，也是我们的责任之本。

Andrew Morton

# 前言

在我刚开始有把自己的内核开发经验集结成册，撰写一本书的念头时，我其实也觉得有点头绪繁多，不知道该从何下手。我实在不想落入传统内核书籍的窠臼，照猫画虎地再写这么一本。不错，前人著述备矣，但我终究是要写出点儿与众不同的东西来，我的书该如何定位，说实话，这确实让人颇费思量。

后来，灵感终于浮现出来，我意识到自己可以从一个全新的视角看待这个主题。开发内核是我的工作，同时也是我的嗜好，内核就是我的挚爱。这些年来，我不断搜集与内核有关的奇闻轶事，不断积攒关键的开发诀窍，依靠这些日积月累的材料，我可以写一本关于开发内核该做什么——更重要的是——不该做什么的书籍。从本质上说，这本书仍旧是描述 Linux 内核是如何设计和实现的，但是写法却另辟蹊径，所提供的信息更倾向于实用。通过本书，你就可以做一些内核开发的工作了——并且是使用正确的方法去做。我是一个注重实效的人，因此，这是一本实践的书，它应当有趣、易读且有用。

我希望读者可以从这本书中领略到更多 Linux 内核的精妙之处（写出来的和没写出来的），也希望读者敢于从阅读本书和读内核代码开始跨越到开始尝试开发可用、可靠且清晰的内核代码。当然如果你仅仅是兴致所至，读书自娱，那也希望你能从中找到乐趣。

从第 1 版到现在，又过了一段时间，我们再次回到本书，修补遗憾。本版比第 1 版和第 2 版内容更丰富：修订、补充并增加了新的内容和章节，使其更加完善。本版融合了第 2 版以来内核的各种变化。更值得一提的是，Linux 内核联盟<sup>①</sup>做出决定，近期内不进行 2.7 版内核的开发，于是，内核开发者打算继续开发并稳定 2.6 版。这个决定意味深长，而本书从中的最大受益就是在 2.6 版上可以稳定相当长时间。随着内核的成熟，内核“快照”才有机会能维持得更久远一些。本书可作为内核开发的规范文档，既认识内核的过去，也着眼于内核的未来。

## 使用这本书

开发 Linux 内核不需要天赋异秉，不需要有什么魔法，连 Unix 开发者普遍长着的络腮胡子都不一定要有。内核虽然有一些有趣并且独特的规则和要求，但是它和其他大型软件项目相比，并没有太大差别。像所有的大型软件开发一样，要学的东西确实不少，但是不同之处在于数量上的积累，而非本质上的区别。

认真阅读源码非常有必要，Linux 系统代码的开放性其实是弥足珍贵的，不要无动于衷地将它搁置一边，浪费了大好资源。实际上就是读了代码还远远不够呢，你应该钻研并尝试着动手改动一些代码。寻找一个 bug 然后去修改它，改进你的硬件设备的驱动程序。增加新功能，即使看

---

<sup>①</sup> 这一决定是在加拿大渥太华 2004 年夏季举办的 Linux 内核年度开发者大会上做出的。

起来微不足道，寻找痛痒之处并解决。只有动手写代码才能真正融会贯通。

## 内核版本

本书基于 Linux 2.6 内核系列。它并不涵盖早期的版本，当然也有一些例外。比如，我们会讨论 2.4 系列内核中的一些子系统是如何实现的，这是因为简单的实现有助于传授知识。特别说明的是，本书介绍的是最新的 Linux 2.6.34 内核版本。尽管内核总在不断更新，任何努力也难以捕获这样一只永不停息的猛兽，但是本书力图适合于新旧内核的开发者和用户。

虽然本书讨论的是 2.6.34 内核，但我也确保了它同样适用于 2.6.32 内核。后一个版本往往被各个 Linux 发行版本奉为“企业版”内核，所以我们可以各种产品线上见到其身影。该版本确实已经开发了数年（类似的“长线”版本还有 2.6.9、2.6.18 和 2.6.27 等）。

## 读者范围

本书是写给那些有志于理解 Linux 内核的软件开发者的。本书并不逐行逐字地注解内核源代码，也不是指导开发驱动程序或是内核 API 的参考手册（如果存在标准的内核 API 的话）。本书的初衷是提供足够多的关于 Linux 内核设计和实现的信息，希望读过本书的程序员能够拥有较为完备的知识，可以真正开始开发内核代码。无论开发内核是为了兴趣还是为了赚钱，我都希望能够带领读者快速走进 Linux 内核世界。本书不但介绍了理论而且也讨论了具体应用，可以满足不同读者的需要。全书紧紧围绕着理论联系实践，并非一味强调理论或是实践。无论你研究 Linux 内核的动机是什么，我都希望这本书都能将内核的设计和实现分析清楚，起到抛砖引玉的作用。

因此，本书覆盖了从核心内核系统的应用到内核设计与实现等各方面的内容。我认为这点很重要，值得花工夫讨论。例如，第 8 章讨论的是所谓的下半部机制。本章分别讨论了内核下半部机制的设计和实现（核心内核开发者或者学者会感兴趣），随即便介绍了如何使用内核提供的接口实现你自己的下半部（这对设备驱动开发者可能很有用处）。其实，我认为上述两部分内容是相得益彰的，虽然核心内核开发者主要关注的问题是内核内部如何工作，但是也应该清楚如何使用接口；同样，如果设备驱动开发者了解了接口背后的实现机制，自然也会受益匪浅。

这好比学习某些库的 API 函数与研究该库的具体实现。初看，好像应用程序开发者仅仅需要理解 API——我们被灌输的思想是，应该像看待黑盒子一样看待接口。另外，库的开发者也只关心库的设计与实现，但是我认为双方都应该花时间相互学习。能深刻了解操作系统本质的应用程序开发者无疑可以更好地利用它。同样，库开发者也决不应该脱离基于此库的应用程序，埋头开发。因此，我既讨论了内核子系统的设计，也讨论了它的用法，希望本书能对核心开发者和应用开发者都有用。

我假设读者已经掌握了 C 语言，而且对 Linux 比较熟悉。如果读者还具有与操作系统设计相关的经验和其他计算机科学的概念就更好了。当然，我也会尽可能多地解释这些概念，但如果你仍然不能理解这些知识的话，请看本书最后参考资料中给出的一些关于操作系统设计方面的经典书籍。

本书很适合在大学中作为介绍操作系统的辅助教材，与介绍操作系统理论的书相搭配。对于



大学高年级课程或者研究生课程来说，可直接使用本书作为教材。

### 第 3 版致谢

与其他作者一样，我决非是一个人躲在山洞里孤苦地写出这本书来的（那也是一件美差，因为有熊相伴）。我能最终完成本书原稿是与无数建议和关怀分不开的。仅仅一页纸无法容纳我的感激，但我还是要衷心地感谢所有给予我鼓励、给予我知识和给予我灵感的朋友和同事。

首先我要对为此书付出辛勤劳作的所有编辑表示感谢，尤其要感谢我的组稿编辑 Mark Taber，他为这一版的出版从头到尾倾注了许多心血。还要特别感谢业务开发编辑 Michael Thurston 和项目组织编辑 Tonya Simps。o。

本版的技术编辑 Robert P. J. Day 也是我需要倍加感谢的人，他独到的洞察力和准确的校对使书稿质量大大提高。尽管他的工作称得上完美，但如果书中仍留有错误，责任由我承担。需要十分感谢的还有 Adam Belay、Zack Brown、Martin Pool 以及 Chris Rivera，他们对第 1 版和第 2 版所做的一切努力我依然记忆犹新。

许多内核开发者为我提供了大力支持，回答了许多问题，还有那些撰写代码的人——本书正是由于有了这些代码，才有了存在的意义。他们是：Andrea Arcangeli、Alan Cox、Greg Kroah-Hartman、Dave Miller、Patrick Mochel、Andrew Morton、Nick Piggin 和 Linus Torvalds。

我要大力感谢我的同事们。他们的创造力和智慧无与伦比，能与他们一起工作其乐无穷。因为篇幅原因，请谅解我不能列出所有人的名字。但不得不提的是：Alan Blount、Jay Crim、Chris Danis、Chris DiBona、Eric Flatt、Mike Lockwood、San Mehat、Brian Rogan、Brian Swetland、Jon Trowbridge 和 Steve Vinter，感谢他们给予我的支持、知识和友谊！

还有许多值得尊敬和热爱的人，他们是：Paul Amici、Mikey Babbitt、Keith Barbag、Jacob Berkman、Nat Friedman、Dustin Hall、Joyce Hawkins、Miguel de Icaza、Jimmy Krehl、Doris Love、Linda Love、Brette Luck、Randy O'Dowd、Sal Ribaud。o 和他了不起的妈妈 Chris Rivera、Carolyn Rodon、Joey Shaw、Sarah Stewart、Jeremy VanDoren 和他的家人、Luis Villa、Steve Weisberg 和他的家人以及 Helen Whisnant 等。

最后，非常感谢我的父母。

内核黑客万岁！

Robert Love  
Boston

# 作者简介

Robert Love 是一位资深的开源软件开发者、讲师和作者，他使用 Linux 和贡献于 Linux 已超过 15 年。目前他是 Google 公司的资深软件工程师，是 Android 移动平台内核开发团队的成员；在去 Google 工作之前，他就职于 Novell 公司，任职 Linux 桌面系统的首席架构师；在去 Novell 之前，他是 MontaVista 和 Ximain 公司的内核开发工程师。

Robert 参与的内核项目包括抢占式内核、进程调度器、内核事件层、通知机制、VM 改进，以及一些设备驱动。

Robert 曾经发表过许多关于 Linux 内核的演讲和文章；他还是《Linux journal》电子杂志的编辑。另外，除了本书，他的著作还包含《Linux System Programming》和《Linux in a Nutshell》。

Robert 在佛罗里达大学获得数学学士学位和计算机理学学士学位。

# 目 录

译者序	2.4.4 不要轻易在内核中使用浮点数	18
序 言	2.4.5 容积小而固定的栈	18
前 言	2.4.6 同步和并发	18
作者简介	2.4.7 可移植性的重要性	19
	2.5 小结	19
第 1 章 Linux 内核简介	第 3 章 进程管理	20
1.1 Unix 的历史	3.1 进程	20
1.2 追寻 Linus 足迹：Linux 简介	3.2 进程描述符及任务结构	21
1.3 操作系统和内核简介	3.2.1 分配进程描述符	22
1.4 Linux 内核和传统 Unix 内核的比较	3.2.2 进程描述符的存放	23
1.5 Linux 内核版本	3.2.3 进程状态	23
1.6 Linux 内核开发者社区	3.2.4 设置当前进程状态	25
1.7 小结	3.2.5 进程上下文	25
第 2 章 从内核出发	3.2.6 进程家族树	25
2.1 获取内核源码	3.3 进程创建	26
2.1.1 使用 Git	3.3.1 写时拷贝	27
2.1.1 安装内核源代码	3.3.2 fork()	27
2.1.3 使用补丁	3.3.3 vfork()	28
2.2 内核源码树	3.4 线程在 Linux 中的实现	28
2.3 编译内核	3.4.1 创建线程	29
2.3.1 配置内核	3.4.2 内核线程	30
2.3.2 减少编译的垃圾信息	3.5 进程终结	31
2.3.3 衍生多个编译作业	3.5.1 删除进程描述符	32
2.3.4 安装新内核	3.5.2 孤儿进程造成的进退维谷	32
2.4 内核开发的特点	3.6 小结	34
2.4.1 无 libc 库抑或无标准头文件	第 4 章 进程调度	35
2.4.2 GNU C	4.1 多任务	35
2.4.3 没有内存保护机制	4.2 Linux 的进程调度	36

4.3 策略 .....	36	5.4.2 参数传递 .....	60
4.3.1 I/O 消耗型和处理器消耗型的 进程 .....	36	5.5 系统调用的实现 .....	61
4.3.2 进程优先级 .....	37	5.5.1 实现系统调用 .....	61
4.3.3 时间片 .....	38	5.5.2 参数验证 .....	62
4.3.4 调度策略的活动 .....	38	5.6 系统调用上下文 .....	64
4.4 Linux 调度算法 .....	39	5.6.1 绑定一个系统调用的最后步骤 .....	65
4.4.1 调度器类 .....	39	5.6.2 从用户空间访问系统调用 .....	67
4.4.2 Unix 系统中的进程调度 .....	40	5.6.3 为什么不通过系统调用的 方式实现 .....	68
4.4.3 公平调度 .....	41	5.7 小结 .....	68
4.5 Linux 调度的实现 .....	42	第 6 章 内核数据结构 .....	69
4.5.1 时间记账 .....	42	6.1 链表 .....	69
4.5.2 进程选择 .....	44	6.1.1 单向链表和双向链表 .....	69
4.5.3 调度器入口 .....	48	6.1.2 环形链表 .....	70
4.5.4 睡眠和唤醒 .....	49	6.1.3 沿链表移动 .....	71
4.6 抢占和上下文切换 .....	51	6.1.4 Linux 内核中的实现 .....	71
4.6.1 用户抢占 .....	53	6.1.5 操作链表 .....	73
4.6.2 内核抢占 .....	53	6.1.6 遍历链表 .....	75
4.7 实时调度策略 .....	54	6.2 队列 .....	78
4.8 与调度相关的系统调用 .....	54	6.2.1 kfifo .....	79
4.8.1 与调度策略和优先级相关的 系统调用 .....	55	6.2.2 创建队列 .....	79
4.8.2 与处理器绑定有关的系统调用 .....	55	6.2.3 推入队列数据 .....	79
4.8.3 放弃处理器时间 .....	56	6.2.4 摘取队列数据 .....	80
4.9 小结 .....	56	6.2.5 获取队列长度 .....	80
第 5 章 系统调用 .....	57	6.2.6 重置和撤销队列 .....	80
5.1 与内核通信 .....	57	6.2.7 队列使用举例 .....	81
5.2 API、POSIX 和 C 库 .....	57	6.3 映射 .....	81
5.3 系统调用 .....	58	6.3.1 初始化一个 idr .....	82
5.3.1 系统调用号 .....	59	6.3.2 分配一个新的 UID .....	82
5.3.2 系统调用的性能 .....	59	6.3.3 查找 UID .....	83
5.4 系统调用处理程序 .....	60	6.3.4 删除 UID .....	84
5.4.1 指定恰当的系统调用 .....	60	6.3.5 撤销 idr .....	84
		6.4 二叉树 .....	84

6.4.1	二叉搜索树	84	8.2	软中断	110
6.4.2	自平衡二叉搜索树	85	8.2.1	软中断的实现	111
6.5	数据结构以及选择	87	8.2.2	使用软中断	113
6.6	算法复杂度	88	8.3	tasklet	114
6.6.1	算法	88	8.3.1	tasklet 的实现	114
6.6.2	大 o 符号	88	8.3.2	使用 tasklet	116
6.6.3	大 $\theta$ 符号	89	8.3.3	老的 BH 机制	119
6.6.4	时间复杂度	89	8.4	工作队列	120
6.7	小结	90	8.4.1	工作队列的实现	121
第 7 章	中断和中断处理	91	8.4.2	使用工作队列	124
7.1	中断	91	8.4.3	老的任务队列机制	126
7.2	中断处理程序	92	8.5	下半部机制的选择	127
7.3	上半部与下半部的对比	93	8.6	在下半部之间加锁	128
7.4	注册中断处理程序	93	8.7	禁止下半部	128
7.4.1	中断处理程序标志	94	8.8	小结	129
7.4.2	一个中断例子	95	第 9 章	内核同步介绍	131
7.4.3	释放中断处理程序	95	9.1	临界区和竞争条件	131
7.5	编写中断处理程序	96	9.1.1	为什么我们需要保护	132
7.5.1	共享的中断处理程序	97	9.1.2	单个变量	133
7.5.2	中断处理程序实例	97	9.2	加锁	134
7.6	中断上下文	99	9.2.1	造成并发执行的原因	135
7.7	中断处理机制的实现	100	9.2.2	了解要保护些什么	136
7.8	/proc/interrupts	102	9.3	死锁	137
7.9	中断控制	103	9.4	争用和扩展性	138
7.9.1	禁止和激活中断	103	9.5	小结	140
7.9.2	禁止指定中断线	105	第 10 章	内核同步方法	141
7.9.3	中断系统的状态	105	10.1	原子操作	141
7.10	小结	106	10.1.1	原子整数操作	142
第 8 章	下半部和推后执行的 工作	107	10.1.2	64 位原子操作	144
8.1	下半部	107	10.1.3	原子位操作	145
8.1.1	为什么要用下半部	108	10.2	自旋锁	147
8.1.2	下半部的环境	108	10.2.1	自旋锁方法	148
			10.2.2	其他针对自旋锁的操作	149

10.2.3 自旋锁和下半部 .....	150	11.7.2 定时器竞争条件 .....	180
10.3 读-写自旋锁 .....	150	11.7.3 实现定时器 .....	180
10.4 信号量 .....	152	11.8 延迟执行 .....	181
10.4.1 计数信号量和二值信号量 .....	153	11.8.1 忙等待 .....	181
10.4.2 创建和初始化信号量 .....	154	11.8.2 短延迟 .....	182
10.4.3 使用信号量 .....	154	11.8.3 schedule_timeout() .....	183
10.5 读-写信号量 .....	155	11.9 小结 .....	185
10.6 互斥体 .....	156	第 12 章 内存管理 .....	186
10.6.1 信号量和互斥体 .....	158	12.1 页 .....	186
10.6.2 自旋锁和互斥体 .....	158	12.2 区 .....	187
10.7 完成变量 .....	158	12.3 获得页 .....	189
10.8 BLK: 大内核锁 .....	159	12.3.1 获得填充为 0 的页 .....	190
10.9 顺序锁 .....	160	12.3.2 释放页 .....	191
10.10 禁止抢占 .....	161	12.4 kmalloc() .....	191
10.11 顺序和屏障 .....	162	12.4.1 gfp_mask 标志 .....	192
10.12 小结 .....	165	12.4.2 kfree() .....	195
第 11 章 定时器和时间管理 .....	166	12.5 vmalloc() .....	196
11.1 内核中的时间概念 .....	166	12.6 slab 层 .....	197
11.2 节拍率: HZ .....	167	12.6.1 slab 层的设计 .....	198
11.2.1 理想的 HZ 值 .....	168	12.6.2 slab 分配器的接口 .....	200
11.2.2 高 HZ 的优势 .....	169	12.7 在栈上的静态分配 .....	203
11.2.3 高 HZ 的劣势 .....	169	12.7.1 单页内核栈 .....	203
11.3 jiffies .....	170	12.7.2 在栈上光明正大地工作 .....	203
11.3.1 jiffies 的内部表示 .....	171	12.8 高端内存的映射 .....	204
11.3.2 jiffies 的回绕 .....	172	12.8.1 永久映射 .....	204
11.3.3 用户空间和 HZ .....	173	12.8.2 临时映射 .....	204
11.4 硬时钟和定时器 .....	174	12.9 每个 CPU 的分配 .....	205
11.4.1 实时时钟 .....	174	12.10 新的每个 CPU 接口 .....	206
11.4.2 系统定时器 .....	174	12.10.1 编译时的每个 CPU 数据 .....	206
11.5 时钟中断处理程序 .....	174	12.10.2 运行时的每个 CPU 数据 .....	207
11.6 实际时间 .....	176	12.11 使用每个 CPU 数据的原因 .....	208
11.7 定时器 .....	178	12.12 分配函数的选择 .....	209
11.7.1 使用定时器 .....	178	12.13 小结 .....	209

第 13 章 虚拟文件系统 .....	210	14.5.7 I/O 调度程序的选择 .....	245
13.1 通用文件系统接口 .....	210	14.6 小结 .....	246
13.2 文件系统抽象层 .....	211	第 15 章 进程地址空间 .....	247
13.3 Unix 文件系统 .....	212	15.1 地址空间 .....	247
13.4 VFS 对象及其数据结构 .....	213	15.2 内存描述符 .....	248
13.5 超级块对象 .....	214	15.2.1 分配内存描述符 .....	249
13.6 超级块操作 .....	215	15.2.2 撤销内存描述符 .....	250
13.7 索引节点对象 .....	217	15.2.3 mm_struct 与内核线程 .....	250
13.8 索引节点操作 .....	219	15.3 虚拟内存区域 .....	251
13.9 目录项对象 .....	222	15.3.1 VMA 标志 .....	251
13.9.1 目录项状态 .....	222	15.3.2 VMA 操作 .....	253
13.9.2 目录项缓存 .....	223	15.3.3 内存区域的树型结构和内存 区域的链表结构 .....	254
13.10 目录项操作 .....	224	15.3.4 实际使用中的内存区域 .....	254
13.11 文件对象 .....	225	15.4 操作内存区域 .....	255
13.12 文件操作 .....	226	15.4.1 find_vma() .....	256
13.13 和文件系统相关的数据结构 .....	230	15.4.2 find_vma_prev() .....	257
13.14 和进程相关的数据结构 .....	232	15.4.3 find_vma_intersection() .....	257
13.15 小结 .....	233	15.5 mmap() 和 do_mmap(): 创建地址 区间 .....	258
第 14 章 块 I/O 层 .....	234	15.6 mmap() 和 do_mmap(): 删除 地址区间 .....	259
14.1 剖析一个块设备 .....	234	15.7 页表 .....	260
14.2 缓冲区和缓冲区头 .....	235	15.8 小结 .....	261
14.3 bio 结构体 .....	237	第 16 章 页高速缓存和页回写 .....	262
14.3.1 I/O 向量 .....	238	16.1 缓存手段 .....	262
14.3.2 新老方法对比 .....	239	16.1.1 写缓存 .....	262
14.4 请求队列 .....	240	16.1.2 缓存回收 .....	263
14.5 I/O 调度程序 .....	240	16.2 Linux 页高速缓存 .....	264
14.5.1 I/O 调度程序的工作 .....	241	16.2.1 address_space 对象 .....	264
14.5.2 Linus 电梯 .....	241	16.2.2 address_space 操作 .....	266
14.5.3 最终期限 I/O 调度程序 .....	242	16.2.3 基树 .....	267
14.5.4 预测 I/O 调度程序 .....	244	16.2.4 以前的页散列表 .....	268
14.5.5 完全公正的排队 I/O 调度 程序 .....	244		
14.5.6 空操作的 I/O 调度程序 .....	245		

16.3	缓冲区高速缓存 .....	268	18.2	内核中的 bug .....	296
16.4	flusher 线程 .....	268	18.3	通过打印来调试 .....	296
16.4.1	膝上型计算机模式 .....	270	18.3.1	健壮性 .....	296
16.4.2	历史上的 bdflush、kupdated 和 pdflush .....	270	18.3.2	日志等级 .....	297
16.4.3	避免拥塞的方法：使用多线程 .....	271	18.3.3	记录缓冲区 .....	298
16.5	小结 .....	271	18.3.4	syslogd 和 klogd .....	298
第 17 章	设备与模块 .....	273	18.3.5	从 printf() 到 printk() 的转换 .....	298
17.1	设备类型 .....	273	18.4	oops .....	298
17.2	模块 .....	274	18.4.1	ksymoops .....	300
17.2.1	Hello, World .....	274	18.4.2	kallsyms .....	300
17.2.2	构建模块 .....	275	18.5	内核调试配置选项 .....	301
17.2.3	安装模块 .....	277	18.6	引发 bug 并打印信息 .....	301
17.2.4	产生模块依赖性 .....	277	18.7	神奇的系统请求键 .....	302
17.2.5	载入模块 .....	278	18.8	内核调试器的传奇 .....	303
17.2.6	管理配置选项 .....	279	18.8.1	gdb .....	303
17.2.7	模块参数 .....	280	18.8.2	kgdb .....	304
17.2.8	导出符号表 .....	282	18.9	探测系统 .....	304
17.3	设备模型 .....	283	18.9.1	用 UID 作为选择条件 .....	304
17.3.1	kobject .....	283	18.9.2	使用条件变量 .....	305
17.3.2	ktype .....	284	18.9.3	使用统计量 .....	305
17.3.3	kset .....	285	18.9.4	重复频率限制 .....	305
17.3.4	kobject、ktype 和 kset 的 相互关系 .....	285	18.10	用二分查找法找出引发罪恶的 变更 .....	306
17.3.5	管理和操作 kobject .....	286	18.11	使用 Git 进行二分搜索 .....	307
17.3.6	引用计数 .....	287	18.12	当所有的努力都失败时：社区 .....	308
17.4	sysfs .....	288	18.13	小结 .....	308
17.4.1	sysfs 中添加和删除 kobject .....	290	第 19 章	可移植性 .....	309
17.4.2	向 sysfs 中添加文件 .....	291	19.1	可移植操作系统 .....	309
17.4.3	内核事件层 .....	293	19.2	Linux 移植史 .....	310
17.5	小结 .....	294	19.3	字长和数据类型 .....	311
第 18 章	调试 .....	295	19.3.1	不透明类型 .....	313
18.1	准备开始 .....	295	19.3.2	指定数据类型 .....	314
			19.3.3	长度明确的类型 .....	314



19.3.4 char 型的符号问题 .....	315	20.2.4 花括号 .....	325
19.4 数据对齐 .....	315	20.2.5 每行代码的长度 .....	326
19.4.1 避免对齐引发的问题 .....	316	20.2.6 命名规范 .....	326
19.4.2 非标准类型的对齐 .....	316	20.2.7 函数 .....	326
19.4.3 结构体填补 .....	316	20.2.8 注释 .....	326
19.5 字节顺序 .....	318	20.2.9 typedef .....	327
19.6 时间 .....	319	20.2.10 多用现成的东西 .....	328
19.7 页长度 .....	320	20.2.11 在源码中减少使用 ifdef .....	328
19.8 处理器排序 .....	320	20.2.12 结构初始化 .....	328
19.9 SMP、内核抢占、高端 内存 .....	321	20.2.13 代码的事后修正 .....	329
19.10 小结 .....	321	20.3 管理系统 .....	329
第 20 章 补丁、开发和社区 .....	322	20.4 提交错误报告 .....	329
20.1 社区 .....	322	20.5 补丁 .....	330
20.2 Linux 编码风格 .....	322	20.5.1 创建补丁 .....	330
20.2.1 缩进 .....	323	20.5.2 用 Git 创建补丁 .....	331
20.2.2 switch 语句 .....	323	20.5.3 提交补丁 .....	331
20.2.3 空格 .....	324	20.6 小结 .....	332
		参考资料 .....	333