

HZ BOOKS
华章教育

计 算 机 科 学 丛 书

CRC Press
Taylor & Francis Group

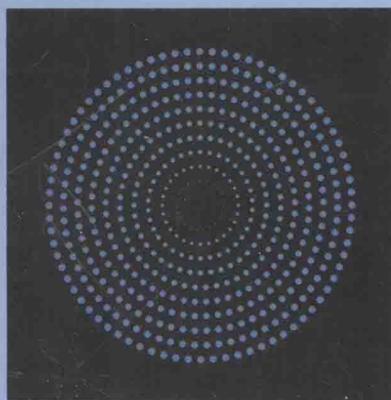
操作系统设计

Xinu方法

(美) Douglas Comer 著 邹恒明 周亮 曹浩 刘正邦 等译
普度大学 上海交通大学

Operating System Design
The Xinu Approach, Linksys Version

Operating System Design
The Xinu Approach
Linksys Version



Douglas Comer



机械工业出版社
China Machine Press

计 算 机 科 学 丛 书

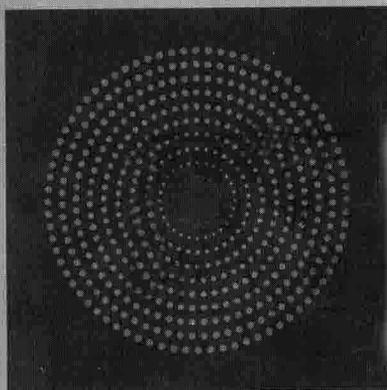
操作系统设计

Xinu方法

(美) Douglas Comer 著 邹恒明 周亮 曹浩 刘正邦 等译
普度大学 上海交通大学

Operating System Design
The Xinu Approach, Linksys Version

Operating System Design
The Xinu Approach
Linksys Version



Douglas Comer



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

操作系统设计: Xinu 方法 / (美) 科默 (Comer, D.) 著; 邹恒明等译. —北京: 机械工业出版社, 2013. 10
(计算机科学丛书)

书名原文: Operating System Design: The Xinu Approach, Linksys Version

ISBN 978-7-111-42826-8

I. 操… II. ①科… ②邹… III. 操作系统—程序设计 IV. TP316

中国版本图书馆 CIP 数据核字 (2013) 第 123969 号

版权所有·侵权必究

封底无防伪标均为盗版

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: 01-2012-4866

Operating System Design: The Xinu Approach (Linksys Version) by Douglas Comer (ISBN: 978-1-4398-8109-5).
Copyright © 2012 by Taylor & Francis Group LLC.

Authorized translation from the English language edition published by CRC Press, part of Taylor & Francis Group LLC. All rights reserved.

China Machine Press is authorized to publish and distribute exclusively the Chinese (Simplified Characters) language edition. This edition is authorized for sale in the People's Republic of China only (excluding Hong Kong, Macao SAR and Taiwan). No part of this publication may be reproduced or distributed in any form or by any means, or stored in a database or retrieval system, without the prior written permission of the publisher.

Copies of this book sold without a Taylor & Francis sticker on the cover are unauthorized and illegal.

本书原版由 Taylor & Francis 出版集团旗下 CRC 出版公司出版, 并经授权翻译出版。版权所有, 侵权必究。

本书中文简体字翻译版授权由机械工业出版社独家出版, 限在中国大陆地区销售。未经出版者书面许可, 不得以任何方式复制或抄袭本书的任何内容。

本书封面贴有 Taylor & Francis 防伪标, 无标者不得销售。

本书对操作系统的内存管理、进程管理、进程协调和同步、进程间通信、实时时钟管理、设备无关的 I/O、设备驱动、网络协议、文件系统等进行了详细的介绍, 并利用分层的设计范式, 以一种有序、易于理解的方式来阐述这些内容。本书以 Xinu 操作系统为系统设计的样板和模式, 从一个裸机开始, 一步一步地设计和实现一个小型但优雅的操作系统。本书的样本代码可以运行在 Linksys E2100L 无线路由器上。

本书适用于高年级的本科生或低年级的研究生, 也适用于那些想了解操作系统的计算机从业人员。学习本书前, 学生需要具备基本的程序设计能力, 应当理解基本的数据结构, 包括链表、栈和队列, 并且应当用 C 语言写过程序。

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: 盛思源

藁城市京瑞印刷有限公司印刷

2014 年 1 月第 1 版第 1 次印刷

185mm × 260mm · 23.5 印张

标准书号: ISBN 978-7-111-42826-8

定 价: 79.00 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

投稿热线: (010) 88379604

购书热线: (010) 68326294 88379649 68995259

读者信箱: hzjsj@hzbook.com

文艺复兴以降，源远流长的科学精神和逐步形成的学术规范，使西方国家在自然科学的各个领域中取得了垄断性的优势；也正是这样的传统，使美国在信息技术发展的六十多年间名家辈出、独领风骚。在商业化的进程中，美国的产业界与教育界越来越紧密地结合，计算机学科中的许多泰山北斗同时身处科研和教学的最前线，由此而产生的经典科学著作，不仅擘划了研究的范畴，还揭示了学术的源变，既遵循学术规范，又自有学者个性，其价值并不会因年月的流逝而减退。

近年，在全球信息化大潮的推动下，我国的计算机产业发展迅猛，对专业人才的需求日益迫切。这对计算机教育界和出版界都既是机遇，也是挑战；而专业教材的建设在教育战略上显得举足轻重。在我国信息技术发展时间较短的现状下，美国等发达国家在其计算机科学发展的几十年间积淀和发展的经典教材仍有许多值得借鉴之处。因此，引进一批国外优秀计算机教材将对我国计算机教育事业的发展起到积极的推动作用，也是与世界接轨、建设真正的世界一流大学的必由之路。

机械工业出版社华章公司较早意识到“出版要为教育服务”。自1998年开始，我们就将工作重点放在了遴选、移译国外优秀教材上。经过多年的不懈努力，我们与 Pearson, McGraw-Hill, Elsevier, MIT, John Wiley & Sons, Cengage 等世界著名出版公司建立了良好的合作关系，从他们现有的数百种教材中甄选出 Andrew S. Tanenbaum, Bjarne Stroustrup, Brian W. Kernighan, Dennis Ritchie, Jim Gray, Alfred V. Aho, John E. Hopcroft, Jeffrey D. Ullman, Abraham Silberschatz, William Stallings, Donald E. Knuth, John L. Hennessy, Larry L. Peterson 等大师名家的一批经典作品，以“计算机科学丛书”为总称出版，供读者学习、研究及珍藏。大理石纹理的封面，也正体现了这套丛书的品位和格调。

“计算机科学丛书”的出版工作得到了国内外学者的鼎力襄助，国内的专家不仅提供了中肯的选题指导，还不辞劳苦地担任了翻译和审校的工作；而原书的作者也相当关注其作品在中国的传播，有的还专程为其书的中译本作序。迄今，“计算机科学丛书”已经出版了近两百个品种，这些书籍在读者中树立了良好的口碑，并被许多高校采用为正式教材和参考书籍。其影印版“经典原版书库”作为姊妹篇也被越来越多实施双语教学的学校所采用。

权威的作者、经典的教材、一流的译者、严格的审校、精细的编辑，这些因素使我们的图书有了质量的保证。随着计算机科学与技术专业学科建设的不断完善和教材改革的逐渐深化，教育界对国外计算机教材的需求和应用都将步入一个新的阶段，我们的目标是尽善尽美，而反馈的意见正是我们达到这一终极目标的重要帮助。华章公司欢迎老师和读者对我们的工作提出建议或给予指正，我们的联系方式如下：

华章网站：www.hzbook.com

电子邮件：hzsj@hzbook.com

联系电话：(010) 88379604

联系地址：北京市西城区百万庄南街1号

邮政编码：100037



华章教育

华章科技图书出版中心

——像小孩那样学习

我实在告诉你们，你们若不回转，变成小孩的样式，断不得进天国。

——马太福音 18: 3

这是耶稣教训门徒所说的话。虽然讲的是天国的事情，但不经意间也指出了一条学习的路径。众所周知，小孩接受新鲜事物的能力非常强，小孩能够随着时间的推移而不知不觉地学会说话和日常的生活规则，而成人学习一门新语言或一种新技能则通常较为困难。

小孩对语言或技能的掌握是在成人所说所行的点滴中先行模仿然后再悟出的，也就是通过模仿从琐细中理出了头绪。小孩的这种学习方法是一种典型的先实际后理解、先细节后全景、从只见树木到顿见森林的转变。这与成人所惯常采用的学习方式有很大不同。成人采用的学习方法通常是先理论后实践、先全貌后细节、先森林后树木式的由高至低逐步推进的方式。这在大学的“操作系统”课程学习中体现得十分明显。一般大学的普遍做法是先介绍操作系统的全貌，然后介绍操作系统的个体组成部分，再附加一些动手实验，所谓通过实验加深对理论的理解（其他课程的学习方式也大同小异）。此种模式是一种典型的自顶向下的演绎模式，由于其由来已久，人们已经习以为常，并没有想过有什么不妥。但操作系统及许多大学课程教学中所呈现的事实是，很多人对其感到枯燥难懂，难以掌握。

既然如此，为什么不试试另一种学习模式，即先动手后理解、先细节后抽象、由树木到森林，像小孩一样来学习操作系统或其他专业课程呢？对于操作系统来说，完全可以采用由底至上的方法，先不讲操作系统理论和模型，而是先来实际设计一个计算机的内存管理系统，在动手实践中逐渐抽象出规律，然后上升为理论和模型。在讨论完内存管理之后再阐述如何调度程序，这样逐次推进，最后搭起整个操作系统。这就像把一大堆积木呈现给读者，先搭起一个长方体，然后搭起几个轮子，再搭上一些窗户，最终，一辆汽车就形成了。

本书所采取的正是由底至上、由动手到理论、由具体到抽象、由细节到整体、由树木到森林的学习模式来阐述操作系统的原理和设计。该书以细节为核心，以设计与实现为手段，通过设计和实现操作系统的各种功能来引入操作系统的原理和模型，并以设计人员在构建一个操作系统时所遵循的工作次序来组织全书。该书从一个裸机开始，一步一步地设计和实现一个小型但优雅的操作系统——Xinu。本书每一章描述设计 Xinu 系统架构里的一个组件，并提供示例软件来演示该层架构所提供的功能。这种与传统模式的截然不同让人一开始有些反感，但在体会了此种方式的美妙后又恍然大悟。

该书的最大特点是实践性强，以与实际情况非常贴切的场景为背景，以动手操作为推手，以实际代码为讨论对象，将操作系统的各种实战原理和模式娓娓道来，易于理解和消化。读者只要顺着书的讲解，并配合运行和分析其附带的代码，即可掌握操作系统的设计和实现。该书在美国普度大学所受到的青睐从一个侧面佐证了这一点。

本书适用于高年级的本科生或低年级的研究生，也适用于那些想了解操作系统的计算机从业人员。鉴于该书的篇幅和所涉及的具体动手实现，建议分配 4 个学分或以上来组织本门课程。另外，考虑到并不是所有人都喜爱先细节后抽象的模式，建议与一本优良的阐述操作系统原理的书结合进行讲授，效果或将更加令人满意。

本书由上海交通大学 2012 年春季高级操作系统课程组集体翻译，因人员众多，就不一一列出。在此特向对本书翻译提供帮助的人员表示感谢。

鉴于译者水平所限，书中纰漏甚至错误在所难免，谨望读者不吝指正。

邹恒明

2013 年 5 月于美国达拉斯

建造计算机操作系统有点像编织锦缎。这两种工作的最终成品都是一个和谐一致、大型、复杂的人造系统。在每一种情况下，最后的人造成品都是由细微但却精巧的步骤所构造。在编织锦缎时，细节是至关重要的，因为一点点不协调的瑕疵都很容易观察到。就像锦缎里的缎面一样，加入到操作系统里的每个新组件都需要与整体的设计相协调。从这个角度看，将不同片段组装起来的机械加工只是整个建造过程中的一小部分，一个大师级的产品必须以某个模式为蓝本，所有参与系统设计的工作人员都必须遵守这种模式。

有讽刺意味的是，现有的操作系统教材或课程很少对底层的模式和原理进行解释，而这些模式和原理正是操作系统构造的基础。在学生看来：操作系统似乎是一个暗箱，而现有的教材则加强了这种误解，因为这些教材所解释的不过是操作系统的功能，其关注的也只是操作系统各种能力的使用。更为重要的是，学生在学习操作系统时采取的是从操作系统外面来察看的方式，从而常常导致这样一种感觉：认为操作系统由一组抽象的界面所组成，这些界面下的功能由一大堆晦涩神秘的代码连接在一起，而这些神秘的代码本身还包含着许多与机器硬件直接相关的、无规律可寻的奇技巧术。

令人惊奇的是，学生一旦从大学毕业，就马上觉得与操作系统有关的工作已经结束，自己不再需要理解或学习操作系统，因为由商业公司和开源社区所构造的现有操作系统足以应付各种需要，没有自己什么事情了。但没有什么比这种想法离真理更远了。有讽刺意味的是，尽管为个人计算机设计传统操作系统的公司数量比以前更少了，但社会和行业对操作系统技能的需求却在增长，许多公司雇佣大学生来从事操作系统方面的工作。社会上这些对操作系统技能的需求来源于更便宜的微处理器，这些便宜的微处理器嵌入在智能手机、视频游戏、iPod、Internet 路由器、线缆和机顶盒以及打印机中。

在与嵌入式系统打交道时，有关原理和结构的知识非常关键，因为程序员可能需要在现有的操作系统内部构造某种或某个新的机制，或者对现有操作系统进行修改以便可以在新的硬件平台上运行。而且，为嵌入式设备编写应用程序时需要理解下层的操作系统。如果不理解操作系统设计的各种细微之处，则不可能充分利用这些小型嵌入式处理器的能力。

本书的目的是揭开操作系统设计中的神秘感，将方方面面的材料整合为一个系统化的整体。本书对操作系统的主要系统组件进行了详细阐述，并以一种层次架构的设计范式来组织这些组件，从而以一种有序、可理解的方式来展开这些内容。与其他评述性书籍不同的是，本书并不尽可能多地提供不同方案，呈现给读者的将是一个基于传统过程的、使用实际的、直截了当的原语所构造的操作系统。本书从一个裸机开始，一步一步地设计和实现一个小型但优雅的操作系统的。这个称为 Xinu 的操作系统将成为系统设计的样板和模式。

虽然 Xinu 操作系统的规模较小，可以完全容纳在本书中，但该系统却包括了构成一个普通操作系统的全部组件：内存管理、进程管理、进程协调和同步、进程间通信、实时时钟管理、设备独立的输入输出、设备驱动、网络协议和一个文件系统。本书将这些组件组织成一个层次架构，使它们之间的相互连接清晰可见、设计过程浅显易懂。尽管规模小，但 Xinu 却拥有大型系统的能力。此外，Xinu 并不是一个玩具系统，它在许多商业产品中得到了应用。使用该系统的厂商包括 Mitsubishi、Lexmark、HP、IBM、Woodward (woodward.com)、Barnard Software 和 Mantissa 公司。学生通过本书可以学到的重要一课是：不管是小型嵌入式系统还是大型系统，好的系统设计都一样重要，一个系统的大部分能力都来自于良好的抽象。

本书所覆盖的所有议题都以一种特定的次序排列，这种次序就是设计人员在构建操作系统时所遵守的工作次序。本书每一章描述设计架构里的一个组件，并提供示例软件来演示由该层架构所提供的功能。使用这种方式具有如下几种优点：第一，每一章所解释的操作系统的功能子集均比上一章所讨论的功能子集更大，这种安排使我们在考虑一层特定架构的设计和实现时不用关心后续层面的实现。

第二，每一章的细节描述在第一次阅读时可以跳过去，读者只需要理解该层所提供的服务即可，而不是这些服务是如何实现的。第三，如果按次序阅读本书，读者可以先理解一个功能，然后在后面看到该功能是如何被后续部分所使用的。第四，有智力挑战的议题（如对并发的支持）出现在书的较前面，高层次的操作系统服务则出现在后面。在本书中，读者将看到大部分核心的功能仅仅只用几行代码就可以完成，这样我们就可以将大部分的代码（网络和文件系统）放到书的较后面，在读者已经做出了充分的思想准备后再进行讲解。

如前所述，与其他关于操作系统的许多书不一样的是，本书并不试图对每个系统组件的每种实现方案进行评估，也不对现有的商业系统进行综述。而是选择对一组使用最广泛的操作系统原语的实现细节进行阐述。例如，在讨论进程协调的一章，我们解释的是信号量（使用最广泛的进程协调原语）原语，而对其他原语（如监视器）的讨论则放到练习里。我们的目的是展示如何将原语在传统的硬件上实现，消除神秘。学生一旦理解了一组特定原语的魔力，其他原语的实现也就容易掌握了。

本书的示例代码可以运行在 Linksys E2100L 无线路由器上，该无线路由器在零售商店里就可以买到。只不过，我们并不是将 Linksys 硬件作为一个无线路由器来使用。我们的做法是，打开 Linksys 设备，将一根串行线连接到其控制端口，使用该串行线来中断 Linksys 正常的启动过程，并通过输入命令来迫使 Linksys 硬件下载和运行一个 Xinu 操作系统副本。也就是说，我们基本上忽略供应商所提供的软件，而是对其底层的硬件进行控制来运行 Xinu。

本书适用于高年级的本科生或者研究生，也适用于那些想了解操作系统的计算机从业人员。在本书所提供的全部材料里，虽然没有任何议题的难度达到不能理解的程度，但学习本书的全部内容可能需要超过一学期的时间。本科生里很少有学生能够熟练地阅读串行程序，而理解运行时环境的细节或机器架构的学生就更少了。因此，必须对学生进行仔细引导，以便使其可以掌握进程管理和进程同步的知识。如果时间有限，我推荐覆盖的内容包括第 1 章 ~ 第 7 章（进程管理）、第 9 章（基本的内存管理）、第 12 章（中断处理）、第 13 章（时钟管理）、第 14 章（设备无关的 I/O）和第 19 章（文件系统）。此外，对于一个完整学期的本科生课程来说，讨论第 20 章的远程文件系统等基本的远程访问议题也很重要。对于研究生课程来说，学生应当完整地阅读整本书，课堂讨论则应该专注于一些微妙的细节、各种折中和不同实现方案的比较。不管是本科生课程还是研究生课程，都应该包括的两个议题是：1) 在初始化阶段，当一个运行中的程序转化为一个进程时所发生的各种改变；2) 当输入行里的字符序列作为一个字符串变量传递给命令进程时，在操作系统壳里所发生的转化。

在所有情况下，如果学生能够在实验室中对系统进行动手实验，则学习的效果将大幅提高。理想的状态下，学生可以在课程的最初几天或几个星期开始使用这个系统，然后再试图理解系统的内部结构。本书第 1 章提供了几个例子和一些能够引起学生兴趣的实验（令人吃惊的是，很多学生在学习过操作系统课程后，却没有写过一个并发程序或使用过操作系统功能）。

如果要在一个学期内覆盖本书的大部分内容，则要求极快的进度，而这在本科生课程里难以达到。此时，选择略去哪些内容将很大程度上取决于选修本课程的学生的背景。在系统课程里，我们需要课堂讲解时间来帮助学生理解动机和细节。如果学生修过的“数据结构”课程里对内存管理和表处理进行过讨论，则本书第 4 章和第 9 章的内容可以略过。如果学生在将来会选修网络方面的课程，则第 17 章的网络协议内容也可以跳过。此外，本书包括一章远程磁盘系统和一章远程文件系统，这两章的内容存在一些相似之处，可以略过一章。相对来说，远程磁盘系统一章的内容可能更加贴切，因为该章引入了磁盘块缓存的议题，而该议题对于许多操作系统来说都非常重要。

在研究生课程里，课堂时间可以用来讨论动机、原理、折中、不同原语集和不同的实现方案比较。学生在本课程学习结束后，应当对进程模型、中断和进程之间的关系有一个深刻的理解，同时也将具备理解、创建和修改系统组件的能力。学生应当在大脑中建立起了整个系统的完整概念模型，并且知道所有的组件之间是如何交互协作的。

我推荐在各个层面上设计程序设计实验。本书的许多练习都推荐对代码进行修改或者测量，或者尝试不同的实现方案。相关的软件可在下面的网站上免费下载，该网站上还列有如何创建一个 Linksys

实验平台的指令：www.xinu.cs.purdue.edu。

因为 Linksys 的硬件非常便宜，所以构建一个实验的成本很低。此外，我们也有用于其他硬件平台的软件版本，这些版本包括 x86 和 ARM 的一个功能有限的版本。

本书中的许多练习都建议进行改进、实验和不同实现，但是也可以设计大型实验项目。可以用于不同硬件平台的大型实验例子包括：虚拟内存系统、不同计算机之间指令执行的同步机制、虚拟网络的设计等。例如，普度大学的一些学生就将 Xinu 操作系统移植到了各种处理器上，或者为各种 I/O 设备编写了设备驱动程序。

学习本书前，学生需要具备基本的程序设计能力。学生应当理解基本的数据结构，这些基本结构包括链表、栈和队列，并且应当用 C 语言写过程序。

最后，我鼓励设计人员尽可能使用高级程序设计语言，仅在必要的情况下才使用汇编语言。根据这种原则，Xinu 操作系统的大部分都是用 C 语言编写的。少数一些与机器相关的功能，如上下文切换和中断分配器的最底层功能，则是用汇编语言写成的。所有的汇编语言代码都附有解释和注释，使读者无需学习汇编语言的细节就可以理解这些代码。此外，我们还提供用于其他平台的 Xinu 版本，这样就可以对在各种平台上实现 Xinu 操作系统的成本进行比较。例如，我们可以对在 MIPS 处理器上实现 Xinu 所需要的代码量和在其他处理器架构（如 x86）上实现 Xinu 所需要的代码量进行比较。

本书的成书要归功于我过去在商业操作系统上所获得的各种经验，这些经验有好也有坏。虽然 Xinu 操作系统与现有的操作系统在内部机制上并不相同，但其基本的思想却并不新颖。另外，虽然 Xinu 系统里的许多概念和名称都来自于 UNIX 系统，但读者应当注意，这两个系统里的许多函数所使用的参数和内部结构有巨大的不同。因此，为一个系统所写的应用程序在未经修改的情况下不能在另一个平台上运行。

我感谢为 Xinu 项目贡献了思想、辛劳和激情的所有人的帮助。在过去的岁月里，普度大学的许多研究生都从事过本系统的工作，他们为 Xinu 进行过移植，写过设备驱动。从原始的系统版本开始，Xinu 到目前已经走过了 30 多年的历程。本书的 Xinu 版本是原始版本的一个完全重写，但却保留了原始设计的优雅。Dennis Brylow 将 Xinu 移植到了 Linksys 平台，并且创建了许多底层的构件，包括启动代码、上下文切换和 Ethernet 驱动。Dennis 还设计了重启机制，并应用在普度大学的实验室里。另外，我特别要感谢我的妻子和我的合作伙伴 Christine，她的仔细编辑和建议让本书改善良多。

Douglas E. Comer

2011 年 8 月

关于作者

Operating System Design: The Xinu Approach, Linksys Version

Douglas Comer 是美国普度大学 (Purdue University) 计算机系的杰出教授, 国际公认的计算机网络、TCP/IP 协议、Internet 和操作系统设计的专家。Comer 发表论文无数, 出版专著多部, 是一位为研究和教育而开发课程体系和实验项目的先驱。

作为一个多产作家, Comer 博士的书被翻译成 16 种语言, 广泛应用于全世界的计算机科学、工程和工商管理管理等院系和行业。Comer 博士划时代的三卷巨著《Internetworking with TCP/IP》对网络和网络教育产生了革命性的影响。他所编写的教科书和富有创意的实验手册已经塑造和继续塑造着研究生和本科生的教学课程体系。

Comer 博士所写书籍的精确和洞见反映出他在计算机系统领域的广泛背景。他的研究横跨硬件和软件。他创建了一个完整的操作系统 Xinu, 编写了设备驱动程序, 为传统计算机和网络处理器实现了网络协议软件。Comer 博士的研究成果已经应用到工业界的各种产品中。

Comer 博士创建和主讲的课程包括“网络协议”、“操作系统”、“计算机体系架构”, 其听众既有大学生和学术界同仁, 也有工业界的工程师。他创新的教育实验让他和他的学生能够设计和实现大型、复杂系统的原型, 并对结果原型的性能进行度量。Comer 博士长期在公司、大学和会议上讲课和演说, 还为工业界提供咨询服务, 以帮助他们设计计算机网络和系统。

二十多年来, Comer 教授担任研究期刊《Software—Practice and Experience》主编。在普度大学停薪留职期间, 他在思科公司 (Cisco) 担任研究副总裁。Comer 博士是 ACM 院士、普度教育学院院士和无数奖项的获得者, 其中包括 Usenix 终身成就奖。

关于 Comer 博士的更多信息可在下面网站找到: www.cs.purdue.edu/people/comer。

关于 Comer 博士所著书籍的更多信息可在下面网站找到: www.comerbooks.com。

出版者的话	
译者序	
前言	
关于作者	
第 1 章 引言和概述	1
1.1 操作系统	1
1.2 本书的研究方法	1
1.3 分层设计	2
1.4 Xinu 操作系统	3
1.5 操作系统不是什么	3
1.6 从外面看操作系统	4
1.7 其他章节概要	4
1.8 观点	5
1.9 总结	5
练习	5
第 2 章 并发执行与操作系统服务	6
2.1 引言	6
2.2 多活动的编程模型	6
2.3 操作系统服务	7
2.4 并发处理的概念和术语	7
2.5 串行程序和并发程序的区别	8
2.6 多进程共享同一段代码	9
2.7 进程退出与进程终止	11
2.8 共享内存、竞争条件和同步	11
2.9 信号量与互斥	14
2.10 Xinu 中的类型命名方法	15
2.11 使用 Kputc 和 Kprintf 进行操作系统 的调试	16
2.12 观点	16
2.13 总结	16
练习	17
第 3 章 硬件和运行时环境概览	18
3.1 引言	18
3.2 E2100L 的物理和逻辑结构	18
3.3 处理器结构和寄存器	19
3.4 总线操作：获取 - 存储范式	19
3.5 直接内存访问	19
3.6 总线地址空间	20
3.7 内核段 KSEG0 和 KSEG1 的内容	20
3.8 总线启动的静态配置	21
3.9 调用约定和运行时栈	21
3.10 中断和中断处理	22
3.11 异常处理	23
3.12 计时器硬件	23
3.13 串行通信	24
3.14 轮询与中断驱动 I/O	24
3.15 内存缓存和 KSEG1	24
3.16 存储布局	24
3.17 内存保护	25
3.18 观点	25
练习	25
第 4 章 链表与队列操作	26
4.1 引言	26
4.2 用于进程链表的统一数据结构	26
4.3 简洁的链表数据结构	27
4.4 队列数据结构的实现	28
4.5 内联队列操作函数	29
4.6 获取链表中进程的基础函数	29
4.7 FIFO 队列操作	30
4.8 优先级队列的操作	32
4.9 链表初始化	33
4.10 观点	34
4.11 总结	34
练习	35
第 5 章 调度和上下文切换	36
5.1 引言	36
5.2 进程表	36
5.3 进程状态	38
5.4 就绪和当前状态	38
5.5 调度策略	38
5.6 调度的实现	39
5.7 上下文切换的实现	41
5.8 内存中保存的状态	41

5.9	在 MIPS 处理器上切换上下文	41	7.16	总结	70
5.10	重新启动进程执行的地址	43	练习		71
5.11	并发执行和 null 进程	44	第 8 章 消息传递		72
5.12	使进程准备执行和调度不变式	44	8.1	引言	72
5.13	推迟重新调度	45	8.2	两种类型的消息传递服务	72
5.14	其他进程调度算法	47	8.3	消息使用资源的限制	72
5.15	观点	47	8.4	消息传递函数和状态转换	73
5.16	总结	47	8.5	send 的实现	73
练习		47	8.6	receive 的实现	74
第 6 章 更多进程管理		49	8.7	非阻塞消息接收的实现	75
6.1	引言	49	8.8	观点	75
6.2	进程挂起和恢复	49	8.9	总结	75
6.3	自我挂起和信息隐藏	49	练习		76
6.4	系统调用的概念	50	第 9 章 基本内存管理		77
6.5	禁止中断和恢复中断	51	9.1	引言	77
6.6	系统调用模板	51	9.2	内存的类型	77
6.7	系统调用返回 SYSERR 和 OK 值	51	9.3	重量级进程的定义	77
6.8	挂起的实现	52	9.4	小型嵌入式系统的内存管理	78
6.9	挂起当前进程	53	9.5	程序段和内存区域	78
6.10	suspend 函数的返回值	53	9.6	嵌入式系统中的动态内存分配	79
6.11	进程终止和进程退出	54	9.7	低层内存管理器的设计	79
6.12	进程创建	56	9.8	分配策略和内存持久性	80
6.13	其他进程管理函数	59	9.9	追踪空闲内存	80
6.14	总结	60	9.10	低层内存管理的实现	80
练习		61	9.11	分配堆存储	82
第 7 章 协调并发进程		62	9.12	分配栈存储	83
7.1	引言	62	9.13	释放堆和栈存储	84
7.2	进程同步的必要性	62	9.14	观点	86
7.3	计数信号量的概念	63	9.15	总结	87
7.4	避免忙等待	63	练习		87
7.5	信号量策略和进程选择	63	第 10 章 高级内存管理和虚拟内存		88
7.6	等待状态	64	10.1	引言	88
7.7	信号量数据结构	64	10.2	分区空间分配	88
7.8	系统调用 wait	65	10.3	缓冲池	88
7.9	系统调用 signal	66	10.4	分配缓冲池	89
7.10	静态和动态信号量分配	66	10.5	将缓冲池返回给缓冲池	90
7.11	动态信号量的实现示例	67	10.6	创建缓冲池	91
7.12	信号量删除	68	10.7	初始化缓冲池表	93
7.13	信号量重置	69	10.8	虚拟内存和内存复用	93
7.14	多核处理器之间的协调	69	10.9	实地址空间和虚地址空间	93
7.15	观点	70	10.10	支持按需换页的硬件	94

10.11	使用页表的地址翻译	95	13.3	实时时钟和计时器硬件	118
10.12	页表项中的元数据	95	13.4	处理实时时钟中断	119
10.13	按需换页以及设计上的问题	95	13.5	延时与抢占	119
10.14	页面替换和全局时钟算法	96	13.6	使用计时器来模拟实时时钟	120
10.15	观点	97	13.7	抢占的实现	120
10.16	总结	97	13.8	使用增量链表对延迟进行有效 管理	120
	练习	97	13.9	增量链表的实现	121
第 11 章	高层消息传递	98	13.10	将进程转入睡眠	122
11.1	引言	98	13.11	定时消息接收	124
11.2	进程间通信端口	98	13.12	唤醒睡眠进程	127
11.3	端口实现	98	13.13	时钟中断处理	127
11.4	端口表初始化	99	13.14	时钟初始化	128
11.5	端口创建	100	13.15	间隔计时器管理	129
11.6	向端口发送消息	101	13.16	观点	130
11.7	从端口接收消息	102	13.17	总结	130
11.8	端口的删除和重置	103		练习	130
11.9	观点	106	第 14 章	设备无关的 I/O	132
11.10	总结	106	14.1	引言	132
	练习	106	14.2	I/O 和设备驱动的概念结构	132
第 12 章	中断处理	107	14.3	接口抽象和驱动抽象	133
12.1	引言	107	14.4	I/O 接口的一个示例	134
12.2	中断的优点	107	14.5	打开 - 读 - 写 - 关闭范式	134
12.3	中断分配	107	14.6	绑定 I/O 操作和设备名	134
12.4	中断向量	107	14.7	Xinu 中的设备名	135
12.5	中断向量号的分配	108	14.8	设备转换表概念	135
12.6	硬件中断	108	14.9	设备和共享驱动的多个副本	136
12.7	中断请求的局限性和中断多路 复用	109	14.10	高层 I/O 操作的实现	138
12.8	中断软件和分配	109	14.11	其他高层 I/O 函数	138
12.9	中断分配器底层部分	110	14.12	打开、关闭和引用计数	141
12.10	中断分配器高层部分	112	14.13	devtab 中的空条目和错误条目	143
12.11	禁止中断	114	14.14	I/O 系统的初始化	143
12.12	函数中中断代码引起的限制	115	14.15	观点	146
12.13	中断过程中重新调度的必要性	115	14.16	总结	147
12.14	中断过程中的重新调度	115		练习	147
12.15	观点	116	第 15 章	设备驱动示例	148
12.16	总结	116	15.1	引言	148
	练习	117	15.2	tty 抽象	148
第 13 章	实时时钟管理	118	15.3	tty 设备驱动的组成	149
13.1	引言	118	15.4	请求队列和缓冲区	149
13.2	定时事件	118	15.5	上半部和下半部的同步	150

15.6	硬件缓冲区和驱动设计	151	17.5	网络数据包的定义	198
15.7	tty 控制块和数据声明	151	17.6	网络输入进程	199
15.8	次设备号	153	17.7	UDP 表的定义	202
15.9	上半部 tty 字符输入 (ttyGetc) ...	153	17.8	UDP 函数	203
15.10	通用上半部 tty 输入 (ttyRead)	154	17.9	互联网控制报文协议	210
15.11	上半部 tty 字符输出 (ttyPutc) ...	155	17.10	动态主机配置协议	211
15.12	开始输出 (ttyKickOut)	156	17.11	观点	214
15.13	上半部 tty 多字符输出 (ttyWrite)	157	17.12	总结	214
15.14	下半部 tty 驱动函数 (ttyInterrupt)	157	练习	214	
15.15	输出中断处理 (ttyInter_out)	159	第 18 章 远程磁盘驱动	215	
15.16	tty 输入处理 (tty Inter-in)	161	18.1	引言	215
15.17	tty 控制块初始化 (ttyInit)	166	18.2	磁盘抽象	215
15.18	设备驱动控制	168	18.3	磁盘操作驱动支持	215
15.19	观点	169	18.4	块传输和高层 I/O 函数	215
15.20	总结	169	18.5	远程磁盘范式	216
练习	169		18.6	磁盘操作的语义	216
第 16 章 DMA 设备和驱动 (以太网)	171		18.7	驱动数据结构的定义	217
16.1	引言	171	18.8	驱动初始化 (rdsInit)	221
16.2	直接内存访问和缓冲区	171	18.9	上半部打开函数 (rdsOpen)	223
16.3	多缓冲区和环	171	18.10	远程通信函数 (rdscomm)	224
16.4	使用 DMA 的以太网驱动例子	172	18.11	上半部写函数 (rdsWrite)	226
16.5	设备的硬件定义和常量	172	18.12	上半部读函数 (rdsRead)	228
16.6	环和内存缓冲区	174	18.13	刷新挂起的请求	231
16.7	以太网控制块的定义	175	18.14	上半部控制函数 (rdsControl)	231
16.8	设备和驱动初始化	177	18.15	分配磁盘缓冲区 (rdsbufalloc) ...	233
16.9	分配输入缓冲区	181	18.16	上半部关闭函数 (rdsClose)	234
16.10	从以太网设备中读取数据包	182	18.17	下半部通信进程 (rdsprocess)	235
16.11	向以太网设备中写入数据包	183	18.18	观点	239
16.12	以太网设备的中断处理	185	18.19	总结	239
16.13	以太网控制函数	187	练习	240	
16.14	观点	189	第 19 章 文件系统	241	
16.15	总结	189	19.1	文件系统是什么	241
练习	189		19.2	文件操作的示例集合	241
第 17 章 最小互联网协议栈	190		19.3	本地文件系统的设计	242
17.1	引言	190	19.4	Xinu 文件系统的数据结构	242
17.2	所需的功能	190	19.5	索引管理器的实现	243
17.3	同步对话、超时和进程	191	19.6	清空索引块 (lfbclear)	246
17.4	ARP 函数	192	19.7	获取索引块 (lfbget)	247
			19.8	存储索引块 (lfbput)	247
			19.9	从空闲链表中分配索引块 (lfballoc)	248

19.10	从空闲链表中分配数据块 (lfdballocc)	249	20.13	写远程文件	286
19.11	使用设备无关的 I/O 函数的文件 操作	250	20.14	远程文件的定位	288
19.12	文件系统的设备设置和函数 名称	251	20.15	远程文件单字符 I/O	288
19.13	本地文件系统打开函数 (lfsOpen)	251	20.16	远程文件系统控制函数	289
19.14	关闭文件伪设备 (lflClose)	256	20.17	初始化远程文件数据结构	292
19.15	刷新磁盘中的数据 (lflflush)	256	20.18	观点	293
19.16	文件的批量传输函数 (lflWrite, lflRead)	257	20.19	总结	293
19.17	在文件中查找一个新位置 (lflSeek)	258	练习	294	
19.18	从文件中提取一个字节 (lflGetc)	259	第 21 章 句法名字空间	295	
19.19	改变文件中的一个字节 (lflPutc)	260	21.1	引言	295
19.20	载入索引块和数据块 (lflsetup)	261	21.2	透明与名字空间的抽象	295
19.21	主文件系统设备的初始化 (lflsInit)	264	21.3	多种命名方案	295
19.22	伪设备的初始化 (lflIInit)	264	21.4	命名系统设计的其他方案	296
19.23	文件截断 (lfltruncate)	265	21.5	基于句法的名字空间	296
19.24	初始文件系统的创建 (lflscreate)	267	21.6	模式和替换	297
19.25	观点	269	21.7	前缀模式	297
19.26	总结	269	21.8	名字空间的实现	297
练习	269	21.9	名字空间的数据结构和常量	297	
第 20 章 远程文件机制	270	21.10	增加名字空间前缀表的映射	298	
20.1	引言	270	21.11	使用前缀表进行名字映射	299
20.2	远程文件访问	270	21.12	打开命名文件	302
20.3	远程文件语义	270	21.13	名字空间初始化	303
20.4	远程文件设计和消息	271	21.14	对前缀表中的项进行排序	305
20.5	远程文件服务器通信	276	21.15	选择一个逻辑名字空间	305
20.6	发送一个基本消息	278	21.16	默认层次和空前缀	305
20.7	网络字节序	279	21.17	额外的对象操作函数	306
20.8	使用设备范式的远程文件系统	279	21.18	名字空间方法的优点和限制	306
20.9	打开远程文件	280	21.19	广义模式	307
20.10	检查文件模式	282	21.20	观点	307
20.11	关闭远程文件	283	21.21	总结	308
20.12	读远程文件	284	练习	308	
			第 22 章 系统初始化	309	
			22.1	引言	309
			22.2	引导程序: 从头开始	309
			22.3	操作系统初始化	309
			22.4	在 E2100L 上启动一个可选的 映像	310
			22.5	Xinu 初始化	310
			22.6	系统启动	312
			22.7	从程序转化为进程	316
			22.8	观点	316

22.9 总结	316	25.2 用户接口	323
练习	316	25.3 命令和设计原则	323
第 23 章 异常处理	317	25.4 一个简化壳的设计决策	324
23.1 引言	317	25.5 壳的组织和操作	324
23.2 异常、陷阱和恶意中断	317	25.6 词法符号的定义	324
23.3 panic 的实现	317	25.7 命令行语法的定义	325
23.4 观点	318	25.8 Xinu 壳的实现	325
23.5 总结	318	25.9 符号的存储	327
练习	318	25.10 词法分析器代码	327
第 24 章 系统配置	319	25.11 命令解释器的核心	330
24.1 引言	319	25.12 命令名查询和内部处理	336
24.2 多重配置的需求	319	25.13 传给命令的参数	336
24.3 Xinu 系统配置	319	25.14 向外部命令传递参数	337
24.4 Xinu 配置文件的内容	320	25.15 I/O 重定向	339
24.5 计算次设备号	321	25.16 示例命令函数 (sleep)	340
24.6 配置 Xinu 系统的步骤	322	25.17 观点	341
24.7 观点	322	25.18 总结	341
24.8 总结	322	练习	342
练习	322	附录 1 操作系统移植	343
第 25 章 一个用户接口例子: Xinu 壳	323	附录 2 Xinu 设计注解	349
25.1 引言	323	索引	352

引言和概述

我们的小小系统也有风光的时刻。

——Alfred, Lord Tennyson

1.1 操作系统

每一个智能设备和计算机系统中都隐藏着这么一类软件，它们控制着处理信息、管理资源以及与显示屏、网络、磁盘和打印机等设备通信的工作。总的来说，这些进行控制和协调工作的代码通常叫做执行器、监视器、任务管理器，或者内核，而我们将使用一个更宽泛的术语操作系统。

计算机操作系统是人类创造的最复杂的物体之一：计算机操作系统允许多个计算进程和用户同时共享一个CPU，保护数据免受未经授权的访问，并保持独立输入/输出(I/O)设备的正确运行。操作系统提供的高级服务都是通过向复杂的硬件发送一系列详细的命令实现的。有趣的是，操作系统并不是从外部控制电脑的独立机制——它还包括一些软件，这些软件由执行应用程序的同一处理器执行。事实上，当处理器运行应用程序的时候，处理器是不能执行操作系统的，反之亦然。

保证操作系统总在应用程序运行结束后重新夺回控制权的安排机制使得操作系统的设计变得非常复杂。操作系统最令人印象深刻的方面来自于服务和硬件之间的不同：操作系统在低级的硬件上提供高级服务。随着本书内容的推进，读者就会理解系统软件处理像串行接口这样简单的设备需要做的事情。而其中的哲学原理很简单：操作系统应该提供让编程更加容易的抽象，而不是反映底层硬件设备的抽象。因此，我们得出结论：

1

设计操作系统时，应该隐藏底层的硬件细节，并创建一个为应用程序提供高级服务的抽象机器。

操作系统的设计并不是人们所熟知的工艺。最初，由于计算机的缺乏和价格的昂贵，只有少数程序员有从事操作系统相关工作的机会。而现在，由于先进的微电子技术降低了制造成本使得微处理器不再昂贵，操作系统便成为一种商品，与此同时也只有少数程序员从事操作系统方面的工作。有趣的是，由于微处理器变得非常便宜，大多数电子设备都是从可编程处理器构建得到，而不是从离散的逻辑构建得到。因此，设计与实现微机和微控制器的软件系统不再是专家的专利，它已成为一个称职的系统程序员必须能胜任的技术。

幸运的是，随着生产新机器的技术的发展，我们对于操作系统的理解也在不断提高。研究人员已经找出了根本问题，制定了设计原则，定义了基本的组件，并设计了组件一起工作的机制。更重要的是，研究人员还定义了一系列的抽象，如文件和当前进程（这些抽象对于所有的操作系统都是相同的），并且已经找到了实现这些抽象的有效方式。最后，我们知道了如何将操作系统的不同组件组织成一个有意义的系统设计与实现。

同早期系统相比，现代操作系统是简洁的、可移植的。设计良好的系统都遵循着将软件分割成一系列基本组件的基本设计模式。因此，现代系统就变得更容易理解和修改，相比早期的系统其处理开销也比较小。

供应商出售的大型商业操作系统通常包括很多额外的软件组件。例如，一个典型的操作系统软件发行版包括编译器、连接器、装载程序、库函数和一系列的应用程序。为了区分这些额外的软件和一个基本的操作系统，我们有时会用内核指代常驻在内存中并且提供诸如并发进程支持等关键性服务的代码。在本书中，操作系统这个术语指的就是内核，而不包括其他附加的功能。一个最小化内核功能的设计有时称为微内核设计。我们的讨论就将集中在微内核上。

2

1.2 本书的研究方法

本书讲解了如何构建、设计并且实现操作系统的内核。书中使用了工程学方法，而不是仅仅罗列

操作系统的特性和抽象地对其进行描述。这种方法向我们展示了每一个抽象是如何建立的，以及如何将这些抽象组织成一个优雅、高效的设计。

这种工程学方法有两个优势。第一，因为本书的内容涵盖操作系统的每一部分，所以读者会看到整个系统如何融合在一起，而不仅仅是一两个部分之间如何交互。第二，由于读者可以得到书中描述的所有部分的源代码，所以任何部分的实现都没有什么神秘的地方——读者可以获得一份系统的副本来检查、修改、工具化、测量、扩展或者将其移植到其他架构。在本书的最后，读者会看到操作系统的每个部分是如何满足设计需求的，以帮助读者理解可选的设计方案。

本书的关注点是实现，这意味着代码是本书的一个重要组成部分。事实上，代码是讨论的核心，必须通过阅读和学习所罗列的程序来欣赏其中的微妙之处和工程中的细节。例子代码都非常精简，这意味着读者可以集中精力在概念的理解上而不需要费力地阅读许多页的代码。但某些练习建议的改进或修改需要读者深入细节或者找到其他方案。熟练的程序员会找到更多方法来改进和扩展我们的系统。

1.3 分层设计

如果设计得好，操作系统的内部可以如最常规的程序一样优雅、简洁。为了达到优雅的目标，本书所描述的设计，将系统功能划分为8个大类，并将这些组件组织成多个层次。系统的每层提供了一个抽象的服务，该服务又通过低层提供的抽象服务实现。该方法的特点很明显：逐渐变大的层次子集可以形成更加强大的系统。我们将会看到如何利用分层方法来帮助设计者降低操作系统的复杂性。

该方法的另一个重要特点体现在运行时的效率上——设计者可以在不引入额外开销的情况下将操作系统的片段构建为一个层次结构。不过，该方法不同于传统的分层系统。在传统的分层系统中， K 层的函数只能调用 $K-1$ 层的函数。在这种多层次的方法中，层次只为设计者提供了一个概念模型——在运行时，高层的函数可以直接调用低层的任何函数。我们将看到，直接调用使得整个系统更有效率。

图1-1说明在本书中所使用的层次结构，给出了我们将讨论的组件的预览，并展示了其中所有片段的结构。

分层结构的核心是计算机硬件。虽然硬件不是操作系统本身的一部分，但现代硬件具有的特性允许其与操作系统紧密集成在一起。因此，我们可以认为硬件是层次结构中的第0层。

从硬件开始往上，操作系统的每个层次都会提供更强大的原语，从而为应用程序屏蔽原始的硬件。内存管理层控制和分配内存。进程管理层是操作系统最基础的组成部分，它包括进程调度和上下文切换。接下来一层的功能包含了进程管理的其余部分内容，包括创建、杀死、挂起和恢复进程。进程管理的上层是进程协调组件，它实现了信号量。实时时钟管理的功能包含在接下来的一个层次中，它允许应用软件在一定的时间内推迟响应。实时时钟上面的一层是与设备无关的I/O程序层，它提供我们所熟悉的服务，如读（read）和写（write）操作。在设备程序之上的一层实现网络通信，更上面的一层则实现了文件系统。

系统的内部组织不应该与系统提供的服务相混淆。虽然将组件组织成不同的层次可以使设计和实

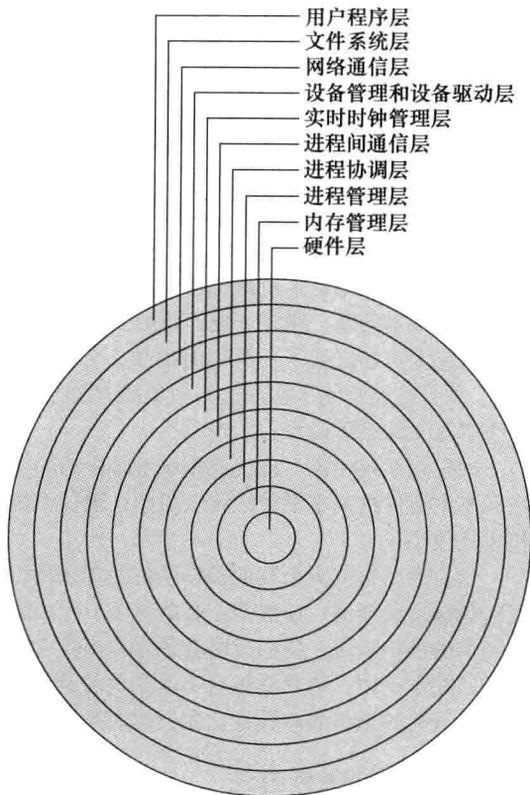


图1-1 在本书中使用的多层次结构