

国外著名高等院校
信息科学与技术优秀教材

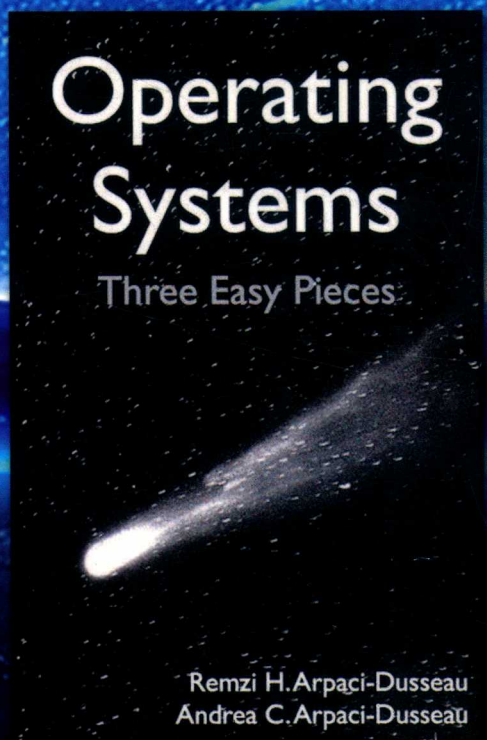
异步图书
www.epubit.com

操作系统导论

[美] 雷姆兹·H.阿帕希杜塞尔 (Remzi H. Arpaci-Dusseau)

著 王海鹏 译

[美] 安德莉亚·C.阿帕希杜塞尔 (Andrea C. Arpaci-Dusseau)



 中国工信出版集团

 人民邮电出版社
POSTS & TELECOM PRESS

国外著名高等院校
信息科学与技术优秀教材

操作系统导论

[美] 雷姆兹·H.阿帕希杜塞尔 (Remzi H. Arpaci-Dusseau)

著

[美] 安德莉亚·C.阿帕希杜塞尔 (Andrea C. Arpaci-Dusseau)

王海鹏 译

人民邮电出版社

北京

图书在版编目 (C I P) 数据

操作系统导论 / (美) 雷姆兹·H. 阿帕希杜塞尔,
(美) 安德莉亚·C. 阿帕希杜塞尔著; 王海鹏译. — 北
京: 人民邮电出版社, 2019.6
国外著名高等院校信息科学与技术优秀教材
ISBN 978-7-115-50823-2

I. ①操… II. ①雷… ②安… ③王… III. ①操作系
统一—高等学校—教材 IV. ①TP316

中国版本图书馆CIP数据核字(2019)第028300号

版权声明

Simplified Chinese translation copyright ©2019 by Posts and Telecommunications Press

ALL RIGHTS RESERVED

Operating Systems: Three Easy Pieces by Remzi H Arpaci-Dusseau, Andrea C Arpaci-Dusseau, ISBN 9781985086593.

Copyright © 2018 by Remzi H Arpaci-Dusseau and Andrea C Arpaci-Dusseau.

本书中文简体版由作者授权人民邮电出版社出版。未经出版者书面许可, 对本书的任何部分不得以任何方式或任何手段复制和传播。

版权所有, 侵权必究。

◆ 著 [美] 雷姆兹·H.阿帕希杜塞尔 (Remzi H. Arpaci-Dusseau)
[美] 安德莉亚·C.阿帕希杜塞尔 (Andrea C. Arpaci-Dusseau)

译 王海鹏

责任编辑 陈冀康

责任印制 焦志炜

◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路11号

邮编 100164 电子邮件 315@ptpress.com.cn

网址 <http://www.ptpress.com.cn>

北京鑫正大印刷有限公司印刷

◆ 开本: 787×1092 1/16

印张: 31.25

字数: 742千字

2019年6月第1版

印数: 1-3000册

2019年6月北京第1次印刷

著作权合同登记号 图字: 01-2016-1423号

定价: 99.00元

读者服务热线: (010)81055410 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京东工商广登字 20170147号

内 容 提 要

这是一本关于现代操作系统的书。全书围绕虚拟化、并发和持久性这 3 个主要概念展开，介绍了所有现代系统的主要组件（包括调度、虚拟内存管理、磁盘和 I/O 子系统、文件系统）。

本书共 50 章，分为 3 个部分，分别讲述虚拟化、并发和持久性的相关内容。本书大部分章节均先提出特定的问题，然后通过书中介绍的技术、算法和思想来解决这些问题。作者以对话形式引入所介绍的主题概念，行文诙谐幽默却又鞭辟入里，力求帮助读者理解操作系统中虚拟化、并发和持久性的原理。

本书内容全面，并给出了真实可运行的代码（而非伪代码），还提供了相应的练习，适合高等院校相关专业教师教学和高校学生自学。

前 言

致本书读者

欢迎阅读本书！我们希望你阅读本书时，就像我们撰写它时一样开心。本书的英文书名为《Operating Systems: Three Easy Pieces》，这显然是向理查德·费曼（Richard Feynman）针对物理学主题创作的、最了不起的一套讲义[F96]致敬。虽然本书不能达到这位著名物理学家设定的高标准，但也许足够让你了解什么是操作系统（以及更一般的系统）。

本书围绕3个主题元素展开讲解：虚拟化（virtualization）、并发（concurrency）和持久性（persistence）。对于这些概念的讨论，最终延伸到讨论操作系统所做的大多数重要事情。希望你在这个过程中体会到一些乐趣。学习新事物很有趣，对吧？

每个主要概念在若干章节中加以阐释，其中大部分章节都提出了一个特定的问题，然后展示了解决它的方法。这些章节很简短，尝试（尽可能地）引用作为这些想法真正来源的源材料。我们写这本书的目的之一就是厘清操作系统的发展脉络，因为我们认为这有助于学生更清楚地理解过去是什么、现在是什么、将来会是什么。在这种情况下，了解香肠的制作方法几乎与了解香肠的优点一样重要。

我们在整本书中采用了几种结构，值得在这里介绍一下。

无论何时，在试图解决问题时，我们首先要说明最重要的问题是什么。我们在书中明确提出关键问题（*crux of the problem*），并希望通过本书其余部分提出的技术、算法和思想来解决。

在许多地方，我们将通过显示一段时间内的行为来解释系统的工作原理。这些时间线（*timeline*）是理解的本质。如果你知道会发生什么，例如，当进程出现页故障时，你就可以真正了解虚拟内存的运行方式。如果你理解日志文件系统将块写入磁盘时发生的情况，就已经迈出了掌握存储系统的第一步。

整本书中有许多“补充”和“提示”，为主线讲解增添了一些趣味性。“补充”倾向于讨论与主要文本相关的内容（但可能不是必要的）；“提示”往往是一般经验，可以应用于所构建的系统。

在整本书中，我们使用最古老的教学方法之一——对话（*dialogue*）。这些对话用于介绍主要的主题概念，并不时地复习这些内容。这也让我们得以用更幽默的方式写作。好吧，你觉得它们是有用还是幽默，完全是另一回事。

在每一个主要部分的开头，我们将首先呈现操作系统提供的抽象（*abstraction*），然后在后续章节中介绍提供抽象所需的机制、策略和其他支持。抽象是计算机科学各个方面的基础，因此它在操作系统中也是必不可少的。

在所有的章节中，我们尝试使用可能的真实代码 (real code)，而非伪代码 (pseudocode)。因此书中几乎所有的示例，你应该能够自己输入并运行它们。在真实系统上运行真实代码是了解操作系统的最佳方式，因此建议你尽可能这样做。

在本书的各个部分，我们提供了一些作业 (homework)，确保你进一步理解书中的内容。其中许多作业都是对操作系统的一些模拟 (simulation) 程序。你应该下载作业，并运行它们，以此来测验自己。作业模拟程序具有以下特征：通过给它们提供不同的随机种子，你可以产生几乎无限的问题，也可以让模拟程序为你解决问题。因此，你可以一次又一次地自测，直至很好地理解了这些知识。

本书最重要的附录是一组项目 (project)，可供你通过设计、测试和实现自己的代码，来了解真实系统的工作原理。所有项目 (以及上面提到的代码示例) 都是使用 C 编程语言 (C programming language) [KR88] 编写的。C 是一种简单而强大的语言，是大多数操作系统的基础，因此值得添加到你的工具箱中。附录中含有两种类型的项目 (请参阅在线附录中的想法)。第一类是系统编程 (system programming) 项目。这些项目非常适合那些不熟悉 C 和 UNIX，并希望学习如何进行底层 C 编程的人。第二类基于在麻省理工学院开发的实际操作系统内核，称为 xv6 [CK+08]。这些项目非常适合已经有一些 C 的经验并希望深入研究操作系统的学生。在威斯康星大学，我们以 3 种不同的方式开课：系统编程、xv6 编程，或两者兼而有之。

致使用本书作为教材的教师

这门课程很适合 15 周的学期，因此授课教师可以在合理的深度范围内讲授大部分主题。如果是 10 周的学期，那么可能需要从每个部分中删除一些细节。还有一些章节是关于虚拟机监视器的，我们通常会在学期的某个时候插入这些章节，或者在虚拟化部分的结尾处，抑或在接近课程结束时作为补充。

本书中的并发主题比较特别。它是许多操作系统书籍中靠前的主题，而在本书中是直到学生了解了 CPU 和内存的虚拟化之后才开始讲解的。根据我们近 15 年来教授本课程的经验，学生很难理解并发问题是如何产生的，或者很难理解人们试图解决它的原因。那是因为他们还不了解地址空间是什么、进程是什么，或者为什么上下文切换可以在任意时间点发生。然而，一旦他们理解了这些概念，那么再引入线程的概念和由此产生的问题就变得相当容易，或者至少比较容易。

我们尽可能使用黑板 (或白板) 来讲课。在着重强调概念的时候，我们会将一些主要的想法和例子带进课堂，并在黑板上展示它们。讲义有助于为学生提供需要解决的具体问题。在着重强调实践的时候，我们就将笔记本电脑连上投影仪，展示实际代码。这种授课风格特别适用于并发的内容以及所有的讨论部分。在这些部分中，教师可以向学生展示与其项目相关的代码。我们通常不使用幻灯片来讲课，但现在我们已经为那些喜欢这种演示风格的人提供了一套教学 PPT。

如果你想要任何这些教学辅助材料，请给 contact@epubit.com.cn 发电子邮件。

最后一个请求：如果你使用免费在线章节，请直接访问作者网站。这有助于我们跟踪

使用情况（过去几年中，本书英文版下载超过 100 万次！），并确保学生获得最新和最好的版本。

致使用本书上课的学生

如果你是读这本书的学生，那么我们很荣幸能够提供一些材料来帮助你学习操作系统的知识。我们至今还能够回想起我们使用过的一些教科书（例如，Hennessy 和 Patterson 的著作 [HP90]，这是一本关于计算机架构的经典著作），并希望这本书能够成为你美好的回忆之一。

你可能已经注意到，这本书英文版的在线版本是免费的，并且可在线获取^①。有一个主要原因：教科书一般都太贵了。我们希望，这本书是新一波免费材料中的第一本（指电子版），以帮助那些寻求知识的人——无论他们来自哪个国家，或者他们愿意花多少钱购买一本书。

我们也希望，在可能的情况下，向你指出书中大部分材料的原始资料——多年来的优秀论文和人物，他们让操作系统领域成为现在的样子。想法不会凭空产生，它们来自聪明勤奋的人（包括众多图灵奖获得者^②），因此如果有可能，我们应该赞美这些想法和人。我们希望这样做能有助于更好地理解已经发生的变革，而不是说好像我们写这本书时那些思想一直就存在一样 [K62]。此外，也许这样的参考文献能够鼓励你深入挖掘，而阅读该领域的著名论文无疑是良好的学习方法之一。

致谢

这里感谢帮助我们编写本书的人。重要的是，你的名字可以出现在这里！但是，你必须提供帮助。请向我们发送一些反馈，帮助完善本书。你可能会出名！或者，至少在某本书中有你的名字。

到目前为止，提供帮助的人包括 Abhirami Senthilkumaran*, Adam Drescher* (WUSTL), Adam Eggum, Aditya Venkataraman, Adriana Iamnitchi and class (USF), Ahmed Fikri*, Ajaykrishna Raghavan, Akiel Khan, Alex Wyler, Ali Razeen (Duke), AmirBehzad Eslami, Anand Mundada, Andrew Valencik (Saint Mary's), Angela Demke Brown (Toronto), B. Brahmananda Reddy (Minnesota), Bala Subrahmanyam Kambala, Benita Bose, Biswajit Mazumder (Clemson), Bobby Jack, Björn Lindberg, Brennan Payne, Brian Gorman, Brian Kroth, Caleb Sumner (Southern Adventist), Cara Lauritzen, Charlotte Kissinger, Chien-Chung Shen (Delaware)*, Christoph Jaeger, Cody Hanson, Dan Soendergaard (U. Aarhus), David Hanle (Grinnell), David Hartman, Deepika Muthukumar, Dheeraj Shetty (North Carolina State), Dorian Arnold (New

① 这里的题外话：我们在这里所说的“免费”并不意味着开源，也不意味着该书没有受到通常保护的版权——它是受到保护的！我们的意思是你可以下载章节，并使用它们来了解操作系统。为什么不是一本开源的书，不像 Linux 一样是一个开源内核？当你阅读它时，这本书应该像一次对话，某人向你解释某事。因此，这就是我们的方法。

② 图灵奖是计算机科学的最高奖项。它就像诺贝尔奖，但你可能从未听说过。

Mexico), Dustin Metzler, Dustin Passofaro, Eduardo Stelmaszczyk, Emad Sadeghi, Emily Jacobson, Emmett Witchel (Texas), Erik Turk, Ernst Biersack (France), Finn Kuusisto*, Glen Granzow (College of Idaho), Guilherme Baptista, Hamid Reza Ghasemi, Hao Chen, Henry Abbey, Hrishikesh Amur, Huanchen Zhang*, Huseyin Sular, Hugo Diaz, Itai Hass (Toronto), Jake Gillberg, Jakob Olandt, James Perry (U. Michigan-Dearborn)*, Jan Reineke (Universität des Saarlandes), Jay Lim, Jerod Weinman (Grinnell), Jiao Dong (Rutgers), Jingxin Li, Joe Jean (NYU), Joel Kuntz (Saint Mary's), Joel Sommers (Colgate), John Brady (Grinnell), Jonathan Perry (MIT), Jun He, Karl Wallinger, Kartik Singhal, Kaushik Kannan, Kevin Liu*, Lei Tian (U. Nebraska-Lincoln), Leslie Schultz, Liang Yin, Lihao Wang, Martha Ferris, Masashi Kishikawa (Sony), Matt Reichoff, Matty Williams, Meng Huang, Michael Walfish (NYU), Mike Griepentrog, Ming Chen (Stonybrook), Mohammed Alali (Delaware), Murugan Kandaswamy, Natasha Eilbert, Nathan Dipiazza, Nathan Sullivan, Neeraj Badlani (N.C. State), Nelson Gomez, Nghia Huynh (Texas), Nick Weinandt, Patricio Jara, Perry Kivolowitz, Radford Smith, Riccardo Mutschlechner, Ripudaman Singh, Robert Ordóñez and class (Southern Adventist), Rohan Das (Toronto)*, Rohan Pasalkar (Minnesota), Ross Aiken, Ruslan Kiselev, Ryland Herrick, Samer Al-Kiswany, Sandeep Ummadi (Minnesota), Satish Chebrolu (NetApp), Satyanarayana Shanmugam*, Seth Pollen, Sharad Punuganti, Shreevatsa R., Sivaraman Sivaraman*, Srinivasan Thirunarayanan*, Suriyahprakash Balaram Sankari, Sy Jin Cheah, Teri Zhao (EMC), Thomas Griebel, Tongxin Zheng, Tony Adkins, Torin Rudeen (Princeton), Tuo Wang, Varun Vats, William Royle (Grinnell), Xiang Peng, Xu Di, Yudong Sun, Yue Zhuo (Texas A&M), Yufui Ren, Zef RosnBrick, Zuyu Zhang。 特别感谢上面标有星号的人，他们的改进建议尤其重要。

此外，衷心感谢 Joe Meehean 教授 (Lynchburg) 为每一章所做的详细注解，感谢 Jerod Weinman 教授 (Grinnell) 和他的全班同学提供的令人难以置信的小册子，感谢 Chien-Chung Shen 教授 (Delaware) 的细致阅读和建议，感谢 Adam Drescher (WUSTL) 的细致阅读和建议，感谢 Glen Granzow (College of Idaho) 提供详细的评论和建议，感谢 Michael Walfish (NYU) 详细的改进建议。上述所有人都给予本书作者巨大的帮助，优化了本书的内容。

另外，非常感谢这些年来参加 537 课程的数百名学生。特别是 2008 年秋季课程的学生，鼓励我们第一次以书面形式写下了这些讲义（他们厌倦了没有任何类型的教科书可读——有进取心的学生！），然后不吝称赞，让我们继续前行（一位同学在那一年的课程评估中喜不自禁地说：“老天爷！你们完全应该写一本教科书！”）。

我们也非常感谢那些参加 xv6 项目实验课程的少数人，这个实验课程大部分现已纳入主要的 537 课程。2009 年春季班的 Justin Cherniak, Patrick Deline, Matt Czech, Tony Gregerson, Michael Griepentrog, Tyler Harter, Ryan Kroiss, Eric Radzikowski, Wesley Reardan, Rajiv Vaidyanathan 和 Christopher Waclawik。2009 年秋季班的 Nick Bearson, Aaron Brown, Alex Bird, David Capel, Keith Gould, Tom Grim, Jeffrey Hugo, Brandon Johnson, John Kjell, Boyan Li, James Loethen, Will McCardell, Ryan Szaroletta, Simon Tso 和 Ben Yule。2010 年春季班的 Patrick Blesi, Aidan Dennis-Oehling, Paras Doshi, Jake Friedman, Benjamin Frisch, Evan Hanson, Pikkili Hemanth, Michael Jeung, Alex Langenfeld, Scott Rick, Mike Treffert, Garret Staus, Brennan Wall, Hans Werner, Soo-Young Yang 和 Carlos Griffin。

虽然没有直接帮助这本书的写作，但我们的研究生教会了我们很多关于系统的知识。我们在威斯康星大学时经常与他们交谈，并且所有真正的工作都是他们做的——通过告诉我们他们在做什么，我们每周都能学习到新事物。下面的列表包括我们已发布论文的当前和以前的学生，带有星号标志的名字是在我们的指导下获得博士学位的人：Abhishek Rajimwale, Andrew Krioukov, Ao Ma, Brian Forney, Chris Dragga, Deepak Ramamurthi, Florentina Popovici *, Haryadi S. Gunawi *, James Nugent, John Bent *, Jun He, Lanyue Lu, Lakshmi Bairavasundaram *, Laxman Visampalli, Leo Arulraj, Meenali Rungta, Muthian Sivathanu *, Nathan Burnett *, Nitin Agrawal *, Ram Alagappan, Sriram Subramanian *, Stephen Todd Jones *, Suli Yang, Swaminathan Sundararaman*, Swetha Krishnan, Thanh Do*, Thanumalayan S. Pillai, Timothy Denehy*, Tyler Harter, Venkat Venkataramani, Vijay Chidambaram, Vijayan Prabhakaran *, Yiyi Zhang *, Yupu Zhang *, Zev Weiss。

最后要感谢 Aaron Brown，他多年前（2009 年春季）首次参加该课程，接着参加了 xv6 实验课程（2009 年秋季），最后还成为了两个课程的研究生助教（从 2010 年秋季到 2012 年春季）。他不知疲倦的工作极大地改善了项目的状态（特别是 xv6 项目），因此有助于改善威斯康星大学无数本科生和研究生的学习体验。正如 Aaron 所说的（以他通常的简洁方式）：“谢谢。”

最后的话

叶芝有一句名言：“教育不是注满一桶水，而是点燃一把火。”他说得既对也错^①。你必须“给桶注一点水”，这本书当然可以帮助你完成这部分的教育。毕竟，当你去 Google 面试时，他们会问你一个关于如何使用信号量的技巧问题，确切地知道信号量是什么感觉真好，对吧？

但是，叶芝的主要观点显而易见：教育的真正要点是让你对某些事情感兴趣，可以独立学习更多关于这个主题的东西，而不仅仅是你需要消化什么才能在某些课程上取得好成绩。正如我们的父亲（雷姆兹的父亲 Vedat Arpacı）曾经说过的，“在课堂以外学习。”

我们编写本书以激发你对操作系统的兴趣，让你能自行阅读有关该主题的更多信息，进而与你的教授讨论该领域正在进行的所有令人兴奋的研究，甚至参与这些研究。这是一个伟大的领域，充满了激动人心和精妙的想法，以深刻而重要的方式塑造了计算历史。虽然我们知道这种火不会为你们所有人点燃，但我们希望这能对许多人，甚至是少数人有所帮助。因为一旦火被点燃，那你就真正有能力做出伟大的事情。因此，教育过程的真正意义在于：前进，学习许多新的和引人入胜的主题，通过学习不断成熟，最重要的是，找到能为你点火的东西。^②

威斯康星大学计算机科学教授 雷姆兹和安德莉亚夫妇

^① 如果他真的说了这句话。与许多名言一样，这句名言的历史也是模糊不清的。

^② 如果这听起来像我们承认过去曾是纵火犯，那你可能理解错了。如果这听起来很俗气，好吧，因为它确实是的，但你必须原谅我们。

参考资料

[CK+08] “The xv6 Operating System” Russ Cox, Frans Kaashoek, Robert Morris, Nikolai Zeldovich.

xv6 是作为原来 UNIX 版本 6 的移植版开发的，它代表了通过一种美观、干净、简单的方式来理解现代操作系统。

[F96] “Six Easy Pieces: Essentials Of Physics Explained By Its Most Brilliant Teacher” Richard P. Feynman
Basic Books, 1996

这本书摘取了 1993 年的《费曼物理学讲义》中 6 个最简单的章节。如果你喜欢物理学，那么就读一读这本很优秀的读物吧。

[HP90] “Computer Architecture a Quantitative Approach” (1st ed.) David A. Patterson and John L. Hennessy
Morgan-Kaufman, 1990

在读本科时，这本书成为了我们去攻读研究生的动力。我们后来都很高兴与 Patterson 合作，他为我们研究事业基础的奠定给予了极大的帮助。

[KR88] “The C Programming Language” Brian Kernighan and Dennis Ritchie Prentice-Hall, April 1988

每个人都应该拥有一本由发明该语言的人编写的 C 编程参考书。

[K62] “The Structure of Scientific Revolutions” Thomas S. Kuhn
University of Chicago Press, 1962

这是关于科学过程基础知识的著名读物，包括科学过程的整理工作、异常、危机和变革。我们要做的是整理工作。

资源与支持

本书由异步社区出品，社区 (<https://www.epubit.com/>) 为你提供相关资源和后续服务。

配套资源

本书为教师提供如下教学辅助资源：

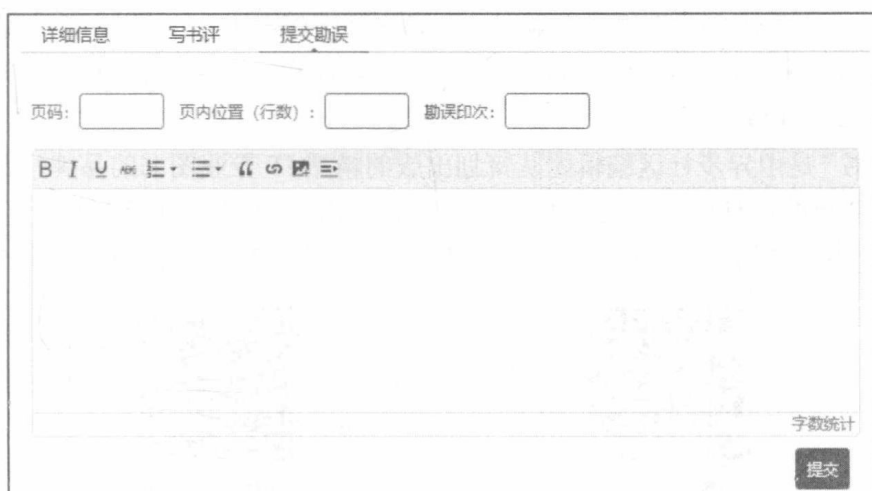
- 教学 PPT 和听课笔记；
- 考试题和参考答案；
- 讨论题和作业；
- 项目说明和指导。

如果您是教师，希望获得教学配套资源，请发邮件到 contact@epubit.com.cn 申请，或者在社区本书页面中直接联系本书的责任编辑。

提交勘误

作者和编辑尽最大努力来确保书中内容的准确性，但难免会存在疏漏。欢迎您将发现的问题反馈给我们，帮助我们提升图书的质量。

当您发现错误时，请登录异步社区，按书名搜索，进入本书页面，单击“提交勘误”，输入勘误信息，单击“提交”按钮即可。本书的作者和编辑会对您提交的勘误进行审核，确认并接受后，您将获赠异步社区的 100 积分。积分可用于在异步社区兑换优惠券、样书或奖品。



The screenshot shows a web interface for reporting errors. At the top, there are three tabs: '详细信息' (Details), '写书评' (Write Review), and '提交勘误' (Submit Error), with the last one being active. Below the tabs, there are three input fields: '页码:' (Page Number), '页内位置 (行数):' (Page Position (Line Number)), and '勘误印次:' (Error Count). Below these fields is a rich text editor with a toolbar containing icons for bold (B), italic (I), underline (U), list (bulleted and numbered), link (S), and other text formatting options. At the bottom right of the form, there is a '字数统计' (Character Count) label and a '提交' (Submit) button.

扫码关注本书

扫描下方二维码，您将会在异步社区微信服务号中看到本书信息及相关的服务提示。



与我们联系

我们的联系邮箱是 contact@epubit.com.cn。

如果您对本书有任何疑问或建议，请您发邮件给我们，请在邮件标题中注明本书书名，以便我们更高效地做出反馈。

如果您有兴趣出版图书、录制教学视频，或者参与图书翻译、技术审校等工作，可以发邮件给我们；有意出版图书的作者也可以到异步社区在线提交投稿（直接访问 www.epubit.com/selfpublish/submission 即可）。

如果您是学校、培训机构或企业，想批量购买本书或异步社区出版的其他图书，也可以发邮件给我们。

如果您在网上发现有针对异步社区出品图书的各种形式的盗版行为，包括对图书全部或部分内容的非授权传播，请您将怀疑有侵权行为的链接发邮件给我们。您的这一举动是对作者权益的保护，也是我们持续为您提供有价值的内容的动力之源。

关于异步社区和异步图书

“异步社区”是人民邮电出版社旗下 IT 专业图书社区，致力于出版精品 IT 技术图书和相关学习产品，为作译者提供优质出版服务。异步社区创办于 2015 年 8 月，提供大量精品 IT 技术图书和电子书，以及高品质技术文章和视频课程。更多详情请访问异步社区官网 <https://www.epubit.com>。

“异步图书”是由异步社区编辑团队策划出版的精品 IT 专业图书的品牌，依托于人民邮电出版社近 30 年的计算机图书出版积累和专业编辑团队，相关图书在封面上印有异步图书的 LOGO。异步图书的出版领域包括软件开发、大数据、AI、测试、前端、网络技术 etc。



异步社区



微信服务号

目 录

第 1 章 关于本书的对话	1	2.4 持久性	9
第 2 章 操作系统介绍	3	2.5 设计目标	11
2.1 虚拟化 CPU	4	2.6 简单历史	12
2.2 虚拟化内存	6	2.7 小结	15
2.3 并发	7	参考资料	15

第 1 部分 虚拟化

第 3 章 关于虚拟化的对话	18	作业 (测量)	47
第 4 章 抽象: 进程	19	第 7 章 进程调度: 介绍	48
4.1 抽象: 进程	20	7.1 工作负载假设	48
4.2 进程 API	20	7.2 调度指标	49
4.3 进程创建: 更多细节	21	7.3 先进先出 (FIFO)	49
4.4 进程状态	22	7.4 最短任务优先 (SJF)	50
4.5 数据结构	24	7.5 最短完成时间优先 (STCF)	51
4.6 小结	25	7.6 新度量指标: 响应时间	52
参考资料	25	7.7 轮转	52
作业	26	7.8 结合 I/O	54
问题	26	7.9 无法预知	54
第 5 章 插叙: 进程 API	28	7.10 小结	55
5.1 fork() 系统调用	28	参考资料	55
5.2 wait() 系统调用	29	作业	56
5.3 最后是 exec() 系统调用	30	问题	56
5.4 为什么这样设计 API	32	第 8 章 调度: 多级反馈队列	57
5.5 其他 API	34	8.1 MLFQ: 基本规则	57
5.6 小结	34	8.2 尝试 1: 如何改变优先级	58
参考资料	34	8.3 尝试 2: 提升优先级	60
作业 (编码)	35	8.4 尝试 3: 更好的计时方式	61
问题	35	8.5 MLFQ 调优及其他问题	61
第 6 章 机制: 受限直接执行	37	8.6 MLFQ: 小结	62
6.1 基本技巧: 受限直接执行	37	参考资料	63
6.2 问题 1: 受限制的操作	38	作业	64
6.3 问题 2: 在进程之间切换	40	问题	64
6.4 担心并发吗	44	第 9 章 调度: 比例份额	65
6.5 小结	45	9.1 基本概念: 彩票数表示份额	65
参考资料	45	9.2 彩票机制	66

9.3 实现.....	67	15.2 一个例子.....	101
9.4 一个例子.....	68	15.3 动态(基于硬件)重定位.....	103
9.5 如何分配彩票.....	68	15.4 硬件支持: 总结.....	105
9.6 为什么不是确定的.....	69	15.5 操作系统的问题.....	105
9.7 小结.....	70	15.6 小结.....	108
参考资料.....	70	参考资料.....	109
作业.....	71	作业.....	110
问题.....	71	问题.....	110
第 10 章 多处理器调度(高级)	73	第 16 章 分段	111
10.1 背景: 多处理器架构.....	73	16.1 分段: 泛化的基址/界限.....	111
10.2 别忘了同步.....	75	16.2 我们引用哪个段.....	113
10.3 最后一个问题: 缓存亲和度.....	76	16.3 栈怎么办.....	114
10.4 单队列调度.....	76	16.4 支持共享.....	114
10.5 多队列调度.....	77	16.5 细粒度与粗粒度的分段.....	115
10.6 Linux 多处理器调度.....	79	16.6 操作系统支持.....	115
10.7 小结.....	79	16.7 小结.....	117
参考资料.....	79	参考资料.....	117
第 11 章 关于 CPU 虚拟化的总结对话	81	作业.....	118
第 12 章 关于内存虚拟化的对话	83	问题.....	119
第 13 章 抽象: 地址空间	85	第 17 章 空闲空间管理	120
13.1 早期系统.....	85	17.1 假设.....	120
13.2 多道程序和时分共享.....	85	17.2 底层机制.....	121
13.3 地址空间.....	86	17.3 基本策略.....	126
13.4 目标.....	87	17.4 其他方式.....	128
13.5 小结.....	89	17.5 小结.....	130
参考资料.....	89	参考资料.....	130
第 14 章 插叙: 内存操作 API	91	作业.....	131
14.1 内存类型.....	91	问题.....	131
14.2 malloc()调用.....	92	第 18 章 分页: 介绍	132
14.3 free()调用.....	93	18.1 一个简单例子.....	132
14.4 常见错误.....	93	18.2 页表存在哪里.....	134
14.5 底层操作系统支持.....	96	18.3 列表中究竟有什么.....	135
14.6 其他调用.....	97	18.4 分页: 也很慢.....	136
14.7 小结.....	97	18.5 内存追踪.....	137
参考资料.....	97	18.6 小结.....	139
作业(编码).....	98	参考资料.....	139
问题.....	98	作业.....	140
第 15 章 机制: 地址转换	100	问题.....	140
15.1 假设.....	101	第 19 章 分页: 快速地址转换(TLB)	142
		19.1 TLB 的基本算法.....	142

19.2 示例: 访问数组.....	143	21.7 小结.....	170
19.3 谁来处理 TLB 未命中.....	145	参考资料.....	171
19.4 TLB 的内容.....	146	第 22 章 超越物理内存: 策略	172
19.5 上下文切换时对 TLB 的处理.....	147	22.1 缓存管理.....	172
19.6 TLB 替换策略.....	149	22.2 最优替换策略.....	173
19.7 实际系统的 TLB 表项.....	149	22.3 简单策略: FIFO.....	175
19.8 小结.....	150	22.4 另一简单策略: 随机.....	176
参考资料.....	151	22.5 利用历史数据: LRU.....	177
作业 (测量).....	152	22.6 工作负载示例.....	178
问题.....	153	22.7 实现基于历史信息的算法.....	180
第 20 章 分页: 较小的表	154	22.8 近似 LRU.....	181
20.1 简单的解决方案: 更大的页.....	154	22.9 考虑脏页.....	182
20.2 混合方法: 分页和分段.....	155	22.10 其他虚拟内存策略.....	182
20.3 多级页表.....	157	22.11 抖动.....	183
20.4 反向页表.....	162	22.12 小结.....	183
20.5 将页表交换到磁盘.....	163	参考资料.....	183
20.6 小结.....	163	作业.....	185
参考资料.....	163	问题.....	185
作业.....	164	第 23 章 VAX/VMS 虚拟内存系统	186
问题.....	164	23.1 背景.....	186
第 21 章 超越物理内存: 机制	165	23.2 内存管理硬件.....	186
21.1 交换空间.....	165	23.3 一个真实的地址空间.....	187
21.2 存在位.....	166	23.4 页替换.....	189
21.3 页错误.....	167	23.5 其他漂亮的虚拟内存技巧.....	190
21.4 内存满了怎么办.....	168	23.6 小结.....	191
21.5 页错误处理流程.....	168	参考资料.....	191
21.6 交换何时真正发生.....	169	第 24 章 内存虚拟化总结对话	193
第 2 部分 并发			
第 25 章 关于并发的对话	196	参考资料.....	207
第 26 章 并发: 介绍	198	作业.....	208
26.1 实例: 线程创建.....	199	问题.....	208
26.2 为什么更糟糕: 共享数据.....	201	第 27 章 插叙: 线程 API	210
26.3 核心问题: 不可控的调度.....	203	27.1 线程创建.....	210
26.4 原子性愿望.....	205	27.2 线程完成.....	211
26.5 还有一个问题: 等待另一个 线程.....	206	27.3 锁.....	214
26.6 小结: 为什么操作系统课要研究 并发.....	207	27.4 条件变量.....	215
		27.5 编译和运行.....	217
		27.6 小结.....	217

参考资料	218	30.3 覆盖条件	260
第 28 章 锁	219	30.4 小结	261
28.1 锁的基本思想	219	参考资料	261
28.2 Pthread 锁	220	第 31 章 信号量	263
28.3 实现一个锁	220	31.1 信号量的定义	263
28.4 评价锁	220	31.2 二值信号量(锁)	264
28.5 控制中断	221	31.3 信号量用作条件变量	266
28.6 测试并设置指令(原子交换)	222	31.4 生产者/消费者(有界缓冲区) 问题	268
28.7 实现可用的自旋锁	223	31.5 读者—写者锁	271
28.8 评价自旋锁	225	31.6 哲学家就餐问题	273
28.9 比较并交换	225	31.7 如何实现信号量	275
28.10 链接的加载和条件式存储指令	226	31.8 小结	276
28.11 获取并增加	228	参考资料	276
28.12 自旋过多:怎么办	229	第 32 章 常见并发问题	279
28.13 简单方法:让出来吧,宝贝	229	32.1 有哪些类型的缺陷	279
28.14 使用队列:休眠替代自旋	230	32.2 非死锁缺陷	280
28.15 不同操作系统,不同实现	232	32.3 死锁缺陷	282
28.16 两阶段锁	233	32.4 小结	288
28.17 小结	233	参考资料	289
参考资料	233	第 33 章 基于事件的并发(进阶)	291
作业	235	33.1 基本想法:事件循环	291
问题	235	33.2 重要 API: select() (或 poll())	292
第 29 章 基于锁的并发数据结构	237	33.3 使用 select()	293
29.1 并发计数器	237	33.4 为何更简单?无须锁	294
29.2 并发链表	241	33.5 一个问题:阻塞系统调用	294
29.3 并发队列	244	33.6 解决方案:异步 I/O	294
29.4 并发散列表	245	33.7 另一个问题:状态管理	296
29.5 小结	246	33.8 什么事情仍然很难	297
参考资料	247	33.9 小结	298
第 30 章 条件变量	249	参考资料	298
30.1 定义和程序	250	第 34 章 并发的总结对话	300
30.2 生产者/消费者(有界缓冲区) 问题	252		
		第 3 部分 持久性	
第 35 章 关于持久性的对话	302	36.3 标准协议	304
第 36 章 I/O 设备	303	36.4 利用中断减少 CPU 开销	305
36.1 系统架构	303	36.5 利用 DMA 进行更高效的数据 传送	306
36.2 标准设备	304		

36.6 设备交互的方法.....	307	39.8 获取文件信息.....	348
36.7 纳入操作系统：设备驱动程序.....	307	39.9 删除文件.....	349
36.8 案例研究：简单的 IDE 磁盘驱动 程序.....	309	39.10 创建目录.....	349
36.9 历史记录.....	311	39.11 读取目录.....	350
36.10 小结.....	311	39.12 删除目录.....	351
参考资料.....	312	39.13 硬链接.....	351
第 37 章 磁盘驱动器	314	39.14 符号链接.....	353
37.1 接口.....	314	39.15 创建并挂载文件系统.....	354
37.2 基本几何形状.....	314	39.16 小结.....	355
37.3 简单的磁盘驱动器.....	315	参考资料.....	355
37.4 I/O 时间：用数学.....	318	作业.....	356
37.5 磁盘调度.....	320	问题.....	356
37.6 小结.....	323	第 40 章 文件系统实现	357
参考资料.....	323	40.1 思考方式.....	357
作业.....	324	40.2 整体组织.....	358
问题.....	324	40.3 文件组织：inode.....	359
第 38 章 廉价冗余磁盘阵列 (RAID)	326	40.4 目录组织.....	363
38.1 接口和 RAID 内部.....	327	40.5 空闲空间管理.....	364
38.2 故障模型.....	327	40.6 访问路径：读取和写入.....	364
38.3 如何评估 RAID.....	328	40.7 缓存和缓冲.....	367
38.4 RAID 0 级：条带化.....	328	40.8 小结.....	369
38.5 RAID 1 级：镜像.....	331	参考资料.....	369
38.6 RAID 4 级：通过奇偶校验节省 空间.....	333	作业.....	370
38.7 RAID 5 级：旋转奇偶校验.....	336	问题.....	371
38.8 RAID 比较：总结.....	337	第 41 章 局部性和快速文件系统	372
38.9 其他有趣的 RAID 问题.....	338	41.1 问题：性能不佳.....	372
38.10 小结.....	338	41.2 FFS：磁盘意识是解决方案.....	373
参考资料.....	339	41.3 组织结构：柱面组.....	373
作业.....	340	41.4 策略：如何分配文件和目录.....	374
问题.....	340	41.5 测量文件的局部性.....	375
第 39 章 插叙：文件和目录	342	41.6 大文件例外.....	376
39.1 文件和目录.....	342	41.7 关于 FFS 的其他几件事.....	377
39.2 文件系统接口.....	343	41.8 小结.....	378
39.3 创建文件.....	343	参考资料.....	378
39.4 读写文件.....	344	第 42 章 崩溃一致性：FSCK 和日志	380
39.5 读取和写入，但不按顺序.....	346	42.1 一个详细的例子.....	380
39.6 用 fsync()立即写入.....	346	42.2 解决方案 1：文件系统检查 程序.....	383
39.7 文件重命名.....	347	42.3 解决方案 2：日志 (或预写日志).....	384