

WILEY

高可用性 系统设计

Evan Marcus
Hal Stern
著
汪青青 卢祖英 译



清华大学出版社

高可用性系统设计

Evan Marcus 著
Hal Stern

汪青青 卢祖英 译

清华大学出版社
北京

Evan Marcus & Hal Stern
Blueprints for High Availability
EISBN: 0-471-43026-9

Copyright © 2003 by Wiley Publishing, Inc., Indianapolis, Indiana.

All Rights Reserved. Authorized translation from the English language edition published by John Wiley & Sons, Inc. Simplified Chinese translation edition is published and distributed exclusively by Tsinghua University Press under the authorization by John Wiley & Sons, Inc., within the territory of the People's Republic of China only (excluding Hong Kong, Macao SAR and Taiwan). Unauthorized export of this edition is a violation of the Copyright Act. Violation of this Law is subject to Civil and Criminal Penalties.

本书中文简体字翻译版由美国 John Wiley & Sons, Inc. 公司授权清华大学出版社在中华人民共和国境内（不包括中国香港、澳门特别行政区和中国台湾地区）独家出版发行。未经许可之出口视为违反著作权法，将受法律之制裁。未经出版者预先书面许可，不得以任何方式复制或抄袭本书的任何部分。

北京市版权局著作权合同登记号 图字：01-2004-1996

版权所有，翻印必究。举报电话：010-62782989 13501256678 13801310933

本书封面贴有清华大学出版社防伪标签，无标签者不得销售。

本书防伪标签采用特殊防伪技术，用户可通过在图案表面涂抹清水，图案消失，水干后图案复现；或将表面膜揭下，放在白纸上用彩笔涂抹，图案在白纸上再现的方法识别真伪。

图书在版编目 (CIP) 数据

高可用性系统设计/(美) 马克思 (Marcus, E.) 等著；汪青青，卢祖英译. —北京：清华大学出版社，2005.7
书名原文：Blueprints for High Availability

ISBN 7-302-10865-X

I . 高… II . ①马… ②汪… ③卢… III . 计算机系统—系统设计 IV . TP302.1

中国版本图书馆 CIP 数据核字 (2005) 第 038000 号

出 版 者：清华大学出版社 地 址：北京清华大学学研大厦
<http://www.tup.com.cn> 邮 编：100084
社 总 机：010-62770175 客户服务：010-62776969

责任编辑：常晓波

封面设计：立日新

印 刷 者：北京市清华园胶印厂

装 订 者：三河市李旗庄少明装订厂

发 行 者：新华书店总店北京发行所

开 本：185×230 印张：27.75 字数：622 千字

版 次：2005 年 7 月第 1 版 2005 年 7 月第 1 次印刷

书 号：ISBN 7-302-10865-X/TP · 7227

印 数：1 ~ 2500

定 价：59.00 元

前　　言

第 2 版

《高可用性系统设计》的第 1 版读者反映良好，令人非常满意。看到我们关于高可用性的著作赢得广大读者的认可，让我们深受鼓舞。我们收到很多关于写作风格的反馈，这些反馈告诉我们如何深入浅出地讲解技术问题。

虽然我们得到的评论基本上都是持肯定态度的，但我们很清楚第 1 版的不足之处。在第 2 版中，我们充实了第 1 版中显得单薄的部分，此次的章节安排凝聚了我们更多的心血。

毫无疑问，我们的“实战传奇”（Tales from the Field）最受读者的好评。读者来信说，他们坐下来后，翻开本书就直接跳到“实战传奇”。这的确能让人欣喜。而另一方面，收集这些故事并娓娓道来，也确实给我们带来了许多乐趣。在本版中，我们又增加了许多新的传奇。大家可以先睹为快！

我们要感谢 John Wiley & Sons 的编辑团队。这次又是 Carol Long 督促着本书的完稿，Carol Long 不允许我们不遵守截稿日期，不允许我们有任何拖延行为。我们必须推出一本与第一版同样受欢迎的书。Carol Long 应该得到充分的肯定。Scott Amerman 是此次团队中的新成员。他热情洋溢的鼓励和始终如一的坚持汇成了一股强大的动力，促使我们最终如期完稿。

Evan Marcus 致辞

Hal 和我完成《高可用性系统设计》第 1 版的写作已经快 4 年了，在这段时间里发生了很多变化。对我个人来说最大的变化是我家有了一名新的成员。在写这本书时，我的儿子 Jonathan 快 3 岁了。过去 4 年更大的变化是电脑变得更便宜、更普及了，电脑使用起来更简单了。Jonathan 经常坐在电脑前，打开电脑，登录，放进一张 CD-ROM，然后开始玩游戏，这都是他自己做的。他还会浏览像 www.pbskids.org 这样的网站。我认为一个不太会穿衣服的孩子用起电脑来却得心应手是很不一般的。

在过去的 4 年（事实上是 4 年多的时间）中，社会发生的最大变化发生在 2001 年 9 月 11 日，那天恐怖分子袭击了纽约和华盛顿特区。我一直住在纽约市的郊区（新泽西北部），那里的人们对失去的朋友、邻居和安全有很深的感触。但对于这本书来说，我会把讨论集中在计算机技术和高可用性受到了什么影响上。

在第 1 版中，有一章专门讨论灾难恢复，在那一章里我们有很多重要的问题没讨论。在第 2 版中，我们把关于灾难恢复的那一章全部重写了一遍（也就是第 20 章“灾难恢复”）。

计划”),该章部分包括了我们在9·11后学到和听到的教训。我们还增加了一章(第21章“弹性企业”),该章讲述了纽约期货交易所最著名的故事,讲述了他们怎样从9·11中恢复业务并在攻击发生后12小时内准备恢复交易。当读到纽约期货交易所的故事时,读者会看到我们没讨论他们使用的恢复技术。我们有意这样安排,因为我们认为最重要的不是技术,而是人们的努力,他们不仅使公司生存下来而且还使公司茁壮成长。

第21章在另一本书中以几乎同样的形式出现。在两版《高可用性系统设计》的写作期间,我是VERITAS公司一本名为*The Resilient Enterprise*内部读物的编者之一和投稿人,这章原本是我为那本书而写。我向Richard Barker、Paul Massiglia及那本书的每一位作者表示感谢,他们同意我在这里再次使用这章的内容。

但有些人从不真正吸取教训。9·11事件过后,假如再发生一次袭击,如何使公司的弹性更强,大家就此事进行了议论。很多人认为如果公司把数据分布在不同的地点,并确保没有单点故障,那么效果会更好。由于经济在袭击中受到重创,没有资金立即投入到保护性措施中来,并且随着时间的流逝,其他需要优先做的事情也出现了,本应用来复制数据和进行异地备份的钱也另作他用。很多需要保护的公司在9·11后无所作为,这是一件令人惋惜的事情。如果再发生一次袭击,那么结果将不仅仅是惋惜了。

当然,技术在过去的4年里发生了变化。我们感到有必要增加一章内容来讨论可用性领域新的和流行的技术。第8章是关于SAN、NAS和存储虚拟化的概述。我们还添加了第22章,该章讨论了一些新兴技术。

尽管社会、技术和家庭发生了变化,但我们在第1版中讨论的高可用性的基本原则没有改变。促使第1本书出版的使命宣言仍然有效:“仅仅安装集群软件然后一走了之是得不到高可用性的”。达到高可用性所需的技术并不自动包含在系统和操作系统厂商的产品中。它仍然是困难、复杂并且昂贵的。

在本版中,我们对高可用性的成本和收益采取了更实际的看法,这使可用性指数模型更详细和优秀。技术章节按照它们在指数图中的顺序安排:前面的章节讨论更基本、更便宜的可用性技术,例如备份和磁盘镜像;后面的章节讨论更为复杂更为昂贵的技术,例如复制和灾难恢复,这些技能能提供最高级别的可用性。

自从第1版出版后,很多事情都发生了变化,我们需要重复一下在前言中的注释:有些读者抱怨书中缺少简单、通用的答案。这有两个原因:第一,每个现场和计算机系统出现的问题都不同,指望适合一个有10 000名员工的全球金融机构的方法也适合10人的律师事务所是不实际的,我们给出选择并让读者决定哪一种更适合他或她的应用环境;第二,我使用计算机工作了15年,我认识到大多数计算机问题的答案都比较令人遗憾:“答案要视情况而定”。

撰写一本这样的书的工作量是很大的,单枪匹马根本就不可能完成。很幸运,有很多优秀的人给予了我帮助和支持。我要再次向我的好妻子Carol,两个漂亮的女儿Hannah和Madeline以及可爱的儿子Jonathan表达永恒的爱和感激,Carol忍受了我所有的荒谬的兴

趣和爱好（例如写书）。没有他们的爱和支持，是不可能有这本书的。还要感谢 Roberta 和 David Marcus 给我的父母的爱，以及我的岳父母 Gladys 和 Herb Laden（他们还没给我那个食谱）。

还要感谢我在 VERITAS 很多的朋友和同事，他们给我提供了各种各样的帮助，他们是 Jason Bloomstein、Steven Cohen、John Colgrove、Roger Cummings、Roger Davis、Oleg Kiselev、Graham Moore、Roger Reich、Jim “El Jefe” Senicka 和 Marty Ward。还要感谢 VERTTAS 公司在纽约和新泽西的 woodbridge 的办公室的所有朋友和同事，多年来他们对我的各种工作提供了大量的帮助，特别感谢 Joseph Hand、Vito Vultaggio、Victor DeBellis、Rich Faille、我的室友 Lowell Shulman 和我们当年的新同事 Phil Carty。

还必须感谢我在写本书我那部分时我在 VERITAS 公司的上司：Richard Barker、Mark Bregman、Fred van den Bosch、Hans van Rietschote 以及 Paul Borrill，感谢他们的帮助、支持，特别是那些星期五。感谢我在 VERITAS 跨产品运作部的同事和好朋友，特别是 Guy Bunker 博士、Chris Chandler、Paul Massiglia 和 Paula Skoe，他们给了我巨大的帮助。

更要感谢这些年来与我共同工作过的同事，他们是 Greg Schulz、Greg Schuweiler、Mindy Anderson、Evan Marks 和 Chuck Yerkes。

再一次，特别感谢纽约期货交易所的 Pat Gambaro 和 Steve Bass，当我把他们的故事写到书中时他们的慷慨和给予的协助令人吃惊，他们还允许我一次次地为了修改和添加信息打扰他们。他们太出色了，他们对自己工作的骄傲是最充分的证明。此外，他们还知道皇后区一些很棒的饭店。

Mark Fitzpatrick 是我多年来出色的朋友和支持者。正是 Mark 在 1996 年把我带入 VERITAS，在此之前他读过我写的一篇关于高可用性的文章，在此书的第 2 版期间他还是我的第一位评论员和指导。非常感谢你，Marky-Mark。

最后，我必须赞扬和我共同执笔的作者。Hal 从我们多年前在 Sun 共事起就是我的同事和好朋友。我在第 1 版中就提到了这一点，现在更是如此。要不是 Hal，这本书只会停留在创意上，Hal 帮我把我的一个永远不可能实现的想法变成了一本真正的书，为此我向他表示我永恒的敬意和感激。

Hal Stern 致辞

如果互联网时代用类似狗年（译注：狗年（dog year）是指对狗来说一年的时间尺度。与人年（human year）相比，过去的标准是 1 人年等于 7 狗年。不过最新的研究认为，人年与狗年的比率可能还与狗的大小有关）那样的标准来衡量的话，那么本书第 1 版出版后的 4 年则代表了半个技术生命周期。我们看到了众多.com 公司的大起大落，我们看到了网络作为整个社会组织的一部分而兴起，无论是我们的孩子互发即时信息亦或我们一边通过无线网络阅读电子邮件一边喝咖啡。我们不再通过电话进行惩罚，在家中我们通过电子化的方式把孩子的行动进行限制（关掉孩子们的 DHCP 服务）。我们的孩子只想让这些东西

工作起来，我们这些这方面的人士有责任确定满足每个人对新社会粘合剂可靠性的需求。

由于信息技术的每个隐藏角落和缝隙都进行了联网，联网系统的可靠性也因此变得更复杂了。在本书的第 2 版中，我们设法在不同的逻辑层化解问题的复杂性，解决问题。虽然很多备受推崇的.com 公司没等欢呼声过去就挺不住了，但有些.com 公司活了下来，并成为真正的实时永远在线的企业：它们是 ebay.com、amazon.com、诸如 orbitz.com 的旅游网站以及实时体育广播网站（比如 mlb.com、Major League Baseball 的网上之家）。在过去的 4 年中我领悟到了这一点：在网络的另一端总会有一个活生生的生活在真实世界中的人存在。那个人忍受不了沙漏鼠标指针，费解的错误信息或不一致的行为。使系统具有高可用性所提出的挑战超越了防止宕机的要求，我们需要考虑预防用户体验的变化。

我要逐一表达我的感谢。在过去的 4 年中，我出色的妻子 Toby、女儿 Elana 和儿子 Benjamin 在容忍我的坏脾气和古怪的同时一直都支持着我。在两版书的写作期间，我转到 Sun 和 AOL-Netscape 联盟的软件公司，与我共事的那些杰出同事负责把软件栈的上层变得更可靠。Daryl Huff、“Big Hal” Jespersen、Sreeram Duvvuru 和 Matt Stevens 投入了很多时间解释状态复制方案和 Web 服务器的可靠性。Sun 首席技术官办公室的 Rick Lytel、Kenny Gross、Larry Votta 和 David Trindade 给我讲解了很多可用性工程方面的数学和科学知识。David 是那些令人吃惊的人之一，很少有人能使应用数学在真实世界中如此有趣。Larry 和 Kenny 正在探索关于软件可用性的新方法，Larry 正在把守旧的远程通信技术思想和 Web 服务结合起来并再次证明强大的基础设计原则历久不衰。

在机构的软件工作方面，我和职棒大联盟和全国曲棍球联合会在他们的网页方面合作，并且很愉快。MLB 高级媒体的首席技术官 Joe Choti 对数百万棒球迷同时传来的（电子）声音的缩放问题有很好的见解。NHL 的 IT 部门集团副总裁 Peter DelGiacco 也生活在实时世界，并且他对媒体、内容和正确性的观察很受赏识。很遗憾 George Spehar（我在写本书第一版时的良师和很多稿件灵感的来源）不幸死于癌症，我非常怀念他。

最后，在过去十年的大部分时间里，Evan Marcus 在技术联系及私下和我的关系方面都很密切。有了 Evan 对材料的组织，重组，修改材料和不知疲倦的热情，这本书的第 2 版才有可能得到迅速的出版。加拿大的电视名人 Scott Russell 曾说过，如果你“告诉我事实，我会忘记；告诉我真相，我学到了东西；告诉我一个故事，我便记住。”感谢你 Evan，感谢你抓住了技术的真相，并把它们编织成一个引人注目的技术故事。

第 1 版前言

技术类书籍包含的内容一般都很广泛，从带有精彩注释的代码清单到讲述晦涩协议非凡特性的枯燥的纯理论大部头，应有尽有。当我们决定写这本书的时候，我们要表达几乎 15 年来的两种体验，这是一个挑战。我们的书里面没有什么代码，这不是程序员手册也不是什么初级的入门书。可用性以及弹性和可预测性的更深层次的概念要求处理时讲求规程和程序。这本书代表了我们在规定如何制定这些规程，定义并细化这些程序，以及可信

地部署系统时所付出的最出色的劳动。在一天结束时，如果一个被设计为应具有高可用性的系统停止运行了，受到影响的是你的名声和你的工程技术。我们的目标是为你的技术提供真实世界的、可行的建议。当你看到书中的“实战传奇”时，你读到的是我们（稍微有些炫耀）所经历的非常好的或非常坏的设计经历。

我们力求在处理材料时寻求平衡。工程学总是在成本和功能之间、面世时间和功能之间、速度优化和安全设计之间寻求平衡。我们把可用性视为端到端网络计算的问题——在此情况下，可用性和性能同样重要。当你读这本书的时候，无论是按顺序读或依兴趣及问题挑着读，请记住你在寻找平衡。成本、复杂性和可用性水平都是你会遇到的问题，我们的任务是指导你在面对某一特定应用程序和环境时，决定各种方法的度。

我们还要感谢 John Wiley & Sons 的全体编辑团队。Carol Long 认为我们的想法可行，然后她指导，游说甚至以美味的午餐诱惑把我们的努力变成了你现在看到的。还要特别感谢 Christina Berry 和 Micheline Frederick，他们所做的编辑和写作工作及提出的建议提高了这本书整体的可读性，使它更流畅。你们是一流的团队，我们感谢你们在过去的 18 个月中对我们的支持。

Evan Marcus 致辞

这本书是 2 年多的准备和写作，7 年多高可用性系统（及人们认为是高可用性的系统）使用经验以及 15 年多的计算机系统综合使用经验的产物。从事了多年为高可用性领域的公司提供咨询和高可用软件产品的商家技术工作后，我发现我每次都在回答着同样的问题。问题都是有关关键系统如何获得最大可能的可用性。系统和其上运行的应用程序也许不同，但有关高可用性的问题都相同。我一直在寻找这方面的书，但一直没找到。

1992 年，我的工作开始与 Fusion System 的 High Availability for Sun 产品紧密相关了。该软件产品被认为是有史以来在 Sun 微系统工作站上运行的第一个高可用性或故障转移软件。软件允许某台事先指定的计算机在另一台计算机失效的时候迅速地、自动地进行干预并接手那台计算机的工作。从事了几年的通用系统的管理咨询工作后，我发现高可用性的概念很令人着迷。实际上，这个产品或工具做了优秀系统管理员工作，并把这项工作提到了另一个高度。优秀的系统管理员努力工作确保系统运行，提供所要求提供的服务，管理员也常常为他们的工作而感到骄傲。但尽管他们努力工作，系统仍旧会瘫痪，数据还是会丢失。这个产品提供的可用性是前所未有的。

High Availability for Sun 这款产品是个工具。和其他工具一样，依使用者知识和经验的不同它的表现可好可坏。我们曾设置过一些很不错的故障转移对。我们也设置过很差劲的配置。成功的案例是由经验丰富的和考虑周到的系统管理员做的，他们了解软件的目的，还明白这只是一个工具而不是什么万灵药。失败的案例是由于顾客没有做磁盘镜像，或把两套系统插到同一个电源上，或运行质量很差的应用程序，这些顾客期望 High Availability for Sun 自动解决系统的所有问题。

成功地实现 High Availability for Sun 的人们明白这个工具并不能代替他们来运行系统。他们知道仍需要大量的管理规程来保证系统以他们所期望的方式运行。他们明白该产品只解决了问题的一部分。

今天，即使那个曾经称为 High Availability for Sun 的产品的名字、公司和代码基址已至少改变了 3 次，有人明白容错软件能做什么不能做什么，而仍有另一些人把它看成是解决所有问题的最重要的因素。现今还有很多经验不丰富的系统管理员，他们并不熟悉所有和转出关键系统有关的问题。经理和制定预算的人认为可以免费获得高可用性系统，或者只需很少或不需要额外的工作。天下可没有免费的晚餐。

使系统具有高可用性（即使不借助故障转移软件）是涉及系统管理各个方面的技能。知道如何设置高可用性系统会使你成为一个更优秀的全面的系统管理员，并使你对你的雇主来说更有价值（即使你实际上从未有机会使用哪怕一个故障转移配置）。

我们希望在本书中指出在所配置过的成千个高可用性关键系统的配置过程中学到的东西。实际情况是我们不可能经历过读者在设置高可用性的过程中遇到的所有问题。然而，我们相信，我们一般的建议对于很多特定的情况也适用。

有些读者抱怨书中没有简单、通用的答案。这有两个原因：第一，每个地点和每个计算机系统出现的问题都不同，指望适合一个有 10 000 名员工的全球金融机构的方法也适合 10 人的律师事务所是不实际的，我们给出选择并让读者决定哪一种更适合他或她的环境；第二，我使用计算机工作了 15 年，我认识到大多数计算机的问题的答案令人遗憾：“视情况而定”。

我们假定读者的技术水平不同。除了极个别的情况，书中材料的技术性不是很强。我不是个喜欢卖弄比特字节的人（虽然 Hal 是那样的人），因此我设法给那些更像我的人写书。编写代码的部分偏重喜欢使用比特和字节的人，但这是例外而不是规则。

当我向朋友和同事描述这个项目时，他们的第一个问题通常是这是不是一本关于 Unix 或 NT 的书。说实话这两者都是。很显然，Hal 和我关于 Unix（特别是 Solaris）的经验都很丰富。但书中的技巧不都和特定的操作系统有关。这些技巧都很通用，其中许多也适用于计算机领域以外。为失效的单元保存备份的思想在下列场合也适用：航空、跳伞（那个讨厌的备用伞）以及其他由于失效可能导致致命、几乎致命或很危险的场合。毕竟，不带上备用胎你是不会开始长途旅行的，不是么？繁忙的交叉路口不会只有一个交通灯，如果灯泡坏了会发生什么？虽然我们的很多例子都以 Sun 和 Solaris 系统为例，但我们在尽可能的情况下都会使用 NT 或其他 Unix 系统下的例子。

整本书中，我们使用适合讨论的厂商的产品为例。我们并不认可它们，我们只是以它们的名字为例。

首先也是最重要的，我要感谢我的家庭。没有我的妻子 Carol 和女儿 Hannah 及 Madeline

的爱、支持和理解（或至少是容忍），我是不可能写出这本书的。我还要特别感谢我的家人和朋友，在我废寝忘食地在家里写作的时候他们假装表示理解。看，这确实是一本书！

还要感谢 Michael Kanaval（我们想你，Mike）提供的启发和一些精彩的例子；感谢 Joseph J. Hand 帮我处理一些 NT 材料；感谢 Michael Zona 和 John Costa 帮我处理备份工作；感谢 Mark Fitzpatrick 和 Bob Zarrow 以前和正在教给我的关于故障转移和通用高可用性的知识；感谢 Mindy Anderson 和 Eric Burgener 在群集和 SAN 方面的工作。也要感谢我的父母 Roberta 和 David Marcus 以及我的岳父母 Gladys 和 Herb Laden（我现在可以得到菜谱了么？）；感谢 Ed Applebaum、Ann Sheridan 和 Dinese Christopher 以及 VERITAS 和其他地方的每一个给这个项目提供建议和对于项目表现出热情和兴趣的每一个人。特别感谢 Mark Fannon 和 Evan Marks 提供的技术检查和一般性的帮助。

感谢数不清的顾客、用户和多年来与我共事的同事，特别感谢 Morgan Stanley Dean Witter、Bear Stearns、Deutsche Bank、J. P. Morgan、Sun Microsystems、VERITAS Software、Open Vision 和 Fusion Systems 的人们。

特别感谢 Hal Stern。1997 年中我终于下定决心写书。因为以前从没写过书，我知道我需要帮助。我给 Hal（他曾写过一本非常成功的书）发了电子邮件，向他寻求最初的指导。他回信问我们是不是有可能合作。我苦苦思考了很长时间（大约 2 纳秒）然后回复了一个充满热情的“是”。Hal 认为我们应在写作之前应做一个幻灯演示。原先的 250 个幻灯迅速增加到 400 个（我们仍然在每年的技术会议上作演示）。通过幻灯演示，我们可以确定缺少哪些内容，问题在哪里出现及内容如何跟随。把这些幻灯片变成你面前的书是相对（非常相对）简单的工作。Hal 起初还联系过 Wiley 的 Carol Long 并在我们的第一次技术会议上就开始了日程。要不是 Hal，这本书将只会是我脑子里的一个想法。

Hal Stern 致辞

我十年以前就接触可靠系统了，当时我在 Foxboro 公司工作，为公司把他们的实时、工业控制系统从专有硬件移植到 Sun 平台上。你从不会去考虑设备驱动器被悬挂或磁盘驱动器出现故障时产生的影响，除非设备驱动器将大型油漆搅拌桶的阀打开了或者磁盘驱动器沿着阿拉斯加石油管道位于几英尺的雪下。由于互联网的流行，可靠性和“正常运行时间工程学”成了我们的家常便饭，因为网上冲浪的人使我们把大多数问题都当做实时系统。作为系统管理员我们必须决定向可用性投入多少资金，努力达到 4 个 9（99.99%）或 5 个 9（99.999%）的正常运行时间，而资方却在谈论硬件变得多么的便宜。不存在正确的答案，每一件事情都是管理、运作、资金、策略、信任和时间的微妙平衡。需要由你决定可以接受的 9 的个数。我希望可以帮你获得足够的信息以便做出选择。

要是没有家庭的爱和支持是不可能有这本书的。我非常感谢我的妻子 Toby 和孩子 Elana 及 Benjamin，我要给他们一个热烈的拥抱，是的，爸爸现在可以不用再研究了。我还要感谢 Sun Microsystems 当前或以前的雇员，他们教给了我很多关于可用性的知识。我

还要感谢他们的点子和鼓励。他们是：Carol Wilhelmy, Jon Simms, Chris Drake, Larry McVoy, Brent Callaghan, Ed Graham, Jim Mauro, Enis Konuk, Peter Marcotte, Gayle Belli, Scott Oaks 和 Wendy Tal-mont。Pete Lega 参加了几次关于复杂性、恢复和自动化的马拉松式的会议，他的贡献也应得到尊重。Chris Kordish 和 Bob Sokol 都是 Sun Microsystems 公司的员工，他们审查了原稿并提出了评论和指导。Larry Bernstein (AT&T 退休的网络运行部的副总裁) 对我提出了质疑，让我多学学“电信运营商级”工程学，能和一位真正的电信先锋讨论问题是我的荣幸。Foxboro 公司的 Avi Nash 和 Randy Rohrbach 给我上了关于故障转移的第一课。Strike Technologies、Bear Stearns、Fidelity Investments、Deutsche Bank、Morgan Stanley Dean Witter 和 State Street Bank 中的很多人都认为书中的想法是可行的。即使保密条款禁止提到你的名字，我还要真诚地感谢你与我分享工程设计的机会。特别感谢 George Spehar (各个方面都是一位真正的绅士)，他提出了很多关于管理和经济决策的睿智的建议。Ed Braginsky (BEA 系统高级技术的副总裁) 是我八年多的好友，他还是优秀的工程师 (时间比 8 年要多)。他对排队系统，事务处理，异步设计及 BEA 共同创始人 Alfred Chuang 的理念的解释对我的帮助是不可估量的。当然，还要感谢爸爸、妈妈，他们教会了我可靠的重要性。

最后，非常非常感谢 Evan Marcus。我们在共同合作的一个客户项目中结识，那个项目要求在早晨很短的时间内找出性能问题。我之前从没见过 Evan，然而他开车带我在新泽西参观并且一路上说个不停。我那时就该意识到他有完成一本书的毅力，有说服我和他一起干的力量。感谢你 Evan，感谢你的耐心，理解以及促使我摆脱写作阻塞症、冬日里的忧郁和极度疲惫的独特才干。和你一起旅行，工作和教书是很愉快的事情。

关于作者

Evan Marcus 是 VERITAS Software 公司的总工程师和数据可用性专家。他从 1992 年起就参与高可用性系统的设计工作，当时他与其他人共同设计了第一个基于 Sun 的商业群集软件的关键部分。他在一家华尔街大型金融机构做过一段时间的股票交易厅系统管理员之后，又在 VERITAS Software 做了 4 年的销售工程师，咨询和撰写包括高可用性，集群和数据恢复等许多不同的问题。他为很多杂志和网站写过文章，包括最近的 TechTarget.com，他还是很多业界事件的很受尊敬的发言人。自从完成了本书第 1 版后，他还成为 *The Resilient Enterprise* (VERITAS 2002 年出版的关于数据恢复的书，这是 VERITAS 出版的第一部包括业界作者合作的书) 的编辑和特约编辑。Evan 拥有利哈伊大学 (Lehigh University) 计算机科学的学士学位，同时还是拉特斯哥大学 (Rutgers University) 的工商管理硕士。

Hal Stern 是 Sun Microsystems 的副总裁和杰出的工程师。他是 Sun Services 的首席技术官，负责高可靠系统和其上运行的联网程序的设计模式。在 Sun 工作的 10 多年的时间里，Hal 是 Sun ONE (iPlanet) 基础设施产品部门的首席技术官，还是 Sun 的美国东北部销售区域的技术专家。Hal 为众多大型金融机构和电子清算网络、两家大型专业运动联盟以及几家最大的电信设备和服务公司做过有关架构、性能和可靠性方面的工作。

Hal 为 *SunWorld* 杂志担任了 5 年的特约编辑，还是 IDG *JavaWorld* 杂志的编辑和编辑咨询委员会成员。在加入 Sun 之前，Hal 为波士顿地区一个刚建立的公司开发了分子建模软件，同时他还是普林斯顿大学 (Princeton University) 的研究人员。他拥有普里斯顿大学工程学的学士学位。不工作时，他会在少年棒球联盟担任教练，打冰球，为新泽西恶魔队加油，而且还拼命地打高尔夫减肥。

目 录

第1章 介绍	1	3.2.2 停机故障的间接损失	28
1.1 为什么需要一本可用性的书	2	3.3 可用性的价值	30
1.2 问题解决方法	2	3.3.1 例子1：双节点群集配置	33
1.3 不包括的内容	3	3.3.2 例子2：未知的停机损失	36
1.4 我们的任务	3	3.4 可用性变化区间	37
1.5 可用性指数	4	3.5 可用性指数图	39
1.6 总结	5	3.6 停机过程	40
1.7 本书的组织结构	5	3.6.1 停机	41
1.8 要点	6	3.6.2 数据丢失	42
第2章 测量数据	7	3.6.3 降级模式	43
2.1 测量可用性	7	3.6.4 预定停机	44
2.1.1 “9”表示法	9	3.7 要点	46
2.1.2 定义停机故障	11		
2.1.3 引起停机故障的原因	11	第4章 可用性政治策略	47
2.1.4 可用性	12	4.1 开始游说	47
2.1.5 平均数	14	4.1.1 从内部着手	47
2.1.6 可接受性	15	4.1.2 然后走出去	48
2.2 故障模式	16	4.1.3 开始行动	50
2.2.1 硬件	16	4.2 你的听众	53
2.2.2 环境和物理故障	17	4.2.1 获得听众	53
2.2.3 网络故障	18	4.2.2 了解听众	53
2.2.4 文件和打印服务器故障	18	4.3 表达信息	53
2.2.5 数据库系统故障	19	4.3.1 幻灯演示	54
2.2.6 网页和应用程序服务器故障	20	4.3.2 报告	54
2.2.7 拒绝服务攻击	21	4.4 传递信息之后	55
2.3 对测量的信心	22	4.5 要点	55
2.3.1 可恢复性	22		
2.3.2 Sigma (σ) 和“9”表示法	23	第5章 20条关键的高可用性设计原则	57
2.4 要点	24	5.1 #20：切勿贪便宜	57
第3章 可用性的价值	25	5.2 #19：不要想当然	58
3.1 高可用性的含义	25	5.3 #18：消除单点故障	59
3.2 停机故障损失	27	5.4 #17：执行安全	60
3.2.1 停机故障直接损失	27	5.5 #16：加强服务器的性能	61
		5.6 #15：留意速度	62
		5.7 #14：实施更改控制	63

5.8 #13: 时时备案	64	恢复所需要的磁盘空间	108
5.9 #12: 采用服务级协议	65	6.9 总结	108
5.10 #11: 超前策划	66	6.10 要点	109
5.11 #10: 尽量多试验	67	第7章 高度可用的数据管理	110
5.12 #9: 隔离你的环境	68	7.1 四个基本原理	111
5.13 #8: 以史为鉴	69	7.1.1 磁盘发生故障的可能性	111
5.14 #7: 设计要留有余地	70	7.1.2 磁带盘上的数据	111
5.15 #6: 选择成熟的软件	70	7.1.3 保护数据	112
5.16 #5: 选择成熟可靠的硬件	72	7.1.4 确保数据的可达性	112
5.17 #4: 重新使用配置	73	7.2 数据存储和管理的六个独立层次	112
5.18 #3: 利用外部资源	74	7.3 磁盘硬件与连通性术语	113
5.19 #2: 一步一个脚印	75	7.3.1 SCSI	113
5.20 #1: 尽量简单化	76	7.3.2 光纤通道	115
5.21 要点	78	7.3.3 多路径	116
第6章 备份与恢复	79	7.3.4 多主机	117
6.1 备份的基本规则	79	7.3.5 磁盘阵列	117
6.2 备份能否真正提供高可用性	81	7.3.6 热交换	117
6.3 需要对什么进行备份	81	7.3.7 逻辑设备(LUN)和卷	118
6.3.1 对备份进行备份	82	7.3.8 JBOD(就是一组磁盘)	118
6.3.2 获得异地备份	82	7.3.9 热备件	118
6.4 备份软件	83	7.3.10 写入高速缓存	118
6.4.1 商业软件还是自主研发	83	7.3.11 存储区域网络 (SAN)	118
6.4.2 商业备份软件实例	83	7.4 RAID 技术	120
6.4.3 商业备份软件的特性	84	7.4.1 RAID 的级别	121
6.5 备份性能	86	7.4.2 其他种类的 RAID	128
6.5.1 提高备份性能: 找出瓶颈	86	7.5 磁盘空间和文件系统	133
6.5.2 解决性能问题	90	7.5.1 大磁盘还是小磁盘	134
6.6 备份类型	93	7.5.2 当 LUN 填满时会出现什么 情况	135
6.6.1 增量备份	93	7.5.3 管理磁盘和卷的可用性	136
6.6.2 数据库增量备份	95	7.5.4 文件系统的恢复	137
6.6.3 缩短备份窗口	96	7.6 要点	137
6.6.4 热备份	96		
6.6.5 数据越少, 越省时间 (和空间)	97		
6.6.6 使用更多的硬件	99		
6.6.7 复杂的软件特征	101		
6.7 处理备份磁带和数据	104		
常规备份安全	106		
6.8 恢复	107		
		第8章 存储区域网络、网络连接存储与存储虚 拟化	139
		8.1 存储区域网络	139
		8.1.1 选用 SAN 的理由	141
		8.1.2 SAN 硬件设备简介	143
		8.2 网络连接存储	144

8.3 SAN 与 NAS 比较	145	10.6 要点	198
8.4 存储虚拟化	149	第 11 章 人与程序	199
8.4.1 选择存储虚拟化的理由	150	11.1 系统管理与修正	199
8.4.2 存储虚拟化的类型	150	11.1.1 维护计划与步骤	200
8.5 要点	153	11.1.2 系统修正	201
第 9 章 组网	154	11.1.3 备用设备方针	203
9.1 网络故障分类	155	11.1.4 预防性维护	204
9.1.1 网络可靠性挑战	155	11.2 供应商管理	204
9.1.2 网络故障模式	156	11.2.1 选择关键的供应商	205
9.1.3 物理设备故障	157	11.2.2 与供应商合作	207
9.1.4 IP 层故障	158	11.2.3 在系统恢复中供应商的角色	208
9.1.5 拥塞引起的故障	160	11.3 安全性	209
9.2 构建冗余网络	162	11.3.1 数据中心的安全	211
9.2.1 虚拟 IP 地址	163	11.3.2 病毒与蠕虫	211
9.2.2 冗余网络连接	164	11.4 文档	212
9.2.3 多重网络的配置	167	11.4.1 文档的使用者	213
9.2.4 IP 路由冗余	170	11.4.2 文档与安全	214
9.2.5 网络恢复模式选择	172	11.4.3 检查文档	214
9.3 负载平衡和网络重定向	173	11.5 系统管理员	215
9.3.1 循环 DNS	173	11.6 内部扩增	217
9.3.2 网络重定向	174	故障标识	219
9.4 动态 IP 地址	176	11.7 要点	219
9.5 网络服务可靠性	176	第 12 章 客户端与用户	220
9.5.1 网络服务依赖性	177	12.1 强化企业客户端	220
9.5.2 强化核心服务	179	12.1.1 客户端备份	221
9.5.3 拒绝服务攻击	180	12.1.2 客户端补给	222
9.6 要点	182	12.1.3 瘦客户端	223
第 10 章 数据中心和本地环境	183	12.2 容许数据服务故障	224
10.1 数据中心	183	12.2.1 文件服务器客户端恢复	224
10.1.1 数据中心机架	185	12.2.2 数据库应用程序恢复(Database Application Recovery)	226
10.1.2 平衡安全性和可访问性	187	12.2.3 Web 客户端恢复(Web Client Recovery)	227
10.1.3 数据中心观光	188	12.3 要点	229
10.1.4 异地主机设施	189	第 13 章 应用程序设计	230
10.2 电	191	13.1 应用程序恢复概览	231
UPS	191	13.1.1 应用程序的故障模式	231
10.3 线缆铺设	193		
10.4 冷却及环境问题	195		
10.5 系统命名惯例	196		

13.1.2 应用程序恢复技术	232	多个实例对比更大的实例	267
13.1.3 更软性的故障	234	14.4 基于 Web 的服务可靠性	268
13.2 从系统故障中进行应用程序恢复	234	14.4.1 Web 服务器群集	268
13.2.1 虚拟内存耗尽	235	14.4.2 应用服务器	270
13.2.2 I/O 错误	236	14.4.3 目录服务器	272
13.2.3 数据库应用程序的重新 连接	236	14.4.4 Web 服务标准	273
13.2.4 网路连通性	237	14.5 要点	274
13.2.5 重启网络服务	238	第 15 章 本地群集和故障转移	276
13.2.6 网络拥塞、重发和超时 设定	239	15.1 群集技术简介	277
13.3 内部应用程序故障	241	15.2 服务器故障和故障转移	279
13.3.1 内存访问错误	241	15.3 逻辑性的以应用为中心的思想	281
13.3.2 内存滥用和恢复	242	15.4 故障转移的要求	282
13.3.3 挂起进程	243	15.4.1 服务器	284
13.4 开发人员“卫生学”	243	15.4.2 服务器间的差异	284
13.4.1 返回值检查	244	15.4.3 网络	286
13.4.2 边界条件检查	245	15.4.4 磁盘	292
13.4.3 基于值的安全	246	15.4.5 应用程序	295
13.4.4 日志支持	247	15.5 大型群集	295
13.5 进程复制	248	15.6 要点	296
13.5.1 冗余服务进程	249	第 16 章 故障转移管理和难题	297
13.5.2 进程状态多路广播	250	16.1 故障转移管理软件	297
13.5.3 检查点技术	251	16.2 部件监控	298
13.6 不做假设，管理一切	252	16.2.1 实施检测的人和关于其他 部件监测的问题	299
13.7 要点	253	16.2.2 当部件检测失败时	300
第 14 章 数据和 Web 服务	254	16.3 进行手工故障转移的时机	301
14.1 网络文件系统服务	254	16.4 自主开发的故障转移软件 还是商业软件	303
14.1.1 检测 RFC 故障	255	16.5 商业故障转移管理软件	304
14.1.2 NFS 服务器的约束	256	16.6 当好的故障转移软件出错时	305
14.1.3 文件锁定	258	16.6.1 脑裂综合症	305
14.1.4 失效文件句柄	260	16.6.2 不受欢迎的故障转移	309
14.2 数据库服务器	261	16.7 验证和检测	310
14.2.1 管理恢复时间	262	16.7.1 状态转换图	310
14.2.2 破坏之中求生存	264	16.7.2 测试作品	312
14.2.3 任何（高）速度下的不安 全状态	264	16.8 管理故障转移	313
14.3 冗余和可用性	266	16.8.1 系统监测	313
		16.8.2 控制台	313

16.8.3 工具	314	第 20 章 灾难恢复计划	361
16.8.4 时间问题	315	20.1 DR 计划的是与非	362
16.9 其他群集话题	315	20.2 DR 计划的 3 个主要目标	362
16.9.1 复制数据群集	315	20.2.1 员工的健康与保护	362
16.9.2 群集之间的距离	317	20.2.2 企业的存活	363
16.9.3 负载均衡群集和故障 转移	317	20.2.3 企业的连续性	363
16.10 要点	318	20.3 良好的 DR 计划	363
第 17 章 故障转移结构	319	20.4 准备构建 DR 计划	364
17.1 双节点故障转移结构	319	20.5 选择 DR 现场	368
17.1.1 “主-从” 故障转移	319	20.5.1 实际位置	368
17.1.2 “主-主” 故障转移	324	20.5.2 DR 现场安全	371
17.1.3 “主-主” 还是 “主-从”	325	20.5.3 停留在 DR 现场的时间	372
17.2 服务组故障转移	326	20.6 分发 DR 计划	372
17.3 更大型的群集系统结构	328	20.6.1 DR 计划内容	372
17.3.1 N 对 1 群集系统	328	20.6.2 分发措施	373
17.3.2 N 加 1 群集系统	329	20.7 计划受众	374
17.4 群集系统的规模应该有多大	331	20.8 时间线	375
17.5 要点	332	20.9 灾难恢复小组任务指派	376
第 18 章 数据复制	333	20.9.1 指派人员	376
18.1 复制概述	333	20.9.2 管理层的角色	377
18.2 进行复制的原因	334	20.10 DR 计划的多与寡	378
18.3 复制类型	334	20.11 共用 DR 现场	379
18.3.1 四类按延迟时间划分的复 制类型	334	20.12 装备 DR 现场	380
18.3.2 五种按启动程序划分的 复制类型	338	20.13 DR 计划的测试	381
18.4 有关复制的其他思想	351	20.13.1 高质量演习的特性	382
18.4.1 SAN: 复制的另一种方式	351	20.13.2 演习计划	383
18.4.2 多个目的地系统	352	20.13.3 演习之后	387
18.4.3 远程应用程序故障转移	354	20.14 三种演习类型	387
18.5 要点	354	20.14.1 全面演练	387
第 19 章 虚拟机和资源管理	355	20.14.2 桌上演练	388
19.1 分区和域: 系统级的 VM	356	20.14.3 电话链演练	388
19.2 容器: 操作系统级的 VM	357	20.15 灾难对人员的影响	389
19.3 资源管理	358	20.15.1 对灾难的典型反应	389
19.4 要点	360	20.15.2 企业应采取的措施	390
		20.16 要点	391
第 21 章 弹性企业	392		
21.1 纽约期货交易所	392		
21.1.1 第一次灾难的发生	394		
21.1.2 大型交易所决不该是 这样的	395		