

世界著名计算机教材精选

PEARSON
Prentice
Hall

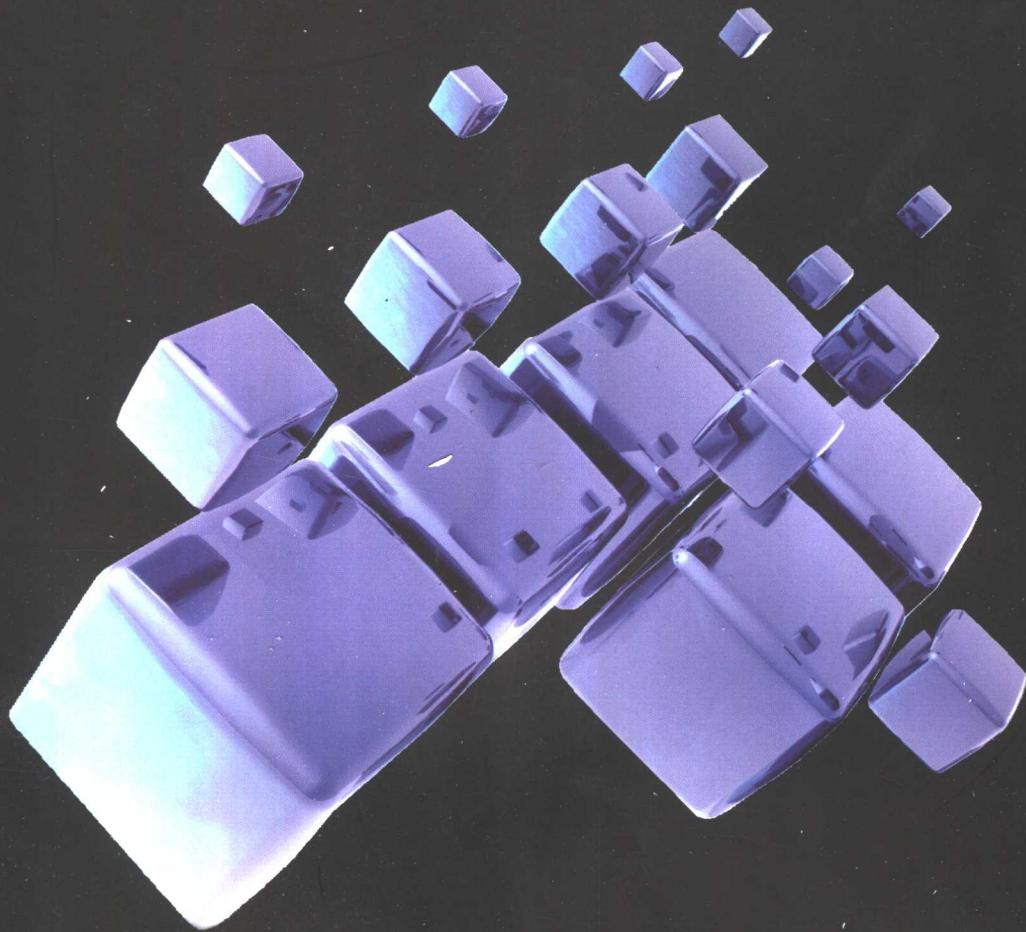
分布式系统

原理与范型

Andrew S. Tanenbaum
Maarten van Steen

杨剑峰 常晓波 李敏

著
译



DISTRIBUTED SYSTEMS

Principles and Paradigms



清华大学出版社

世界著名计算机教材精选

分布式系统原理与范型

Andrew S. Tanenbaum, Maarten van Steen 著

杨剑峰 常晓波 李 敏 译

清华大学出版社
北京

Simplified Chinese edition copyright © 2004 by PEARSON EDUCATION ASIA LIMITED and TSINGHUA UNIVERSITY PRESS.

Original English language title from Proprietor's edition of the Work.

Original English language title: Distributed Systems: Principles and Paradigms, by Andrew S. Tanenbaum and Maarten van Steen. Copyright © 2002

EISBN: 0-13-088893-1

All Rights Reserved.

Published by arrangement with the original publisher, Pearson Education, Inc., publishing as Prentice Hall.

This edition is authorized for sale only in the People's Republic of China (excluding the Special Administrative Region of Hong Kong and Macao).

本书中文简体翻译版由 Prentice Hall 授权给清华大学出版社在中国境内(不包括中国香港、澳门特别行政区)出版发行。

北京市版权局著作权合同登记号 图字: 01-2003-0557

版权所有,翻印必究。举报电话: 010-62782989 13901104297 13801310933

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。

图书在版编目(CIP)数据

分布式系统原理与范型/特南鲍姆(Tanenbaum, A. S.),范施特恩(van Steen, M.)著;杨剑峰等译. —北京: 清华大学出版社, 2004. 9

(世界著名计算机教材精选)

书名原文: Distributed Systems: Principles and Paradigms

ISBN 7-302-08961-2

I. 分… II. ①特… ②范… ③杨… III. 分布式操作系统—教材 IV. TP316.4

中国版本图书馆 CIP 数据核字(2004)第 063337 号

出版者: 清华大学出版社

<http://www.tup.com.cn>

社总机: 010-62770175

地址: 北京清华大学学研大厦

邮 编: 100084

客户服务: 010-62776969

责任编辑: 袁勤勇

印 刷 者: 北京四季青印刷厂

装 订 者: 北京市密云县京文制本装订厂

发 行 者: 新华书店总店北京发行所

开 本: 185×260 印张: 39.75 字数: 915 千字

版 次: 2004 年 9 月第 1 版 2004 年 9 月第 1 次印刷

书 号: ISBN 7-302-08961-2/TP · 6340

印 数: 1~4000

定 价: 68.00 元

本书如存在文字不清、漏印以及缺页、倒页、脱页等印装质量问题,请与清华大学出版社出版部联系调换。联系电话: (010)62770175-3103 或 (010)62795704

译者序

随着计算机网络,特别是 Internet 的迅猛发展,传统的信息系统概念发生了巨大的变化,这些变化突出地表现在信息的存储、传递、发布以及获取方式所发生的革命性变革。与此同时,基于网络的分布式信息系统在各个领域得到了广泛的应用,在整个社会生活中正发挥着日益突出的作用。Internet 已经越来越多地成为构建信息系统的一个关键组成部分。如何在更为广域和异构的计算环境中有效地发布和获取信息,已成为亟待解决的问题。分布式系统正是解决了上述问题。现在分布式系统的研究、应用日益广泛深入,分布式系统的学习也成为计算机及相关专业必不可少的教学环节。

本书是 Tanenbaum 先生在所著的《分布式操作系统》的基础上,总结了分布式系统方面的最新进展,重新撰写的力作,是分布式系统的权威教材。本书循序渐进地、全面地、深入地讲解了分布式系统的原理,并列出了大量的范型。本书的结构分为两部分:原理和范型。第一部分(第 1~8 章)详细讨论了分布式系统的原理、概念和技术,其中包括通信、进程、命名、同步、一致性和复制、容错以及安全。第二部分(第 9~12 章)给出了一些实际的分布式系统,即基于对象的分布式系统、分布式文件系统、基于文档的分布式系统以及基于协作的分布式系统,介绍了一些实际系统的设计思想和实现技术。全书结构清晰,内容全面经典,系统性与先进性并茂。

本书的目标读者是计算机及相关专业的高年级学生或研究生。从事分布式计算研究和工程应用的科研人员和工程技术人员也会从本书中受益匪浅。

本书是多人共同努力的成果,参与本书翻译、审稿、录排的人员包括:杨剑峰、常晓波、梁金昆、张丽萍、汪青青、朱志博、李敏、李静、李娟、张颖、朱剑平、刘颖、吴东升、杨战伟、郭宁宁、李楠、聂晶、刘恒、刘敏、刘洋、吕喜熹、马睿倩等。全书由杨剑峰、常晓波和李敏负责统稿。

限于译者水平,难免有错误和疏漏之处,恳请读者不吝指正。希望这本书能成为您工作的好帮手。

杨剑峰 常晓波
2004 年 5 月

前　　言

本书的出发点是对 Distributed Operating Systems 一书进行再版修订,但笔者很快就发现自 1995 年以来很多技术发生了改变,要完全体现出这些变化,仅仅对该书进行修订是不够的,而是需要写一本全新的书。因此,这本新书有了一个新的标题:《分布式系统原理和范型》。标题的改变体现了对重点的调整。虽然我们仍然讨论一些操作系统的问题,但现在这本书还从更广泛的意义上研究分布式系统。例如,WWW 作为已建立的最大的分布式系统,在 Distributed Operating Systems 一书中完全没有提到,因为它并不是一个操作系统。而在本书中,它几乎占去整整一章。

本书分为两部分:原理和范型。第 1 章是对主题的总体介绍。接下来的第 2~8 章分别讨论我们认为最重要的原理:通信、进程、命名、同步、一致性和复制、容错以及安全性。

实际的分布式系统通常围绕一些范型来组织的,例如“所有事物都是文件”。接下来的第 9~12 章分别介绍一个不同的范型,并描述使用该范型的一些重要系统。涉及到的范型包括基于对象的系统、分布式文件系统、基于文档的系统以及基于协作的系统。

第 13 章包含一份附有说明的参考书目,可供该主题的进一步学习使用,还包含本书中引用的著作列表。

本书是作为计算机科学的大学高年级学生或研究生课程而编写的。因此,本书有一个 Web 站点,站点中以各种格式放置了本书中用到的 PowerPoint 表和图。要访问该站点,在 <http://www.prenhall.com/tanenbaum> 页面上点击本书标题即可。将本书作为教材使用的教授可以通过联系当地的 Prentice Hall 代理机构得到一本习题解答手册。当然,本书也十分适合希望更多地了解这一重要主题的社会人士。

许多人以多种方式对本书作出了贡献。我们尤其要感谢 Arno Bakker、Gerco Ballintijn、Brent Callaghan、Scott Cannon、Sandra Cornelissen、Mike Dahlin、Mark Darbyshire、Guy Eddon、Amr el Abbadi、Vincent Freeh、Chandana Gamage、Ben Gras、Bob Gray、Michael van Hartskamp、Philip Homburg、Andrew Kitchen、Ladislav Kohout、Bob Kutter、Jussipekka Leiwo、Leah McTaggart、Eli Messenger、Donald Miller、Shivakant Mishra、Jim Mooney、Matt Mutka、Rob Pike、Krithi Ramamirtham、Shmuel Rotenstreich、Sol Shatz、Gurdip Singh、Aditya Shivram、Vladimir Sukonnik、Boleslaw Szymanski、Laurent Therond 和 Leendert van Doorn,感谢他们阅读了部分书稿并提出了宝贵意见。

最后,我们还要感谢我们的家庭。Suzanne 已经经历过很多次这样的情况了。她从未说过“我受够了”,尽管这个念头肯定在她脑海里出现过。谢谢你!Barbara 和 Marvin 现在对教授们为谋生所做的工作有了更好的了解,并且认识到好教材和坏教材之间的差别。现在,对我来说他们是我努力创作出更多好教材的动力所在。

本书使用指南

我们使用本书中的材料已经很多年了,主要是用作大学高年级学生和研究生的教材。而且,这些材料还曾经作为为时1~2天的有关分布式系统和中间件的研讨会的基本资料,参加这些研讨会的人包括ICT专家(技术上的)。下面是我们根据经验对本书使用方式提出的一些建议。

大学高年级学生和研究生教材

如果作为大学高年级学生和研究生的教材,本书通常可以在12~15周内完成教学。我们发现,在大多数学生看来,分布式系统由很多似乎彼此紧密结合的主题所组成。在本书的组织上,我们按照不同的原理介绍这些主题,分别讲授各个原理,这对学生领会重点内容有很大帮助。这样安排的效果是当第一部分(第1~8章)结束时,即在讨论范型之前,学生已经对本书主题在整体上有了一个相当好的把握。

然而,分布式系统的领域涵盖许多不同的主题,其中一些主题在初次学习时很难理解。因此,我们强烈建议学生们随着课程的进展学习适当的章节。从Web站点(<http://www.prenhall.com/tanenbaum>)可以获得所有PowerPoint表,将它们预先分发下去,以便学生在课堂中能够积极参与讨论。这种方法非常成功,并得到了学生们的高度评价。

所有的材料都包括在一个为时15周的课程中。大多数时间花费在讲授分布式系统的原理,也就是前8章所包括的材料上。在讨论范型时,我们的经验是:只需要介绍要点。直接从书中学习每个案例的详细内容比在课堂上听授更加容易。例如,尽管书中有基于对象的系统的内容达80页之多,但我们只用一周的时间讲授这类系统。下面是一个课程进度安排建议表(表0.1),其中包括每次讲座中包括的主题。

表0.1 课程进度安排建议

周	主题	章	讲授内容
1	绪论	1	全部
2	通信	2	2.1~2.3
3	通信	2	2.4~2.5
4	进程	3	全部
5	命名	4	4.1~4.2
6	命名	4	4.3
6	同步	5	5.1~5.2

续表

周	主题	章	讲授内容
7	同步	5	5.3~5.6
8	一致性和复制	6	6.1~6.4
9	一致性和复制	6	6.5~6.6
9	容错	7	7.1~7.3
10	容错	7	7.4~7.6
11	安全性	8	8.1~8.2
12	安全性	8	8.3~8.7
13	基于对象的系统	9	全部
14	文件系统	10	全部
15	基于文档的系统	11	全部
15	基于协作的系统	12	全部

并不是所有材料都需要在课堂上讲授;我们希望学生能够自学特定的部分,尤其是细节部分。在讲授时间少于 15 周的情况下,我们建议跳过有关范型的章节,让感兴趣的学生自己学习这些部分。

如果用于低年级的课程,我们推荐将本书的学习延长至两个学期,并增加实验作业。例如,可以通过让学生修改一些组件,使这些组件具有容错性、处理多播 RPC 等功能来使学生理解简单的分布式系统。

行业的专业研讨会

在 1~2 天的研讨会上,通常将本书作为主要的背景材料使用。然而,如果跳过所有细节,仅将重点放在分布式系统的本质上,则有可能在两天内讲完整本书。此外,要使内容的表达更加生动实用,有必要重新安排章节的顺序,以提早说明原理是如何得到应用的。对于研究生来说,一般是在了解原理的应用之前(有时甚至根本不了解原理的具体应用)先对原理进行为期 10 周的学习,但专业人士如果能了解这些原理的实际应用,就会有更大的学习动力。下面是一个为期 2 天课程的试验性进度表(表 0.2),该表按照逻辑单元进行划分。

表 0.2 按逻辑单元划分的课程进度

第 1 天				
单元	时间(分)	主 题	章	重 点
1	90	绪论	1	客户-服务器体系结构
2	60	通信	2	RPC/RMI 和消息传递
3	60	基于协作的系统	12	消息传递问题
4	60	进程	3	移动代码和代理
5	30	命名	4	位置跟踪
6	90	基于对象的系统	9	CORBA

续表

第 2 天				
单元	时间(分)	主 题	章	重 点
1	90	一致性和复制	6	模型和协议
2	60	基于文档的系统	11	Web 缓存/复制
3	60	容错	7	进程组与 2PC
4	90	安全性	8	基本思想
5	60	分布式文件系统	10	NFS v3 和 v4

个人学习

本书同样也适用于个人学习。如果具有足够的时间和动力,建议读者仔细阅读整本书。

如果没有足够的时间仔细阅读所有材料,我们建议只集中学习最重要的主题。下面的表格中列举一些章节,我们认为这些章节涵盖了关于分布式系统的最重要的主题(表 0.3)。

表 0.3 自学内容

章	主 题	小 节
1	绪论	1.1、1.2、1.4.3、1.5
2	通信	2.2、2.3、2.4
3	进程	3.3、3.4、3.5
4	命名	4.1、4.2
5	同步	5.2、5.3、5.6
6	一致性和复制	6.1、6.2.2、6.2.5、6.4、6.5
7	容错	7.1、7.2.1、7.2.2、7.3、7.4.1、7.4.3、7.5.1
8	安全性	8.1、8.2.1、8.2.2、8.3、8.4
9	基于对象的系统	9.1、9.2、9.4
10	分布式文件系统	10.1、10.4
11	基于文档的系统	11.1
12	基于协作的系统	12.1、12.2 或 12.3

比较好的做法是对学习这些建议的材料需要花费的时间进行估算,但这在很大程度上取决于读者的背景知识,对各种背景的读者很难做一个一般性的估计。然而,如果一个具有全职工作的人抽出晚上的时间阅读本书,则可能至少花费几周时间。

目 录

第 1 章 绪论		1
1. 1	分布式系统的定义	1
1. 2	目标	3
1. 2. 1	让用户连接到资源	3
1. 2. 2	透明性	4
1. 2. 3	开放性	6
1. 2. 4	可扩展性	7
1. 3	分布式系统的硬件	12
1. 3. 1	多处理器系统	13
1. 3. 2	同构式多计算机系统	15
1. 3. 3	异构式多计算机系统	16
1. 4	分布式系统的软件	17
1. 4. 1	分布式操作系统	18
1. 4. 2	网络操作系统	26
1. 4. 3	中间件	28
1. 5	客户-服务器模型	33
1. 5. 1	客户与服务器	33
1. 5. 2	应用程序的分层	38
1. 5. 3	客户-服务器体系结构	40
1. 6	小结	43
	习题	43
第 2 章 通信	45	
2. 1	分层协议	45
2. 1. 1	低层协议	48
2. 1. 2	传输协议	50
2. 1. 3	高层协议	52
2. 2	远程过程调用	54
2. 2. 1	基本的 RPC 操作	55
2. 2. 2	参数传递	58
2. 2. 3	扩展的 RPC 模型	61
2. 2. 4	实例: DCE RPC	64
2. 3	远程对象调用	68
2. 3. 1	分布式对象	68

2.3.2 将客户绑定到对象	70
2.3.3 静态远程方法调用与动态远程方法调用	72
2.3.4 参数传递	73
2.3.5 实例 1: DCE 远程对象	74
2.3.6 实例 2: Java RMI	76
2.4 面向消息的通信	79
2.4.1 通信中的持久性和同步性	79
2.4.2 面向消息的暂时通信	83
2.4.3 面向消息的持久通信	86
2.4.4 示例: IBM MQSeries	91
2.5 面向流的通信	95
2.5.1 为连续媒体提供支持	95
2.5.2 流与服务质量	98
2.5.3 流同步	101
2.6 小结	103
习题	104
 第 3 章 进程	107
3.1 线程	107
3.1.1 线程简介	107
3.1.2 分布式系统中的线程	112
3.2 客户	114
3.2.1 用户界面	114
3.2.2 客户端软件与分布透明性	116
3.3 服务器	117
3.3.1 设计上常见的 important 问题	117
3.3.2 对象服务器	120
3.4 代码迁移	125
3.4.1 代码迁移方案	125
3.4.2 迁移与本地资源	128
3.4.3 异构系统中的代码迁移	131
3.4.4 实例: D'Agents	132
3.5 软件代理	136
3.5.1 分布式系统中的软件代理	136
3.5.2 代理技术	138
3.6 小结	140
习题	141
 第 4 章 命名	144
4.1 实体的命名	144
4.1.1 名称、标识符和地址	144

4.1.2	名称解析	148
4.1.3	名称空间的实现	152
4.1.4	示例：域名系统	158
4.1.5	示例：X.500	161
4.2	移动实体的定位	165
4.2.1	实体命名与定位	165
4.2.2	简单方法	167
4.2.3	基于起始位置的方法	169
4.2.4	分层方法	171
4.3	删除无引用的实体	176
4.3.1	无引用对象的问题	177
4.3.2	引用计数	178
4.3.3	引用列表	181
4.3.4	标识不可到达实体	182
4.4	小结	187
	习题	188
第 5 章	同步	190
5.1	时钟同步	190
5.1.1	物理时钟	191
5.1.2	时钟同步算法	194
5.1.3	使用同步时钟	197
5.2	逻辑时钟	198
5.2.1	Lamport 时间戳	199
5.2.2	向量时间戳	201
5.3	全局状态	203
5.4	选举算法	206
5.4.1	欺负(Bully)算法	206
5.4.2	环算法	207
5.5	互斥	208
5.5.1	集中式算法	208
5.5.2	分布式算法	209
5.5.3	令牌环算法	211
5.5.4	三个算法的比较	212
5.6	分布式事务	213
5.6.1	事务模型	213
5.6.2	事务的分类	216
5.6.3	实现	218
5.6.4	并发控制	220
5.7	小结	226
	习题	227

第6章 一致性和复制	229
6.1 简介	229
6.1.1 复制的目的	230
6.1.2 对象复制	230
6.1.3 作为扩展技术的复制	232
6.2 以数据为中心的一致性模型	233
6.2.1 严格一致性	234
6.2.2 线性化和顺序一致性	236
6.2.3 因果一致性	239
6.2.4 FIFO一致性	240
6.2.5 弱一致性	242
6.2.6 释放一致性	244
6.2.7 入口一致性	245
6.2.8 一致性模型小结	247
6.3 以客户为中心的一致性模型	248
6.3.1 最终一致性	249
6.3.2 单调读	250
6.3.3 单调写	251
6.3.4 写后读	252
6.3.5 读后写	253
6.3.6 实现	254
6.4 分发协议	256
6.4.1 副本放置	256
6.4.2 更新传播	259
6.4.3 epidemic 协议	262
6.5 一致性协议	264
6.5.1 基于主备份的协议	264
6.5.2 复制的写协议	267
6.5.3 高速缓存相关性协议	270
6.6 实例	271
6.6.1 Orca	272
6.6.2 因果一致的懒惰复制	276
6.7 小结	279
习题	280
第7章 容错性	283
7.1 容错性简介	283
7.1.1 基本概念	283
7.1.2 典型故障	285
7.1.3 使用冗余来掩盖故障	287
7.2 进程恢复	288

7.2.1	设计问题	288
7.2.2	故障掩盖和复制	290
7.2.3	故障系统的协议	290
7.3	可靠的客户-服务器通信	293
7.3.1	点到点通信	293
7.3.2	出现失败时的 RPC 语义	293
7.4	可靠的组通信	298
7.4.1	基本的可靠多播方法	298
7.4.2	可靠多播中的可扩展性	299
7.4.3	原子多播	301
7.5	分布式提交	307
7.5.1	两阶段提交	307
7.5.2	三阶段提交	312
7.6	恢复	313
7.6.1	简介	314
7.6.2	检查点	316
7.6.3	消息日志	318
7.7	小结	320
	习题	321

第 8 章	安全性	323
8.1	安全性介绍	323
8.1.1	安全威胁、策略和机制	323
8.1.2	设计问题	328
8.1.3	加密	331
8.2	安全通道	337
8.2.1	身份验证	338
8.2.2	消息完整性和机密性	344
8.2.3	安全组通信	346
8.3	访问控制	349
8.3.1	访问控制中的一般问题	349
8.3.2	防火墙	352
8.3.3	保护移动代码	354
8.4	安全管理	359
8.4.1	密钥管理	359
8.4.2	安全组管理	363
8.4.3	授权管理	364
8.5	实例: KERBEROS	368
8.6	实例: SESAME	370
8.6.1	SESAME 组件	370
8.6.2	PAC	372
8.7	实例: 电子付费系统	373

8.7.1 电子付费系统	373
8.7.2 电子付费系统中的安全性	375
8.7.3 协议实例	377
8.8 小结	381
习题	382
第9章 基于对象的分布式系统	384
9.1 CORBA	384
9.1.1 CORBA 概述	385
9.1.2 通信	390
9.1.3 进程	395
9.1.4 命名	399
9.1.5 同步	402
9.1.6 缓存与复制	403
9.1.7 容错性	404
9.1.8 安全性	406
9.2 分布式组件对象模型(DCOM)	408
9.2.1 DCOM 概述	408
9.2.2 通信	413
9.2.3 进程	415
9.2.4 命名	417
9.2.5 同步	420
9.2.6 复制	420
9.2.7 容错性	420
9.2.8 安全性	421
9.3 Globe	423
9.3.1 Globe 概述	423
9.3.2 通信	430
9.3.3 进程	430
9.3.4 命名	432
9.3.5 同步	435
9.3.6 复制	435
9.3.7 容错性	437
9.3.8 安全性	438
9.4 CORBA、DCOM 和 Globe 的比较	439
9.4.1 基本原理	439
9.4.2 通信	440
9.4.3 进程	441
9.4.4 命名	441
9.4.5 同步	442
9.4.6 缓存与复制	442
9.4.7 容错性	442

9.4.8 安全性	442
9.5 小结	444
习题	444
第 10 章 分布式文件系统	446
10.1 SUN 网络文件系统	446
10.1.1 NFS 概述	447
10.1.2 通信	450
10.1.3 进程	451
10.1.4 命名	452
10.1.5 同步	458
10.1.6 缓存和复制	462
10.1.7 容错性	464
10.1.8 安全性	466
10.2 Coda 文件系统	469
10.2.1 Coda 概述	469
10.2.2 通信	471
10.2.3 进程	472
10.2.4 命名	473
10.2.5 同步	474
10.2.6 缓存和复制	477
10.2.7 容错性	480
10.2.8 安全性	482
10.3 其他分布式文件系统	484
10.3.1 Plan 9: 资源统一为文件	485
10.3.2 xFS: 无服务器的文件系统	489
10.3.3 SFS: 可扩展的安全性	494
10.4 分布式文件系统的比较	496
10.4.1 设计理念	497
10.4.2 通信	497
10.4.3 进程	497
10.4.4 命名	498
10.4.5 同步	499
10.4.6 缓存和复制	499
10.4.7 容错性	499
10.4.8 安全性	500
10.5 小结	501
习题	501
第 11 章 基于文档的分布式系统	503
11.1 WWW	503

11.1.1	WWW 概述	504
11.1.2	通信	511
11.1.3	进程	515
11.1.4	命名	520
11.1.5	同步	522
11.1.6	缓存和复制	522
11.1.7	容错性	526
11.1.8	安全性	526
11.2	Lotus Notes	527
11.2.1	Lotus Notes 概述	527
11.2.2	通信	529
11.2.3	进程	530
11.2.4	命名	531
11.2.5	同步	533
11.2.6	复制	533
11.2.7	容错性	535
11.2.8	安全性	535
11.3	WWW 和 Lotus Notes 的比较	538
11.4	小结	542
	习题	542
第 12 章	基于协作的分布式系统	544
12.1	协作模型介绍	544
12.2	TIB/Rendezvous	546
12.2.1	TIB/Rendezvous 概述	546
12.2.2	通信	548
12.2.3	进程	551
12.2.4	命名	551
12.2.5	同步	553
12.2.6	缓存和复制	554
12.2.7	容错性	554
12.2.8	安全性	556
12.3	Jini	557
12.3.1	Jini 概述	558
12.3.2	通信	560
12.3.3	进程	561
12.3.4	命名	563
12.3.5	同步	565
12.3.6	缓存和复制	567
12.3.7	容错性	567
12.3.8	安全性	567
12.4	TIB/Rendezvous 和 Jini 的比较	568

12.5 小结	571
习题	571

第 13 章 阅读材料和参考书目 573

13.1 对进一步阅读的建议	573
13.1.1 介绍性和综述性的著作	573
13.1.2 通信	574
13.1.3 进程	575
13.1.4 命名	576
13.1.5 同步	576
13.1.6 一致性与复制	577
13.1.7 容错性	578
13.1.8 安全性	579
13.1.9 面向对象的分布式系统	580
13.1.10 分布式文件系统	581
13.1.11 基于文档的分布式系统	582
13.1.12 基于协作的分布式系统	583
13.2 参考书目列表	583