

中国计算机学会学术著作丛书

对等网络：

结构、应用与设计

Peer-to-Peer Network

Structure, Application and Design

陈贵海 李振华 著
倪明选 刘云浩 校

清华大学出版社
北京

内 容 简 介

对等网络是分布式系统和计算机网络相结合的产物,它在应用领域和学术界获得了广泛的重视和成功,被称为“改变 Internet 的新一代网络技术。”本书由浅入深,全面系统地介绍了对等网络的各个方面内容,重点在于对三代 P2P 网络、P2P 网络结构和设计机制的讲解。首先介绍对等网络的概念、历史、特点,重点介绍三代 P2P 网络,同时讲述各种 P2P 应用系统和软件;然后深入讨论 P2P 网络设计所要考虑的核心机制、优化网络性能的增强机制和 P2P 模拟器的设计;最后分析了 P2P 的现状和发展趋势。

本书是有关对等网络的学术专著,包含对等网络产生以来在学术界、商业领域有价值、有影响力的成果,内容编写力求前沿、具体、实用。本书既可作为大学计算机、电子信息等专业教材,也可作为研究所、商业公司中对等网络研究者的重要参考书。

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。

版权所有,侵权必究。侵权举报电话:010-62782989 13501256678 13801310933

图书在版编目(CIP)数据

责任编辑:

责任校对:

责任印制:

出版发行:清华大学出版社

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

c-service@tup.tsinghua.edu.cn

社总机:010-62770175

投稿咨询:010-62772015

地 址:北京清华大学学研大厦 A 座

邮 编:100084

邮购热线:010-62786544

客户服务:010-62776969

印刷者:

装订者:

经 销:全国新华书店

开 本: 印张: 字数: 千字

版 次:2007 年 月第 2 版 印次:2007 年 月第 1 次印刷

印 数:1~ 000

定 价: .00 元

本书如存在文字不清、漏印、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:010-62770177 转 3103 产品编号:

中国计算机学会学术著作丛书

- | 名誉主任委员：张效祥
- | 主任委员：唐泽圣
- | 副主任委员：陆汝铃
- | 委员：(以姓氏笔画为序)
王 珊 李晓明 吕 建
林惠民 罗军舟 郑纬民
施伯乐 焦金生 谭铁牛





序

Preface

第

一台电子计算机诞生于 20 世纪 40 年代,到目前为止,计算机的发展已远远超出了其创始者的想象。计算机的处理能力越来越强,应用面越来越广,应用领域也从单纯的科学计算渗透到社会生活的方方面面:从工业、国防、医疗、教育、娱乐直至人们的日常生活,计算机的影响可谓无处不在。

计算机之所以能取得上述地位并成为全球最具活力的产业,原因在于其高速的计算能力、庞大的存储能力以及友好灵活的用户界面。而这些新技术及其应用有赖于研究人员多年不懈的努力。学术研究是应用研究的基础,也是技术发展的动力。

自 1992 年起,清华大学出版社与广西科学技术出版社为促进我国计算机科学技术与产业的发展,推动计算机科技著作的出版,设立了“计算机学术著作出版基金”,并将资助出版的著作列为中国计算机学会的学术著作丛书。时至今日,本套丛书已出版学术专著近 50 种,产生了很好的社会影响,有的专著具有很高的学术水平,有的则奠定了一类学术研究的基础。中国计算机学会一直将学术著作的出版作为学会的一项主要工作。本届理事会将秉承这一传统,继续大力支持本套丛书的

出版,鼓励科技工作者写出更多的优秀学术著作,多出好书,多出精品,为提高我国的知识创新和技术创新能力,促进计算机科学技术的发展和进步做出更大的贡献。

中国计算机学会

2002年6月14日



前 言

Foreword

对

等网络(Peer-to-Peer Network,简称 P2P 网络)是分布式系统和计算机网络相结合的产物,它在网络协议的应用层,打破过去的“客户/服务器”模式,让所有网络成员享有“自由、平等、互联”的功能,不再有客户、服务器之分,任何两个网络结点之间都能共享文件、传递消息。对等网络起源于 1999 年风行一时的音乐文件共享软件 Napster,随后则是一系列人们耳熟能详的网络软件:Gnutella*,KaZaA,BitTorrent,eDonkey/eMule,Skype,等等。虽然从产生到现在只有短短几年的历史,但是对等网络在应用领域和学术界获得了广泛的重视和成功,并占据了当前 Internet 超过一半的带宽资源,被称为“改变 Internet 的新一代网络技术”。

在应用领域,1999 年出现的第一个 P2P 网络 Napster,在美国 6 个月内即拥有 5000 万注册用户,成为网络时代的一个奇迹。此后许多计算机公司开始投资 P2P 网络的开发,如著名的无结构 P2P 网络 Gnutella、KaZaA 和 eDonkey/eMule。2002 年,家喻户晓的 BitTorrent 网络出现,其简称“BT”已成为“自由下载、文件

* 注意 Gnutella 的发音为[nju:ˈtelə],同 new-tella,字母 G 不发音。

共享”的代名词。BT 在中国的应用尤为广泛，2006 年发布的中国互联网统计报告显示：中国 1.11 亿网民中有 27.8% 的人使用过 BT 软件（超过 3000 万人），其流行程度由此可见一斑。除了文件共享，P2P 软件在其他多个应用领域也带来了冲击性的变革：网络语音电话软件 Skype 直接威胁到电话公司的利益，网络电视软件如 PPLive、TvAnts 等使得人们能通过 Internet 流畅地观看电视直播，诸如此类的应用还有很多。

在学术界，虽然 P2P 的思想起源很早并一直有人提及，但直到 2001 年“P2P 网络”这一概念才获得人们的共识并成为研究热点。计算机领域许多重要会议（如 SIGCOMM, SPAA, PODC, ICNP, INFOCOM, ICDCS 等）和刊物（如 IEEE/ACM Transactions on Networking, IEEE Transactions on Parallel and Distributed Systems, IEEE Journal on Selected Areas in Communications 等）从 2001 年开始发布和刊登 P2P 领域的论文，同年出现了 P2P 领域最著名的结构化网络模型：Chord, CAN, Tapestry 和 Pastry。P2P 领域自身的两个专业会议 IEEE P2P 和 IPTPS 分别于 2001 年、2002 年举行。人们提出了许多新颖、独特的 P2P 网络模型，开发了许多用以实现、完善 P2P 网络的核心机制及增强机制，其中不少理论成果已经转化为 Internet 上的实用软件。

到目前为止，P2P 还是一项比较新的技术，其中还存在很多问题，同时也存在着各种挑战和机遇。虽然有太多的问题和障碍，但 P2P 从无到有，从一个民间小软件发展到改变 Internet 面貌、改变人们交流方式的一项新技术，这一过程是值得人们回味的，其中蕴藏的力量和意义，值得我们付出勇气、热情、耐心去努力研究并发掘。

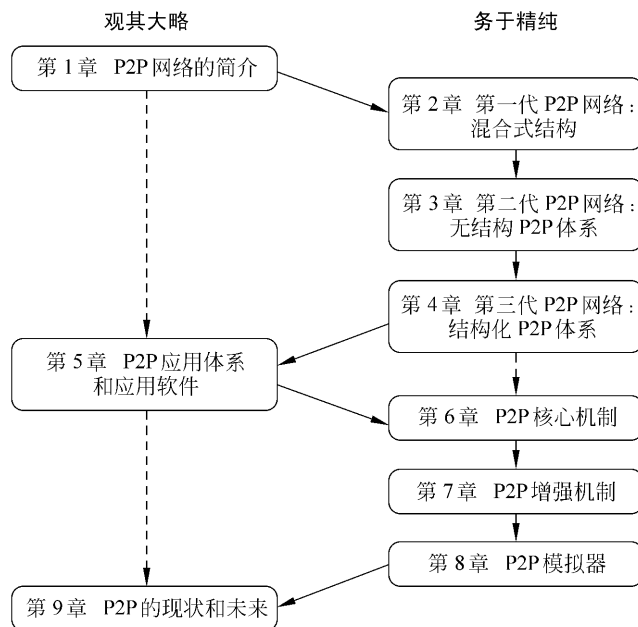
阅读本书的方法

追比圣贤，原是读书人的本意。《三国演义》用这样一段话来描述诸葛亮读书：“孔明与博陵崔州平、颍川石广元、汝南孟公威与徐元直四人为密友。此四人务于精纯，惟孔明独观其大略。尝抱膝长吟，而指四人曰：‘公等仕进可至刺史、郡守。’众问孔明之志若何，孔明但笑而不答。每常自比管仲、乐毅，其才不可量也。”

我们以为这段话，正说出了读书的要义。做学问不必执着于一事一物，对于大多数事情，只要能“观其大略”，明白是怎么回事就足够了；生命有限，青春短暂，我们要做的，是关注那些真正对自己有价值的东西，将精力集中于真正值得去研究的问题，努力“务于精纯”。

这也正是我们希望的读者阅读本书的方法。我们真心希望您付出时间之

后能有所得，所以建议读者按照如下线索来阅读本书：



给读者的鼓励

第一次读这本书时你或许不能完全理解，但不用担心，实际上作者在撰写本书时很多相关论文第一次也没看懂，比如 Freenet 的论文就至少看了 5 遍！读者完全可以跳过暂时看不懂的内容，因为我们在写作时已努力使各章各节尽量独立，当你往后看得更多时，前面感到困惑的内容很可能就豁然开朗了。

“事理因人言而悟者，有悟还有迷，总不如自悟之了了。”在阅读本书的时候，千万不要忘记您自己的批判和思考。无论是批评之言，还是建议之语，希不吝赐教。我们的信箱是：gchen@nju.edu.cn 或 lizhenhua@dislab.nju.edu.cn。

或许世上只有两种人能成为作家：痛苦的或快乐的，因为只有不一样的心情才能承受写作的漫漫孤寂。在作者心中，只能靠想象去构思他的读者空间，所以我们希望能看到、听到您任何形式的反馈，感受到您阅读时的欢欣、困惑与体会。只有这样，我们所构思的空间才能变得真实、饱满，让我们觉得写作时所有的单调、枯燥与艰辛，都是值得付出的。

愿本书带给您的，不仅仅是技术和知识！

参考文献来源

本书所引用的参考文献绝大多数来自于网络通信领域的重要国际会议和刊物,文献发表时间集中于 2001—2007 年之间,以保证内容的权威性和时效性;除此之外,我们也参考了大量 Web 资源,尤其是 P2P 相关网站上的资料。我们所关注的会议、刊物列举如下:

缩写名称	英文全称	中文名称
SIGCOMM	ACM Annual Conference of the Special Interest Group on Data Communication	ACM 数据通信特别兴趣组年会
SPAA	ACM Symposium on Parallel Algorithms and Architectures	ACM 并行算法和体系讨论会
PODC	ACM Symposium on Principles of Distributed Computing	ACM 分布式计算原理讨论会
ICNP	International Conference on Network Protocols	国际网络协议会议
INFOCOM	IEEE Conference on Computer Communications	IEEE 计算机通信会议
NOSSDAV	ACM International Workshop on Network and Operating System Support for Digital Audio and Video	ACM 数字音频和视频的网络和操作系统支持讨论会
SOSP	ACM Symposium on Operating Systems Principles	ACM 操作系统原理讨论会
OSDI	USENIX Symposium on Operating Systems Design and Implementation	USENIX 操作系统设计和实现讨论会
USITS	USENIX Symposium on Internet Technologies and Systems	USENIX 因特网技术和系统讨论会
ASPLOS	International Conference on Architectural Support for Programming Languages and Operating Systems	国际程序设计语言和操作系统的体系支持会议
ICDCS	International Conference on Distributed Computing Systems	国际分布式计算系统会议
ICPP	International Conference on Parallel Processing	国际并行处理会议
IPDPS	International Parallel and Distributed Processing Symposium	国际并行和分布式处理讨论会
IPTPS	International Workshop on Peer-to-Peer Systems	国际 P2P 系统讨论会
IEEE P2P	IEEE International Conference on Peer-to-Peer Computing	IEEE 国际 P2P 计算会议
刊物名一般不缩写	IEEE/ACM Transactions on Networking	IEEE/ACM 网络学报
	IEEE Transactions on Parallel and Distributed Systems	IEEE 并行和分布式系统学报
	IEEE Journal on Selected Areas in Communications	IEEE 通信精选领域杂志

致谢

本书的写作、出版得到两大项目的支持：一是 2005 年国家自然科学基金项目“实用化对等网络技术的研究”；二是 2005 年江苏省自然科学基金前期预研项目“新型 P2P 计算技术的基础研究”。香港科技大学倪明选教授和刘云浩教授对本书的写作给予了有益的指导，华中科技大学金海教授和国防科技大学肖侬教授通读了本书的初稿，并提出了许多有价值的建议。邱彤庆、谢军峰、袁瑞峰、叶懋、刘之育、陈欢、于南南等同学对部分章节提出了修改建议。在此谨表示衷心的感谢！

陈贵海 李振华

2007 年 3 月

于南京大学鼓楼校区



目 录

Contents

前言	I
第 1 章 P2P 网络简介	1
1.1 什么是 P2P 网络?	1
1.2 P2P 网络的发展历程	6
1.3 为什么需要 P2P 网络	11
1.4 P2P 网络本质的特点	14
1.5 P2P 网络的各种应用	19
第 2 章 第一代 P2P 网络：混合式 P2P 体系	26
2.1 Napster——P2P 网络的先驱	26
2.1.1 Napster 出现的背景和它创造的 奇迹	26
2.1.2 Napster 网络的工作原理	27
2.1.3 Napster 的性能分析	28
2.1.4 Napster 的陨落和它的现状	29
2.1.5 Napster 的缺陷和新的混合式 P2P 网络	32
2.2 BitTorrent——分片优化的新一代混合式 P2P 网络	33
2.2.1 BitTorrent 的曲折历史	33
2.2.2 BitTorrent 体系原理	35

2.2.3	BitTorrent 分片机制	37
2.2.4	BitTorrent 阻塞算法	39
2.2.5	BitTorrent 性能分析	40
2.2.6	BitTorrent 体系总结	42
2.2.7	关于 BT 的一个重要事实澄清：BT 伤硬盘吗？	43
2.3	第一代 P2P 网络的特点	44
2.3.1	拓扑结构	44
2.3.2	查询与路由	44
2.3.3	容错、自适应和匿名性	44
2.3.4	增强机制	44
第 3 章	第二代 P2P 网络：无结构 P2P 体系	45
3.1	Gnutella——纯分布式无结构 P2P 网络	45
3.1.1	Gnutella 出现的背景	45
3.1.2	Gnutella 体系的工作原理	46
3.1.3	Gnutella 网络的性能分析	49
3.1.4	Napster 与 Gnutella 的比较	51
3.1.5	Gnutella 协议 0.6 版	52
3.2	KaZaA——基于超结点的无结构 P2P 网络	53
3.2.1	KaZaA 和 FastTrack 介绍	53
3.2.2	KaZaA 的工作原理	54
3.2.3	KaZaA 协议语法和语义	55
3.2.4	KaZaA 技术细节	56
3.2.5	KaZaA 性能分析	58
3.2.6	KaZaA 网络总结	58
3.3	eDonkey/eMule——分块下载的双层无结构 P2P 网络	59
3.3.1	eDonkey、eMule 和 Overnet 介绍	59
3.3.2	eDonkey 工作原理	60
3.3.3	eDonkey 文件分块	61
3.3.4	eDonkey 性能分析	62
3.3.5	eDonkey 网络总结	63
3.4	Freenet——自由、安全、匿名的无结构 P2P 网络	63
3.4.1	Freenet 出现的背景和发展历史	63
3.4.2	Freenet 的密码学基础	65
3.4.3	Freenet 中数据的查询与获取	68

3.4.4	Freenet 中数据的存储与管理	69
3.4.5	Freenet 网络新结点加入	70
3.4.6	Freenet 协议细节	71
3.4.7	Freenet 性能分析	72
3.4.8	Freenet 的安全性和匿名性分析	74
3.4.9	Freenet 体系总结	75
3.5	无结构 P2P 网络的特点	76
3.5.1	覆盖网拓扑结构	76
3.5.2	路由和定位方法	77
3.5.3	容错性与自适应	80
3.5.4	可扩展性	80
3.5.5	安全性与匿名性	80
3.5.6	增强机制——复制	80
3.5.7	优势和缺陷	81
第 4 章	第三代 P2P 网络：结构化 P2P 体系	82
4.1	Chord 与 CFS——简单、精确的环形 P2P 网络	83
4.1.1	Chord 介绍	83
4.1.2	Chord 基础工作原理	84
4.1.3	Chord 对象定位算法	87
4.1.4	Chord 结点加入算法	88
4.1.5	Chord 自适应算法	90
4.1.6	Chord 容错性和复制、缓存	91
4.1.7	Chord 实验分析	91
4.1.8	Chord 总结	93
4.1.9	CFS 介绍	94
4.1.10	CFS 文件系统结构	95
4.1.11	CFS 对 Chord 的改进：前驱列表定位、服务器选择和 结点 ID 认证	96
4.1.12	CFS 中的复制、缓存和负载均衡	97
4.1.13	CFS 总结	99
4.2	CAN——简单、容错的多维空间 P2P 网络	99
4.2.1	CAN 介绍	99
4.2.2	CAN 网络构建	100
4.2.3	CAN 增强机制：多维、多空间、多散列	102

4.2.4	CAN 的“区域超载”	103
4.2.5	CAN 中的复制与缓存	104
4.2.6	CAN 总结	104
4.3	Tapestry 与 OceanStore——广域的超立方体结构 P2P 网络	104
4.3.1	Tapestry 简介	104
4.3.2	Tapestry 路由和定位	106
4.3.3	Tapestry 动态结点算法	108
4.3.4	Tapestry 体系架构	110
4.3.5	Tapestry 总结	111
4.3.6	OceanStore 简介	111
4.3.7	OceanStore 的命名机制和存取控制	113
4.3.8	OceanStore 的路由和定位算法	116
4.3.9	OceanStore 的更新模型	117
4.3.10	OceanStore 的深度归档存储	119
4.3.11	OceanStore 的内省优化	119
4.3.12	OceanStore 总结	121
4.4	Pastry 与 PAST——容错的混合式超立方体结构 P2P 网络	122
4.4.1	Pastry 简介	122
4.4.2	Pastry 路由	124
4.4.3	Pastry 自组织和自适应	126
4.4.4	Pastry 的局部性	128
4.4.5	Pastry 实验分析	129
4.4.6	Pastry 总结	130
4.4.7	PAST 简介	131
4.4.8	PAST 操作	131
4.4.9	PAST 安全机制	132
4.4.10	PAST 存储管理	133
4.4.11	PAST 副本转移和文件转移	133
4.4.12	PAST 缓存管理	134
4.4.13	PAST 总结	134
4.5	其他著名结构化 P2P 网络——Kademlia、SkipNet 等	135
4.5.1	其他著名结构化 P2P 网络简介	135
4.5.2	Kademlia——基于异或度量的 P2P 信息系统	136
4.5.3	SkipNet——基于跳表、提供显式局部性的 P2P 模型	139
4.6	常数度 P2P 模型——Viceroy、Koorde 和 Cycloid 等	147

4.6.1	常数度 P2P 模型概要	147
4.6.2	Viceroy——基于蝴蝶结构的常数度 P2P 模型	148
4.6.3	Koorde——整合 Chord、de Bruijn 图的常数度 P2P 模型	151
4.6.4	Cycloid——基于 CCC 的常数度 P2P 模型	154
4.7	结构化 P2P 网络的特点与分析	159
4.7.1	覆盖网拓扑结构	159
4.7.2	分布式散列表(DHT)	160
4.7.3	路由和定位	160
4.7.4	动态结点算法(自组织、自适应)	163
4.7.5	容错性与安全性	163
4.7.6	局部性	164
4.7.7	增强机制:复制、缓存和分片	165
4.7.8	P2P 网络各项属性总结	165
第 5 章	P2P 应用体系和应用软件	167
5.1	P2P 应用清单	167
5.2	文件共享	174
5.2.1	BitTorrent 的使用	174
5.2.2	eDonkey 的使用	180
5.2.3	百宝——优秀的国产 P2P 音乐共享软件	184
5.2.4	Maze 文件共享系统	185
5.2.5	国产 P2P 文件共享软件评点	186
5.3	多媒体传输	188
5.3.1	Skype——优秀的网络语音传输工具	188
5.3.2	PPLive——不错的国产 P2P 网络电视软件	190
5.3.3	TvAnts——支持内网的 P2P 电视蚂蚁	191
5.3.4	AnySee 视频直播系统	193
5.4	实时通信和协同工作	193
5.4.1	P2P 实时通信软件评点	193
5.4.2	Groove 虚拟办公室——优秀的 P2P 协同工作空间	195
5.5	分布式数据存取	196
5.5.1	Granary 广域存储服务系统	197
5.6	分布式计算	198
5.6.1	GPU——Gnutella 全球处理单元	198
5.6.2	SETI@Home——分布式计算缘起	199

5.7	P2P 搜索引擎	201
5.8	其他应用介绍	201
5.8.1	TinyP2P——15 行代码的 P2P 软件	202
5.8.2	JXTA——开放式 P2P 开发平台	203
第 6 章	P2P 核心机制	207
6.1	覆盖网拓扑结构	209
6.2	分布式散列表	211
6.2.1	散列函数	213
6.2.2	安全散列函数	213
6.2.3	一致性散列函数	214
6.3	路由和定位	215
6.3.1	混合式 P2P 网络的路由和定位方法	215
6.3.2	无结构 P2P 网络的路由和定位方法	216
6.3.3	结构化 P2P 网络的路由和定位方法	217
6.3.4	P2P 网络定位至少需要多少跳?	217
6.3.5	结点度和网络直径的折中关系对路由算法的影响	218
6.4	查询和搜索	219
6.4.1	路由索引	220
6.4.2	基于 DHT 的 P2P 网络复合查询	222
6.4.3	前缀散列树	224
6.5	动态结点算法	228
6.5.1	混合式 P2P 网络的动态结点算法	228
6.5.2	无结构 P2P 网络的动态结点算法	228
6.5.3	结构化 P2P 网络的动态结点算法	230
6.6	容错性	231
6.6.1	为了保证容错, 结点度至少需要多少?	232
6.6.2	容错的传统方法——冗余	233
6.6.3	容错性的分类	233
6.6.4	网络动态性的分类	234
6.6.5	容错性的重要参数——崩溃点	235
6.6.6	P2P 覆盖网分割问题	238
第 7 章	P2P 增强机制	242
7.1	P2P 系统的性能	245

7.2	复制与缓存	246
7.2.1	关于复制的一个简单而有用的结论	246
7.2.2	无结构 P2P 网络复制的理论模型	247
7.2.3	结构化 P2P 网络中的复制和缓存	249
7.3	分片	251
7.3.1	实用的传统分片技术	251
7.3.2	结构化 P2P 网络中的层次化数据分块	252
7.3.3	冗余编码的数学原理	252
7.3.4	冗余编码的优点	255
7.3.5	冗余编码 vs. 复制	256
7.4	负载均衡、异构性与热点问题	258
7.4.1	负载均衡	259
7.4.2	异构性	259
7.4.3	Gia——异构性自适应的无结构 P2P 网络模型	260
7.4.4	热点问题	262
7.4.5	有用的引申——用 P2P 技术来解决 C/S 热点问题	263
7.5	拓扑意识和一致性问题	266
7.5.1	处理一致性问题的三种传统方法	266
7.5.2	希尔伯特编码——邻近 ID 选择	267
7.5.3	CFS 的下一跳选择——邻近路由选择	268
7.5.4	Pastry 的路由表构造——邻近邻居选择	269
7.5.5	LTM——动态邻近邻居选择	270
7.6	匿名、声誉和信任	272
7.6.1	匿名的各种方法	272
7.6.2	Tarzan——P2P 匿名网络层	273
7.6.3	P2P 声誉、信任涉及的问题和解决方法	275
7.6.4	EigenTrust 算法——完备的 P2P 声誉管理	277
7.7	P2P 安全问题	280
7.7.1	P2P 网络中的攻击方式和安全对策	280
7.7.2	P2P 面临的非技术性问题	283
第 8 章	P2P 模拟与仿真	285
8.1	P2P 模拟器的设计意义和准则	286
8.1.1	P2P 模拟器的设计意义	286
8.1.2	P2P 模拟器的设计准则	287

8.2	经典的网络模拟器与拓扑产生器	288
8.2.1	经典的网络模拟器 NS-2	288
8.2.2	Transit-Stub 模型与 GT-ITM 拓扑产生器	289
8.2.3	通用拓扑产生器 BRITE	291
8.3	P2P 模拟器	292
8.3.1	通用 P2P 模拟器 p2psim	292
8.3.2	三层 P2P 模拟器 3LS	297
8.3.3	数据包层 Gnutella 模拟器 GnutellaSim	298
8.3.4	PeerSim 模拟器简介	300
8.4	全球网络服务仿真平台 PlanetLab	300
第 9 章	P2P 的现状和未来	302
9.1	P2P 的主要研究组织	303
9.1.1	世界计算机领域最有影响力的组织	303
9.1.2	研究 P2P 的著名高校	305
9.1.3	研究 P2P 的著名公司	307
9.1.4	研究 P2P 的其他组织	308
9.2	P2P 的重要国际刊物和会议	308
9.2.1	P2P 的重要国际会议	309
9.2.2	P2P 的重要国际刊物	311
9.3	P2P 的主要商业模式	311
9.4	P2P 与其他领域的融合	315
9.4.1	P2P 与无线网络的融合	315
9.4.2	P2P 与网格计算的融合	318
9.5	P2P 的未来	319
9.5.1	P2P 未来的商业应用趋势	319
9.5.2	P2P 未来的学术研究趋势	320
9.5.3	结束语：探讨 P2P 的未来	320
	参考文献	322
	索引	330