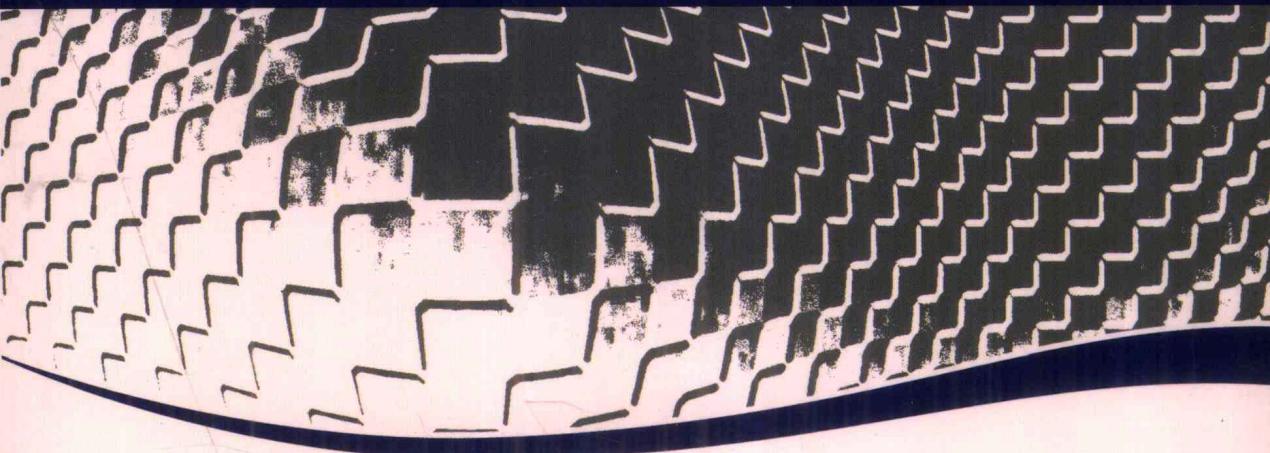


可扩展计算学术著作丛书



陈汉华 金海 著

# 大规模分布式内容检索技术



科学出版社

可扩展计算学术著作丛书

# 大规模分布式内容检索技术

陈汉华 金海著

科学出版社

北京

## 内 容 简 介

大规模分布式内容检索是近年来分布式系统方向的一个热点研究领域。本书全面地阐述了各种体系结构的分布式大规模内容检索系统的关键技术和核心理论，并对各项技术和理论的来龙去脉进行了详细深入的分析。本书通过丰富的文献资料和研究成果，从研究者的视角对大规模分布式内容检索技术进行了深入剖析，是分布式处理系统领域的学术专著。

本书可供高等院校计算机科学与技术相关专业的高年级本科生、研究生、教师、研究人员及工程技术人员阅读参考，也可作为相关专业的研究生教材。

### 图书在版编目(CIP)数据

大规模分布式内容检索技术/陈汉华,金海著. —北京:科学出版社,2011  
(可扩展计算学术著作丛书)

ISBN 978-7-03-031417-8

I. 大… II. ①陈…②金… III. 机器检索:情报检索 IV. G354.4

中国版本图书馆 CIP 数据核字(2011)第 106910 号

责任编辑: 魏英杰 杨向萍 / 责任校对: 林青梅

责任印制: 赵 博 / 封面设计: 鑫联必升

科学出版社出版

北京东黄城根北街 16 号

邮政编码: 100717

<http://www.sciencep.com>

新蕾印刷厂 印刷

科学出版社发行 各地新华书店经销

\*

2011 年 5 月第 一 版 开本: B5(720×1000)

2011 年 5 月第一次印刷 印张: 19 1/2

印数: 1—3 000 字数: 372 000

**定价: 70.00 元**

(如有印装质量问题, 我社负责调换)

# 《可扩展计算学术著作丛书》编委会

名誉主编：陈国良

主编：金海

编委：（按姓氏汉语拼音排序）

安 虹(中国科技大学)	戴国忠(中科院软件研究所)
董守斌(华南理工大学)	董小社(西安交通大学)
过敏意(上海交通大学)	何炎祥(武汉大学)
胡 煊(兰州大学)	金 海(华中科技大学)
李明禄(上海交通大学)	李晓明(北京大学)
刘 克(国家自然科学基金委员会)	刘云浩(清华大学)
罗军舟(东南大学)	钱德沛(北京航空航天大学)
孙凝辉(中科院计算技术研究所)	王兴伟(东北大学)
王志英(国防科技大学)	魏英杰(科学出版社)
吴朝晖(浙江大学)	肖 依(国防科技大学)
谢向辉(江南计算技术研究所)	徐志伟(中科院计算技术研究所)
杨天若(华中科技大学)	臧斌宇(复旦大学)
郑庆华(西安交通大学)	郑纬民(清华大学)
周兴社(西北工业大学)	

## 序

从万维网的辉煌到对等系统的风靡,再到以在线社会网络应用为成功代表的云应用的崛起,互联网用户之间的信息共享不断向更便捷、更高效的形式发展。近二十年来,分布式计算系统的体系结构分分合合,网络应用层新概念、新事物层出不穷,技术飞速进步归根结底归功于摩尔定律所起的作用。在互联网兴起的初期,大多数连接到互联网上的普通用户由于计算机性能、资源等因素的限制,并没有能力对外提供服务。因而,互联网逐步形成了以少数服务器为中心的状况——客户机/服务器(C/S)模式成为Web和FTP等互联网应用的主要架构。时至今日,即使是廉价的普通机器也具备了相当的计算、存储和通信能力,甚至超过了老一代曾经风云一时的超级计算机。对等模型顺应了技术的发展而迅速兴起,并为开发大规模分布式内容共享系统奠定了核心基础。

分布式内容检索是构建大规模分布式应用系统的关键支撑技术。由于大规模分布式系统中往往不适合或难于搭建中心索引服务,因此大规模分布式内容检索的核心问题是如何建立高效的分布式索引,并基于分布式索引开发能透明支持复杂查询处理的协议和算法。《大规模分布式内容检索技术》一书围绕此核心问题展开,通过结合对等分布式模型的概念、结构、资源描述和组织、资源定位和路由选择、结果融合及排序方法等,从结构化对等网络、非结构化对等网络和混合式对等网络各自的特点出发,全面系统地阐述了各种环境下进行分布式大规模内容检索的策略。在每一类方法的讨论中,作者系统地对各项关键技术进行了深入的介绍。同时,对各项关键技术的来龙去脉进行了详细地剖析,为读者还原了各项创新性技术从思想萌芽到成熟应用的技术纵深。

该书内容结合了作者学术研究和系统开发的实践,介绍的许多技术源于作者在香港科技大学计算机科学与工程系访问期间的科研成果。

该书内容丰富、鞭辟入里,是分布式处理系统领域一本难得的学术专著,故向读者郑重推荐。

倪明选

香港科技大学讲座教授

2011年4月于香港九龙清水湾

## 前　　言

随着网络技术的迅猛发展和网络应用的迅速普及,互联网日益形成一个巨大的分布式信息库。互联网应用产生的超大规模信息对现有的网络数据管理基础设施提出了新的严峻挑战。互联网信息库的无限扩张性和与生俱来的分布式特性使非集中式的数据管理和共享机制的兴起成为一种必然趋势。大规模分布式内容检索技术具有重要的学术价值和应用价值。大规模分布式内容检索系统的核心问题是如何建立高效的分布式索引以支持大规模网络环境下的复杂内容检索。本书从此核心问题出发,通过扩展传统对等模式的概念、结构、资源描述与组织、资源发现与路由、结果融合与排序等,全面、深入、系统地论述了利用对等模型构建大规模分布式内容检索系统的解决方案和关键技术。

全书共 9 章,其中各章的主要内容安排如下:

第 1 章为绪论。本章介绍了大规模分布式内容检索面临的挑战。从对等网络的概念和特性入手,对现有大规模分布式内容检索技术进行了分类介绍。

第 2 章为分布式哈希表及单关键字全局索引。本章介绍了分布式定位的经典算法——分布式哈希表。首先介绍了构造分布式哈希表的理论基础——一致性哈希;然后介绍了四种经典的分布式哈希表算法,包括 Chord、CAN、Pastry 和 Tapestry,并对各种算法进行了对比分析;最后分别介绍了基于分布式哈希表构建的典型内容检索系统 eSearch 和 Minerva 的基本原理。

第 3 章为布隆滤波。本章系统介绍了布隆滤波这种时空高效的集合哈希编码技术。首先从传统的哈希编码技术入手,分析通过引入一定的误判率来实现高效的集合编码;然后介绍了布隆滤波的几种扩展形式,包括计数布隆滤波、压缩布隆滤波、动态布隆滤波;最后介绍了布隆滤波在网络和分布式系统中的应用。

第 4 章为基于分布式哈希表单关键字索引的搜索。本章介绍了如何在基于分布式哈希表的单关键字索引上进行多关键字复杂搜索。首先详细分析比较了几种分布式 Top-k 算法,然后重点介绍了一种利用优化布隆滤波技术来实现多关键字搜索的 PWEB 系统的设计。

第 5 章为多关键字全局索引及搜索。本章介绍了如何通过扩展标准分布式哈希表实现关键字集索引。首先介绍了经典关键字赋权模型,然后重点介绍了两个基于分布式哈希表关键字集索引的系统 HDK 和 TSS。

第 6 章为基于复制的联邦式对等搜索策略。本章介绍了如何通过复制策略来提高联邦式搜索系统。首先对各种复制策略进行了理论分析,然后重点介绍了

WP、BubbleStorm 和 BloomCast 三种有效的复制策略。

第 7 章为基于内容路由的联邦式搜索策略。本章介绍了如何重构路由表以提高联邦式搜索引擎的检索效率。分别介绍了基于语义小世界的模型、基于兴趣局部性的路由策略和一种基于语义覆盖网的搜索系统 SemreX。

第 8 章为混合式对等搜索策略。本章介绍了如何结合分布式哈希表和洪泛协议各自的特点设计混合搜索策略。分别详细介绍了基于预先探测的混合搜索策略 SimpleHybrid、基于 Gossip 的混合搜索策略 GAB 和基于查询难度感知的混合搜索策略 QRank。

第 9 章为大规模在线社会网络搜索。本章通过分析大规模社会网络的特点，结合传统分布式检索中的成熟技术，对大规模社会网络搜索中的关键技术进行了展望。

本书的出版得到了国家重点基础研究发展计划(973 计划)课题“基于语义网格的语义关联存贮模型及管理和通信平台(2003CB317003)”和国家自然科学基金委员会与香港研究资助局(NSFC/RGC)联合科研基金项目“因特网上基于对等网络的大规模实时视频系统：理论和实践(60731160630)”的资助。

作 者

2010 年 12 月于瑜伽山

# 目 录

## 序

## 前言

<b>第1章 绪论</b>	1
1.1 对等网络概述	1
1.2 基于对等模式的大规模分布式文本内容检索	5
1.3 大规模分布式文本内容检索研究面临的挑战	6
1.4 大规模分布式文本内容检索技术分类	7
1.4.1 基于结构化分布式哈希表的分布式全局倒排索引	7
1.4.2 基于非结构化对等网络的联邦式搜索网络	9
1.4.3 混合对等网络搜索引擎	11
1.5 本书内容	12
参考文献	15
<b>第2章 分布式哈希表及单关键字全局索引</b>	20
2.1 分布式哈希表	21
2.1.1 Chord: 基于二分查找的环状对等结构	24
2.1.2 CAN: 基于多维空间划分的对等结构	28
2.1.3 Pastry: 基于多分查找的前缀匹配对等结构	34
2.1.4 Tapestry: 基于多分查找的对等结构	40
2.2 现有分布式哈希表算法的比较	45
2.3 利用分布式哈希表构建单关键字全局索引	46
2.3.1 eSearch: 基于分布式哈希表的水平索引	47
2.3.2 Minerva: 在查询中挖掘关联关键字	48
2.3.3 局限性	52
参考文献	53
<b>第3章 布隆滤波</b>	54
3.1 哈希编码的时间/空间权衡	54
3.1.1 一种经典的哈希编码方法	55
3.1.2 两种存在误判率的哈希编码方法	56
3.1.3 计算因子	56
3.1.4 三种哈希编码方法的数学分析	58

3.1.5 时空性能比较 .....	61
3.2 布隆滤波的基本理论 .....	62
3.2.1 布隆滤波概念 .....	62
3.2.2 位向量长度的下界 .....	64
3.2.3 布隆滤波与集合运算 .....	65
3.3 布隆滤波的扩展形式 .....	66
3.3.1 计数布隆滤波 .....	66
3.3.2 压缩布隆滤波 .....	71
3.3.3 动态布隆滤波 .....	74
3.4 布隆滤波的应用 .....	87
3.4.1 早期应用 .....	87
3.4.2 分布式缓存 .....	88
3.4.3 P2P 网络 .....	88
3.4.4 资源路由 .....	89
3.4.5 数据包路由 .....	90
3.4.6 基础设施测量 .....	91
参考文献 .....	91
<b>第 4 章 基于分布式哈希表单关键字索引的搜索 .....</b>	<b>94</b>
4.1 结构化对等网多关键字检索面临的挑战 .....	94
4.2 Top-k 查询策略 .....	95
4.2.1 倒排索引 .....	95
4.2.2 Top-k 裁剪算法 .....	97
4.2.3 性能评估 .....	102
4.3 PWEB 系统 .....	104
4.3.1 PWEB 网络结构 .....	105
4.3.2 多关键字搜索通信开销优化策略 .....	106
4.3.3 扩展性算法 .....	112
4.3.4 分布式交集运算执行顺序优化策略 .....	114
4.3.5 搜集关键字全局统计信息 .....	115
4.3.6 模拟仿真方法 .....	117
4.3.7 性能评估 .....	121
4.4 小结 .....	132
参考文献 .....	132
<b>第 5 章 多关键字全局索引及搜索 .....</b>	<b>135</b>
5.1 分布式关键字集索引面临的挑战 .....	135

5.2 文本检索中的关键字权重方法 .....	136
5.2.1 关关键字权重模型 $TF \times IDF$ .....	136
5.2.2 理解逆文档频率 .....	140
5.2.3 用逆向总关键字频率替换逆文档频率的尝试 .....	144
5.2.4 词频在相关权重模型中的探索 .....	144
5.3 HDK: 基于高区分关键字集的索引技术 .....	146
5.3.1 关关键字集倒排索引 .....	147
5.3.2 高区分关键字集索引 .....	147
5.3.3 基于高区分关键字集索引的搜索 .....	148
5.3.4 扩展性分析 .....	148
5.3.5 性能评估 .....	150
5.4 TSS: 基于关键字集索引的 P2P 搜索系统 .....	153
5.4.1 TSS 系统结构 .....	153
5.4.2 分布式关键字集索引 .....	155
5.4.3 模拟测试方法 .....	159
5.4.4 性能评估 .....	161
参考文献 .....	167
<b>第 6 章 基于复制的联邦式对等搜索策略 .....</b>	<b>169</b>
6.1 理论分析 .....	169
6.1.1 模型建立 .....	170
6.1.2 均匀复制策略和比例复制策略 .....	170
6.1.3 平方根复制策略 .....	171
6.1.4 混合复制策略 .....	173
6.1.5 分布式复制算法的实现 .....	174
6.2 基于随机游走的随机复制策略 .....	176
6.2.1 生日悖论和理论下界 .....	177
6.2.2 随机游走复制策略和搜索协议 .....	179
6.2.3 性能评估 .....	183
6.3 BubbleStorm: 基于随机多图的概率穷尽搜索策略 .....	187
6.3.1 副本数量的确定 .....	188
6.3.2 网络大小的测量 .....	188
6.3.3 随机多图与随机采样 .....	189
6.3.4 洪泛和随机游走的完美结合 .....	189
6.3.5 系统分析 .....	190
6.3.6 性能评估 .....	191

6.4 BloomCast: 基于轻量级分布式哈希表的随机采样.....	194
6.4.1 BloomCast 网络结构 .....	194
6.4.2 网络结点数量估计 .....	195
6.4.3 随机结点采样 .....	197
6.4.4 基于布隆滤波的复制算法 .....	198
6.4.5 多关键字搜索 .....	198
6.4.6 性能评估 .....	199
6.5 PlanetP: 基于全局摘要索引的复制策略 .....	205
6.5.1 全局目录索引复制 .....	206
6.5.2 结点排序模型 .....	207
6.5.3 查询处理算法 .....	208
6.5.4 性能评估 .....	208
参考文献.....	211
<b>第7章 基于内容路由的联邦式搜索策略.....</b>	<b>213</b>
7.1 基于语言模型的路由选择 .....	213
7.1.1 联邦式搜索引擎的两层结构 .....	213
7.1.2 语言模型 .....	214
7.1.3 相对熵 .....	216
7.1.4 搜索算法 .....	216
7.1.5 性能评估 .....	219
7.2 基于语义小世界模型的联邦式对等搜索 .....	221
7.2.1 语义空间和向量 .....	221
7.2.2 构造语义小世界 .....	223
7.2.3 降低语义小世界的维度 .....	224
7.2.4 基于语义小世界的搜索 .....	226
7.2.5 性能评估 .....	226
7.3 基于兴趣局部性的路由 .....	229
7.3.1 兴趣局部性 .....	229
7.3.2 基于兴趣局部性的拓扑和路由 .....	229
7.3.3 性能评估 .....	230
7.4 SemreX 系统 .....	232
7.4.1 SemreX 系统模型 .....	232
7.4.2 语义覆盖网 .....	236
7.4.3 基于语义覆盖网的查询搜索算法 .....	243
7.4.4 性能评估 .....	246

---

参考文献.....	252
<b>第8章 混合式对等搜索策略.....</b>	<b>254</b>
8.1 混合对等搜索面临的挑战 .....	254
8.2 基于预先探测的混合策略 .....	256
8.2.1 Boon Thau Loo 的 Gnutella 实验 .....	256
8.2.2 SimpleHybrid 混合 P2P 搜索策略 .....	260
8.2.3 性能评估 .....	261
8.3 基于 Gossip 的混合搜索选择.....	262
8.3.1 收集全局统计信息 .....	262
8.3.2 使用全局信息进行搜索选择 .....	264
8.3.3 洪泛阈值的调节 .....	264
8.3.4 性能评估 .....	265
8.4 难度感知的混合式搜索策略 .....	268
8.4.1 很多复本≠很多结点 .....	268
8.4.2 QRank 设计 .....	269
8.4.3 用 QRank 进行混合查询.....	273
8.4.4 自适应混合查询 .....	274
8.4.5 QRank 仿真器设计 .....	275
8.4.6 性能评估 .....	276
参考文献.....	285
<b>第9章 大规模在线社会网络搜索.....</b>	<b>287</b>
9.1 大规模在线社会网络搜索面临的挑战 .....	287
9.2 在线社会网络系统研究现状 .....	288
9.3 流行在线社会网络的数据划分与定位 .....	289
9.4 大规模在线社会网络内容搜索关键技术 .....	290
9.4.1 流式文本摘要技术 .....	291
9.4.2 基于摘要索引的排序算法 .....	292
9.4.3 多跳邻居摘要聚合技术 .....	292
9.4.4 基于社区局部性降低摘要索引开销 .....	293
参考文献.....	294

# 第 1 章 绪 论

随着网络技术的迅猛发展和网络应用的迅速普及,互联网日益形成一个巨大的分布式信息库。互联网应用产生的超大规模信息对现有的网络数据管理基础设施提出了新的更为严峻的挑战。互联网信息库的无限扩张性和与生俱来的分布式特性使非集中式的数据管理和共享机制成为近年来的研究热点<sup>[1]</sup>。大规模分布式内容检索研究具有重要的学术价值和应用价值。

对等计算技术是分布式系统和计算机网络结合的产物,它在网络协议的应用层打破了传统的客户机/服务器(C/S)模式,以自主、平等的原则将处于网络边缘的计算、存储、通信、信息等各种资源有效地共享起来,形成协作网络<sup>[2,3]</sup>。自诞生以来的短短几年时间里,对等文件共享应用,如 Gnutella<sup>[4]</sup>、KaZaA<sup>[5]</sup>、BitTorrent<sup>[6]</sup>等都取得了极大的成功,并占据了当前互联网一半以上的网络流量。对等模式因其可扩展性、鲁棒性和动态自适应性等优点,在大规模互联网应用的数据管理和搜索领域日益展现出巨大的潜力。近年来,基于对等模式的大规模分布式数据管理和内容搜索系统如雨后春笋般涌现出来<sup>[7~13]</sup>。

本书将围绕大规模分布式内容检索应用展开,全面系统地阐述大规模分布式内容检索系统的关键理论和支撑技术。

## 1.1 对等网络概述

在互联网产生的初期,大多数连接到互联网上的普通用户由于计算机性能、资源等因素的限制,并没有能力对外提供服务。互联网逐步形成了以少数服务器为中心的状况,C/S 模式成为互联网应用的主要模式。C/S 架构对客户机的资源要求非常少,满足了早期互联网用户希望以非常低廉的成本接入互联网的需求。最近二十年来,计算机性能按照摩尔定律增长,采用廉价的计算机也能够搭建分布式高性能数据中心,提供互联网规模的服务。时至今日,即使是处于网络边缘的个人电脑也具备了相当的计算、存储和通信能力,甚至超过十多年前的超级计算机。这使得人们在互联网上进行更大规模、更广泛、更自主的信息共享成为可能。

1999 年,18 岁的美国波士顿东北大学的一年级新生 Fanning 在校园里编写了一个共享软件 Napster。Napster 能够很方便地让用户在网上交换最新的 MP3 音乐文件。与传统的提供音乐下载的网站不同,Napster 的服务器中没有一首音乐,它只维护一个轻量级的用户音乐文件索引提供查询功能,而音乐文件实际上是在

用户之间直接交换的，而且系统中的某用户在从网络中别的用户那里获取 MP3 文件的同时，也为其他用户提供下载服务。同年，Napster 程序被放到了互联网上供网友下载。Napster 一夜成名，突然成为人们争相转告的互联网“杀手锏”应用，它就像一个拥有无限货架的虚拟音像店<sup>[14]</sup>，令无数散布在互联网上的音乐爱好者美梦成真。短短一年半之间，Napster 就吸引了超过 5000 万的下载用户<sup>[15]</sup>，一度成为互联网上增长最快的应用。Napster<sup>[16]</sup>采用的对等模式随后引发了一场互联网地震，其背后隐藏的巨大的 P2P 世界也从此浮出水面。

在传统的基于 C/S 模式的互联网信息共享应用中，如 Web、FTP 文件服务等，用户之间如果要共享数据，必须经过服务器交换。在 FTP 应用中，用户先将文件上传到某个 FTP 服务器上，然后其他用户再到相应站点下载需要的文件。这种方式不但共享效率低、不可扩展，而且对用户而言极不方便。相比较而言，采用对等模式，数据可以自主分散在网络中，不像以往的 C/S 模式那样把数据存放在中心服务器上。采用对等模式，网络结点(Peer)既可以获取其他结点的资源或服务，又可以提供资源或服务，即兼具客户机和服务器的双重身份，这种特征使得互联网上的任意两台计算机之间直接共享文件成为可能<sup>[15]</sup>。Napster 无意间采用的对等模式打破了 C/S 模式的束缚，充分挖掘了蕴藏在互联网中的巨大的计算、存储和通信资源，高度共享，形成协作网络。对等计算作为新型的分布式计算模型，正深刻地影响着互联网用户发布、共享和获取信息的方式，并被称为改变互联网的新一代网络技术<sup>[2]</sup>。

一般来说，对等网络具有以下基本特征：

① 自主平等性。对等网络中的资源分散在所有结点上，信息交换直接在结点之间自主进行，无需中间环节的介入，避免了可能的性能瓶颈，提高了资源共享的效率。对等结点之间交换哪些信息，可由结点本身自主决定，而不必受限于其他网络结点。在对等模式中，虽然网络结点能力不同，但所有成员在功能、地位上都是平等的，没有谁拥有超越其他结点的特权，没有谁能控制或者限制其他结点。对网络模型而言，平等指的是打破传统的 C/S 模式，取消服务器这一特权结点的存在，让所有网络结点之间平等地交流信息。平等性是对等网络的工作基础，对等网络对网络带宽的高效利用、对网络结点潜力的充分开发以及可扩展性等，都是基于其平等性的<sup>[15]</sup>。

② 可扩展性。在传统的 C/S 架构中，系统能够容纳的用户数量和提供服务的能力主要受服务器的资源限制，为了支持互联网上的大量用户，需要在服务器端使用高性能的计算机和高带宽的网络。一般来说，当系统网络结点数目增加时，服务器的负载也随之线性增长，如果原来的服务器不堪重负，则必须购买增加新的服务器来分担原有服务器的负载。此外，C/S 架构往往需要集中式服务器之间的同步、协同处理，也会产生大量的开销，使系统难以扩展。相比之下，对等网络具有非常

高的可扩展性,当结点数目大量增加时,随之增加的存储、处理、通信开销被更多的结点分担,系统结点的平均负载不会增加太多。因为网络用户规模的增大也带来了系统整体资源和服务能力的同步扩充,所以系统无需增加额外的设备。对等网络这种良好的可扩展性优点已经在一些大规模实例系统中得到了证明<sup>[2]</sup>。

③ 魯棒性和动态自适应性。互联网上随时可能出现异常情况,如网络中断、网络拥塞、结点失效等,都会影响系统的可用性和稳定性。在传统的集中式服务模式中,集中式服务器成为整个系统的要害所在,一旦发生异常,将影响到所有用户使用系统服务。相比较而言,在对等网络中,服务是由分散的各个结点协同提供的,部分结点或网络遭到破坏不会造成整个协同网络的瘫痪。对等网络是动态变化的,结点加入或离开,数据随时动态更新。对等网络协议一般针对这种动态性提供自适应拓扑调整、更新功能,保持网络的连通性和服务的可用性<sup>[2]</sup>。

④ 隐私性。随着互联网的普及和计算、存储及通信能力的飞速增长,收集隐私信息正在变得越来越容易。例如,目前的互联网攻击者往往通过监控用户的流量特征,获得IP地址,进而追踪到个人用户。实践中解决互联网隐私问题主要采用中继转发的方法,将通信双方隐藏在众多的网络实体中,实现匿名通信。在传统的匿名通信系统中,实现这一机制依赖于某些中继服务器结点。而在对等网络中,信息的传输分散在各结点之间进行,所有参与者都可以提供中继转发的功能,信息传输无需经过某个特定的集中环节,从而大大提高了匿名通信的灵活性和可靠性,能够为用户提供更好的隐私保护。另外,对等网络也采用分布式哈希表(DHT)<sup>[18~21]</sup>等技术保障其匿名性<sup>[22]</sup>。分布式哈希表采用一致性哈希<sup>[23]</sup>将用户信息(如IP地址等)、数据对象信息映射到一个没有任何意义的数值标识(ID),用来唯一地标识用户和数据对象。一致性哈希具有单向不可逆性和抗冲突性等特性,这使得网络攻击者很难从ID破解出它实际所代表的真实信息<sup>[15]</sup>。

采用对等模式,网络结点在网络体系结构的应用层中建立虚拟连接,从而整个分布式系统中的所有结点互连组成了一个逻辑虚拟网络,简称覆盖网络(overlay network)。覆盖网络构建于底层物理网络之上,结点间通信依赖于底层物理网络通信协议的支持<sup>[24]</sup>。

根据覆盖网络的结构,对等网络可以分为集中索引对等网络、结构化对等网络和非结构化对等网络<sup>[24]</sup>。

### (1) 集中索引对等网络

集中索引对等网络维护一个集中式的目录服务器用于注册和查找网络中的数据。用户访问数据时,首先搜索请求被发送到集中目录服务器进行资源定位,获取数据的地址,然后直接从地址所在结点获取数据,而不需要经过集中服务器。集中索引服务器对等模式,虽然避免了数据在服务器上集中放置,但集中式的索引服务器仍然是系统的核心。这自然使此类系统保留了集中结构的一些缺点:集中索引

服务器是系统潜在的性能瓶颈<sup>[25]</sup>,随着系统规模的增长,搜索请求和目录更新会造成巨大的负载,保障系统的可扩展性是这类系统面临的一个主要问题;集中索引也是系统潜在的单一失效点,索引服务器故障或者网络攻击,都有可能造成整个系统的崩溃。集中索引对等系统的典型代表是 Napster 和 Maze<sup>[26]</sup>。

## (2) 结构化对等网络

结构化对等网络<sup>[27~29]</sup>一般使用分布式哈希表<sup>[18~21]</sup>技术对数据的放置进行严格控制。分布式哈希表以分布的方式维护逻辑上统一的哈希表,具有良好的可扩展性。另外,哈希函数的映射原理使结构化对等网络只提供精确查找功能,这虽然能有效保证系统查全率(recall)<sup>[30,31]</sup>,但不直接支持灵活的查询处理,如多关键字搜索和区间查询<sup>[32]</sup>等。研究人员在结构化对等网络方面进行了大量卓有成效的研究工作,下面简要介绍基于分布式哈希表的一些流行算法。

① Chord<sup>[18,33]</sup>是麻省理工学院和加州大学伯克利分校设计的一种基于分布式哈希表的查找和路由协议。Chord 采用统一的哈希函数将网络对象(key, value)中的 key 和网络结点信息(如 IP 地址)映射到相同的值空间,然后按照哈希值临近原则将对象(key, value)存储在最接近 key 的哈希值的结点上。Chord 协议采用了类似二分查找的方法,每次查找发送的消息数为  $O(\log_2 n)$ ,其中  $n$  为网络结点数。结点离开或加入网络时,拓扑维护所需网络消息开销为  $O(\log_2^2 n)$ ,因此一般认为当网络结点动态性很强时,Chord 维护成本比较高。总体来说,Chord 算法精妙,并具有扩展性好、查询效率高、负载均衡、自适应性拓扑维护等优点。然而,由于采用哈希映射,Chord 仅支持精确的对象查询,不直接支持复杂查询应用,这也是其他几种分布式哈希表算法的共同问题。

② CAN<sup>[19]</sup>是加州大学伯克利分校和 AT&T 设计的一种分布式查找和路由算法。CAN 将网络中所有结点映射到一个  $n$  维的笛卡儿空间中,并为每个结点尽可能均匀地分配一块区域。CAN 采用的哈希函数通过对(key, value)对中的 key 进行哈希运算,得到笛卡儿空间中的一个点,并将(key, value)对存储到该点所在区域对应的结点上。CAN 采用的路由算法简单有效,当前结点知道目标点的坐标后,就利用贪心算法的思想,将请求传给其四邻中坐标最接近目标点的结点,依次路由,直至查询达到目标结点。CAN 具有良好的可扩展性,其路由路径长度为  $O(\sqrt{d}n)$ ,结点状态信息和网络规模无关,为  $O(d)$ ,其中  $n$  为结点数,  $d$  为系统维数。

③ Tapestry<sup>[20]</sup>是加州大学伯克利分校 Zhao 等设计的面向广域网分布式数据存取的结构化对等系统。Tapestry 结点通过逐位匹配的方式路由查找请求,所以结点同样匹配自己标识符的每一个前缀相同而下一位不同的邻居信息。逐位匹配的思想类似于 Chord 的二分查找思想,有效提高了路由效率。系统中的每个结点有  $O(\log_2 n)$  个邻居,由于每跳匹配一位,故路由路径为  $O(\log_2 n)$  跳,其中  $n$  为结点数。网络中的结点能够通过选择最佳邻居结点保证最小化平均路径延迟。

Tapestry在构建拓扑时考虑了覆盖网络和底层网络的匹配问题,有效提高了系统效率。

④ Pastry<sup>[21]</sup>是微软研究院提出的一种可扩展的分布式查找路由协议,可用于构建大规模的对等系统。Pastry 源于 Plaxton 路由机制<sup>[34]</sup>,Plaxton 是第一个基于分布式哈希表的算法,但设计目标不是用于对等系统,所以只假设了静态环境,并且没有提供很有效的路由查找。在 Pastry 中每个结点分配一个 128 位的结点标识符号(nodeID),所有的结点标识符被映射到环形的 nodeID 空间,范围从  $0 \sim 2^{128} - 1$ ,结点加入系统时通过哈希结点 IP 地址在 128 位 nodeID 空间中随机分配。Pastry 采用了类似 Tapestry 的逐位匹配思想,算法具有指数收敛性。在 Pastry 中路由一个消息需要  $O(\log_2 n)$  步,每个结点需要维护  $O(\log_2 n)$  个人口,其中  $n$  为网络结点数。

### (3) 非结构化对等网络

非结构化对等网络<sup>[35~37]</sup>对覆盖网络拓扑结构以及系统中资源的放置没有严格的限制<sup>[38]</sup>,因此具有较好的鲁棒性,能够很好地应对网络中结点的高度动态变化。该类系统中的搜索主要是通过结点间的消息洪泛转发,协议简单且易于支持各类查询语法。然而,采用基于洪泛思想的算法,搜索消息数会在网络上指数增长并导致大量冗余的通信负载,从而限制了非结构化对等网络的可扩展性<sup>[4]</sup>。早期的非结构化对等网络不考虑结点之间的区别,后期的对等网络协议考虑了网络中结点之间存在的性能差异,一些具有较高网络带宽和较长在线时间的网络结点以超级结点的身份加入对等网络,每个普通结点至少与一个超级结点连接。超级结点在搜索过程中代理与之相连的普通结点执行搜索请求。研究表明,对于广泛分布的热点数据,非结构对等协议搜索性能很高,而对于分布非常稀少的数据,其搜索性能却很差<sup>[30,39,40]</sup>。非结构对等系统的典型代表是 Gnutella 和 KaZaA 等。

## 1.2 基于对等模式的大规模分布式文本内容检索

随着网络技术的迅猛发展和网络应用的迅速普及,互联网日益成为一个巨大的分布式信息库。特别是以 FaceBook、Wikipedia、MySpace、YouTube、Google Document 和 Flickr 等为代表的云计算应用兴起以来,越来越多的用户正以更广泛、更直接和更高效的方式参与到互联网信息的创造和共享行列中来。互联网信息库的无限扩张性和与生俱来的分布式特性使研究非集中式的数据管理和共享机制成为一种必然趋势。基于分布式模型的新一代内容检索技术已成为近年来国内外学术界和工业界的新热点。对等计算模型凭借其分布式、易扩展、容错性高等优点,日益在互联网信息搜索方面显示出巨大的潜力。

① 互联网信息的无限增长和集中式搜索引擎数据中心处理能力的局限性之