

P2P对等网络

原理与应用

蔡康 唐宏 丁圣勇 郑贵锋 编著



科学出版社

P2P 对等网络原理与应用

蔡 康 唐 宏 丁圣勇 郑贵锋 编著

科学出版社

北 京

内 容 简 介

本书较为系统地介绍了 P2P 的理论基础,对 P2P 的基础路由,如 DHT 算法、DHT 性能作了深入介绍,对最新的理论成果网络编码也作了深入浅出的分析。同时,本书对 P2P 传送过程进行了建模,从模型的高度抽象提取了传送过程中的几个核心参数,并分析核心参数之间的依赖关系,为 P2P 传送优化和播放器缓存设计提供了有价值的理论参考。在此基础上,本书列举了大量 P2P 的应用实例,为读者理解 P2P 应用方法提供了丰富的参考。此外,本书还前瞻性地提出了 P2P 在 IPv6 网络环境下的问题,指出在 IPv6 与 IPv4 共存的环境下,P2P 必须依赖自身的算法来自适应不同的网络环境,并提出了具体的解决方案。这些方案的有效性已经通过实践验证,为 P2P 向 IPv6 发展提供了重要参考。

在结构上,本书按照从理论到实践、从抽象到具体、从简单到深入的顺序安排内容,主要面向希望全面掌握 P2P 知识的初级读者和 P2P 软件的开发者。通过阅读本书,读者能够快速掌握 P2P 的基础原理,并循序渐进地深入理解 P2P 的核心理论和应用技术。

图书在版编目(CIP)数据

P2P 对等网络原理与应用/蔡康等编著. —北京:科学出版社,2011

ISBN 978-7-03-031582-3

I. ①P… II. ①蔡… III. ①因特网-基本知识 IV. ①TP393.4

中国版本图书馆 CIP 数据核字(2011)第 113230 号

责任编辑:裴 育 于 红 / 责任校对:陈玉凤
责任印制:赵 博 / 封面设计:耕者设计工作室

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>

丽源印刷厂印刷

科学出版社发行 各地新华书店经销

*

2011 年 6 月第 一 版 开本:B5(720×1000)

2011 年 6 月第一次印刷 印张:19 1/4

印数:1—4 000 字数:372 000

定价:60.00 元

(如有印装质量问题,我社负责调换)

前 言

凭借天然的去中心化对等通信理念，P2P 无论是在学术领域还是在应用领域都获得了极大的关注，在理论和应用上都获得了显著的发展，如经典的 DHT 算法、风靡全球的 BitTorrent 软件等。随着用户对文件下载速率和流媒体质量要求的不断提升，P2P 几乎已经成为大型通信系统中不可缺少的技术，甚至在传统的集中式解决方案中，在服务器之间采用 P2P 技术实现自动负载均衡也逐渐开始流行。基于这样的广泛需求，编写一本能够将理论与实践结合起来的 P2P 书籍显得尤为重要。

事实上，作者参与了 P2P 重大项目，负责整个 P2P 系统的设计与代码开发工作，在研发过程中，深感 P2P 领域包含的内容非常丰富。要灵活应用 P2P 技术实现真实的系统，就必须对 P2P 的理论和实现等各方面的内容有全面掌握。从实践经验来看，掌握这些技术需要查阅大量的理论资料并进行相关实验。基于此种需求，作者将实际开发过程中涉及的各种理论和应用技术整理并编写成册，尽量以通俗易懂的方式对各种复杂的算法进行解释，同时又不失深度，以期给有兴趣深入学习 P2P 或有志于开发 P2P 的读者提供相对完备的参考资料。

此外，作者对电信运营网络有一手的经验数据，并参与制定 IPv6 的发展规划。作者结合这些经验对 IPv6 网络下的 P2P(P6P)进行了分析，提出了 P6P 在 IPv4 和 IPv6 混杂网络下的解决方案，这些内容为关注未来 P2P 发展的读者提供了重要参考。

本书的具体内容主要包括 P2P 综述、P2P 路由技术、P2P 传送技术、P2P 主要应用、P2P 运营管理、P4P 优化、IPv6 P2P 技术等。其中，综述部分介绍 P2P 的定义、发展历史与现状、目前主流的应用等；路由技术部分介绍集中式路由算法、分布式路由算法、混合式路由算法等；传送技术部分介绍传统的推拉模式以及基于网络编码的传送算法，并结合缓存管理详细解释 P2P 传送算法设计时需要考虑的各种参数；主要应用部分介绍典型的文件下载、流媒体、及时通信等具体的 P2P 应用原理；运营管理部分介绍 P2P 对网络产生的冲击以及运营商采用的运营管理方法；P4P 优化部分主要介绍目前流行的 P2P 流量优化方法；IPv6 P2P 技术部分介绍 IPv6 环境下 P2P 面临的主要问题以及相应的解决方案。

目 录

前言

| | |
|--------------------------------------|----|
| 1 P2P 简介 | 1 |
| 1.1 P2P 定义 | 1 |
| 1.2 P2P 特点 | 2 |
| 1.3 P2P 发展历史与现状 | 4 |
| 1.3.1 P2P 发展的四个阶段 | 4 |
| 1.3.2 国外 P2P 技术的研究现状 | 6 |
| 1.3.3 国内 P2P 技术的研究现状 | 8 |
| 1.3.4 P2P 的网络流量 | 8 |
| 1.4 P2P 的主要应用领域与代表软件 | 10 |
| 1.4.1 下载 | 10 |
| 1.4.2 流媒体 | 11 |
| 1.4.3 即时通信 | 14 |
| 1.4.4 其他领域 | 16 |
| 1.5 P2P 产业 | 17 |
| 1.5.1 P2P 产业链 | 17 |
| 1.5.2 版权问题 | 20 |
| 1.5.3 P2P 与电信网络运营 | 21 |
| 1.6 本章总结 | 22 |
| 2 P2P 网络核心技术——拓扑结构与内容路由 | 23 |
| 2.1 P2P 网络基本概念 | 23 |
| 2.2 集中式 P2P 网络 | 25 |
| 2.3 纯分布式 P2P 网络 | 27 |
| 2.3.1 小世界模型 | 27 |
| 2.3.2 纯分布式 P2P 网络的网络拓扑与内容路由 | 30 |
| 2.4 混合式 P2P 网络 | 33 |
| 2.5 结构化 P2P 网络 | 35 |
| 2.5.1 DHT 算法概述 | 37 |
| 2.5.2 Chord 算法 | 38 |
| 2.5.3 Pastry 算法 | 43 |

| | | |
|----------|-------------------------|------------|
| 2.5.4 | CAN 算法 | 47 |
| 2.6 | 本章总结 | 48 |
| 3 | P2P 网络核心技术——内容传送 | 49 |
| 3.1 | 非实时内容传送技术 | 49 |
| 3.1.1 | 基本传送技术 | 50 |
| 3.1.2 | 基于网络编码的模式 | 52 |
| 3.2 | 实时内容传送技术 | 55 |
| 3.3 | NAT 穿越 | 62 |
| 3.4 | 本章总结 | 65 |
| 4 | P2P 开发平台 | 66 |
| 4.1 | JXTA | 67 |
| 4.1.1 | JXTA 介绍 | 67 |
| 4.1.2 | JXTA 层次结构 | 67 |
| 4.1.3 | JXTA 协议 | 68 |
| 4.1.4 | JXTA 相关概念 | 69 |
| 4.1.5 | 开发实例 | 75 |
| 4.2 | Python | 90 |
| 4.2.1 | Python 介绍 | 90 |
| 4.2.2 | Python 的基本语法和结构 | 91 |
| 4.2.3 | 开发实例 | 93 |
| 4.3 | 本章总结 | 98 |
| 5 | P2P 文件共享应用 | 99 |
| 5.1 | P2P 文件共享应用系统 | 99 |
| 5.2 | BitTorrent 下载系统 | 99 |
| 5.2.1 | BT 系统结构 | 100 |
| 5.2.2 | BT 网络协议分析 | 101 |
| 5.2.3 | CTorrent 程序源码分析 | 110 |
| 5.3 | eMule 下载系统 | 114 |
| 5.3.1 | eMule 系统结构 | 115 |
| 5.3.2 | eMule 网络协议分析 | 118 |
| 5.3.3 | eMule 源代码分析 | 124 |
| 5.4 | 本章总结 | 139 |
| 6 | P2P 网络流媒体应用 | 140 |
| 6.1 | 流媒体系统概述 | 140 |
| 6.1.1 | 流媒体系统架构 | 140 |

| | |
|-------------------------|-----|
| 6.1.2 P2P 流媒体系统 | 142 |
| 6.2 PeerCast 流媒体传输系统 | 143 |
| 6.2.1 PeerCast 系统结构 | 144 |
| 6.2.2 PeerCast 网络协议 | 144 |
| 6.2.3 频道组织结构 | 145 |
| 6.2.4 工作流程 | 145 |
| 6.2.5 算法原理 | 148 |
| 6.2.6 PeerCast 源代码分析 | 151 |
| 6.3 本章总结 | 168 |
| 7 P2P 网络即时通信应用 | 169 |
| 7.1 即时通信 | 169 |
| 7.2 Skype 通信系统 | 169 |
| 7.2.1 Skype 简介 | 169 |
| 7.2.2 Skype 系统结构 | 171 |
| 7.2.3 Skype 协议分析 | 173 |
| 7.3 本章总结 | 187 |
| 8 P2P 网络搜索应用 | 188 |
| 8.1 P2P 搜索原理及算法 | 188 |
| 8.1.1 非结构化 P2P 网络搜索算法 | 188 |
| 8.1.2 结构化 P2P 网络搜索算法 | 191 |
| 8.1.3 其他搜索算法 | 193 |
| 8.1.4 算法对比分析 | 194 |
| 8.2 典型应用 | 195 |
| 8.2.1 搜索引擎工作原理 | 195 |
| 8.2.2 YaCy 搜索引擎系统 | 195 |
| 8.3 本章总结 | 198 |
| 9 P2P 网络运营系统体系架构 | 199 |
| 9.1 终端呈现 | 200 |
| 9.2 P2P 业务封装 | 200 |
| 9.2.1 子系统功能 | 200 |
| 9.2.2 子系统接口 | 201 |
| 9.3 P2P 基础服务 | 204 |
| 9.3.1 子系统功能 | 204 |
| 9.3.2 子系统接口 | 205 |
| 9.4 内容提供 | 206 |

| | | |
|-----------|-------------------------|------------|
| 9.5 | 发布管理 | 206 |
| 9.5.1 | 子系统功能 | 206 |
| 9.5.2 | 子系统接口 | 207 |
| 9.6 | 认证/计费管理 | 207 |
| 9.6.1 | 子系统功能 | 207 |
| 9.6.2 | 子系统接口 | 207 |
| 9.7 | 本章总结 | 208 |
| 10 | P2P 网络监控 | 209 |
| 10.1 | P2P 网络监控的意义 | 209 |
| 10.1.1 | P2P 网络监控概念 | 209 |
| 10.1.2 | P2P 监控现状 | 212 |
| 10.1.3 | P2P 监控意义 | 214 |
| 10.2 | P2P 网络监测手段 | 216 |
| 10.2.1 | 传统 P2P 监测手段 | 216 |
| 10.2.2 | 基于 DPI 技术的 P2P 监测 | 219 |
| 10.2.3 | P2P 监测手段小结 | 223 |
| 10.3 | P2P 网络控制手段 | 224 |
| 10.3.1 | 法律政策手段 | 224 |
| 10.3.2 | 经济手段 | 227 |
| 10.3.3 | 技术手段 | 232 |
| 10.3.4 | P2P 控制手段小结 | 242 |
| 10.4 | P2P 网络监控系统 | 243 |
| 10.4.1 | DPI 系统的实现 | 243 |
| 10.4.2 | DPI 系统流量识别过程 | 247 |
| 10.4.3 | 旁路部署式 DPI 系统 | 250 |
| 10.4.4 | 串接部署式 DPI 系统 | 251 |
| 10.4.5 | 集成式 DPI 系统 | 253 |
| 10.4.6 | DPI 系统综合比较 | 256 |
| 10.4.7 | DPI 系统功能和性能要求 | 258 |
| 10.4.8 | P2P 网络监控发展趋势 | 266 |
| 10.5 | 本章总结 | 269 |
| 11 | P2P 网络未来趋势 | 270 |
| 11.1 | 综合平台 | 272 |
| 11.2 | 协议标准化 | 275 |
| 11.3 | 终端统一化 | 276 |

| | |
|--------------------------|------------|
| 11.4 从 P2P 到 P4P | 278 |
| 11.5 从 IPv4 到 IPv6 | 283 |
| 11.6 P2P 和云计算 | 288 |
| 11.7 本章总结 | 293 |
| 主要参考文献 | 294 |

1 P2P 简介

1.1 P2P 定义

P2P 是“peer-to-peer”的缩写,peer 在英语里有“(地位、能力等)同等者”、“同事”和“伙伴”等含义,因此 P2P 通常被称为对等网,网络中的各个节点被称为对等体(peer)。在 P2P 网络中,每个节点的地位是对等的,它既能充当网络服务的请求者,又能对其他计算机的请求作出响应,提供资源和服务。P2P 网络利用客户端的处理能力,实现了通信与服务端的无关性,改变了目前互联网以服务器为中心的状态,重返“非中心化”。P2P 网络的本质思想实质上打破了互联网中传统的客户端/服务器(client/server,C/S)结构,令各对等体具有自由、平等通信的能力,体现了互联网自由、平等的本质。

目前,在学术界和工业界对 P2P 没有一个统一的定义,不同的研究学者和机构分别从不同的角度给出了 P2P 的定义。这些定义之间并不矛盾,均从不同侧面反映了 P2P 的内在特点。

Graham 通过三个关键条件对 P2P 进行定义:

- (1) 具有服务器质量的可运行计算机;
- (2) 具有独立于 DNS 的寻址系统;
- (3) 具有与可变连接合作的能力。

Abere 通过描述 P2P 系统的七个特征来定义 P2P 系统,这些特征如下:

- (1) 没有集中的协调中心;
- (2) 没有集中的数据库;
- (3) P2P 节点没有整个系统的全局视图;
- (4) 全局的行为依靠局部的相互作用;
- (5) 所有存在的数据和服务都是可访问的;
- (6) P2P 节点是自治的;
- (7) P2P 节点及其相互之间的连接都是不可靠的。

惠普实验室(Hewlett-Packard Laboratories)的 Milojevic 将 P2P 系统定义为一类采取分布式方式、利用分布式资源完成关键功能的系统。分布式资源包括计算能力、存储空间、数据、网络带宽,以及各种存在的可用资源。关键功能可以是分布式计算、数据内容共享、通信与协作或平台服务。分布式方式可以应用到算法、

数据、元数据或所有方面,但并不排除在系统或应用程序的某些部分保留集中式的方式。典型的 P2P 系统主要应用在互联网边缘或 Ad-Hoc 网络环境中。

Shirky 将 P2P 系统定义为一种利用互联网边缘的各种可用资源(如存储空间、计算能力、媒体内容、人力资源等)的应用程序。因为访问这些分散的边缘资源意味着要在连接不稳定和 IP 地址不可预见的环境下工作,所以 P2P 节点必须能够独立于 DNS 系统且拥有独立于集中服务器的完全的自治。

Intel 公司的 P2P 工作组(P2P Working Group, P2P WG)将 P2P 系统定义为“通过在系统之间直接进行交换来共享计算机资源和服务的系统”,这些资源与服务包括信息交换、处理器时钟、缓存和磁盘空间等。

IBM 对 P2P 的定义则更为广泛,认为 P2P 是由若干互联协作的计算机构成的系统,系统具备以下特征:

(1) 系统依存于边缘化(非中央式服务器)设备的主动协作,每个成员直接从其他成员而不是从服务器的参与中受益;

(2) 系统中成员同时扮演服务器与客户端的角色;

(3) 系统应用的用户能够意识到彼此的存在而构成一个虚拟或实际的群体。

从研究的角度看,P2P 包含三个层面的含义。

(1) P2P 实现技术:实现 P2P 应用系统时所用到的技术,包括相关协议(如 Gnutella、FastTrack 等)。

(2) P2P 通信模式:与传统的 C/S 模式不同,每个通信方都具有相同的逻辑能力,并且每个通信方都有能力发起一个通信过程。

(3) P2P 网络:由 P2P 节点、附属管理设备(如索引服务器等)及其相关应用等组成的可实现 P2P 功能的网络。它是一种运行在互联网上动态变化的逻辑网络。每个 P2P 系统都对应一个 P2P 网络。P2P 网络是一种具有较高扩展性的分布式系统结构,其对等概念是指网络中的物理节点在逻辑上具有相同的地位,而非处理能力的对等。

简单地说,P2P 可以让不同计算机用户之间不经过中继设备直接交换数据或服务。在 P2P 网络中,每个节点的地位都是相同的,具备客户端和服务端双重特性,可以同时作为服务使用者和服务提供者。由于 P2P 的飞速发展,互联网的存储模式将由目前的“内容位于中心”模式转变为“内容位于边缘”模式,改变互联网目前以大网站为中心的状态,重返“非中心化”。本书后继章节将主要讲述 P2P 网络层面的原理和应用,同时涵盖部分 P2P 实现技术和 P2P 通信模式层面的内容。

1.2 P2P 特点

P2P 其实并非一个全新的事物。事实上自互联网诞生之日起,P2P 就已经存

在,它是互联网的起源和基础。P2P 改变了目前互联网中占主导地位的客户/服务器结构中信息在消费者和生产者之间的不平衡。如图 1-1 所示,由于 P2P 网络没有中心节点(中心服务器),网络中的每个节点具有信息消费者和信息提供者的双重身份,同时具有信息通信方面的功能,因此 P2P 应用的实现扩展性强,实现方式灵活多样,部署成本低,给互联网的发布和共享带来了巨大的空间。

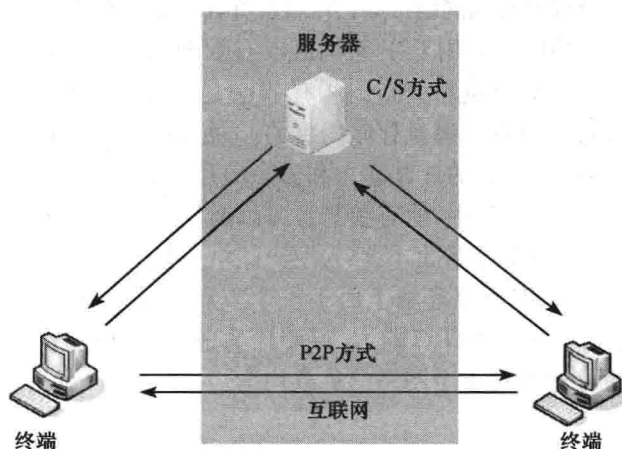


图 1-1 P2P 方式与 C/S 方式的区别

总体来说,P2P 具有以下特点。

- (1) P2P 是动态的:动态地提供信息和服务。
- (2) P2P 是双向的:切实实现信息和服务的交换与共享。
- (3) P2P 是直接的:无中介、等级和格式限制,直接交换信息和服务。
- (4) P2P 是平等的:生产者与消费者地位平等,角色合二为一。
- (5) P2P 是及时的:无服务器参与空间分配,可提供实时、可升级的信息。
- (6) P2P 是有效的:可充分利用个人计算机的硬件设备、传输信息和服务时目标确定。

与 C/S 结构相对比,P2P 的优势体现在非中心化、可扩展性、健壮性、高性能/价格比、隐私保护和负载均衡这几个方面。

(1) 非中心化(decentralization)。网络中的资源和服务分散在所有节点上,信息传输和服务的实现都直接在节点之间进行,可以无需中间环节和服务器的介入,避免了可能的瓶颈。

(2) 可扩展性(scalability)。在 P2P 网络中,节点在获取其他节点资源的同时也为其他节点服务。随着用户的加入,不仅服务的需求增加了,系统整体的资源和服务能力也在同步扩充,始终能较便捷地满足用户的需要。这就使 P2P 系统的服务能力能够随需求的增长而自然增长,具有“与生俱来”的可扩展性,能够解决传统

C/S 架构容易产生“热点效应”的问题。

(3) 健壮性(robustness)。P2P 网络通常都是以自组织的方式建立起来的,并允许节点自由地加入和离开。P2P 网络天生具有耐攻击、高容错的优点。由于服务是分散在各个节点之间进行的,当部分节点或网络遭到破坏时,对其他部分的影响很小。P2P 网络一般在部分节点失效时能够自动调整网络的整体拓扑,保持与其他节点的连通性。P2P 网络还能根据网络带宽、节点数、负载等变化不断地作自适应的调整。

(4) 高性能/价格比。采用 P2P 架构可以有效地利用互联网中散布的大量普通节点,将计算任务或存储资料分布到所有节点上。利用其中空置的计算能力或存储空间,达到高性能计算和海量存储的目的。通过利用网络中的大量空闲资源,可以用更低的成本提供更强的计算能力和更高的存储能力。

(5) 隐私保护。在 P2P 中,所有参与者可以提供中继转发的功能,因此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。

(6) 负载均衡。P2P 网络环境下由于每个节点既是服务器又是客户端,减少了对传统 C/S 结构中服务器计算能力、存储能力的要求,同时由于资源分布在多个节点,因此能更好地实现整个网络的负载均衡。

1.3 P2P 发展历史与现状

1.3.1 P2P 发展的四个阶段

P2P 并不是一项新的技术,实际上,在互联网诞生之初基本的 P2P 技术就已经出现了。在早期的互联网中,P2P 实际是占主导地位的网络结构。当时,互联网仅用于学术研究,由专业的用户操作和维护,并使用当时最强大的计算机。这些计算机既是网络客户端又是服务器。每一台主机都拥有固定的 IP 地址及域名,允许其他主机在需要时与它们通信。

1979 年,Truscott 和 Ellis 开发了早期基于 P2P 的典型应用:新闻讨论组(uses network,USENET)。它是互联网上信息传播的一个重要组成部分,也是互联网上一种高效率的交流方式。它通过由个人或公司负责维护的新闻服务器提供服务,并可管理成千上万个新闻组。USENET 通过电话线成批地进行文件交换,这通常是在夜间进行的,因为此时的长途话费更低。另一个早期 P2P 的典型应用 FidoNet 是由 Jennings 在 1984 年开发的。该软件用来在不同的通告系统(BBS)或电子消息中心的用户之间进行信息交换,通常服务于一些兴趣组并通过调制器(modem)进行访问。

1993 年,第一个流行的 Web 浏览器 Mosaic 诞生了,互联网也进入了第二个发展阶段。Mosaic 可以在一个页面上同时显示图片和文本,这使得互联网比以往更具有可访问性。万维网(world wide web,WWW)风靡一时,使互联网逐步形成

了以少数服务器为中心的客户/服务器结构。在客户/服务器结构下,对客户机的资源要求非常少,因此可以使用户以非常低廉的成本方便地连接互联网,推动了互联网的快速普及。

但是,随着互联网的逐渐普及并深入人们的日常生活,人们需要更直接、更广泛的信息交流,可以实现更多的资源和服务共享。普通用户希望能够更全面地参与到互联网的信息交流中,而计算机和网络性能的提升也使其具有了现实的可能性。在此背景下,P2P 再一次受到了广泛关注。

1999年,Napster 软件将 P2P 重新带回网络世界,P2P 应用才真正地迅速流行起来。Napster 允许对等的用户不受任何干涉地进行上传和下载。通过 Napster 提供的软件,乐迷可以共享自己硬盘上的音乐文件,同时可以搜索并下载其他用户共享出来的音乐文件。与提供免费音乐下载的 MP3.com 不同,Napster 并不提供 MP3 音乐资源,只提供动态刷新的 MP3 目录服务。Napster 在短时间里吸引了 5000 万用户,它的成功使人们意识到 P2P 扩展到整个互联网领域的可能性。

到目前为止,P2P 技术的发展经历了四个阶段:集中式、纯分布式、混合式和结构化模型。每个阶段都代表着一种 P2P 技术的网络模型,其中,混合式模型是当前最为成熟、有效的实现方式。

集中式 P2P 又被称为第一代 P2P 应用。一般来说,每种模式均需要有一个中心服务器来负责记录共享信息以及回答对这些信息的查询。每个对等体对它要共享的信息及进行的通信负责,根据需要下载它所需要的其他对等体上的信息,因此这种模式被称做“集中式 P2P”。人们通常把肖恩·范宁(Shawn Fanning)于 1998 年编写的 Napster 程序视为典型意义上 P2P 应用的开始。Napster 程序的运作模式是通过将安装 Napster 音乐共享软件的某使用者的个人计算机连接到 Napster 服务器中的一台,将该使用者可用来交换的音乐文件索引“上传”到服务器中建立资料库。资料库在其中扮演媒介的角色,它将告知需要该音乐文件的其他使用者在何处有他想要的文件。当其他使用者确定要从该使用者那里取得音乐文件时,两个用户间即开始建立连接进行下载。由于 Napster 这类软件需依靠中央服务器来管理集中的文件列表,因此这一代的 P2P 文件交换系统的生命力十分脆弱,只要关闭服务器,交换就停止。

纯分布式 P2P 被称为第二代 P2P 应用,它的出现是为了从根本上改进第一代 P2P 软件的缺陷。纯分布式 P2P 网络中不存在中枢服务器,所有的服务及其相关信息完全散布于各个 P2P 节点中,因此它最显著的特点就是“完全的去中心化”。纯分布式 P2P 采用了随机图的组织方式来形成一个松散的网络,在网络中一般采用泛洪的方式查询和定位资源,并通过 TTL 等机制进行限制。Gnutella 是第二代 P2P 应用的典型代表。Gnutella 采用完全分布式的模式,不需要中央服务器的支持,其技术原理简单来说就是在计算机(节点)间相互发送搜索请求,找到文件后再

将信息传回搜索者的计算机(节点)。该种方式需要占用大量网络资源,同时,局部性能较差的节点还将造成 Gnutella 网络分片,从而导致整个网络可用性变差;并且由于网络直径不可控,使得这种网络的可扩展性较差。

混合式 P2P 网络拓扑结构由多个簇构成,每个簇由一个超级节点所代表,不同簇之间通过分布式的 P2P 网络拓扑结构将超节点连接起来。例如,eDonkey 将网络分为服务器层和客户层。客户加入 eDonkey 网络可以获取或者共享文件,服务器的作用是提供文件索引信息和服务器列表,它不传递实际的文件数据,每个客户都可能成为服务器。每个客户连接到一个服务器上进行文件查询或者服务器列表更新,服务器之间自组织成服务器层网络以交换文件索引和服务器列表信息。总体来说,基于超级节点的混合式 P2P 网络比以往有很大程度改进,但超级节点本身的脆弱性也可能导致其簇内的节点处于孤立状态。因此,这种局部索引方法仍存在一定的局限性。

结构化 P2P 应用具有自组织和负载均衡等特点,可消除混合式 P2P 应用中的单节点失效问题,增强了网络扩展性,被称为第四代 P2P 应用。结构化 P2P 具有严格的拓扑结构,一般采用基于分布式散列表(distributed Hash table,DHT)的分布式哈希算法。网络中各节点并不需要维护整个网络的信息,只在节点中存储其邻近的后继节点信息,输入的信息将通过分布式散列函数唯一地映射到某个节点上,通过一些特定的路由算法和该节点建立连接。这种方式有效地减少了节点信息的发送数量,增强了网络的扩展性,避免了存在类似 Napster 的中央处理器,以及像 Gnutella 那样进行泛洪方式查询。目前,基于 DHT 的第四代 P2P 应用有 Chord、Pastry 和 CAN 等网络模型。在第 2 章中将详细介绍这几类 P2P 网络。

1.3.2 国外 P2P 技术的研究现状

21 世纪初,国外大多数著名的学术团体和技术组织均成立或者完善了专门的 P2P 研究组,其中较知名的有:MIT 的 Chord/CFS 研究组(承担 IRIS 计划),UC Berkeley 的 Tapestry/OceanStore 研究组,Microsoft Research 和 Rice University 的 Pastry/PAST 研究组,Intel 公司成立的 P2P 工作组和全球网络论坛(global grid forum,GGF)等。

2003 年,麻省理工学院(MIT)Kaashoek 教授领衔的容错的互联网系统架构项目(Infrastructure for Resilient Internet System,IRIS),用 P2P 的方法去研究并建立新一代互连网络架构,得到 2003 年美国自然科学基金(NSF)在 IT 领域最大的一项基金资助。IRIS 计划的核心技术是 DHT 技术,MIT 的 Chord 项目已经实现基于 DHT 的查询。

2000 年 3 月,加利福尼亚大学(UC)Berkeley 分校以 Zhao 等为首的研究者开始 Tapestry 的开发工作。Tapestry 是一个基于 Plaxton Mesh 的覆盖网结构、面

向广域分布式数据存取、容错的超立方体结构 P2P 模型,它在构建网络时就考虑了拓扑一致性。2003 年,该校的 Tapestry 研究组发布了使用 Java 编写的面向 Linux 操作系统的 Tapestry 2.0 版本软件,该软件可以提供分布式应用的底层覆盖网络或者作为其分布式散列表。OceanStore 是一个基于 Tapestry 的分布式数据存取系统,其目标是提供全球范围广域、持久性的数据存取服务。2004 年 6 月, OceanStore 在 SourceForge 上发布了原型软件及其源代码。

Microsoft Research 的 Rowstron 和 Rice University 的 Druschel 等在 2000 年开始 Pastry 的设计及其应用的开发。Pastry 是一个容错、高效、可扩展的混合式结构 P2P 网络,结合了环形结构与超立方体结构的优点,在互联网上构造了一个分布式、自组织、容错的覆盖网,提供高效的路由、确定性的对象定位和独立于具体应用的负载平衡。Pastry 有着广泛的应用,其中著名的 PAST 归档存储系统就是以 Pastry 作为底层构建的,在互联网上提供安全、高可用、持久性的数据存取服务。

2000 年 8 月, Intel 公司宣布成立 P2P 工作组,正式开展 P2P 研究。该工作组成立的主要目的是加速 P2P 基础设施的建立和相应的标准化工作。P2P 工作组成立后,对 P2P 的术语进行了统一,也形成了相关的草案,但在标准化工作方面进展缓慢。目前, P2P 工作组已经和全球网络论坛合并,由该论坛管理 P2P 相关的工作。同时,工作组成立以后积极与应用开发商合作开发 P2P 应用平台。2002 年, Intel 公司发布了基于 .NET 基础架构之上的 P2P Accelerator Kit(P2P 加速工具包)和 P2P 安全 API 软件包,从而使得利用 Microsoft .NET 开发软件的人员能够迅速建立 P2P 安全 Web 应用。IBM 和 HP 等公司也是 Intel 公司成立的 P2P 工作组的成员,这两家公司在 2000 年 9 月共同推出了一种开放存储技术。这一存储技术利用了 P2P 技术,可以方便地从用户的硬盘向服务器上复制数据。HP 公司还把 P2P 的立足点放在打印技术上,该公司新推出的网络打印技术可以使用户通过 P2P 网络共享打印机。

Sun 公司推出了基于 Java 的开源 P2P 项目——JXTA。JXTA 主要包括一个独立于编程语言、系统平台和网络平台的协议集, JXTA 将建立核心的网络计算技术,提供支持在任何平台、任何地方、任何时间实现 P2P 计算的一整套简单、小巧和灵活的机制。JXTA 项目采用开放源码的方式,因此吸引了大量业界人士参与到 JXTA 技术的研究与应用当中。JXTA 提供了基础性的机制解决当前分布计算应用中面临的问题,实现新一代统一、安全、互操作以及异构的应用。将来 JATX 将不受到内存的限制而支持更多小型移动设备。JXTA 通过 Java 技术和 XML 数据表达的结合,提供了强大的功能使得垂直应用得以交互,并且可以克服目前 P2P 软件中的限制。同时,通过小型、简单、便于开发的构造模块, JXTA 将使开发者从建立各自框架的复杂工作中得以解放,可以潜心关注于建设各类新颖、创造性的分布式计算应用。

1.3.3 国内 P2P 技术的研究现状

国内一些高校和研究机构在 2003 年前后也纷纷开展了 P2P 的研究开发工作,其中,较知名的有北京大学网络实验室开发的 Maze 网络文件系统,清华大学高性能计算研究所开发的 Granary 对等计算存储服务系统,华中科技大学集群与网格计算实验室开发的 AnySee 视频直播系统,广州数联软件技术有限公司开发的国内最大的 P2P 用户分享平台 POCO,深圳市点石软件技术有限公司开发的 P2P 网络娱乐平台 OP 等。

Maze 是北京大学网络实验室开发的一个中心控制与对等连接相融合的对等计算文件共享系统,在结构上类似于 Napster,对等计算搜索方法类似于 Gnutella。网络上的一台计算机,不论是在内网还是外网,都可以通过安装运行 Maze 的客户端软件自由加入和退出 Maze 系统。每个节点可以将自己的一个或多个目录下的文件共享给系统的其他成员,也可以分享其他成员的资源。Maze 支持基于关键字的资源检索,也可以通过好友关系直接获得。

Granary 是清华大学高性能计算研究所自主开发的对等计算存储服务系统。它以对象格式存储数据。另外,Granary 设计了专门的节点信息收集算法 Peer-window 和结构化覆盖网络路由协议 Tourist。

AnySee 是华中科技大学集群与网格计算实验室开发的视频直播系统。它采用了一对多的服务模式,支持部分 NAT 和防火墙的穿越,提高了视频直播系统的可扩展性;同时,它利用近播原则、分域调度的思想,使用 Landmark 路标算法直接建树的方式构建应用层上的组播树,克服了 ESM 等一对多模式系统由连接图的构造和维护带来的负载影响。

广州数联软件技术有限公司开发的 POCO 是中国最大的 P2P 用户分享平台。它采用有安全、流量控制力的且无中心服务器的第三代 P2P 资源交换平台,也是世界范围内少有的盈利的 P2P 平台。目前已经形成了 2600 万海量用户,平均在线 58.5 万,在线峰值突破 71 万,并且全部是宽带用户的用户群。POCO 现已成为中国地区第一的 P2P 分享平台。

深圳市点石软件技术有限公司开发的网络娱乐内容平台 OP(OPenext Media Desktop),可以最直接的方式找到用户想要的音乐、影视、软件、游戏、图片、书籍,以及各种文档,随时在线共享文件容量数以亿计的影视、音乐、图片。OP 整合了 Internet Explorer、Windows Media Player、RealOne Player 和 ACDSec,是国内基于 P2P 的网络娱乐内容平台。

1.3.4 P2P 的网络流量

目前,全球的宽带内容,也就是大流量内容的提供商,几乎都在试图借助 P2P