

# 数据中心

## 的网络互联结构和 流量协同传输管理

郭得科 陈涛 罗来龙 李妍 著

清华大学出版社



A background network diagram consisting of several dark grey circular nodes of varying sizes connected by thin, light grey lines. The nodes are scattered across the page, with a larger node at the top right and several smaller ones elsewhere. The overall structure is a complex, interconnected web.

# 数据中心

## 的网络互联结构和 流量协同传输管理

郭得科 陈涛 罗来龙 李妍 著

清华大学出版社  
北京

## 内 容 简 介

本书对数据中心及其网络互联结构的现状和发展趋势进行了深度剖析,深入介绍了一些新型网络互联结构的设计与优化方法,力求满足数据中心网络的高带宽、高容错、高可扩展性等方面的需求,并通过引入数据中心内关联性流量的协同传输机制,实现对数据中心现有传输能力的高效利用。

本书可用作高等学校计算机专业、软件工程专业、信息工程专业以及其他相近专业的教材或教学参考书,也可供这些专业的研究人员和工程技术人员阅读。

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 13701121933

### 图书在版编目(CIP)数据

数据中心的网络互联结构和流量协同传输管理/郭得科等著. —北京:清华大学出版社,2016

ISBN 978-7-302-44810-5

I. ①数… II. ①郭… III. ①计算机网络—网络结构—研究 ②计算机网络—流量—研究 IV. ①TP393

中国版本图书馆 CIP 数据核字(2016)第 200801 号

责任编辑:薛 慧

封面设计:何凤霞

责任校对:王淑云

责任印制:王静怡

出版发行:清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址:北京清华大学学研大厦 A 座 邮 编:100084

社总机:010-62770175 邮 购:010-62786544

投稿与读者服务:010-62776969, [c-service@tup.tsinghua.edu.cn](mailto:c-service@tup.tsinghua.edu.cn)

质量反馈:010-62772015, [zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

印 装 者:三河市金元印装有限公司

经 销:全国新华书店

开 本:170mm×240mm 印 张:17.75 插 页:1 字 数:318千字

版 次:2016年11月第1版 印 次:2016年11月第1次印刷

印 数:1~1500

定 价:59.00元

---

产品编号:062164-01

# 前言

## 1. 背景

数据中心的出现源于人们对海量数据的高效存储和处理需求。互联网的蓬勃发展和社会的数字化变革,导致网络上的数据呈爆炸式增长,出现了越来越多需要进行大规模数据存储和处理的应用需求。数据中心的规模和应用领域不断扩展,已经渗透到经济、科技、军事以及人们日常生活等各个方面。总体而言,数据中心旨在依据特定网络结构互联大规模服务器和网络设施等硬件资源,形成计算、存储、网络等资源的规模效应和整体优势,进而面向各类上层应用提供网络化存储、网络化计算等弹性服务。

现代社会信息量的爆炸式增长、资源复用技术的成熟和宽带网络的普及,共同促进了云计算的诞生和发展。数据中心可为云计算提供大规模可扩展的基础物理资源,并在云平台的辅助下为用户提供多种类型的云服务。云计算模式的出现拓展了数据中心的使用方式,而数据中心的建设和发展也为云计算的推广和应用奠定了坚实的基础。另外,随着大数据时代的来临,如何从类型多样、规模巨大的大数据中快速提取有价值的信息成为关键。数据中心可为大数据应用提供基础平台,并在解决大数据的存储、大规模分析处理等难题方面具有天然的优势。

云计算和大数据等新技术和应用推动了现代数据中心的快速发展,并使其成为国家和IT企业的核心信

息基础设施。数据中心具有巨大的商业价值和社会效益,其应用领域非常广泛,涉及信息社会的诸多行业和领域。数据中心的网络化存储和网络化计算为很多大规模数据处理模式提供基础服务,而这类数据处理模式在大数据领域、物联网、科学应用等领域也得到了广泛的重视和应用。同时,数据中心也面临着不断的变革,新一代数据中心由数以万计的服务器组成,并通过特定的网络结构互联为一个整体,共同形成一个分布式的计算和存储网络。新一代数据中心因为内在的高可扩展性、高容错性等优势得到了业界的高度关注,并向着虚拟化、软件定义、模块化、绿色节能以及自动化运行维护等方向迈进。

数据中心研究的一个重要理论基础是数据中心网络,其作为基本要素在新一代数据中心中具有至关重要的地位。数据中心网络所考虑的不仅是设备之间的通信协议,更主要的是把交换机和服务器作为一个整体进行拓扑互联、性能优化、资源管理和能耗控制,形成数据中心基础设施在网络化计算能力、网络化存储能力和网络通信能力等方面的综合优势。换句话说,数据中心网络不仅是连接大规模服务器的桥梁,而且是承载网络化存储和网络化计算的基础。数据中心支持的业务往往伴随着服务器之间密集的数据交互,网络资源已成为影响数据中心服务质量的瓶颈,且直接关系到各类用户对数据中心的使用体验。因此,迫切需要开展数据中心网络方面的基础理论研究,从而推动数据中心及相关应用领域的发展。另外,数据中心网络是互联网的重要组成部分,该领域的研究进展也会对下一代互联网的发展产生一定的推动作用。与互联网相比,数据中心具有集中管控等鲜明特点,这有利于开展网络创新技术的探索。数据中心运营商为了提高服务性能和收益,会根据应用需求定制网络架构和协议,并进行创新网络技术的部署。

当前,数据中心的可持续发展面临着一系列关键基础性问题,而数据中心网络是其中一个重要的研究方向。数据中心网络的研究面临很多基础理论和关键技术方面的挑战,例如网络功能的灵活定制、横向可扩展的互联结构、网络资源的高效复用、网络虚拟化、关联性流量的协同传输以及网络能耗的协同控制等。

(1) 网络功能的灵活定制。传统网络设备的控制平面和数据平面紧密耦合,不具有动态性和灵活性,致使传统的数据中心网络只能提供有限且已知的网络功能和服务。如今,数据中心的网络应用需求变得日益丰富和灵活多样,为了提高数据中心的服务质量,需要针对不同的应用需求对网络功能进行灵活配置。而数据中心网络流量难以预测、网络设备可靠性低等环

境特征,也对网络的可动态灵活配置功能提出了新的需求,致使数据中心的传统网络架构面临着严峻的挑战。

(2) 横向可扩展的网络互联结构。数据中心必须具备计算、存储、网络资源按需扩展的能力。有效互联上万台甚至更大规模的服务器,是数据中心提供网络化计算和网络化存储的前提。依靠扩充交换机端口数量或提升端口速率的纵向扩展方法已经远远不能满足数据中心的规模扩展需求。因此,迫切需要对数据中心的规模扩展方式进行改进,探索各种横向扩展模式,进而连接更多的交换机和服务器以实现计算性能和存储容量的按需扩展。

(3) 数据中心网络资源的高效复用。数据中心所采用的网络协议基本源自广域网环境,致使数据中心网络资源的利用率很低。在软件定义的可定制网络框架下,如何设计新型路由和传输协议,以提高数据中心网络的资源利用率并进而提升上层应用的性能,是非常具有挑战性的问题。此外,软件定义的数据中心网络架构,为网络资源、计算资源和存储资源的联合优化提供了新的发展机会。

(4) 数据中心的网络虚拟化。在大量用户竞争使用网络资源的数据中心环境下,要实现网络资源的有效共享和安全隔离。虚拟化是保障数据中心安全和实现资源复用的重要技术,每个用户租用的多个虚拟机之间形成了虚拟数据中心网络。不同用户的虚拟网络之间竞争使用实际的物理网络,而且从安全考虑需要被有效隔离。当前,数据中心普遍采用“尽力而为”的方式共享网络资源,不能很好地支持虚拟数据中心网络的流量隔离和带宽保障需求。

(5) 关联性流量的协同传输问题。数据中心中密集的数据交互行为产生了庞大的“东西向流量”。Multicast、Incast 以及 shuffle 传输是“东西向流量”的主要组成部分。此外,在组成一个 multicast、shuffle、incast 的众多数据流之间存在很大的数据关联性,进而存在非常大的数据流聚合增益。通过和上层应用的联合设计优化,可在不影响应用效果的前提下大幅降低关联性流量造成的网络传输开销,进而降低对数据中心稀缺网络带宽的消耗。

(6) 数据中心网络能耗的协同控制。对数据中心进行高效的能耗管理具有重要的经济效益和社会影响。为此,需要从底层硬件节点、上层协议运行、外围供能系统等环节进行能耗控制的联合优化。其中,数据中心网络层面的能耗控制必不可少,且其能耗控制策略也直接影响到计算设备层面的

能耗控制策略。为了实现多维度的协同能耗控制,数据中心网络需要研究如何实时感知网络能耗、如何在不影响网络性能和可靠性的前提下实现节能流量工程,并尽可能地使用清洁能源。

我们在数据中心的横向可扩展的网络互联结构和关联性流量的协同传输领域进行了一系列深入而系统的研究工作。本书以数据中心的可扩展网络互联结构为基础,深入探讨了一些新型网络互联结构的设计与优化方法,力求满足数据中心网络的高带宽、高容错、高可扩展性等方面的需求,并通过引入数据中心内关联性流量的协同传输机制,实现对数据中心现有传输能力的高效利用。本书绝大部分内容取材于我们近期在国内外重要学术期刊和会议上发表的论文,全面系统地展示了相关领域的很多新的研究成果和进展。

## 2. 内容安排

本书共 10 章,从结构上可分为 3 个部分。

第 1 部分是对数据中心及数据中心网络互联结构发展现状的介绍,包括第 1 章和第 2 章。

第 1 章首先介绍了数据中心的起源和发展、云计算和大数据对数据中心的内在需求以及新一代数据中心的发展趋势。在此基础上,从应用角度详述了数据中心在网络化存储、网络化计算以及数据分析处理等领域的应用现状。最后,从网络功能的灵活定制、横向可扩展的互联结构、网络资源的高效复用、网络虚拟化、关联性流量的协同传输、网络能耗的协同控制等角度探讨了数据中心网络发展所面临的重要挑战。

第 2 章对当前数据中心网络的最新互联结构进行了综述,从构建规则、路由算法、网络性能等方面进行了对比分析。同时,将当前数据中心的网络互联结构按照 5 种类型进行归类,以揭示数据中心网络互联结构设计理念的发展和变化,分别是交换机为核心的互联结构、服务器为核心的互联结构、模块化数据中心的互联结构、随机型数据中心的互联结构以及无线数据中心的互联结构。最后,对未来数据中心网络互联结构的发展趋势提出了一些观点。

第 2 部分介绍数据中心的新型网络互联结构,包括第 3~6 章。

第 3 章介绍以服务器为核心的数据中心互联结构 HCN 和 BCN。在服务器网络端口数目固定不变的情况下,设计了由这类同构服务器互联而成的数据中心网络互联结构。首先采用复合图(compound graph)理论设计

以服务器为核心的网络互联结构 HCN,为只具有两个网络端口的普通服务器和多端口的低成本网络交换机提供高效互联,使其兼具无损可扩展和持续可扩展能力。DCell 和 BCube 等同类数据中心的网络互联结构中每台服务器的网卡和连线会依据网络规模的扩大而增长,其最大规模被每台服务器的网卡数目所限定,不具备持续可扩展能力。在此基础上,在相同的服务器网卡配置和网络直径前提下,构造出规模尽可能大的数据中心网络互联结构 BCN。随后介绍了这两种网络互联结构的高效及容错路由机制,并通过数学分析和综合仿真模拟结果验证了 HCN 和 BCN 的良好拓扑特征。

第 4 章介绍模块化数据中心互联结构 DCube。构建大规模数据中心有两种截然不同的趋势:第 1 种趋势是通过类似 DCell、BCube、HCN 等扩展性网络互联结构,构造出单体大规模数据中心;第 2 种趋势是在大量单体数据中心的基础上,通过模块化的互联结构构造大规模数据中心。本章介绍为模块化数据中心设计的一组模块内网络互联结构 DCube,包括 H-DCube 和 M-DCube,每个 DCube 互联大量配备双网卡的服务器和低成本交换机。大量 DCube 互联结构的数据中心模块进一步互联可形成全新的模块化数据中心。DCube( $n, k$ )由  $k$  个互联的子网络构成,每个子网络都由许多基本构建模块和标准超级立方体结构(或其变种结构 1-möbius)按照复合图理论构成。在此基础上,分析了 H-DCube 和 M-DCube 的路由机制,并与 BCube 和 Fat-Tree 进行了性能分析比较。如果采用诸如 Twisted cube、Flip MCube 和 Fastcube 等标准超级立方体结构的其他变种结构,本章提出的设计方法仍然适用。

第 5 章介绍一种数据中心网络的混合互联结构设计方法,并讨论了一种具体的混合互联结构 R3。如第 2 章所分析,当前的数据中心网络互联结构普遍采用两种设计思路,分别是完全规则的互联结构设计和完全随机的互联结构设计。虽然这两种类型的网络互联结构具有特定的优势,但是也存在内在的缺点和不足。本章论述了一种基于复合图理论的混合互联结构设计方法,可兼容当前的规则互联结构和随机互联结构。进一步地,还提出了一种混合互联结构 R3。该结构采用随机正则图作为基本单元,并采用通用超级立方体这种规则结构将这些基本单元进行互联。该混合结构兼具随机正则图和通用超级立方体的拓扑优势,并可有效避免二者的拓扑缺陷。

第 6 章介绍在数据中心的网络结构中额外引入可见光通信(visible light communication, VLC)链路后,设计无线链路和有线链路混合的网络互联结构 VLCcube,从而提升数据中心的网络性能。具体而言,在 Fat-Tree 这一具有代表性的数据中心网络结构基础之上,在每个机架顶部安装 4 个

VLC 收发装置,即可提供 4 条 10Gbps 左右的无线链路,全体机架上的无线链路组网成为无线 Torus 结构。本章重点介绍了混合拓扑的构建规则、路由策略、批处理流量和在线流量的拥塞感知调度策略,并开展了相关实验验证工作。因为 Fat-Tree 中很多原本 4 跳的数据流可切换到无线 Torus 结构进行短距离传播,VLCcube 取得了比 Fat-Tree 更好的网络性能,而且设计的拥塞感知调度策略可使 VLCcube 的性能进一步得到提升。VLCcube 仅是数据中心中利用 VLC 链路的一种可选方案,未来供应商可基于不同的有线网络互联结构设计完全不同的混合网络互联结构。VLC 链路的引入不仅能和已有的数据中心网络良好地兼容,而且可有效提升数据中心的网络性能和网络设计的灵活性。

第 3 部分是数据中心内关联性流量的协同传输管理。虽然新型互联结构的研究不断提高数据中心的网络传输能力,但是对数据中心现有传输能力的高效利用同样重要。数据中心支持多种分布式计算框架,这些计算框架普遍采用流式计算模型,相邻处理阶段间的大量数据交互产生了严重的东西向流量,其中 multicast、incast、shuffle 等关联性流量占相当大的比重,进而严重影响到上层应用的性能。优化和管理这些关联性流量对高效使用数据中心的网络资源和提升处理作业的性能至关重要。这部分内容包括第 7~10 章。

第 7 章介绍关联性流量 incast 的协同传输管理问题,包括网内聚合和协同传输两个环节。当前很多上层应用决定了 incast 的全体数据流之间存在数据相关性,并在相同的接收端被执行聚合操作。这就促使我们考虑在这些关联性流量的网内传输阶段应尽可能早地而不是仅在流量的接收端进行数据聚合。本章首先以新型数据中心网络结构为背景讨论关联性流量之间数据聚合的可行性和增益,随后探讨实现该网内聚合所必须的基于 incast 树的协同传输方法。为最大化网内聚合的增益,我们为 incast 传输建立最小代价树模型,并设计了两种近似的 incast 树构造方法,其能够仅基于 incast 成员的位置和数据中心拓扑结构生成一棵有效的 incast 树。本章进一步介绍了 incast 树面临的多种动态和容错问题,最后通过实验发现我们所提出的网内聚合方法能大幅度降低 incast 流量造成的传输开销,从而节约了数据中心的网络资源。本章虽然选用 BCube 这种以服务器为核心的数据中心网络互联结构为研究背景,但是提出的关联性流量网内聚合理念也适用于其他类型的数据中心网络互联结构。本章工作是第 8 章和第 9 章的前提。

第 8 章介绍关联性流量 shuffle 的协同传输管理问题,包括网内聚合

和协同传输两个环节。受关联性流量 incast 网内聚合的启发,本章介绍了如何将关联性流量 shuffle 原本在诸多接收端执行的流量聚合操作推送到网络传输环节中执行,通过降低网络内的流量传输从而高效地利用网络资源。首先,针对新型数据中心结构 BCube 中的 shuffle 流量进行网内聚合问题的建模,并提出两种近似方法来高效地构建 shuffle 聚合子图,依据该结构进行流量的协同传输可有效实现预期的网内聚合。本章还介绍了基于布鲁姆滤波器(Bloom filters)的可扩展流量转发模式,从而为大量并存的 shuffle 传输实现各自预期的网内聚合效果。尽管本章选用以 BCube 为依托的网络互联结构,但是提出的关联性流量 shuffle 的网内聚合理念也适用于其他类型的数据中心网络互联结构。

第 9 章介绍不确定关联性流量 incast 的协同传输管理。当上层应用在数据中心内产生关联性 incast 流量时,其内部诸多数据流的发送端和接收端已经确定。第 7 章针对这类关联性流量介绍了如何实施网内聚合和协同传输。但是,很多数据中心应用面临计算节点和存储节点选择的多样性,不同的选择方案会导致对应的 incast 表现出不同的发送端和接收端。这类关联性流量被定义为不确定性 incast,与确定性 incast 流量相比,其具有更多机会来获取更大的数据流网内聚合增益。本章首先深入剖析了不确定 incast 流量的网内聚合问题,并设计了对应的协同传输方法以获取尽可能大的数据流网内聚合增益,包括为 incast 的各个数据流初始化发送端和构造 incast 聚合树两个环节。数据流发送端初始化的目标是令初始化后的全体发送端形成最少数目的群组,每个群组输出的数据流在传输的下一跳网络设备上即可被全部聚合为一个新数据流。为了充分利用初始化环节产生的这种优势,本章提出了两种 incast 聚合树的构建算法。实验结果表明,从减少网络流量和节省网络资源的角度来看,不确定性 incast 传输要优于确定性 incast 传输。

第 10 章介绍关联性流量 multicast 的协同传输管理。多播协议的出发点是从一个发送端将相同的内容传输给一组接收端,进而有效节约网络带宽并降低发送端的负载。数据中心的分布式文件系统为每个数据块提供多个副本,此时传统的 multicast 面临发送端的多样性问题,不再依赖于某个唯一选定的发送端,每个接收端只需从其中一个发送端获得发送内容。本章关注如何使这种发送端不确定的多播造成的网络传输代价尽可能地小,提出了对应的链路代价最小多播森林(minimum cost forest, MCF)模型。针对确定性 multicast 的方法不适用于 MCF 这个 NP 难问题,为此本章提出了两种高效的 MCF 近似算法,即 P-MCF 和 E-MCF。本章在 3 种类

型的数据中心互联结构(随机网络、随机正则网络以及无标度网络)下对MCF问题进行了评估。结果显示:3种网络互联结构下不确定性multicast的MCF都比确定性multicast的最小斯坦纳树占用更少的网络链路资源。

我们的研究工作得到了国家自然科学基金优秀青年科学基金项目(No.61422214)、国家重点基础研究发展计划(“973”计划)青年科学家项目(No.2014CB347800)、湖南省自然科学杰出青年基金项目(No.2016JJ1002)、教育部新世纪优秀人才计划项目以及国防科学技术大学杰出青年基金项目(No.JQ14-05-02)的资助。国防科学技术大学的研究生赵亚威参与了本书1.1节和1.2节内容的撰写工作,研究生胡智尧、任棒棒、史良等同学参与了本书的排版、图表绘制、整理等工作,在此一并表示感谢。

由于作者水平所限,加之数据中心网络的互联结构和内部流量管理的研究仍处于快速发展和变化之中,书中错误和不足之处在所难免,恳请专家、读者予以指正。

郭得科

2016年2月于长沙

# 目录

## 第 1 部分 基础知识

第 1 章 数据中心简介 .....	3
1.1 起源与发展 .....	3
1.1.1 数据中心的概念及分类 .....	3
1.1.2 云计算对数据中心的需求 .....	6
1.1.3 大数据对数据中心的需求 .....	8
1.1.4 新一代数据中心的发展 .....	9
1.2 数据中心的应用领域 .....	14
1.2.1 基于数据中心的网络化存储 .....	15
1.2.2 基于数据中心的网络化计算 .....	16
1.2.3 基于数据中心的大数据应用 .....	17
1.3 数据中心网络面临的挑战 .....	19
1.3.1 功能可灵活定制的数据中心网络 .....	20
1.3.2 横向可扩展的数据中心网络 .....	21
1.3.3 数据中心网络资源的高效复用 .....	22
1.3.4 数据中心的网络虚拟化 .....	23
1.3.5 关联性流量的协同传输问题 .....	24
1.3.6 数据中心网络能耗的协同控制 .....	25
参考文献 .....	27
第 2 章 数据中心网络互联结构的研究现状 .....	29
2.1 引言 .....	29
2.2 以交换机为核心的网络互联结构 .....	31

2.2.1	树型互联结构	32
2.2.2	扁平化互联结构	34
2.2.3	光交换互联结构	37
2.3	以服务器为核心的网络互联结构	39
2.3.1	基于复合图的互联结构	40
2.3.2	基于非复合图的互联结构	42
2.4	模块化数据中心的互联结构	44
2.4.1	模块内的互联结构	45
2.4.2	模块间的互联结构	45
2.5	随机型数据中心的网络互联结构	47
2.5.1	基于小世界模型的数据中心互联结构	48
2.5.2	基于随机正则图的数据中心互联结构	49
2.5.3	基于无标度网络的数据中心互联结构	49
2.6	无线数据中心的网络互联结构	50
2.6.1	基于 60GHz 通信技术的混合互联结构	51
2.6.2	基于 60GHz 通信技术的全无线互联结构	52
2.6.3	基于自由空间通信技术的互联结构	52
2.6.4	基于可见光通信的互联结构	53
2.7	互联结构设计方法的演进和趋势	53
2.7.1	数据中心互联结构设计方法的演进	53
2.7.2	数据中心网络互联结构的发展趋势	55
	参考文献	57

## 第 2 部分 数据中心的新型网络互联结构

第 3 章	以服务器为核心的数据中心互联结构 HCN	63
3.1	引言	63
3.2	HCN 互联结构	65
3.2.1	复合图的基本理论	65
3.2.2	HCN 互联结构的构建方法	67
3.3	BCN 互联结构	69
3.3.1	BCN 互联结构的描述	69
3.3.2	BCN 互联结构的构建方法	71
3.4	BCN 互联结构的路由机制	75

3.4.1	单播通信的单路径路由	75
3.4.2	单播通信的多路径路由	78
3.4.3	容错路由模式	79
3.5	性能评估	82
3.5.1	网络规模	83
3.5.2	网络直径和节点度	85
3.5.3	连通性和路由路径多样性	85
3.5.4	路径长度	87
3.6	相关讨论	88
3.6.1	扩展至多端口服务器	88
3.6.2	位置关联的任务部署	89
3.6.3	服务器路由的影响	89
	参考文献	89
<b>第4章</b>	<b>模块化数据中心互联结构 DCube</b>	<b>91</b>
4.1	引言	91
4.2	DCube 互联结构	93
4.2.1	DCube 互联结构的设计思想	93
4.2.2	H-DCube 互联结构	95
4.2.3	M-DCube 互联结构	95
4.3	DCube 的单播单径路由	97
4.3.1	H-DCube 的单路径路由方法	97
4.3.2	M-DCube 的单路径路由方法	99
4.4	DCube 的单播多径路由及组播传输	102
4.4.1	H-DCube 的多路径路由方法	103
4.4.2	M-DCube 的多路径路由方法	104
4.4.3	组播传输的速率提升	106
4.5	性能评估	107
4.5.1	单播和组播的传输加速能力	108
4.5.2	聚合瓶颈吞吐量	109
4.5.3	成本与布线复杂度的量化比较	110
4.5.4	评估小结	112
4.6	相关问题讨论	114
4.6.1	任务的局部性部署	114

4.6.2	服务器配备更多的 NIC 端口 .....	114
4.6.3	服务器参与路由决策的影响 .....	115
	参考文献 .....	115
<b>第 5 章</b>	<b>数据中心的混合互联结构 R3 .....</b>	<b>117</b>
5.1	引言 .....	117
5.2	混合互联结构的设计方法 .....	118
5.2.1	混合互联结构概述 .....	119
5.2.2	R3: 基于复合图的互联结构 .....	120
5.2.3	混合互联结构数据中心的部署策略 .....	123
5.3	R3 的路由方法 .....	123
5.3.1	基于边着色的标识符分配方法 .....	125
5.3.2	基于标识符的路由方法 .....	125
5.4	R3 的拓扑优化 .....	128
5.4.1	R3 结构设计的影响因素 .....	128
5.4.2	R3 拓扑优化策略 .....	128
5.5	R3 的规模渐进扩展问题 .....	130
5.5.1	现有单元簇中添加节点 .....	131
5.5.2	额外添加新的单元簇 .....	131
5.6	性能评估 .....	133
5.6.1	路由方法的时间开销 .....	133
5.6.2	布线成本比较 .....	134
5.6.3	网络性能比较 .....	135
5.6.4	相关问题讨论 .....	136
	参考文献 .....	137
<b>第 6 章</b>	<b>基于可见光通信的数据中心无线互联结构 .....</b>	<b>139</b>
6.1	问题背景 .....	140
6.1.1	研究动机 .....	140
6.1.2	相关工作 .....	141
6.2	无线互联结构 VLCcube 的设计 .....	142
6.2.1	数据中心内引入 VLC 链路的可行性 .....	143
6.2.2	VLC 信号的干扰问题 .....	143

6.2.3	VLCcube 互联结构设计	145
6.3	VLCcube 的路由设计和流调度策略	149
6.3.1	VLCcube 中的混合路由算法	149
6.3.2	数据流调度问题的建模	150
6.3.3	数据流的批量调度方法	151
6.3.4	数据流的在线调度方法	154
6.4	VLCcube 的性能评估	155
6.4.1	实验设置与实验方法	155
6.4.2	VLCcube 的拓扑性质	156
6.4.3	Trace 流量下的网络性能	158
6.4.4	Stride-2k 流量下的网络性能	158
6.4.5	随机流量下的网络性能	159
6.4.6	拥塞感知的流调度算法评估	160
6.5	相关问题讨论	162
	参考文献	163

### 第 3 部分 数据中心的流量协同传输管理

第 7 章	关联性流量 Incast 的协同传输管理	167
7.1	引言	167
7.2	Incast 传输的网内数据聚合	169
7.2.1	Incast 数据流间网内聚合的可行性分析	170
7.2.2	Incast 最小代价树的建模	171
7.2.3	基于 Incast 树的流间数据聚合实现	172
7.3	高效 Incast 聚合树的构建	174
7.3.1	Incast 树的构建方法	174
7.3.2	Incast 最小代价树的构造方法	177
7.3.3	发送端动态行为的处理方法	178
7.3.4	接收端动态行为的处理方法	179
7.4	相关问题讨论	180
7.4.1	通用 Incast 传输模式	180
7.4.2	其他数据中心结构下的 Incast 传输模式	180
7.4.3	作业特征对 Incast 网内聚合性质的影响	181
7.5	性能评估	182

7.5.1	原型实现 .....	182
7.5.2	数据中心规模对聚合增益的影响 .....	183
7.5.3	Incast 传输规模对聚合增益的影响 .....	185
7.5.4	聚合率对聚合增益的影响 .....	186
7.5.5	Incast 成员分布对聚合增益的影响 .....	187
	参考文献 .....	188
<b>第 8 章 关联性流量 Shuffle 的协同传输管理 .....</b>		<b>193</b>
8.1	引言 .....	193
8.2	Shuffle 传输网内聚合 .....	194
8.2.1	问题建模 .....	194
8.2.2	Incast 聚合树的构造方法 .....	196
8.2.3	Shuffle 聚合子图的构造方法 .....	200
8.2.4	Shffule 聚合子图的容错性能 .....	203
8.3	支持数据流网内聚合的可扩展转发策略 .....	204
8.3.1	通用的转发模式 .....	204
8.3.2	基于交换机内 Bloom 滤波的转发模式 .....	206
8.3.3	基于数据包内 Bloom 滤波的转发模式 .....	207
8.4	性能评估 .....	209
8.4.1	原型实现 .....	209
8.4.2	数据中心规模对聚合增益的影响 .....	210
8.4.3	Shuffle 传输规模对聚合增益的影响 .....	212
8.4.4	聚合率对聚合增益的影响 .....	212
8.4.5	数据包中 Bloom 滤波的大小 .....	215
	参考文献 .....	215
<b>第 9 章 不确定关联性 Incast 的协同传输管理 .....</b>		<b>217</b>
9.1	引言 .....	218
9.2	不确定性 Incast 传输的网内聚合问题 .....	220
9.2.1	不确定性 Incast 传输问题 .....	220
9.2.2	确定性 Incast 传输中的数据流网内聚合 .....	221
9.2.3	不确定性 Incast 传输中的数据流网内聚合 .....	222
9.3	不确定性 Incast 传输的聚合树构造方法 .....	223
9.3.1	网内聚合增益的多样性 .....	223