

大数据，云计算
日益庞大的电子商务信息
迅速增多的智能终端设备

这一切，都需要更强大的数据传输能力
如何避免数据“堵车”？

大数据时代下的 通信需求

——TCP传输原理与优化

徐永士 王新华 编著



中国工信出版集团



电子工业出版社
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY
<http://www.phei.com.cn>

大数据时代下的通信需求

——TCP 传输原理与优化

徐永士 王新华 编著



电子工业出版社

Publishing House of Electronics Industry

北京 · BEIJING

内 容 简 介

本书内容分为两部分。上半部分包含第 1~5 章，偏重介绍相关的互联网算法分类、理论，其中第 1 章介绍了背景知识及数据链路层，第 2 章整体介绍了 TCP 传输的原理与拥塞控制，第 3、4 章介绍了相关算法分类及主要算法，第 5 章介绍了相关理论模型。下半部分包含第 6~9 章，偏重介绍 Linux 系统上的具体实现及测量模拟技术，其中第 6 章从数据流动的角度逐层介绍了 Linux 系统如何实现网络协议 TCP/IP 协议族的各层，第 7 章具体介绍了如何书写一个拥塞控制模块及 Linux 系统自带的主要算法，第 8 章介绍了网络模拟器 NS2 及其他性能测量工具，第 9 章介绍了移动网络和软件定义网络 SDN。

本书侧重 Linux 系统上的 TCP 网络协议实现，但不限于 Linux 系统，其他操作系统及智能终端系统也可以参考。本书可以作为高等院校计算机专业、通信专业的参考用书，也可以作为大型网络中心、云计算服务技术人员的参考用书。

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。

版权所有，侵权必究。

图书在版编目 (CIP) 数据

大数据时代下的通信需求：TCP 传输原理与优化 / 徐永士，王新华编著. —北京：电子工业出版社，2015.8
ISBN 978-7-121-26602-7

I . ①大… II . ①徐… ②王… III . ①计算机网络—通信协议 IV . ①TN915.04

中国版本图书馆 CIP 数据核字 (2015) 第 156341 号

策划编辑：陈韦凯

责任编辑：毕军志

印 刷：涿州市京南印刷厂

装 订：涿州市京南印刷厂

出版发行：电子工业出版社

北京市海淀区万寿路 173 信箱 邮编 100036

开 本：787×1092 1/16 印张：17 字数：435.2 千字

版 次：2015 年 8 月第 1 版

印 次：2015 年 8 月第 1 次印刷

印 数：3 500 册 定价：49.00 元

凡所购买电子工业出版社图书有缺损问题，请向购买书店调换。若书店售缺，请与本社发行部联系，联系及邮购电话：(010) 88254888。

质量投诉请发邮件至 zlts@phei.com.cn，盗版侵权举报请发邮件至 dbqq@phei.com.cn。

服务热线：(010) 88258888。

前　　言

“云计算”、“大数据”这些概念和词汇在网络中随处可见。基于这些技术的一些服务也日益成为人们生活的一部分。如果我们深入想一想就会发现，推动这些改变的因素有哪些呢？首当其冲的当属人们的需求。快捷而廉价的沟通方式永远都是技术进步的动力。于是，一些传统的内容服务提供商之间就展开了一系列竞争，“做大”、“做强”变得不可避免，换句话说，也就是如何更加充分地利用这些资源。

那么，这些服务商都有哪些资源呢？“计算能力”、“存储能力”、“数据传输”，这些就是它们的基本要素。虽然现在硬件的购买成本越来越低，但是这并不意味着拥有成本或运维成本会下降。规模效应成了一个强劲的驱动力。对于“计算能力”和“存储能力”来说，也要更加充分地利用这些硬件资源。例如，对于代表“计算能力”的CPU，一个策略是不让CPU空转，另一个就是干脆关掉那些暂时空闲的机器。在“数据传输”方面，尽量减少需要传输的数据量是一个策略，如数据压缩、缓冲等，但是更根本的策略是提高数据传输的能力。网络设备性能的提升只是其中的一个侧面，如何更加充分地利用网络传输能力是另一个重要的侧面。

当前的互联网是一个无中心的自治网络。从网络传输的角度来看，一个瓶颈因素就是关键链路，主干网上的交换机的性能和链路的电气特性决定着网络传输的能力。这犹如现实生活中的高速公路，或者城市里的主马路，早、晚高峰的时候，每个人都想快速到达自己的目的地，但是在某些情况下，我们看到的是水泄不通的马路和焦急等待的目光。与现实生活中的道路拥堵不同的是，互联网上没有主动调度的“警察”。也就是说，在互联网上缺乏宏观的中心调控。当然这种中心调控的策略也非常不现实，因此，侧重微观调控的拥塞控制算法，成为影响网络性能的一个重要因素。

有关拥塞控制算法的研究由来已久。拥塞控制算法有几个显著的特点。首先，受限于框架，能够被用于调节拥塞的内核变量较少，主要有当前时间、往返时延等，以及丢包事件。因此，如何估计有效带宽、增减的幅度等都需要在算法中仔细考量。网络有延迟，当前的网络状况，影响的实际上是已经发出去的包。这就带来另外的难度。其次，网络（尤其是主干网络）是公共资源，如何公平有效地使用，这个问题不容忽视。网络的用户特性是短连接占有很大的比例，但是长连接的影响更大，也就是说，在资源长时间使用的情况下，如果短连接几乎得不到占有的机会，实际上就是长连接独占了资源，这就会带来极大的不公平。当然衡量公平性的方式有很多，这也是算法需要考虑的地方。在网络快速传递数据包的情况下，高速的计算非常必要，对应于代码，Linux系统已经在内核中实现了网络的大部分功能，因此留给算法的“空间”就很有限。从编程的角度来看，要求代码短小精炼，任何小的错误都可能导致内核停止、崩溃，因此要非常谨慎。

如何准确、有效地确定网络的各种参数，也是随之而来的一个重要问题。现实中使用的网络具有动态性，可重复性低，因此为了能更好地研究网络，不仅需要测试真实网络的带宽、时延、丢包率等，还需要网络模拟工具，尤其是在研究公平性方面。

网络的发展速度非常快，时下兴起的“云计算”、“大数据”、“电子商务”等，更是将网络再次推向了一个高速发展期。从用户的角度来看，数量巨大的手持智能终端是一个巨大的业务来源，众多商家都希望在这个产业链条上有所斩获。另一方面，大型计算中心内部在巨大的压力下，也在重新思考网络资源，大而统一的分配模式终将让位于精细化、个性化的定

制网络，软件定义的网络（SDN）是发展的另一个方向。这些新兴的概念和实践在本书的最后也有所讨论。

本书第1~5章侧重理论，第6~8章较深入地讨论了相关理论在Linux操作系统上的实现，作为本书的最后一章，第9章讨论了当前的热点和方向。请读者根据需求自行选择对照阅读。

本书得到了中科院高能物理研究所计算中心的大力支持。该数据中心在机房建设方面积累了丰富的实践经验，在此特别感谢其为本书付出的辛勤工作。

编著者

目 录

第1章 概述	(1)
1.1 快速发展的互联网	(1)
1.1.1 互联网的发展规模	(1)
1.1.2 争相建设的下一代互联网	(2)
1.1.3 永无止境的带宽需求	(3)
1.1.4 网络传输还需要加速	(3)
1.2 网络互联的基础——网络协议	(4)
1.2.1 OSI 参考模型与 TCP/IP 参考模型之争	(4)
1.2.2 OSI 模型	(5)
1.2.3 “阿帕网”（ARPANET）与 TCP/IP 协议族	(9)
1.2.4 TCP/IP 参考模型与协议族组成	(10)
1.2.5 数据链路层	(15)
1.3 大数据时代带来数据传输的巨大需求	(23)
1.3.1 大数据时代的到来	(23)
1.3.2 “万能”的广域网加速技术	(23)
1.3.3 技术选择	(25)
第2章 TCP 传输的原理与拥塞控制	(26)
2.1 TCP 传输原理	(26)
2.1.1 OSI 参考模型和 TCP/IP 参考模型	(26)
2.1.2 TCP 协议简介	(27)
2.1.3 TCP 数据报的传输	(27)
2.2 传输控制协议 TCP 有限状态机模型	(31)
2.2.1 客户端流程图	(33)
2.2.2 服务器端流程图	(34)
2.3 拥塞控制与 AIMD	(34)
2.3.1 拥塞的定义与发生的原因	(34)
2.3.2 拥塞控制原理 AIMD	(35)
2.3.3 现阶段的 TCP 拥塞	(38)
2.4 糊涂窗口综合症	(39)
2.4.1 发送端产生的症状	(40)
2.4.2 接收端产生的症状	(40)
2.5 其他杂项问题	(41)

第3章 主要的TCP拥塞控制算法	(44)
3.1 概述	(44)
3.1.1 从“第一次”拥塞说算法改进	(44)
3.1.2 “宏观”的解决方案——传输加速	(45)
3.1.3 新的“应用场景”	(45)
3.1.4 拥塞成因概述	(46)
3.1.5 拥塞算法设计的基本要求	(47)
3.2 基本概念与术语	(48)
3.3 TCP拥塞控制算法的演进	(49)
3.3.1 早期的TCP实现	(49)
3.3.2 TCP Tahoe	(49)
3.3.3 TCP Reno	(50)
3.3.4 TCP NewReno	(51)
3.3.5 TCP SACK	(51)
3.3.6 TCP Vegas	(52)
3.3.7 TCP Venetian	(52)
3.3.8 TCP BIC	(54)
3.3.9 TCP CUBIC	(54)
3.3.10 FAST TCP	(54)
3.3.11 Compound TCP	(55)
3.4 讨论	(55)
第4章 TCP传输加速与主要解决方案	(57)
4.1 TCP传输加速概述	(57)
4.2 解决方案分类	(57)
4.2.1 以部署方式分类	(57)
4.2.2 以实施位置分类	(58)
4.2.3 以拥塞反馈信号分类	(60)
4.2.4 基于应用层的改进方案	(61)
4.2.5 典型的隐式拥塞反馈方案	(61)
4.2.6 典型的显式拥塞反馈方案	(62)
4.2.7 基于带宽测量的改进	(64)
4.3 主要的拥塞控制算法	(65)
4.3.1 Scalable TCP	(65)
4.3.2 High Speed TCP	(65)
4.3.3 TCP Vegas	(67)
4.3.4 TCP BIC与TCP CUBIC	(69)
4.3.5 小结	(70)

第 5 章 TCP 传输性能分析与模型	(72)
5.1 端到端的可靠传输	(72)
5.1.1 差错控制过程	(72)
5.1.2 流量控制机制	(72)
5.2 传输时延	(73)
5.2.1 测量方法	(73)
5.2.2 RTT 测量的程序实现	(77)
5.3 分析模型	(84)
5.3.1 概述与进展	(84)
5.3.2 分类	(85)
5.3.3 Jacobson 管道模型	(86)
5.3.4 TCP 吞吐量分析模型	(87)
5.3.5 流体流模型	(89)
5.3.6 其他场景模型	(89)
5.3.7 传输速率上限	(90)
5.3.8 仿真实验	(90)
5.4 性能分析	(91)
5.4.1 链路利用率	(91)
5.4.2 公平性	(91)
5.5 Padhye 吞吐量模型简介	(92)
5.5.1 发送窗口表达式	(92)
5.5.2 吞吐率	(93)
5.5.3 $E[W]$ 和 $E[X]$ 的推导	(93)
5.5.4 $E[W]$ 和 $E[Y]$ 的推导	(94)
5.5.5 丢包概率 p	(95)
5.5.6 $E[A]$ 的推导	(95)
5.5.7 吞吐率表达式	(96)
第 6 章 Linux 网络协议栈	(97)
6.1 网络协议栈与层次结构	(97)
6.1.1 Linux 网络协议栈特点	(97)
6.1.2 标准 TCP/IP 协议与 Linux 网络协议栈具体设计的对比	(98)
6.2 数据结构	(99)
6.2.1 数据包结构	(100)
6.2.2 基本数据结构	(100)
6.3 协议栈的初始化	(117)
6.3.1 <code>sock_init</code> 函数	(118)
6.3.2 <code>net_dev_init</code> 函数	(121)

6.3.3	inet_init 函数	(122)
6.4	Linux 系统网络设备驱动程序	(128)
6.4.1	网络驱动程序的结构	(128)
6.4.2	数据包发送	(128)
6.4.3	数据包接收	(129)
6.5	网络协议层	(131)
6.5.1	数据接收	(131)
6.5.2	数据发送	(133)
6.6	传输层——TCP 协议处理	(135)
6.6.1	TCP 协议的数据接收	(135)
6.6.2	TCP 协议的数据发送	(138)
6.6.3	拥塞控制的事件处理	(139)
第 7 章	拥塞控制模块编程实践	(142)
7.1	拥塞控制模块的调用	(142)
7.1.1	模块的初始化	(142)
7.1.2	主要窗口的计算	(144)
7.1.3	拥塞状态机	(144)
7.1.4	状态处理函数	(147)
7.1.5	“成员函数”的调用关系	(148)
7.2	模块编程基础	(149)
7.2.1	基本数据结构	(149)
7.2.2	内核函数介绍	(151)
7.2.3	编译	(154)
7.3	主要算法介绍	(154)
7.3.1	BIC 算法	(154)
7.3.2	CUBIC 算法	(160)
7.3.3	Vegas 算法	(168)
7.3.4	High Speed TCP 算法	(173)
7.3.5	H-TCP 算法	(175)
7.3.6	Scalable TCP 算法	(182)
7.3.7	Westwood 算法	(183)
7.3.8	Reno 算法	(189)
7.3.9	代码中常见的修饰符	(192)
7.4	用户态获取当前拥塞窗口值编程示例	(194)
7.5	实践举例	(195)
第 8 章	仿真与测量	(196)
8.1	网络仿真	(196)

8.1.1 软件仿真	(196)
8.1.2 ns-2 简介	(197)
8.1.3 OPNET 简介	(223)
8.1.4 硬件模拟	(224)
8.2 性能测量方法	(225)
8.2.1 网络带宽测量	(226)
8.2.2 网络延迟测量	(234)
第 9 章 新的数据传输场景——研究热点	(235)
9.1 无线传输新场景——移动网络	(235)
9.1.1 移动智能终端逐步普及	(235)
9.1.2 需求催生的“新成员”	(236)
9.1.3 移动网络的特性	(237)
9.1.4 移动设备操作平台	(243)
9.1.5 主要的研究进展	(245)
9.2 数据中心内部传输遇到的新问题——TCP Incast	(247)
9.2.1 数据中心的网络架构	(247)
9.2.2 MapReduce 新业务与发展	(248)
9.2.3 TCP Incast 的发生	(249)
9.2.4 国内外研究现状	(251)
9.2.5 主要的解决方案介绍	(252)
9.3 网络发展的新趋势——软件定义网络 SDN 与大二层结构	(254)
9.3.1 软件定义网络	(254)
9.3.2 大二层结构	(254)
9.3.3 虚拟机迁移与数据中心二层网络的变化	(255)
9.3.4 大二层网络需要有多大	(256)
参考文献	(258)

第1章 概述

1.1 快速发展的互联网

1.1.1 互联网的发展规模

众所周知，互联网（万维网，Internet）是在美国较早的军用计算机网络 ARPANET 的基础上，经过不断发展变化而形成的一个国际性的计算机通信网络集合体。它集现代通信技术和现代计算机技术于一身，将各种各样的物理网络联合起来，构成一个整体，实现全球范围内广泛的信息交流和资源共享。

以互联网为代表的信息网络已经逐渐渗透到当今社会的各个领域，成为国家发展和社会进步的重要支柱，以及知识经济的基础载体和支撑环境。它的重要性就如同铁路和高速公路的蓬勃发展给工业社会带来了广泛而深远的影响一样，必将成为 21 世纪全球最重要的基础设施之一。如图 1.1 所示为部分互联网的路由路径的可视化图，右下角的小图为某部分末端的放大图，可以想象互联网的规模和复杂程度。

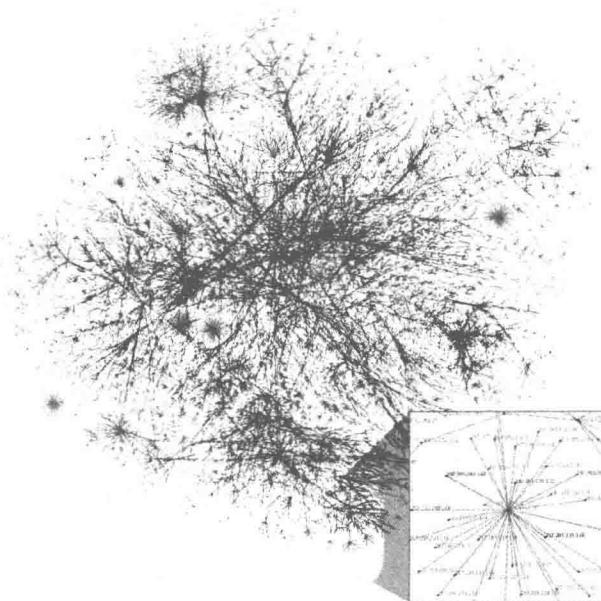


图 1.1 互联网的路由路径可视化图（部分）

据中国互联网信息中心 CNNIC 调查报告，截至 2014 年 6 月底，我国网民数量已达到 6.32 亿，互联网普及率为 46.9%，手机网民规模达到 5.27 亿，手机上网比例增长至 83.4%，超过了计算机。截至 2014 年 6 月，中国互联网网站数约有 273 万。预计 2015 年，我国互联网服务业收入将突破 6000 亿元，直接吸纳就业超过 230 万人，网民数超过 8 亿人，骨干网总带宽较“十一五”期末增长 10 倍。

人们不断提高的信息服务要求，推动着互联网的蓬勃发展。以互联网为基础的信息技术响应人们的需求，呈现出多样化的发展趋势。一方面，为了满足人们摆脱环境与设备束缚，随时、随地获取信息服务的需求，3G、LTE、WiMax 等新型无线移动通信技术开始逐步普及，并成为重要的互联网接入方式；另一方面，针对数据中心、云计算集群等海量数据交换的应用需求，千兆网络、光纤通信等超高速数据传输技术也在互联网中广泛应用，并构成了互联网核心网络的联通脉络。

不难预见，超高速的光通信技术、高速无线通信技术、网格计算和生物计算技术的研究进展会使网络技术在未来十年内产生新的飞跃，朝向以“更大、更快、更及时、更方便、更安全、更可管理和更有效”为标志的下一代互联网发展。

1.1.2 争相建设的下一代互联网

基于这样的共识，下一代互联网络及其应用的研究引起了普遍关注，世界各国相继启动了一系列的重大研究计划与项目，例如，美国的 VBNS (Very high-speed Backbone Network Service)、Internet2 和 NGI (Next Generation Internet) 三大研究计划，就是为了向美国的教育和科研机构提供世界最先进的信息基础设施，并保持美国在高速计算机网络及其应用领域的技术优势，分别在不同的技术和应用层面上共同打造一个全新概念的下一代互联网络。

与此同时，其他国家和地区也相继开展了下一代高速互联网络的研究，英、德、法、日、加等发达国家目前除了拥有政府投资建设和运行的大规模教育和科研网络以外，也都建立了研究高速计算机网络及其典型应用技术的高速网络试验床，例如，加拿大的 CANet4、欧盟的下一代互联试验网主干网 GEANT (Gigabit European Academic Network) 和亚太高级网络联盟 APAN (Asia-Pacific Advanced Network) 等。中国相关机构也在积极开展下一代互联网发展战略研究，其中，下一代互联网 NSFCNET 在北京建立了连接六个节点的 2.5~10Gbps 高速计算机互联研究试验网，分别以 1Gbps 速率连接我国的学术网络 CERNET 和 CSTNET，同时连接国际下一代互联网络交换中心 STARTAP 和亚太地区高速网 APAN 交换中心 Tokyo-XP，完成了与国际下一代互联主干网 Abilene 和 VBNS 的互联，是我国第一个与 Internet2 实现互联的计算机网络。

就目前的现状和未来的发展而言，下一代互联网的骨干带宽必将呈现指数增长的趋势。自 2002 年以来，美国的下一代互联试验网主干网 Internet2 和欧盟的下一代互联试验网主干网 GEANT 不仅在带宽方面不断升级，还在 2002 年完成了 5Gbps 的高速互联；2004 年 2 月，Internet2 的独立高速试验床 Abilene 的骨干带宽从 2.5Gbps 全面升级到 10Gbps；亚太高级网络联盟(APAN)也发起了 GTRN (Global Terabit Research Network) 计划，旨在推动骨干带宽的升级和实现全球互联。下一代互联网建设与发展的各种趋势表明：大规模的高速网络试验环境已经形成，未来几年内，互联网骨干将全面升级到支持近 10Gbps 的高速链路，而且很有可能持续增长。

1.1.3 永无止境的带宽需求

与此同时，各种新型应用的产生对网络的数据传输需求也在不断提高。现在已经有越来越多的研究人员开始经常利用这些高速网络传输 $10\text{Gbps} \sim 1\text{Tbps}$ 的数据，代表性的应用有量子物理学、地球观测、生物信息科学和射电天文学等方面的各种数据密集型的网格应用，以及 Web 站点的镜像和基于 push 的 Web 高速缓存更新等应用。

于是，虽然下一代互联网的骨干带宽呈现指数性的增长，实践中上述海量数据传输业务的用户却并没有切身感受到网络带宽剧增所带来的好处，于是人们开始怀疑高速网络中传输协议的性能。据统计，当前在所有因特网的数据包中，大约有 95% 的数据包传送使用了 TCP 协议，因此针对 TCP 协议相关机制的研究很有实际意义。

为了澄清事实，加州理工大学的 Sylvain Ravot 通过试验手段分析和评价了 Internet 上流行的 TCP Reno 协议的传输性能。在互联 GEANT 和 Internet2 的 WaveTriangle 试验床上，研究人员持续监测芝加哥超级计算中心与 CERN（欧洲核子研究中心）之间 1Gbps 的链路，测量结果表明端到端的有效吞吐量（Goodput）甚至还达不到 400Mbps 。之后，美国北卡罗来纳州立大学的 Lisong Xu 借助仿真实验也证实：TCP 协议在高速网络中确实存在效率问题。

进一步的细致分析将问题症结锁定在 TCP 拥塞控制中“加性增加乘性减小”（AIMD，Additive Increase Multiplicative Decrease）的调节机制及其相关系数上。为了论述方便，我们举例说明：假设高速链路带宽是 10Gbps ，分组大小为 1500 字节（byte），回路延时 RTT（Round Trip Time）为 100ms ，则发送端达到 10Gbps 吞吐量时，发送端拥塞窗口大小应为 83 333 个分组，依据 AIMD 窗口调整规则，TCP 拥塞避免阶段所经历的时间为 4167s ，约 1.2h ，这意味着丰富的带宽资源在长时间内都无法得到充分利用。于是，实践中网络数据传输效率低下便成为必然。究其根本原因是传统的 AIMD 拥塞控制算法在高速网络中适应性不强，效率低下，无法适应高速网络环境。

1.1.4 网络传输还需要加速

上述工程实践中发现的问题引起了众多网络研究者的关注，研究适应于高速网络的拥塞控制算法成为网络研究的新热点。在较短的时间内，研究者已经相继提出了若干新的改进算法。总的来说，高速网络拥塞控制的研究从最初单纯解决 TCP 的低效问题，到围绕公平性、稳定性及收敛性等方面开展了一系列更深入的研究。

传输层（Transport Layer）是互联网分层网络模型中的第四层，其主要任务是屏蔽网络的底层细节，为上层的网络应用提供端到端的总体数据传输控制。拥塞控制（Congestion Control）作为传输层的一项重要的网络传输控制任务，其主要功能是解决端到端网络连接在数据传输过程中存在的数据流量与网络传输能力适配问题。优秀的拥塞控制算法可以在充分、公平地利用网络传输资源的同时，避免网络链路因数据过载而造成的丢包、延迟增加、传输速率下降等拥塞崩溃问题。作为底层网络与上层应用的衔接环节，传输层中的拥塞控制算法在整个互联网网络体系结构中扮演着举足轻重的角色。

随着计算机通信技术的飞速发展，新型网络通信技术和应用服务形式不断涌现，互联网的

异构性和复杂性也随之日益增强，现在的互联网无论是在底层网络性质方面还是在上层应用需求方面，都较设计之初产生了根本性的变化。传统的传输层拥塞控制技术已经难以适应互联网日益复杂的网络结构与应用需求，并逐步成为了整个互联网系统的性能瓶颈。

但是到目前为止，在该研究领域仍然存在很多开放性问题。目前的多数研究没有充分强调模型分析的重要性，缺乏总结性结论和定律的归纳与描述，同时在拥塞控制机制和算法的设计上，过分依赖基于经验的启发式设计结合典型、有限和局部仿真试验验证的设计方法，得到的算法往往是静态和准静态的，不能适应快速变化的动态网络化环境。

高速网络环境下拥塞控制算法的优化设计还存在很大的研究空间。不仅如此，设计适应当前互联网网络结构与应用需求的新型传输方案，改进网络设备的性能，研究复杂异构网络中传输控制的关键技术与核心理论问题，都已经成为目前网络研究领域的重点和热点。

1.2 网络互联的基础——网络协议

网络协议代表着标准化，规定了计算机信息交换中消息格式和意义，是通信双方都必须遵循的一系列规则。在计算机网络中要做到有条不紊地交换数据，就必须遵守一些事先约定好的规则。这些规则明确规定了所交换的数据的格式及相关的同步问题。由此可见，网络协议是计算机网络不可缺少的组成部分。

为了简化网络设计的复杂性，通信协议采用分层的结构，各层协议之间既相互独立又高效地协调工作。对于复杂的通信协议，其结构应该有层次。分层的协议可以带来很多便利。

1. 降低了问题的复杂度，易于实现和维护

各层不需要知道它的下一层是如何实现的，而仅仅需要知道下一层通过层间的接口所提供的服务。由于每一层只实现一种相对独立的功能，因而可将一个难以处理的复杂问题分解为若干个较容易处理的更小一些的问题。这样，整个问题的复杂度就降低了。

这种结构使得实现和调试一个庞大而又复杂的系统变得容易，因为整个系统已经被分解为若干个相对独立的子系统。

2. 各层之间相互独立，灵活性好

某一层发生变化时，只要层间接口关系保持不变，则在该层以上或以下各层均不受影响。此外，对某一层提供的服务还可进行修改。当某层提供的服务不再需要时，甚至可以将该层取消。

3. 促进标准化工作

层级结构下，每一层的功能及其所提供的服务都需要精确的说明。目前广泛使用的有两种标准化模型，分别是改进后的OSI模型和TCP/IP模型，以下详细介绍。

1.2.1 OSI参考模型与TCP/IP参考模型之争

为了使不同体系结构的计算机网络都能互联，国际标准化组织ISO于1977年成立专门机

构来研究这个问题，不久即推出了一个试图使各种计算机在全世界范围内互联成网的标准框架，即著名的开放系统互联基本参考模型 OSIRM (Open Systems Interconnection Reference Model)，简称 OSI。OSI 试图达到一种理想境界，即全世界的计算机网络都遵循这个标准，使得全球所有的计算机都能够很方便地进行互联和交换数据。在 20 世纪 80 年代，许多大公司甚至一些政府机构都纷纷表示支持 OSI。当时看来似乎在不久的将来，全世界一定都会按照 OSI 制定的标准来构造自己的计算机网络。然而到了 90 年代，虽然整套的 OSI 国际标准都已经制定出来了，但由于因特网已抢先在全世界覆盖了相当大的范围，而与此同时却几乎找不到有什么厂家生产出符合 OSI 标准的商业产品。而因特网使用的体系结构是 TCP/IP。法律上的国际标准 OSI 并没有得到市场的认可。非国际标准 TCP/IP 现在获得了最广泛的应用，所以 TCP/IP 被称为事实上的国际标准。

两种主要的体系结构对应层次分布图及 TCP/IP 的三个服务层次如图 1.2 所示。

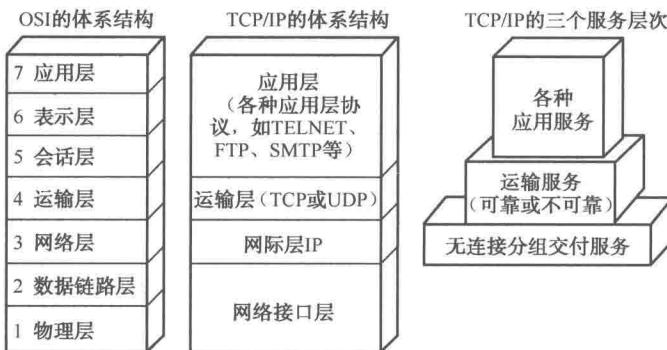


图 1.2 OSI 和 TCP/IP 体系

目前一般认为，OSI 只获得了一些理论研究的成果，但在市场化方面 OSI 则事与愿违地失败了。现今规模最大的、覆盖全世界的计算机网络——因特网使用的是 TCP/IP 体系。OSI/RM 失败的原因可归纳为以下几方面。

- (1) 网络功能在各层的分配差异大，链路层和网络层过于繁重，表示层和会话层又太轻。为此，一般只采用五层模型。
- (2) OSIRM 有关协议和服务定义太复杂且冗余，很难且没有必要在一个网络中全部实现。例如，流量控制、差错控制、寻址在很多层重复。
- (3) 高层的标准化工作唯一性太差，某些功能究竟在哪层不明确。

1.2.2 OSI 模型

历史上，在制定计算机网络标准方面起着重大作用的两大国际组织分别是国际电报与电话咨询委员会 (CCITT) 和国际标准化组织 (ISO)。虽然它们工作领域不同，但随着科学技术的发展，通信与信息处理之间的界限开始变得比较模糊，这个领域也就成了 CCITT 和 ISO 共同关心的领域。1983 年，ISO 发布了著名的 ISO/IEC 7498 标准。该标准定义了网络互联的 7 层框架，也就是开放式系统互联参考模型。OSI 是一个定义良好的协议规范集，并有许多可选具体协议来完成类似的任务。OSI 将计算机网络通信抽象成如图 1.3 所示的模型。

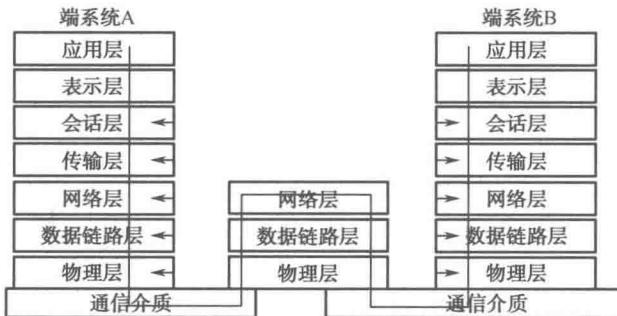


图 1.3 OSI 模型中的数据传输

OSI 将计算机网络体系结构划分为以下七层，各层的名字和基本功能如表 1.1 所示。

表 1.1 OSI 模型各层功能简介

英 文	中 文	功 能 简 介
7 Application	7 应用层	应用层能与应用程序界面沟通，以达到展示给用户的目的。常见的协议有 HTTP、HTTPS、FTP、Telnet、SSH、SMTP、POP3 等
6 Presentation	6 表示层	表示层能为不同的客户端提供数据和信息的语法转换内码，使系统能解读成正确的数据。同时，也能提供压缩解压、加密解密
5 Session	5 会话层	会话层用于为通信双方制定通信方式，并创建、注销会话（双方通信）
4 Transport	4 传输层	传输层用于控制数据流量，并且进行调试及错误处理，以确保通信顺利。而传送端的传输层会为分组加上序号，方便接收端把分组重组为有用的数据或文件
3 Network	3 网络层	网络层的作用是决定如何将发送方的数据传到接收方。该层通过考虑网络拥塞程度、服务质量、发送优先权、每次路由的耗费来决定节点 X 到节点 Y 的最佳路径。我们熟知的路由器就在这一层工作，通过不断地接收与传送数据使得网络变得相互连通
2 Data layer	2 数据链路层	首先，数据链路层的功能在于管理第一层的比特数据，并且将正确的数据传送到没有传输错误的路线中。创建并辨认数据开始及退出的位置同时予以标记。另外，处理由数据受损、丢失甚至重复传输错误的问题，使后续的层级不会受到影响，所以在该层进行数据的调试、重传或修正，决定设备何时进行传输。设备有 Bridge（桥接器）、Switch（交换器）等
1 Physical	1 物理层	物理层定义了所有电子及物理设备的规范。其中特别定义了设备与物理媒介之间的关系，这包括了针脚、电压、线缆规范、集线器、中继器、网卡、主机适配器（在 SAN 中使用的主机适配器），以及其他设备的设计定义

从表 1.1 可以看出来，第一层物理层涉及的是纯电气特性，与软件的关系不大。因为物理层传送的是原始的比特数据流，即设计的目的是为了保证当发送时的信号为二进制“1”时，对方接收到的也是二进制“1”而不是二进制“0”。因而就需要定义哪个设备有几个针脚，其中哪个针脚发送的多少电压代表二进制“1”或二进制“0”。还有诸如一个 bit 需要持续几微妙、传输信号是否在双向同时进行、最初的连接如何创建和最终如何终止等问题。

物理层的主要功能和提供的服务如下。

(1) 在设备与传输媒介之间创建及终止连接。

(2) 参与通信过程使得资源可以在共享的多用户中有效分配，例如，冲突解决机制和流量控制。

(3) 对信号进行调制或转换，使得用户设备中的数字信号定义能与信道上实际传送的数字信号相匹配；这些信号可以经由物理线缆（如铜缆和光缆）或无线信道传送。

OSI 作为一个框架来协调和组织各层所提供的服务，它定义了开放系统的层次结构、层次之间的相互关系，以及各层所包括的可能的任务。但是 OSI 参考模型并没有提供一个可以实现的方法，而是描述了一些概念，用来协调进程间通信标准的制定，即 OSI 参考模型并不是一个标准，而是一个在制定标准时所使用的概念性框架。在 ISO/OSI 模型框架里，可以兼容很多具体的协议，如表 1.2 所示。

表 1.2 协议示例

名字	OSI 协议	TCP/IP 协议	Signaling System 7	AppleTalk	IPX	SNA	UMTS	其他示例
应用层	FTAM, X.400, X.500, DAP, ROSE, RTSE, ACSE, CMIP	NNTP, SIP, SSI, DNS, FTP, Gopher, HTTP, NFS, NTP, DHCP, SMPP, SMTP, SNMP, Telnet, BGP, FCIP	INAP, MAP, TCAP, ISUP, TUP	AFP, ZIP, RTMP, NBP	RIP, SAP	APPN		HL7, Modbus
表示层	ISO/IEC 8823, X.226, ISO/IEC 9576-1, X.236	MIME, SSL, TLS, XDR		AFP				TDI, ASCII, EBCDIC, MIDI, MPEG
会话层	ISO/IEC 8327, X.225, ISO/IEC 9548-1, X.235	Sockets. Session establishment in TCP, RTP, PPTP		ASP, ADSP, PAP	NWLink	DLC		Named pipes, NetBIOS, SAP, half duplex, full duplex, simplex, RPC, SOCKS
传输层	ISO/IEC 8073, TP0, TP1, TP2, TP3, TP4 (X.224), ISO/IEC 8602, X.234	TCP, UDP, SCTP, DCCP			DDP, SPX			NBF
网络层	ISO/IEC 8208, X.25 (Packet-L4 ISO/IEC 8878, X.223, ISO/IEC 8473-1, CLNP X.233)	IP, IPsec, ICMP, IGMP, OSPF, RIP	SCCP, MTP	ATP (TokenTalk or EtherTalk)	IPX		RRC (Radio Resource Control) and BMC (Broadcast/ Multicast Control)	NBF, Q.931, NDP, ARP (maps layer 3 to layer 2 address), IS-IS