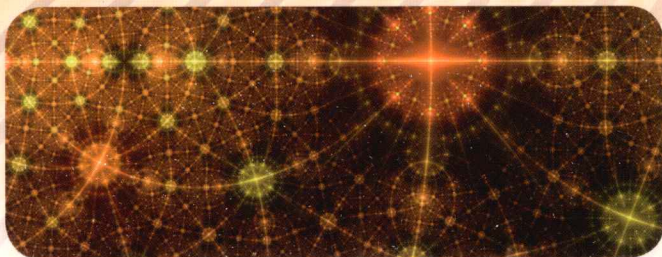


高等学校网络工程系列教材

网络协议分析与实现

Network Protocols Analysis
and Implementation

胡维华 胡昔祥 张 祯 侯宏元 编著



高等教育出版社
HIGHER EDUCATION PRESS

高等学校网络工程系列教材

网络协议分析与实现

Wangluo Xieyi Fenxi yu Shixian

胡维华 胡昔祥 张 楨 侯宏元 编著



高等教育出版社·北京
HIGHER EDUCATION PRESS · BEIJING

内容提要

本书是高等学校网络工程专业系列教材之一,是浙江省重点教材建设项目。本书采用自底向上的方法来分析 TCP/IP 协议栈的核心源代码,主要内容包括 TCP/IP 概述、底层技术、ARP、IP、ICMP、IGMP、RIP、UDP、TCP、网络应用编程接口等,最后通过网络应用编程实例来介绍典型网络应用程序的设计思想和开发步骤,加深学生对网络理论的理解,提高学生基于网络内核的网络编程与软件开发能力。

本书既可作为高等学校网络工程、计算机科学与技术、软件工程、通信工程等电气信息类专业相关课程教材,也可作为网络工程从业人员学习相关技术的高级教程。

图书在版编目(CIP)数据

网络协议分析与实现/胡维华等编著. -- 北京:
高等教育出版社,2012.6

ISBN 978-7-04-034736-4

I. ①网… II. ①胡… III. ①通信协议-
高等学校-教材 IV. ①TN915.04

中国版本图书馆 CIP 数据核字(2012)第 069428 号

策划编辑 刘 艳
插图绘制 尹 莉

责任编辑 刘 艳
责任校对 杨凤玲

封面设计 杨立新
责任印制 韩 刚

版式设计 于 婕

出版发行 高等教育出版社
社 址 北京市西城区德外大街 4 号
邮政编码 100120
印 刷 三河市杨庄长鸣印刷装订厂
开 本 787 mm × 1092 mm 1/16
印 张 23.75
字 数 580 千字
购书热线 010-58581118

咨询电话 400-810-0598
网 址 <http://www.hep.edu.cn>
<http://www.hep.com.cn>
网上订购 <http://www.landradio.com>
<http://www.landradio.com.cn>
版 次 2012 年 6 月第 1 版
印 次 2012 年 6 月第 1 次印刷
定 价 42.00 元(含光盘)

本书如有缺页、倒页、脱页等质量问题,请到所购图书销售部门联系调换
版权所有 侵权必究
物 料 号 34736-00

前 言

对于网络工程、计算机科学与技术、软件工程、通信工程等专业的本科生和研究生,如何让他们在完成“计算机网络”这一专业基础课后,能够较深入地理解抽象的网络协议体系的内在机制和实现机理,进而具有较强的网络协议研究和网络应用软件开发能力,是长期困扰计算机网络教学的一个难题。

传统的网络协议分析教材往往是首先罗列作为网络事实标准的 TCP/IP 协议栈中广泛使用的主要协议,然后分别对它们的功能需求、报文格式和基本方法做概念性介绍,而一般不深入到协议的内在机制和实现机理。有的教材使用了抓包软件,给出了协议软件的部分运行结果,使读者能够直观地了解协议软件内在的运行情况,但这毕竟仅仅是一种实验或验证,不能让读者清晰地了解协议的工作原理、数据结构、模块结构、主要算法等实质性的内容。

解决问题的关键是要阅读和理解源代码。在某种意义上,源代码本身既是最准确的说明书,也是最权威的教科书,因为它所构成的系统是切切实实在运行的。一些抽象的概念、原理和流程最终都体现在源代码上,只有读懂了源代码,才能真正理解协议的工作原理,清楚地知道网络是如何工作的。另外,对于从事系统设计或实现的工程技术人员来说,源代码的阅读和理解是一项重要的基本功。就像写小说,写小说的人大多是读了许多名著和文学评论,而不是读了“小说概论”、“小说原理”之后才学到写作技巧的。我们学习程序设计的人也是如此,尽管这是一个艰苦的过程。

本书就是旨在为可能处于这个艰苦过程中的读者提供帮助与指导的。面对庞大、复杂、抽象、深度交互的 TCP/IP 协议栈的源代码系统,我们在整体分析而又深入研究的基础上理清代码结构,精选出核心的有机联系的协议、模块、函数和数据;对每一个协议软件都从顶层勾画出一个系统的模块(函数)结构图,清晰地展示出它们之间的调用关系;对每一个函数都画出详细的程序流程图,给出源代码,并对代码的关键段落、语句和数据结构分别做出注解。我们相信,读者在学习本书的时候不会感到太庞大、太复杂,经过几个迂回,攻克几个难关,读完全书以后,一定会很好地理解每一个协议的实现方法、内部结构和一些精妙的细节,清晰地了解众多协议及其函数之间的交互作用,以及这些协议及其函数是如何被集成为一个简单而有效的软件系统的。

我们采用自底向上的方法来分析 TCP/IP 协议栈的核心源代码,全书共分成 11 章。第 1 章为 TCP/IP 概述,主要讲述用 IP 实现异构网互联、TCP/IP 分层模型、TCP/IP 协议栈及其分析的基本概念,为全书打下基础,并引出后续章节。第 2 章为底层技术,分别讲述 HDLC 协议、PPP、IEEE 802.3 标准、IEEE 802.11 标准的帧格式和协议流程,以及网络接口层的数据结构、处理流程和部分源代码。第 3 章为 ARP,给出 ARP 的基本原理、报文格式和封装、软件整体结构,具体分析了 ARP 输入处理函数、请求发送函数、缓冲区维护函数等的处理流程和实现代码,讨论了 ARP 攻击问题。第 4 章为 IP,该章是全书的重点章节之一,因为我们认为 IP 是整个 TCP/IP 协议栈的中心环节。在讲述 IP 基本原理、数据报格式、软件整体结构的基础上,给出 IP 进程实现函

数、定向广播函数、IP 数据报发送函数、向上层协议交付 IP 数据报函数、IP 数据报分片函数、IP 数据报重组函数、IP 路由获取函数、路由重定向函数等的处理流程及源代码。第 5 章为 ICMP, 在讲述 ICMP 基本原理、报文格式、软件整体结构的基础上, 着重分析了 ICMP 输入处理函数、重定向报文处理函数、发送 ICMP 报文函数等的处理流程和源代码, 并列举了一个 ICMP 的典型应用——PING 程序的实现。第 6 章为 IGMP, 给出 IGMP 软件整体结构及各函数间的调用关系, 详细分析了主机端的 IGMP 输入处理函数、IGMP 事件处理函数、主机群表的建立与维护函数等的实现代码, 并对路由器端的 IGMP 实现思路做了概述。第 7 章为 RIP, 在简要介绍 RIP 软件整体结构的基础上, 对 RIP 输入处理函数、RIP 输出处理函数、路由表的更新与维护函数等的实现代码做了详细剖析。第 8 章为 UDP, 给出 UDP 软件整体结构, 以及 UDP 输入处理函数、输出处理函数的处理流程与代码, 并列举了一个应用 UDP 与 ICMP 的典型程序 Traceroute 的实现方法。第 9 章为 TCP, 这也是全书的重点章节之一。由于它功能复杂, 我们用 3 幅图描述其软件整体结构, 并给出 TCP 实现中的一个重要的数据结构——TCP 控制块。然后, 对 TCP 输入处理等 18 个 TCP 实现的重要函数分别做了详细的分析。第 10 章为网络应用编程接口, 介绍了 3 个高级网络应用编程技术, 分别是 Socket 编程、多线程编程和 I/O 编程。该章与第 11 章配合起来, 旨在提高读者基于网络内核的编程能力, 这也是本书的重要目标之一。第 11 章为网络应用编程实例, 精选了一批网络应用实例分别说明 Socket 编程、多线程编程和 I/O 编程的具体实现方式, 在介绍时力图清楚地表述带有共性的网络编程的原理和方法, 每个实例都给出了主要的数据结构、函数处理流程和源代码。附录 A、B、C 分别给出图表索引、相关函数和宏参考表以及实验环境搭建方法。本书还配有一张光盘, 光盘中包括可在 PC 上编译后运行的完整 XINU 操作系统源代码以及第 11 章讲述的全部网络应用程序的源代码以及一个用于获取网络报文的软件 wireshark。读者通过阅读完整源代码, 修改其中部分实现程序并编译运行, 可以加深对 XINU 操作系统及运行在其之上的 TCP/IP 协议栈的理解。此外, 光盘中还给出本书各章所有习题解答的参考思路。

本书的读者对象主要是高等学校网络工程、计算机科学与技术、软件工程、通信工程等专业的本科高年级学生和研究生, 应学过“操作系统”、“数据结构”、“计算机网络”和“C 语言程序设计”这 4 门基础课程。对于本科生来说, 要能够抓住书中介绍的协议的实现机理和主要技术, 并能对每个协议软件的主要数据结构、模块或函数的调用关系、主要函数的处理流程, 以及相关协议软件之间的交互等进行描述, 实现中的一些细节和难点可以不做要求。对于研究生来说, 必须对协议的许多重要细节有充分的理解, 在此基础上能够分析其局限性, 思考一些优化与改进的策略, 并能将 TCP/IP 协议栈精妙的实现机理和一些重要结论应用到自己的研究工作中去。本书适用于所有希望了解 TCP/IP 是如何实现、因特网是如何工作的人, 当然他们已具有计算机网络、操作系统、数据结构的一些基本知识, 又粗通 C 语言。例如, 设计 TCP/IP 网络应用的程序设计师、负责维护基于 TCP/IP 协议栈的计算机系统与计算机网络的系统管理员, 以及想深入研究大规模 TCP/IP 协议栈的研究人员和工程师们。

本书中的绝大部分源代码取自于 XINU 操作系统内核, 该系统是开放的, 没有版权与费用问题, 在 PC 上也都能运行。这些软件与商用应用中的同等协议软件的核心思想无太大差别, 有的功能甚至更强。这些软件是公开的, 可以从很多地方获得。采用开源协议软件平台, 还有助于课程实验的开展。

我于 1993 年开始为杭州电子科技大学计算机科学与技术学科的研究生们开设“网络协议分

析”课程,积累了一定的心得与讲稿。本书的编写由我拟定基本思路与大纲,第1、2章由张祯执笔编写,第3章至第6章由侯宏元执笔编写,第7章至第9章由我执笔编写,第10、11章由胡昔祥执笔编写,最后由我统阅、修改并定稿。TCP/IP协议栈是一个十分庞大而又复杂的系统,仅TCP部分就有10多万行的C语言代码,所以协议与函数的选择是本书编写的一个难点;另外,就像软件免不了有错一样,对软件的理解和诠释也免不了会有不当甚至错误,本书一定也是如此。敬请各位读者批评、赐教与指正,以利再版时改进。

本书的编写得到了浙江大学何钦铭教授和陈天洲教授、浙江工业大学陈庆章教授、杭州电子科技大学万健教授和谢红标老师,以及高等教育出版社的支持与指导,在此表示衷心的感谢!我的多位研究生:汤利平、李源洲、贾琳、杨汉赞、徐晶、王科、许华强、冯宗伟、冯伟等帮助我们进行资料整理、流程图绘制和文稿录入等,特别是汤利平在内容选择、流程图绘制和实验环境搭建方面做了许多富有成效的工作,在此也一并表示感谢!

胡维华
2012年4月

郑重声明

高等教育出版社依法对本书享有专有出版权。任何未经许可的复制、销售行为均违反《中华人民共和国著作权法》，其行为人将承担相应的民事责任和行政责任；构成犯罪的，将被依法追究刑事责任。为了维护市场秩序，保护读者的合法权益，避免读者误用盗版书造成不良后果，我社将配合行政执法部门和司法机关对违法犯罪的单位和个人进行严厉打击。社会各界人士如发现上述侵权行为，希望及时举报，本社将奖励举报有功人员。

反盗版举报电话 (010)58581897 58582371 58581879

反盗版举报传真 (010)82086060

反盗版举报邮箱 dd@hep.com.cn

通信地址 北京市西城区德外大街4号 高等教育出版社法务部

邮政编码 100120

目 录

第 1 章 TCP/IP 概述	1	2.2.1 SLIP 与 PPP	18
1.1 网络互联与 TCP/IP	1	2.2.2 PPP 组件	19
1.1.1 计算机网络	1	2.2.3 PPP 的帧格式	19
1.1.2 网络互联	2	2.2.4 PPP 工作流程	20
1.1.3 TCP/IP 协议栈	3	2.2.5 PPP 应用	20
1.2 网络协议的分层	4	2.3 以太网及 IEEE 802.3	21
1.2.1 分层的网络体系结构	4	2.3.1 以太网技术	21
1.2.2 TCP/IP 模型	7	2.3.2 CSMA/CD	21
1.2.3 多路复用和分解	7	2.3.3 帧格式	21
1.2.4 TCP/IP 编址	8	2.4 无线局域网及 IEEE 802.11	22
1.3 TCP/IP 协议栈及其分析	9	2.4.1 无线局域网简介	22
1.3.1 TCP/IP 协议栈	9	2.4.2 CSMA/CA	23
1.3.2 协议栈的处理流程	10	2.4.3 IEEE 802.11 帧格式	24
1.3.3 设备驱动程序和输入输出 程序	11	2.5 网卡驱动和网络接口层的实现	25
1.3.4 网络接口层处理程序	12	2.5.1 以太网接口数据结构	25
1.3.5 IP 层处理程序	12	2.5.2 以太网网卡驱动程序	26
1.3.6 传输层处理程序	13	2.5.3 网络接口层数据结构	28
1.3.7 应用编程接口	13	2.5.4 网络接口层处理流程	31
1.4 本书的代码组织	14	2.5.5 网络接口层的多路分解	32
1.4.1 研究代码的重要性	14	2.5.6 网络初始化	34
1.4.2 XINU 的 TCP/IP 协议栈源 代码	14	习题	39
1.4.3 应用层示例源代码	15	第 3 章 ARP	41
习题	15	3.1 ARP 的基本原理	41
第 2 章 底层技术	16	3.2 ARP 报文格式和封装	41
2.1 HDLC 协议	16	3.2.1 ARP 报文格式	41
2.1.1 HDLC 协议介绍	16	3.2.2 ARP 报文结构的实现	42
2.1.2 帧格式	16	3.3 ARP 软件整体结构	43
2.1.3 帧类型和 HDLC 操作	17	3.4 ARP 输入处理	44
2.1.4 HDLC 协议的应用	18	3.5 发送 ARP 请求报文	48
2.2 PPP	18	3.6 ARP 缓冲区的管理	51
		3.6.1 ARP 缓冲区结构的实现	51
		3.6.2 ARP 缓冲区维护函数	52

3.7 ARP 攻击	55	第 6 章 IGMP	109
习题	55	6.1 IGMP 的基本原理	109
第 4 章 IP	56	6.2 IGMP 报文格式	109
4.1 IP 的基本原理	56	6.2.1 IGMP 报文格式	109
4.2 IP 数据报格式	56	6.2.2 IGMP 报文结构的实现	110
4.2.1 IP 数据报格式	56	6.3 IGMP 软件整体结构	110
4.2.2 IP 数据报结构的实现	57	6.4 主机端的 IGMP 输入处理	111
4.3 IP 软件整体结构	59	6.5 IGMP 事件处理进程	117
4.4 IP 输入处理	61	6.6 主机群表的建立与维护	118
4.4.1 IP 进程的实现	61	6.7 路由器端 IGMP 实现概述	123
4.4.2 IP 定向广播	66	习题	125
4.5 IP 输出处理	70	第 7 章 RIP	126
4.5.1 IP 输出处理	70	7.1 RIP 的基本原理	126
4.5.2 将 IP 数据报交付上层		7.2 RIP 报文格式	127
协议	74	7.2.1 RIP 报文格式	127
4.6 IP 数据报的分片与重组	76	7.2.2 RIP 报文结构的实现	127
4.6.1 分片结构的实现	76	7.3 RIP 软件整体结构	129
4.6.2 IP 数据报分片	77	7.4 RIP 输入处理	129
4.6.3 IP 数据报重组	79	7.4.1 RIP 输入进程的实现	129
4.7 IP 选路	85	7.4.2 RIP 通告报文的处理	132
4.7.1 路由表基本结构	86	7.4.3 RIP 请求报文的处理	134
4.7.2 路由选择	87	7.5 RIP 报文输出进程	136
4.7.3 路由重定向的实现	89	7.5.1 RIP 接口输出控制结构	136
4.8 IPv6	91	7.5.2 RIP 输出进程的实现	137
习题	92	7.5.3 发送一个 RIP 通告报文	139
第 5 章 ICMP	93	7.5.4 构造 RIP 通告报文	141
5.1 ICMP 的基本原理	93	7.5.5 为 RIP 通告报文计算路由	
5.2 ICMP 报文格式	93	度量值	143
5.2.1 ICMP 报文格式	93	7.6 路由表的更新与维护	144
5.2.2 ICMP 报文结构的实现	94	7.6.1 路由表项的添加	144
5.3 ICMP 软件整体结构	96	7.6.2 路由表项的定时维护	150
5.4 ICMP 输入处理	96	7.7 OSPF	153
5.4.1 ICMP 输入处理函数	96	习题	154
5.4.2 ICMP 重定向报文的处理	101	第 8 章 UDP	155
5.5 发送 ICMP 报文	103	8.1 UDP 的基本原理	155
5.6 PING 程序的实现	107	8.2 UDP 数据报格式	155
5.7 ICMPv6	108	8.2.1 UDP 数据报格式	155
习题	108	8.2.2 UDP 数据报结构的实现	156

8.3	UDP 软件整体结构	157	10.3.4	Socket 编程模型	230
8.4	UDP 输入处理	157	10.4	多线程编程	231
8.5	UDP 输出处理	160	10.4.1	线程创建	231
8.6	Traceroute 程序的实现	163	10.4.2	线程属性	231
	习题	163	10.4.3	线程同步	232
第 9 章	TCP	164	10.4.4	线程取消	235
9.1	TCP 的基本原理	164	10.4.5	线程终止	236
9.2	TCP 报文段格式	165	10.5	I/O 编程	236
9.2.1	TCP 报文段格式	165	10.5.1	I/O 模型	236
9.2.2	TCP 报文段结构的实现	166	10.5.2	高级 I/O 函数	239
9.3	TCP 软件整体结构	167	10.5.3	非阻塞 I/O	241
9.4	TCP 控制块结构	169	10.5.4	I/O 复用	242
9.5	TCP 状态机的实现	173		习题	244
9.6	TCP 输入进程及输入状态机	175	第 11 章	网络应用编程实例	245
9.6.1	TCP 输入进程的实现	175	11.1	UDP 套接字的简单应用——TIME	
9.6.2	TCP 输入状态机	178		的实现	245
9.6.3	TCP 对输入数据的处理	189	11.1.1	TIME 简介	245
9.6.4	输入 ACK 报文段的处理	192	11.1.2	TIME 客户端程序	245
9.7	TCP 输出进程及输出状态机	195	11.1.3	connectUDP() 函数	248
9.7.1	TCP 输出进程	195	11.1.4	connectsock() 函数	248
9.7.2	TCP 输出状态机的实现	197	11.1.5	TIME 服务器端程序	251
9.7.3	发送 TCP 报文段	204	11.1.6	passivesock() 函数	253
9.8	流量控制和拥塞控制	210	11.2	TCP 套接字的简单应用——	
9.8.1	设定发送方窗口通告值	210		DAYTIME 的实现	255
9.8.2	设定接收方窗口通告值	212	11.2.1	DAYTIME 简介	255
9.8.3	估算往返时延并设定重传 时间和拥塞窗口	214	11.2.2	DAYTIME 客户端程序	256
9.9	TCP 定时器管理	216	11.2.3	DAYTIME 服务器端 程序	257
9.9.1	TCP 定时结构	217	11.3	利用多路 I/O 和多线程编程——	
9.9.2	TCP 定时进程的实现	217		ECHO 的实现	260
	习题	220	11.3.1	ECHO 简介	260
第 10 章	网络应用编程接口	221	11.3.2	ECHO 客户端程序	260
10.1	网络应用程序	221	11.3.3	ECHO 服务器端程序	263
10.2	网络应用模式	221	11.3.4	tcp_listen() 函数	266
10.3	网络编程接口	223	11.4	网络客户端实例(1)——TELNET	
10.3.1	Socket 编程基本概念	223		客户端的实现	268
10.3.2	Socket 地址结构	227	11.4.1	TELNET 简介	268
10.3.3	Socket 接口函数	227	11.4.2	TELNET 协议数据结构	269

11.4.3	TELNET 客户端程序	273	11.6.2	PING 程序主函数的处理 流程	301
11.4.4	telnet() 函数	274	11.6.3	catcher() 函数	311
11.4.5	connectTCP() 函数	278	11.6.4	pinger() 函数	313
11.4.6	ttwrite() 函数	278	11.6.5	pr_pack() 函数	315
11.4.7	sowrite() 函数	279	11.6.6	in_cksum() 函数	322
11.5	网络客户端实例(2)——TFTP		11.6.7	pr_icmph() 函数	323
	客户端的实现	280	11.6.8	pr_iph() 函数	326
11.5.1	TFTP 简介	280	11.6.9	pr_addr() 函数	327
11.5.2	TFTP 报文首部数据 结构	282	11.6.10	pr_retip() 函数	328
11.5.3	TFTP 客户端程序	283	11.6.11	finish() 函数	329
11.5.4	command() 函数	287	11.6.12	其他函数	330
11.5.5	recvfile() 函数	288		习题	330
11.5.6	sendfile() 函数	293	附录 A	图表索引	331
11.5.7	makerequest() 函数	298	附录 B	相关函数和宏参考表	335
11.5.8	nak() 函数	298	附录 C	实验环境搭建方法	359
11.6	原始套接字的应用——PING		参考文献		368
	的实现	300			
11.6.1	PING 简介	300			

第 1 章

TCP/IP 概述

1.1 网络互联与 TCP/IP

1.1.1 计算机网络

一般而言,计算机网络是指将地理位置不同、具有独立功能的多个计算机终端通过通信线路和网络设备连接起来,在网络软件的管理和协调下,实现资源共享和信息传递的系统。具有独立功能的多个计算机终端广义上可以是计算机、服务器,也可以是目前流行的智能终端(智能手机、智能家电等);通信线路既可以是同轴电缆、双绞线、光纤等有线传输介质,也可以是像无线电波那样的无线传输介质;常见的网络设备主要有集线器、交换机和路由器等;网络软件包括运行在计算机和网络设备上的驱动程序、协议栈和网络应用程序等。

计算机网络在逻辑功能上可以划分为资源子网和通信子网两大部分。资源子网负责对信息进行处理和加工,包括访问网络和处理数据的硬件、软件设施,如联网的计算机、智能终端等;通信子网负责网络信息流的传递、交换、控制以及信号的转换等工作,包括网卡、通信线路、集线器、交换机、路由器、调制解调器,以及通信软件等设施。

计算机网络的种类很多,根据网络覆盖范围的不同,可以分为局域网、城域网、广域网;根据数据的交换方式不同,可以分为电路交换、报文交换、分组交换;根据传输介质不同,可以分为有线网络(同轴电缆、双绞线、光纤等),无线网络(蓝牙、无线局域网、蜂窝式无线网等)。

计算机网络是复杂的,其复杂性体现在由于应用范围的不同,存在着多种网络技术,每一种技术所包含的硬件组成和软件构件都各不相同,其信道访问方式和数据传输方式也都存在差异。这些不同的网络技术有着各自不同的特点:高速以太网价格低廉,但会受到地理范围的限制;同步光纤网能够提供较高的带宽,适合骨干网使用,但价格相对较高;卫星通信网适合远距离、大容量的网络通信,但网络延迟较大;点对点网络使用户得以远距离接入网络;无线局域网适合移动用户在固定的区域内访问网络;蜂窝式无线网可以使用户在快速移动中访问网络。

这些采用不同技术所组建的网络,在网络内部之间可以直接通信,但是从用户的角度看,人们还希望这些采用不同技术的异构网络之间能够相互进行数据传递和资源共享,组成一个更大的网络(网络的网络),这就是网络互联。举例来说,用户 A 可以通过接入使用以太网技术的校园网,与另外一个使用电话点对点拨号上网的用户 B 之间进行邮件通信,同时还和一个坐在时速 300 公里的高铁上的使用 WCDMA 手机进行 3G 上网的用户 C 进行 QQ 聊天。但问题的关键

在于,这些采用不同技术的异构网络之间存在着很大的差异:它们的信道访问方式和数据传输方式不同,其帧格式和物理地址形式也各不相同。那么,如何将这些异构的网络互联在一起呢?

1.1.2 网络互联

网络互联的根本问题是解决网络技术和应用所带来的网络异构性问题,用来解决该问题的技术叫做网络互联技术。网络互联技术通过提供异构网络互联的方法及一组通信约定,使其可以容纳多种不同的硬件技术,屏蔽网络底层的技术细节,从而允许计算机独立于它们的物理网络连接来进行通信。使用网络互联技术将不同的网络连接而成的网络叫做互联网。

目前世界上最大的互联网是因特网(Internet,一般所说的互联网特指因特网),因特网在底层网络和网络应用之间添加了一个 IP(Internet Protocol,网际协议)层,定义了标准的 IP 数据报格式以及标准的 IP 地址格式,向下屏蔽底层物理网络的差异,向上提供一致性的访问接口,各种网络应用通过将数据封装成统一的 IP 数据报格式来进行数据传送,通过统一的 IP 地址来进行寻址,由此解决了异构网络的互联问题。图 1-1 是利用 IP 实现异构网络互联的示意图。



图 1-1 利用 IP 解决异构网络互联问题

借助于 IP 层,网络应用使用 IP 数据报封装数据,使用 IP 地址标识目的主机,屏蔽底层物理网络的差异,从而实现透明的网络数据传输。但是对于 IP 层而言,IP 数据报必须转换成具体物理网络所使用的数据帧进行传输。在上面用户 A、B、C 邮件通信和 QQ 聊天的例子中,以太网中传递的是以太帧,电话点对点拨号网络传递的是 PPP 帧,3G 网络依据采用技术的不同会使用不同的帧格式(WCDMA、CDMA-2000、TD-SCDMA 等),同步光纤网传递的是 SDH 帧。要想实现这些异构网络之间的互联,必须有一个中间的转换设备,将一种帧格式转换为另一种帧格式;此外,为了使 IP 数据报能够顺利到达目的主机,必须有 IP 数据报的选路和转发设备。因特网中完成这两个任务的设备就是路由器。路由器在网络互联中的作用如图 1-2 所示。

图 1-2 显示的是用户 A、B、C 邮件通信和 QQ 聊天的例子。以太网中的用户 A、电话点对点拨号网络中的用户 B 和 3G 网络中的用户 C 分别通过路由器 R1、R2、R3 连接到因特网主干网,路由器 R1、R2、R3 都是采用同步光纤的方式接入因特网主干网。下面以用户 A 发送邮件数据给用户 B 为例来说明路由器在网络互联中的作用。

- 用户 A 的主机将发送的邮件数据先封装到 IP 数据报中,再封装到以太帧中,然后发送到其接入的以太网中,到达路由器 R1。

- 路由器 R1 从以太帧中提取出 IP 数据报,根据目标 IP 地址选择合适的路径,再将其封装成 SDH 帧,转发到因特网主干网中,经过因特网主干网中若干路由器的选路和转发,到达路由器 R2。

- 路由器 R2 从 SDH 帧中提取出 IP 数据报,转换成 PPP 帧,发送到电话网络中,到达用户 B 的主机。

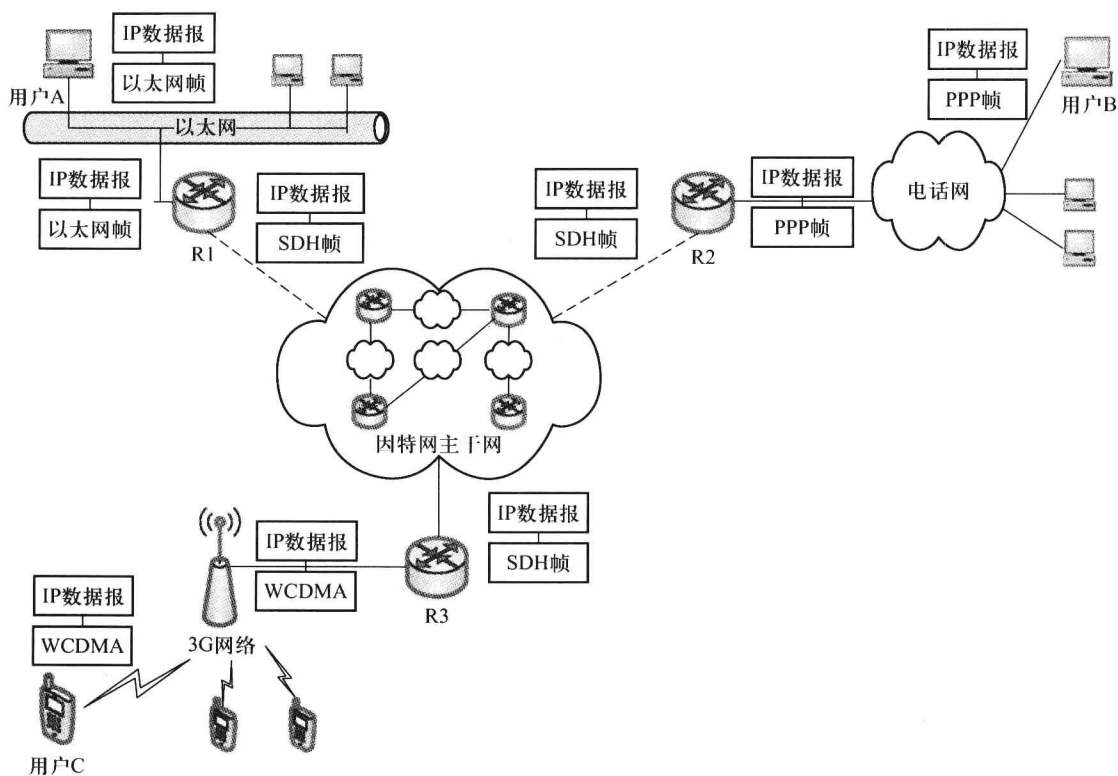


图 1-2 路由器在网络互联中的作用

- 用户 B 的主机提取出 IP 数据报,最终交付上层的邮件应用程序,显示给用户 B。

在这个通信的过程中,用户 A 和用户 B 的邮件通信应用是一个对等实体,它们都直接使用 IP 数据报进行交互通信,虽然 IP 数据报在传输的过程中被多次封装成不同的帧格式,经过了多个底层物理网络,但是这些都是由主机或路由器的 IP 层来完成的。路由器一般连接两个以上的物理网络,在 IP 层进行 IP 数据报的选路和数据转发,转发的同时也完成了帧格式的转换。

综上所述,因特网通过引入 IP 层,解决了异构网络的互联问题。在实现上,路由器是实现网络互联的核心设备,整个因特网就是由无数个用路由器互联起来的物理网络所构成的网络。这些物理网络,无论规模大小、作用如何,不管是使用以太网技术的局域网,还是用做主干网的广域网,或者是连接两台计算机的点到点链路,其地位都是一样的,都可以看成一个网络。从用户的角度看,因特网就是一个统一的虚拟网络,只要用户通过某种物理网络连上了因特网,它就可以与因特网中的任何一台主机通信,不管它们之间间隔了多少个路由器和物理网络。

1.1.3 TCP/IP 协议栈

IP 层的引入解决了异构网络的互联问题,但是要确保一个庞大的、有着众多的技术特征不同的异构网络组成的系统能够正确、高效地运转,仍然有很多其他问题需要解决。

- 路由器要完成 IP 数据报的选路工作,必须对它所连接的网络拓扑结构有一个准确的了解,因此需要有路由协议来帮助其维护路由表。

• IP 对底层网络的通信能力没有做过多的要求,底层网络可能会出现报文出错、丢失等问题;另外,路由器在转发报文时也可能由于输入和输出数据的速率不同而导致报文丢弃,这就需要由差错报告机制来报告差错。

• IP 提供了统一的 IP 地址,但是在具体物理网络中数据传输仍然要使用物理地址,因此需要由地址转换协议来完成 IP 地址和物理地址的转换。

• 因特网支持多播,因此要提供多播组管理的一套流程及其配套协议。

• IP 提供了传输服务,但并不保证数据传输的可靠性和正确性,这对于网络应用而言,还是不够的,很多网络应用需要数据能够可靠、正确地传输到目的主机。

• 根据 IP 地址能够寻址到主机,但是一台主机可能会有多个网络应用,因此还需要其他地址来标识具体的网络应用。

如果上述问题都通过 IP 来解决,将使 IP 过于庞大。为了保证网络各个功能的相对独立性,因特网采用模块化的思想,针对上述每个问题,引入了专门的协议来解决,以便于网络的实现和维护。例如,RIP(Routing Information Protocol,路由信息协议)、OSPF(Open Shortest Path First,开发式最短路径优先)和 BGP(Border Gateway Protocol,边界网关协议)等协议用于维护路由表;ICMP(Internet Control Message Protocol,Internet 报文控制协议)用于传递网络控制和差错信息;ARP(Address Resolution Protocol,地址解析协议)和 RARP(Reverse Address Resolution Protocol,反向地址解析协议)用于实现 IP 地址和物理地址的相互转换;IGMP(Internet Group Management Protocol,Internet 组管理协议)用于实现多播组的管理;UDP(User Datagram Protocol,用户数据报协议)为网络应用提供了无连接的、不可靠的数据报传输服务;TCP(Transmission Control Protocol,传输控制协议)为网络应用提供了面向连接的、基于流的可靠数据传输服务;SCTP(Stream Control Transmission Protocol,流控制传输协议)提供支持多媒体业务的可靠的数据报传输协议。此外,因特网还为一些常见的网络应用制定了标准的协议,例如用于远程登录的 TELNET 协议,用于域名解析的 DNS(Domain Name System,域名系统)协议,用于 IP 地址分配的 DHCP(Dynamic Host Configuration Protocol,动态主机配置协议),用于文件传输的 FTP(File Transfer Protocol,文件传输协议),用于邮件通信发送的 SMTP(Simple Mail Transfer Protocol,简单邮件传输协议),用于邮件接收的 POP3(Post Office Protocol 3,邮局协议版本 3),用于超文本传输的 HTTP(HyperText Transfer Protocol,超文本传输协议),用于网络管理的 SNMP(Simple Network Management Protocol,简单网络管理协议)等。

上述协议和 IP 组成一组协议,构成了因特网中网络互联的基础。在这些协议当中,TCP 和 IP 是最重要的两个协议,故用它们名字组合作为这一组协议的名称,TCP/IP 协议栈(TCP/IP Internet Protocol Suite),简称 TCP/IP。

1.2 网络协议的分层

1.2.1 分层的网络体系结构

网络协议就是通信双方共同遵守的规则和约定的集合,而协议的实现是由具体的硬件和软件模块来完成的,在网络中将这种实现特定功能的模块称为实体。

网络互联是一个复杂的系统工程,包括很多的软件、硬件和相关协议,为了简化系统,设计者采用分而治之的思想来设计和描述网络体系结构,将网络协议的实现按层次进行组织,每一层使用下层提供的服务完成一定的功能,并向上层提供服务。这样,协议栈中的协议就具有上下层次关系。

如图 1-3 所示,采用分层结构后,协议及其实体根据功能被划分到不同的层次上,这样两个结点间的通信体现为两个结点对等层相应实体之间的通信。

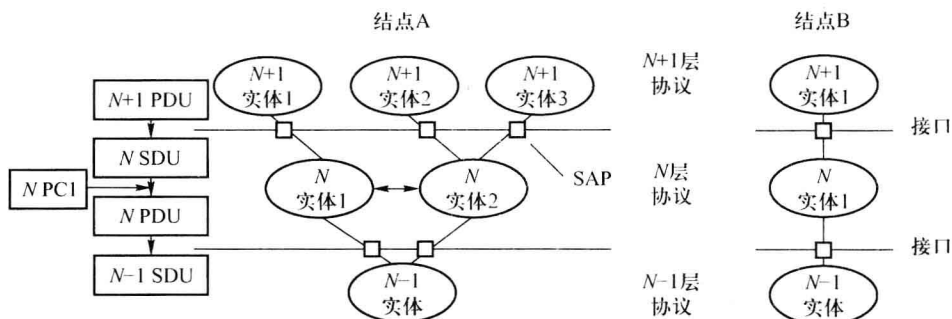


图 1-3 协议分层思想

通信双方对等层中完成相同协议的功能的实体称为对等实体。如图 1-3 中结点 A 的 $N+1$ 实体 1 和结点 B 中的 $N+1$ 实体 1 就是对等实体,它们之间的通信体现了一种横向的通信关系。

对等层实体的通信是通过下一层实体来完成的,同时它又为上一层实体通信提供服务。层与层之间的实体通信是通过相邻层间的服务接口来完成的,因此这个接口称为访问服务点 (Service Access Point, SAP)。发送方的 $N+1$ 层实体根据协议将其封装成协议数据单元 (Protocol Data Unit), 将其和一些控制信息通过 SAP 传递给 N 层实体, N 层实体从 $N+1$ 层实体中得到的报文称为服务数据单元 (Service Data Unit, SDU), N 层实体无需理解 SDU 的含义, 只将其当做要服务的数据。为了传送 SDU, N 层实体可能会将 SDU 拆成几段, 或者将多个 SDU 合并为一段, 在每一段的首部加上该对等实体间约定的协议控制信息 (Protocol Control Information, PCI), 作为 N 层 PDU, 再向下传送。

网络协议的分层有利于简化协议设计, 每一层协议通信只在本层上完成, 协议设计者能够把注意力集中到某一层上而不必担心较低层的通信情况; 其次, 分层还有利于网络互联, 每一层协议仅和与它相邻的上下两层协议进行数据交换; 此外, 分层还可以屏蔽下层协议的变化, 新的底层技术的引入不会对上层的应用协议产生影响。

目前使用最普遍的网络分层模型是 ISO (国际标准化组织) 制定的 OSI (Open System Interconnection, 开放式系统互联) 参考模型。OSI 模型并不是一个标准, 而只是一个在制定标准时所使用的概念性框架, 并没有提供一个可以实现的方法。如图 1-4 所示, OSI 模型把网络通信的工作划分为 7 层, 自底向上分别是物理层、数据链路层 (简称链路层)、网络层、传输层、会话层、表示层和应用层。两个端点实现了整个协议栈, 可以进行端到端的数据通信, 而中间结点可能仅完成数据的交换和转发, 并不需要实现所有的层次。

OSI 各层的功能如下。

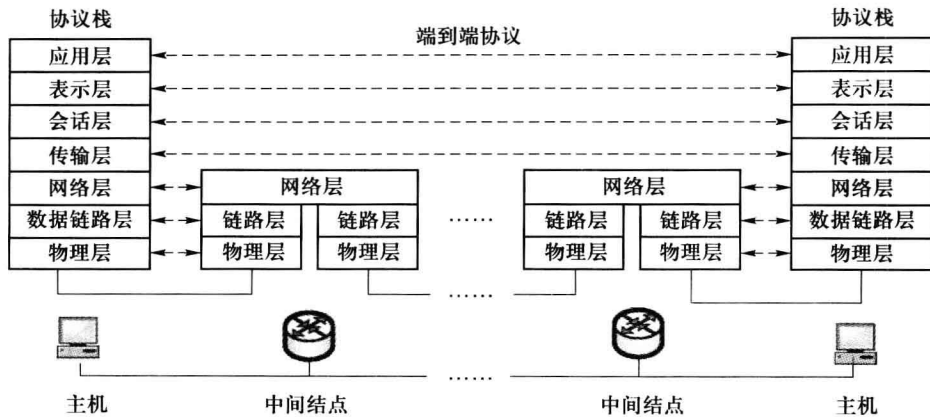


图 1-4 OSI 模型

1. 物理层

物理层为数据链路层提供服务,其功能就是通过传输介质发送比特流,传输的数据单元是比特。该层定义了与物理链路的建立、维护和拆除有关的机械、电气、功能和规范等特性,包括信号线的功能、介质的物理特性、传输速率、位同步、传输模式和连接规格等。

2. 数据链路层

数据链路层为网络层提供数据传输服务,完成两个相邻结点之间的通信问题。该层传送的协议数据单元称为帧(Frame),具体功能包括数据成帧、介质访问控制、物理寻址、差错控制和流量控制等。

3. 网络层

网络层为传输层提供服务,负责从源主机到目的主机的端到端的数据通信,传输的协议数据单元称为数据包(Packet)。具体功能包括网络逻辑地址编址、路由选择、报文转发等。

4. 传输层

传输层为上一层协议提供进程到进程的、有序的、可靠和透明的数据传输服务,传输的协议数据单元称为报文段(Segment)。具体功能包括差错控制、流量控制、拥塞控制、报文的分段与重组和进程寻址等。

5. 会话层

会话层负责会话的控制,具体功能包括会话的建立、维护、同步和会话的终止。

6. 表示层

表示层负责信息的表示和转换,具体功能包括数据的加密和解密、压缩和解压缩、数据的格式转换等。

7. 应用层

应用层是用户和网络的接口,通过网络应用程序完成用户的网络应用需求,例如,电子邮件、文件传输、远程登录等。

归纳起来,7层结构可以分为3个部分。下面3层是通信支持层,完成端到端的通信;上面3层是应用支持层,完成具体的网络应用;中间的传输层起着隔离作用,使网络应用可以和具体的网络无关。