

HZ BOOKS
华章教育

计 算 机 科 学 丛 书

Pearson

TCP/IP详解

卷3: TCP事务协议、HTTP、 NNTP和UNIX域协议

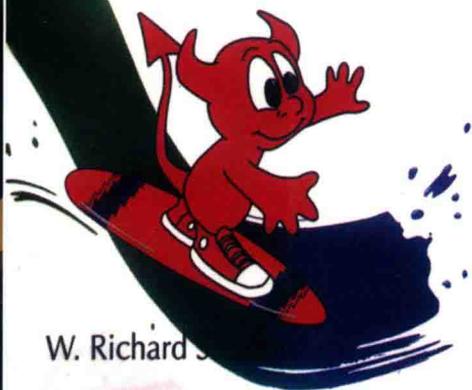
[美] W. 理查德·史蒂文斯 (W. Richard Stevens) 著
胡谷雨 吴礼发 等译
谢希仁 校

TCP/IP Illustrated

Volume 3: TCP for Transactions, HTTP, NNTP, and the UNIX Domain Protocols

TCP/IP Illustrated, Volume 3

TCP for Transactions, HTTP, NNTP,
and the UNIX® Domain Protocols



W. Richard S

ADDITION-WESLEY PROFESSIONAL COMPUTING SERIES

非
外
借



机械工业出版社
China Machine Press

计 算 机 科 学 丛 书

TCP/IP详解

卷3: TCP事务协议、HTTP、 NNTP和UNIX域协议

[美] W. 理查德·史蒂文斯 (W. Richard Stevens) 著
胡谷雨 吴礼发 等译
谢希仁 校

TCP/IP Illustrated

Volume 3: TCP for Transactions, HTTP, NNTP, and the UNIX Domain Protocols

TCP/IP Illustrated, Volume 3

TCP for Transactions, HTTP, NNTP,
and the UNIX® Domain Protocols



W. Rich



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

TCP/IP 详解 卷 3: TCP 事务协议、HTTP、NNTP 和 UNIX 域协议 / (美) W. 理查德·史蒂文斯 (W. Richard Stevens) 著; 胡谷雨等译. —北京: 机械工业出版社, 2019.3

(计算机科学丛书)

书名原文: TCP/IP Illustrated, Volume 3: TCP for Transactions, HTTP, NNTP, and the UNIX Domain Protocols

ISBN 978-7-111-61777-8

I. T… II. ①W… ②胡… III. 计算机网络 - 通信协议 IV. TN915.04

中国版本图书馆 CIP 数据核字 (2019) 第 007421 号

本书版权登记号: 图字 01-2018-7883

Authorized translation from the English language edition, entitled TCP/IP Illustrated, Volume 3: TCP for Transactions, HTTP, NNTP, and the UNIX Domain Protocols, ISBN: 9780201634952, by W. Richard Stevens, published by Pearson Education, Inc., Copyright ©1996 by Addison Wesley.

All rights reserved. No part of this book may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording or by any information storage retrieval system, without permission from Pearson Education, Inc.

Chinese simplified language edition published by China Machine Press, Copyright © 2019.

本书中文简体字版由 Pearson Education (培生教育出版集团) 授权机械工业出版社在中华人民共和国境内 (不包括香港、澳门特别行政区及台湾地区) 独家出版发行。未经出版者书面许可, 不得以任何方式抄袭、复制或节录本书中的任何部分。

本书封底贴有 Pearson Education (培生教育出版集团) 激光防伪标签, 无标签者不得销售。

本书是三卷本套书《TCP/IP 详解》的第 3 卷, 主要内容包括: TCP 事务协议, 即 T/TCP, 它是对 TCP 的扩展, 使客户 - 服务器事务更快、更高效和更可靠; TCP/IP 应用, 主要是 HTTP 和 NNTP; Unix 域协议, 这些协议提供了一种进程之间通信的手段, 当客户与服务器进程在同一台主机上时, Unix 域协议通常要比 TCP/IP 快 1 倍。本书同样给出了大量的实例和实现细节, 并参考引用了卷 2 中的大量源程序。

本书适用于希望理解 TCP/IP 工作原理的读者, 包括编写网络应用程序的程序员以及利用 TCP/IP 维护计算机网络的系统管理员。

出版发行: 机械工业出版社 (北京市西城区百万庄大街 22 号 邮政编码: 100037)

责任编辑: 吴 怡

责任校对: 李秋荣

印 刷: 北京市兆成印刷有限责任公司

版 次: 2019 年 3 月第 1 版第 1 次印刷

开 本: 185mm × 260mm 1/16

印 张: 16.75

书 号: ISBN 978-7-111-61777-8

定 价: 59.00 元

凡购本书, 如有缺页、倒页、脱页, 由本社发行部调换

客服热线: (010) 88378991 88361066

投稿热线: (010) 88379604

购书热线: (010) 68326294 88379649 68995259

读者信箱: hzjsj@hzbook.com

版权所有·侵权必究

封底无防伪标均为盗版

本书法律顾问: 北京大成律师事务所 韩光 / 邹晓东

书中介绍的结构、函数/宏

结构	卷2	卷3	结构	卷2	卷3
cmsghdr		217	sockaddr	58	
ifnet	51		sockaddr_in	127	
in_ifaddr	128		sockaddr_un		182
inpcb	574		socket	350	
mbuf	29		tcpcb	643	74
radix_node	461		tcphdr	641	
radix_node_head	460	60	tcpihdr	643	
rmxp_tao		61	tcpopt		97
route	174	85	timeval	83	
rtentry	464	61	unixdomain		181
rt_metrics	465	61	unixsw		182
rtqk_arg		65	unpcb		183

函数/宏	卷2	卷3	函数/宏	卷2	卷3
CC_INC		73	m_freem	41	
dtom	36		mtod	36	
in_addroute		62	pipe		202
in_clsroute		64	recvit	404	
in_inithead		62	recvmsg	403	
in_localaddr	144		rmx_tao		61
in_matroute		63	rtalloc	482	
in_pcbbind	585		rtallocl	483	
in_pcbconnect	590	72	rtrequest	487	
in_pcbdisconnect	594		sbappend	383	
in_pcbldaddr		71	sbappendaddr	384	
in_pcblookup	582		sbappendcontrol	384	
in_rtqkill		67	sbreserve	384	
in_rtqtime		65	sendit	391	
m_copy	41		sendmsg	388	
m_free	41				

(续)

函数/宏	卷2	卷3	函数/宏	卷2	卷3
SEQ_GEQ	649		tcp_mssrcvd		92
SEQ_GT	649		tcp_msssend		90
SEQ_LEQ	649		tcp_newtcpcb	667	84
SEQ_LT	649		tcp_output	681	77
socantrcvmore	353		tcp_rcvseqinit	756	
socantsendmore	353		tcp_reass	728	98
sockargs	362		tcp_rtlookup		85
socketpair		199	tcp_sendseqinit	756	
soclose	378		tcp_slowtimo	658	75
soconnect2		202	tcp_sysctl		126
socreate	359		tcp_template	707	
sofree	379		TCPT_RANGESET	656	
soisconnected	371		tcp_usrclosed	817	125
soisconnecting	353		tcp_usrreq	815	120
soisdisconnected	353		uipc_usrreq		185
soreceive	410		unp_attach		186
soreserve	384		unp_bind		189
sorflush	377		unp_connect		192
sorwakeup	383		unp_connect2		196
sosend	394	57	unp_detach		188
sowakeup	383		unp_discard		222
splnet	18		unp_disconnect		204
splx	18		unp_dispose		222
tcp_canceltimer	657		unp_drop		206
tcp_close	715	89	unp_externalize		221
tcp_connect		121	unp_gc		225
tcp_dooptions	745	97	unp_internalize		220
tcp_drop	713		unp_mark		230
tcp_gettaocache		86	unp_scan		223
tcp_init	651	75	unp_shutdown		206
tcp_input	739	103			

出版者的话

文艺复兴以来，源远流长的科学精神和逐步形成的学术规范，使西方国家在自然科学的各个领域取得了垄断性的优势；也正是这样的优势，使美国在信息技术发展的六十多年间名家辈出、独领风骚。在商业化的进程中，美国的产业界与教育界越来越紧密地结合，计算机学科中的许多泰山北斗同时身处科研和教学的最前线，由此而产生的经典科学著作，不仅擘划了研究的范畴，还揭示了学术的源变，既遵循学术规范，又自有学者个性，其价值并不会因年月的流逝而减退。

近年，在全球信息化大潮的推动下，我国的计算机产业发展迅猛，对专业人才的需求日益迫切。这对计算机教育界和出版界都既是机遇，也是挑战；而专业教材的建设在教育战略上显得举足轻重。在我国信息技术发展时间较短的现状下，美国等发达国家在其计算机科学发展的几十年间积淀和发展的经典教材仍有许多值得借鉴之处。因此，引进一批国外优秀计算机教材将对我国计算机教育事业的发展起到积极的推动作用，也是与世界接轨、建设真正的世界一流大学的必由之路。

机械工业出版社华章公司较早意识到“出版要为教育服务”。自1998年开始，我们就将工作重点放在了遴选、移译国外优秀教材上。经过多年的不懈努力，我们与 Pearson、McGraw-Hill、Elsevier、MIT、John Wiley & Sons、Cengage 等世界著名出版公司建立了良好的合作关系，从它们现有的数百种教材中甄选出 Andrew S. Tanenbaum、Bjarne Stroustrup、Brian W. Kernighan、Dennis Ritchie、Jim Gray、Afred V. Aho、John E. Hopcroft、Jeffrey D. Ullman、Abraham Silberschatz、William Stallings、Donald E. Knuth、John L. Hennessy、Larry L. Peterson 等大师名家的一批经典作品，以“计算机科学丛书”为总称出版，供读者学习、研究及珍藏。大理石纹理的封面，也正体现了这套丛书的品位和格调。

“计算机科学丛书”的出版工作得到了国内外学者的鼎力相助，国内的专家不仅提供了中肯的选题指导，还不辞劳苦地担任了翻译和审校的工作；而原书的作者也相当关注其作品在中国的传播，有的还专门为其书的中译本作序。迄今，“计算机科学丛书”已经出版了近500个品种，这些书籍在读者中树立了良好的口碑，并被许多高校采用为正式教材和参考书籍。其影印版“经典原版书库”作为姊妹篇也被越来越多实施双语教学的学校所采用。

权威的作者、经典的教材、一流的译者、严格的审校、精细的编辑，这些因素使我们的图书有了质量的保证。随着计算机科学与技术专业学科建设的不断完善和教材改革的逐渐深化，教育界对国外计算机教材的需求和应用都将步入一个新的阶段，我们的目标是尽善尽美，而反馈的意见正是我们达到这一终极目标的重要帮助。华章公司欢迎老师和读者对我们的工作提出建议或给予指正，我们的联系方式如下：

华章网站：www.hzbook.com

电子邮件：hzjsj@hzbook.com

联系电话：(010) 88379604

联系地址：北京市西城区百万庄南街1号

邮政编码：100037



华章科技图书出版中心

本书赞誉

“绝对值得一读！它说明了如何将科学的思想方法和分析方法应用于实际的技术问题……它体现了技术写作和思考的最高水平。”

——Marcus J. Ranum, 防火墙设计师

“是继既清楚又准确的系列卓越标准之后的又一杰出力作。该书的内容覆盖了T/TCP和HTTP，并剖析了WWW，特别及时。”

——Vern Paxson, 劳伦斯伯克利国家实验室网络研究小组

“对需要理解Web服务器行为细节的任何人来说，该书对HTTP的介绍都是无价之宝。”

——Jeffrey Mogul, 数字设备公司

“卷3是对前两卷的自然补充，包括了Web服务中的网络技术和TCP事务传输的深入介绍。”

——Pete Haverlock, 程序管理员, IBM

“在《TCP/IP详解》的最后一卷中，Rich Stevens保持了他在前两卷中给自己设定的高标准：清楚的表达和准确的技术细节。”

——Andras Olah, Twente大学

“这一卷保持了这套书前几卷中的极高质量，在新的方向上扩充了对网络实现技术的深入介绍。对于渴望了解当今Internet工作原理的任何人来说，这套书不可不读。”

——Ian Lance Taylor, 《GNU/Talyor UUCP》的作者

译者序

我们愿意向广大的读者推荐W. Richard Stevens关于TCP/IP的经典著作(共3卷)的中译本。本书是其中的第3卷——《TCP/IP详解 卷3: TCP事务协议、HTTP、NNTP和UNIX域协议》。

大家知道, TCP/IP已成为计算机网络事实上的标准。在关于TCP/IP的论著中, 最有影响的两部著作是: Douglas E. Comer的《用TCP/IP进行网际互连》(一套共3卷), 以及Stevens写的这3卷书。这两套巨著都很有名, 各有其特点。无论是从事计算机网络教学的教师还是进行计算机网络科研的技术人员, 这两套书都应当是必读的。

这套书的特点是内容丰富, 概念清楚且准确, 讲解详细, 例子很多。作者在书中举出的所有例子均在作者安装的计算机网络上做过实际验证, 而且书后还给出了许多经典的参考文献, 并一一写出评注。

第3卷是第1、2卷的继续和深入。读者在学习这一卷时, 应当先具备第1卷和第2卷所阐述的TCP/IP的基本知识和实现知识。本卷仍然采用大量的源代码来讲述协议及其应用的实现, 并且本卷使用的一部分源代码是对第1卷和第2卷中有关源代码的修改, 需要对照参考。这些内容对于编写TCP/IP网络应用程序的程序员和研究TCP/IP的计算机网络研究人员是非常有用的。

本卷的前言由胡谷雨翻译, 第1~5章由胡谷雨、马春华翻译, 第6~12章由胡谷雨、张晖翻译, 第13~15章由吴礼发、李旺翻译, 第16~18章由吴礼发、金凤林翻译, 附录由胡谷雨翻译。全书由谢希仁进行校阅。

限于水平, 翻译中不妥或错误之处在所难免, 敬请广大读者批评指正。

前 言

引言和本书的组织

本书是套书《TCP/IP详解》的第3卷，这套书的卷1是[Stevens 1994]，卷2是[Wright and Stevens 1995]。本书分成三个部分，每个部分覆盖了不同的内容。

- 1) TCP事务协议，通常叫作T/TCP。这是对TCP的扩展，其设计目的是使客户-服务器事务更快、更高效和更可靠。这个目标的实现省略了连接开始时TCP的三次握手，并缩短了连接结束时TIME_WAIT状态的持续时间。我们将会看到，在客户-服务器事务中，T/TCP的性能与UDP相当，而且T/TCP具有可靠性和适应性，这两点相对UDP来说都是很大的改进。

事务是这样定义的：一个客户向服务器发出请求，接下来是服务器给出响应(这里的名词“事务”(transaction)并非数据库中的事务处理，数据库中的事务处理有封锁、两步提交和回退)。

- 2) TCP/IP应用，特别是HTTP(超文本传输协议，WWW的基础)和NNTP(网络新闻传输协议Usenet新闻系统的基础)。

- 3) Unix域协议。这些协议是所有Unix的TCP/IP实现中都提供的，在许多非Unix的实现中也有提供。这些协议提供了一种进程之间通信(IPC)的手段，采用了与TCP/IP中一样的插口[⊖]接口。当客户与服务器进程在同一主机上时，Unix域协议通常要比TCP/IP快1倍。

第一部分是介绍T/TCP，又分成两个小部分。第1~4章介绍协议，并给出了大量实例来说明它们是怎样工作的。这些材料主要是对卷1中24.7节的补充，在那里对T/TCP只是做了简单的介绍。第5~12章介绍T/TCP在4.4BSD-Lite网络代码(即卷2中给出的代码)中的确切实现。由于最早的T/TCP实现迟至1994年9月才发布，已经是本书卷1出版一年以后了，那时卷2也快完成了，因此T/TCP的详细叙述，包括诸多实例和所有的实现细节都只好放在本系列书的卷3中了。

第二部分介绍HTTP和NNTP应用，是卷1的第25~30章中介绍的TCP/IP应用的延续。在卷1出版后的两年里，随着Internet的发展，HTTP得到了极大的流行，而NNTP的使用则在最近的10多年中每年增长了大约75%。T/TCP对HTTP来说也是非常好的，可以这样来用TCP：在少量数据传输中缩短连接时间，因为这种时候连接的建立和拆除时间往往占总时间的大头。在繁忙的Web服务器上，成千上万个不同而且不断变化的客户对HTTP(因此也对TCP)的高负荷使用，也提供了唯一可以对服务器上确切的分组进行考察的机会(第14章)，可以回顾卷1和卷2中给出的TCP/IP的许多特性。

第三部分中的Unix域协议原本是准备在卷2中介绍的，但由于卷2已多达1200页[⊕]而删

⊖ 插口对应的原文是socket，现更常译为“套接字”。——编辑注

⊕ 指原书英文版。——编辑注

去了。在名为《TCP/IP详解》这样的套书中夹杂着TCP/IP以外的协议不免令人奇怪，但Unix域协议几乎15年前就已经伴随着BSD版TCP/IP的实现在4.2BSD中发布了。今天，它们在任何一个从伯克利衍生而来的内核中都在频繁地使用，但它们的使用往往“被掩盖在后台”，大多数用户不知道它们的存在。除了在从伯克利衍生而来的内核中充当Unix管道的基础外，它们的另一个大用户是当客户程序和服务器程序在同一主机(典型的情况是工作站)上时的X Window系统。Unix域的插口也用于进程之间传递描述符，是进程之间通信的一个强大工具。由于Unix域协议所用的插口API(应用编程接口)与TCP/IP所用的插口API几乎是相同的，Unix域协议以最小的代码变化提供了一个简单的手段来增强本地应用的性能。

以上三个部分的每个部分都可以独立阅读。

读者

与这套书的前两卷一样，这一卷是为所有想要理解TCP/IP如何工作的人写的：编写网络应用的程序员，负责维护采用TCP/IP的计算机网络的系统管理员，以及在日常工作中经常与TCP/IP应用程序打交道的用户。

第一和第二部分是理解TCP/IP工作原理的基础。不熟悉TCP/IP的读者应该看看这套书的卷1，见[Stevens 1994]，以便对TCP/IP协议集有一个全面的了解。第一部分的前半部分(第1~4章，TCP/IP中的概念和例子)与卷2无关，可以直接阅读。但后半部分(第5~12章，T/TCP的实现)则需要先熟悉4.4 BSD-Lite网络程序，这些内容在卷2中介绍。

在整本书中有大量的向前和向后参考索引，这些参考索引是针对本书的两个主题，以及对卷1和卷2的内容，为想要了解更详细内容的读者提供的。在本书最后有书中用到的所有缩略语，书中介绍的所有结构、函数和宏(以字母顺序排列)及其介绍起始页码的交叉索引。如果本书引用了卷2中的定义，则该交叉索引也列出了卷2中的定义。

源码版权

本书中引自4.4BSD-Lite版的所有源码(源程序)都包括下面这样的版权声明：

```
/*
 * Copyright (c) 1982, 1986, 1988, 1990, 1993, 1994
 *   The Regents of the University of California.  All rights reserved.
 *
 * Redistribution and use in source and binary forms, with or without
 * modification, are permitted provided that the following conditions
 * are met:
 * 1. Redistributions of source code must retain the above copyright
 *   notice, this list of conditions and the following disclaimer.
 * 2. Redistributions in binary form must reproduce the above copyright
 *   notice, this list of conditions and the following disclaimer in the
 *   documentation and/or other materials provided with the distribution.
 * 3. All advertising materials mentioning features or use of this software
 *   must display the following acknowledgement:
 *   This product includes software developed by the University of
 *   California, Berkeley and its contributors.
 * 4. Neither the name of the University nor the names of its contributors
```

```

*   may be used to endorse or promote products derived from this software
*   without specific prior written permission.
*
* THIS SOFTWARE IS PROVIDED BY THE REGENTS AND CONTRIBUTORS ``AS IS'' AND
* ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE
* IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE
* ARE DISCLAIMED. IN NO EVENT SHALL THE REGENTS OR CONTRIBUTORS BE LIABLE
* FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL
* DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS
* OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION)
* HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT
* LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY
* OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF
* SUCH DAMAGE.
*/

```

第6章路由表的源码则包括下面这样的版权声明：

```

/*
* Copyright 1994, 1995 Massachusetts Institute of Technology
*
* Permission to use, copy, modify, and distribute this software and
* its documentation for any purpose and without fee is hereby
* granted, provided that both the above copyright notice and this
* permission notice appear in all copies, that both the above
* copyright notice and this permission notice appear in all
* supporting documentation, and that the name of M.I.T. not be used
* in advertising or publicity pertaining to distribution of the
* software without specific, written prior permission. M.I.T. makes
* no representations about the suitability of this software for any
* purpose. It is provided "as is" without express or implied
* warranty.
*
* THIS SOFTWARE IS PROVIDED BY M.I.T. ``AS IS''. M.I.T. DISCLAIMS
* ALL EXPRESS OR IMPLIED WARRANTIES WITH REGARD TO THIS SOFTWARE,
* INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF
* MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. IN NO EVENT
* SHALL M.I.T. BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL,
* SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT
* LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF
* USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND
* ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY,
* OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT
* OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF
* SUCH DAMAGE.
*/

```

印刷惯例

当需要显示交互的输入和输出信息时，将用黑体表示键盘输入，而计算机输出则用 Courier 体，并用中文宋体做注释。

```

sun % telnet www.aw.com 80  连接到HTTP服务器
Trying 192.207.117.2...      本行和下一行由Telnet服务器输出
Connected to aw.com.

```

书中总是把系统名作为命令解释程序提示符的一部分(例如sun)，以说明命令是在哪个主机上执行的。在正文中引用的程序名通常都是首字母大写(如Telnet和Tcpdump)，以避免过多的字体形式。

在整本书中，我们会使用这种缩进格式的附加说明来描述实现细节或历史观点。

W. Richard Stevens

图森，亚利桑那

1995年11月

`rstevens@noao.edu`

`http://www.noao.edu/~rstevens`

目 录

出版者的话
本书赞誉
译者序
前言

第一部分 TCP事务协议

第1章 T/TCP概述	1	第4章 T/TCP协议 (续)	43
1.1 概述	1	4.1 概述	43
1.2 UDP上的客户-服务器	1	4.2 客户的端口号和TIME_WAIT状态	43
1.3 TCP上的客户-服务器	6	4.3 设置TIME_WAIT状态的目的	45
1.4 T/TCP上的客户-服务器	12	4.4 TIME_WAIT状态的截断	48
1.5 测试网络	15	4.5 利用TAO跳过三次握手	51
1.6 时间测量程序	15	4.6 小结	55
1.7 应用	17	第5章 T/TCP实现: 插口层	56
1.8 历史	19	5.1 概述	56
1.9 实现	20	5.2 常量	56
1.10 小结	21	5.3 sosend函数	56
第2章 T/TCP协议	23	5.4 小结	58
2.1 概述	23	第6章 T/TCP实现: 路由表	59
2.2 T/TCP中的新TCP选项	23	6.1 概述	59
2.3 T/TCP实现所需变量	25	6.2 代码介绍	59
2.4 状态变迁图	27	6.3 radix_node_head结构	60
2.5 T/TCP的扩展状态	28	6.4 rtentry结构	61
2.6 小结	30	6.5 rt_metrics结构	61
第3章 T/TCP使用举例	31	6.6 in_inithead函数	61
3.1 概述	31	6.7 in_addroute函数	62
3.2 客户重新启动	31	6.8 in_matroute函数	63
3.3 常规的T/TCP事务	33	6.9 in_clsroute函数	63
3.4 服务器收到过时的重复SYN	34	6.10 in_rtqtime函数	64
3.5 服务器重启动	35	6.11 in_rtqkill函数	66
3.6 请求或应答超出报文段最大长度	36	6.12 小结	69
3.7 向后兼容性	39	第7章 T/TCP实现: 协议控制块	70
3.8 小结	41	7.1 概述	70
		7.2 in_pcbaddr函数	71
		7.3 in_pcbconnect函数	71
		7.4 小结	72
		第8章 T/TCP实现: TCP概要	73
		8.1 概述	73
		8.2 代码介绍	73
		8.3 TCP的protosw结构	74

16.4 编码举例	177	17.16 PRU_SHUTDOWN请求和unp_shutdown 函数	205
16.5 小结	179	17.17 PRU_ABORT请求和unp_drop函数.....	206
第17章 Unix域协议: 实现	180	17.18 其他各种请求	207
17.1 概述	180	17.19 小结	209
17.2 代码介绍	180	第18章 Unix域协议: I/O和描述符的传递.....	210
17.3 Unix domain和protosw结构	181	18.1 概述	210
17.4 Unix域插口地址结构	182	18.2 PRU_SEND和PRU_RCVD请求.....	210
17.5 Unix域协议控制块	183	18.3 描述符的传递	214
17.6 uipc_usrreq函数	185	18.4 unp_internalize函数	218
17.7 PRU_ATTACH请求和unp_attach函数.....	186	18.5 unp_externalize函数	220
17.8 PRU_DETACH请求和unp_detach函数.....	187	18.6 unp_discard函数	221
17.9 PRU_BIND请求和unp_bind函数.....	189	18.7 unp_dispose函数	222
17.10 PRU_CONNECT请求和unp_connect 函数	191	18.8 unp_scan函数	222
17.11 PRU_CONNECT2请求和unp_connect2 函数	195	18.9 unp_gc函数	223
17.12 socketpair系统调用	198	18.10 unp_mark函数.....	230
17.13 pipe系统调用	202	18.11 性能(再讨论)	231
17.14 PRU_ACCEPT请求	203	18.12 小结	231
17.15 PRU_DISCONNECT请求和 unp_disconnect函数	204	附录A 测量网络时间.....	232
		附录B 编写T/TCP应用程序	242
		参考文献	246
		缩略语	251

第一部分 TCP事务协议

第1章 T/TCP概述

1.1 概述

本章首先介绍客户-服务器事务概念。我们从使用UDP的客户-服务器应用开始，这是最简单的情形。接着我们编写使用TCP的客户和服务器程序，并由此考察两台主机间交互的TCP/IP分组。然后我们使用T/TCP，证明利用T/TCP可以减少分组数，并给出为利用T/TCP需要对两端的源代码所做的最少改动。

接下来介绍了运行书中示例程序的测试网络，并对分别使用UDP、TCP和T/TCP的客户-服务器应用程序进行了简单的时间耗费比较。我们考察了一些使用TCP的典型Internet应用程序，看看如果两端都支持T/TCP，将需要做哪些修改。紧接着，简要介绍了Internet协议族中事务协议的发展历史，概略叙述了现有的T/TCP实现。

本书全文以及有关T/TCP的文献中，事务一词的含义都是指客户向服务器发出一个请求，然后服务器对该请求做出应答。Internet中最常见的一个例子是，客户向域名服务器(DNS)发出请求，查询域名对应的IP地址，然后域名服务器给出响应。本书中的事务这个术语并没有数据库中的事务那样的含义：加锁、两步提交、回退，等等。

1.2 UDP上的客户-服务器

我们先来看一个简单的UDP客户-服务器应用程序的例子，其客户程序源代码如图1-1所示。在这个例子中，客户向服务器发出一个请求，服务器处理该请求，然后发回一个应答。

```
1 #include "cliserv.h" udplici.c
2 int
3 main(int argc, char *argv[])
4 { /* simple UDP client */
5     struct sockaddr_in serv;
6     char request[REQUEST], reply[REPLY];
7     int sockfd, n;
8
9     if (argc != 2)
10        err_quit("usage: udpcli <IP address of server>");
11
12    if ((sockfd = socket(PF_INET, SOCK_DGRAM, 0)) < 0)
13        err_sys("socket error");
14
15    memset(&serv, 0, sizeof(serv));
16    serv.sin_family = AF_INET;
17    serv.sin_addr.s_addr = inet_addr(argv[1]);
18    serv.sin_port = htons(UDP_SERV_PORT);
```

图1-1 UDP上的简单客户程序

```

16  /* form request[] ... */
17  if (sendto(sockfd, request, REQUEST, 0,
18          (SA) &serv, sizeof(serv)) != REQUEST)
19      err_sys("sendto error");

20  if ((n = recvfrom(sockfd, reply, REPLY, 0,
21          (SA) NULL, (int *) NULL)) < 0)
22      err_sys("recvfrom error");

23  /* process "n" bytes of reply[] ... */

24  exit(0);
25 }

```

udcli.c

图1-1 (续)

本书中所有源代码的格式都是这样。每一非空行前面都标有行号。正文中叙述某段源代码时，这段源代码的起始和结束行号标记于正文段落的左边，如下面的正文所示。有时这些段落前面会有一小段说明，对所描述的源代码进行概要说明。源代码段开头和结尾处的水平线标明源代码段所在的文件名。这些文件名通常都是指我们在1.9节中将介绍的4.4版BSD-Lite中发布的文件。

我们来讨论这个程序的一些有关特性，但不详细描述插口函数，因为我们假设读者对这些函数有一些基本的认识。关于插口函数的细节在参考书[Stevens 1990]的第6章中可以找到。图1-2给出了头文件cliserv.h。

1. 创建UDP插口

10-11 socket函数用于创建一个UDP插口，并将一个非负的插口描述符返回给调用进程。出错处理函数err_sys参见参考书[Stevens 1992]的附录B.2。这个函数可以接受任意数目的参数，但要用vsprintf函数对它们格式化，然后这个函数会打印出系统调用所返回的errno值所对应的Unix出错信息，然后终止进程。

2. 填写服务器地址

12-15 首先用memset函数将Internet插口地址结构清零，然后填入服务器的IP地址和端口号。为简明起见，我们要求用户在程序运行中通过命令行输入一个点分十进制数形式的IP地址(argv[1])。服务器端口号(UDP_SERV_PORT)在头文件cliserv.h中用#define定义，在本章的所有程序首部中都包含了该头文件。这样做是为了使程序简洁，并避免使调用gethostbyname和getservbyname函数的源代码复杂化。

3. 构造并向服务器发送请求

16-19 客户程序构造一个请求(只用一行注释来表示)，并用sendto函数将其发出，这样就有一个UDP数据报发往服务器。同样是为了简明起见，我们假设请求(REQUEST)和应答(REPLY)的报文长度为固定值。实用的程序应当按照请求和应答的最大长度来分配缓存空间，但实际的请求和应答报文长度是变化的，而且一般都比较小。

4. 读取和处理服务器的应答

20-23 调用recvfrom函数将使进程阻塞(即置为睡眠状态)，直至收到一个数据报。接着客户进程处理应答(用一行注释来表示)，然后进程终止。

由于recvfrom函数中没有超时机制，请求报文或应答报文中任何一个丢失都将造成该进程永久挂起。事实上，UDP客户-服务器应用的一个基本问题就是对现实世界中的此类错误缺少健壮性。在本节的末尾将对这个问题做更详细的讨论。