

TCP/IP 详解

卷3: TCP事务协议, HTTP, NNTP和UNIX域协议

TCP/IP Illustrated

Volume 3:

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

(美) W. Richard Stevens 著

胡谷雨 吴礼发 等译

谢希仁 校



机械工业出版社
China Machine Press



Addison-Wesley

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. Richard Stevens 著

胡谷雨 吴礼发 等译

谢希仁 校



机械工业出版社
China Machine Press

本书是“TCP/IP详解系列”的延续。主要内容包括：TCP事务协议，即T/TCP，这是对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，这些协议提供了进程之间通信的一种手段。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快一倍。本书同样采用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如何工作的人，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。

W. Richard Stevens: TCP/IP Illustrated, Volume 3: TCP for Transactions, HTTP, NNTP, and the UNIX Domain Protocols.

Original edition copyright © 1996 by Addison Wesley.

Chinese edition published by arrangement with Addison Wesley Longman, Inc

All rights reserved.

本书中文简体字版由美国Addison Wesley公司授权机械工业出版社独家出版。未经出版者书面许可，不得以任何方式复制或抄袭本书内容。

版权所有，侵权必究。

本书版权登记号：图字：01-1999-2857

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)史蒂文斯(Stevens, W. R.)著; 胡谷雨等译. -北京: 机械工业出版社, 2000.9

(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 7-111-07568-4

I. T… II. ①史… ②胡… III. 计算机网络-传输控制协议 IV. TN915.04

中国版本图书馆CIP数据核字(2000)第36105号

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: 陈贤舜

北京忠信诚胶印厂印刷·新华书店北京发行所发行

2000年9月第1版第1次印刷

787mm × 1092mm 1/16 · 16.75印张

印数: 0 001-10 000册

定价: 35.00元

凡购本书，如有倒页、脱页、缺页，由本社发行部调换

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。这里是其中的第3卷:《TCP/IP详解 卷3: TCP事务协议、HTTP、NNTP、UNIX域协议》。

大家知道, TCP/IP已成为计算机网络的事实上的标准。在关于TCP/IP的论著中, 最有影响的就是两部著作。一部是Douglas E. Comer的《用TCP/IP进行网际互连》, 一套共3卷(中译本已由电子工业出版社于1998年出版), 而另一部就是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的教师还是进行科研的技术人员, 这两套书都应当是必读的。

本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经过实际验证。在本书的最后, 作者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、关于TCP/IP的基本知识和实现知识。本卷仍然采用了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。

本卷书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。

译者

2000年5月于解放军理工大学, 南京

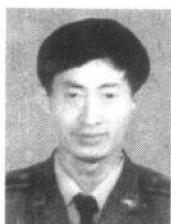
译、校者介绍



谢希仁，中国人民解放军理工大学(南京)计算机系教授，全军网络技术研究中心主任，博士研究生导师，1952年毕业于清华大学电机系电信专业。所编写的《计算机网络》于1992年获全国优秀教材奖。1999年再版的《计算机网络》第2版为普通高等教育“九五”国家级重点教材。近来还主持翻译了Comer写的《TCP/IP网际互联》计算机网络经典教材一套三卷本(电子工业出版社1998年出版)，Harnedy写的《简单网络管理协议教程》(电子工业出版社1999年出版)。



胡谷雨，1963年出生，现为中国人民解放军理工大学(南京)指挥自动化学院计算机系教授、博士研究生导师，全军网络技术研究中心副主任。1983年于浙江大学无线电系获学士学位，并分别于1989和1992年在中国人民解放军通信工程学院计算机网络专业获硕士和博士学位。曾参加和负责了多项军队和863计划的重点科研项目，独立编著出版了《现代通信网和计算机网管理》(电子工业出版社，1996)，在《电子学报》、《通信学报》等刊物上发表论文多篇。目前主要从事计算机网络、通信网管理等方面的教学和研究工作。



吴礼发，1968年出生，现为中国人民解放军理工大学(南京)博士后。分别于1991和1995年在中国人民解放军通信工程学院获学士和硕士学位。1998年10月在南京大学获博士学位。曾作为主要成员参加过多项军队和地方重要科研项目，在《中国科学》、《计算机学报》、《通信学报》和《软件学报》等刊物上发表论文10余篇。目前主要从事计算机网络管理、VSAT卫星通信系统的网控系统和软件工程等方面的教学和研究工作。

对《TCP/IP详解 卷3: TCP事务协议、HTTP、 NNTP和UNIX域协议》的评语

“是绝对值得一读的实例，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又一个杰出后继。该书的内容覆盖了T/TCP和HTTP，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

“该书内容中对HTTP的介绍对需要理解Web服务器行为细节的任何人都是无价之宝。”

——Jeffrey Mogul, 数字设备公司

“卷3是对该书系列的自然补充，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

“在这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

“这一卷保持了该系列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。该书整个系列对于渴望了解当今Internet如何工作的任何人都是不可不读的。”

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

前 言

引言和本书的组织

本书是“TCP/IP详解系列”的延续：此系列的卷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，通常叫做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(超文本传送协议，WWW的基础)和NNTP(网络新闻传送协议，Usenet新闻系统的基础)。

(3) Unix域协议。这些协议是所有Unix中的TCP/IP实现中都提供的，在许多非Unix的实现中也提供。这些协议提供了进程之间通信(IPC)的一种手段，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快一倍。

第一部分，即对T/TCP的介绍，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2小部分，即第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，即HTTP和NNTP应用，是卷1的第25~30章中介绍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查的机会(第14章)，可以观察卷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

[⊖] 指原书英文版。——编者注

去了。在书名为《TCP/IP详解》这样的系列书中夹杂着TCP/IP以外的协议不免令人生奇，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，也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

与本系列书的前两卷一样，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的计算机和网络的系统管理员，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的读者应该看看本系列书的卷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，封底背面则有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参考列表。如果本书引用了卷2中的定义，则该交叉参考列表也列出了卷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)都包括有下面这样的版权说明：

```
/*
 * 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
```


VIII

```
*   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说明:

```
/*
 * 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, 并用中文宋体作注释。

```
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。

在整个这本书中，我们将用像上述这段文字一样的缩进、附加的注释来说明实现细节中的历史阶段。

W. Richard Stevens

Tucson, Arizona

1995年11月

rstevens@noao.edu

<http://www.noao.edu/~rstevens>

目 录

译者序
前言

第一部分 TCP事务协议

第1章 T/TCP概述	1	4.4 TIME_WAIT状态的截断	48
1.1 概述	1	4.5 利用TAO跳过三次握手	51
1.2 UDP上的客户-服务器	1	4.6 小结	55
1.3 TCP上的客户-服务器	6	第5章 T/TCP协议的实现: 插口层	56
1.4 T/TCP上的客户-服务器	12	5.1 概述	56
1.5 测试网络	15	5.2 常量	56
1.6 时间测量程序	15	5.3 sosend函数	56
1.7 应用	17	5.4 小结	58
1.8 历史	19	第6章 T/TCP的实现: 路由表	59
1.9 实现	20	6.1 概述	59
1.10 小结	21	6.2 代码介绍	59
第2章 T/TCP协议	23	6.3 radix_node_head结构	60
2.1 概述	23	6.4 rtable结构	61
2.2 T/TCP中的新TCP选项	23	6.5 rt_metrics结构	61
2.3 T/TCP实现所需变量	25	6.6 in_inithead函数	61
2.4 状态变迁图	27	6.7 in_addroute函数	62
2.5 T/TCP的扩展状态	28	6.8 in_matroute函数	63
2.6 小结	30	6.9 in_clsroute函数	63
第3章 T/TCP使用举例	31	6.10 in_rtqtime函数	64
3.1 概述	31	6.11 in_rtqkill函数	66
3.2 客户重新启动	31	6.12 小结	69
3.3 常规的T/TCP事务	33	第7章 T/TCP实现: 协议控制块	70
3.4 服务器收到过时的重复SYN	34	7.1 概述	70
3.5 服务器重启动	35	7.2 in_pcbladdr函数	71
3.6 请求或应答超出报文段最大长度MSS	36	7.3 in_pcbconnect函数	71
3.7 向后兼容性	39	7.4 小结	72
3.8 小结	41	第8章 T/TCP实现: TCP概要	73
第4章 T/TCP协议 (续)	43	8.1 概述	73
4.1 概述	43	8.2 代码介绍	73
4.2 客户的端口号和TIME_WAIT状态	43	8.3 TCP的protosw结构	74
4.3 设置TIME_WAIT状态的目的	45	8.4 TCP控制块	74
		8.5 tcp_init函数	75
		8.6 tcp_slowtime函数	75
		8.7 小结	76
		第9章 T/TCP实现: TCP输出	77

9.1 概述77

9.2 tcp_output函数77

 9.2.1 新的自动变量77

 9.2.2 增加隐藏的状态标志77

 9.2.3 在SYN_SENT状态不要重传SYN78

 9.2.4 发送器的糊涂窗口避免机制78

 9.2.5 有RST或SYN标志时强制发送报文段79

 9.2.6 发送MSS选项80

 9.2.7 是否发送时间戳选项80

 9.2.8 发送T/TCP的CC选项80

 9.2.9 根据TCP选项调整数据长度83

9.3 小结83

第10章 T/TCP实现: TCP函数84

10.1 概述84

10.2 tcp_newtcpcb函数84

10.3 tcp_rtlookup函数85

10.4 tcp_gettaocache函数86

10.5 重传超时间隔的计算86

10.6 tcp_close函数89

10.7 tcp_msssend函数90

10.8 tcp_mssrcvd函数91

10.9 tcp_dooptions函数96

10.10 tcp_reass函数98

10.11 小结99

第11章 T/TCP实现: TCP输入101

11.1 概述101

11.2 预处理103

11.3 首部预测104

11.4 被动打开的启动105

11.5 主动打开的启动108

11.6 PAWS: 防止序号重复114

11.7 ACK处理115

11.8 完成被动打开和同时打开115

11.9 ACK处理 (续)116

11.10 FIN处理118

11.11 小结119

第12章 T/TCP实现: TCP用户请求120

12.1 概述120

12.2 PRU_CONNECT请求120

12.3 tcp_connect函数120

12.4 PRU_SEND和PRU_SEND_EOF请求124

12.5 tcp_usrclosed函数125

12.6 tcp_sysctl函数126

12.7 T/TCP的前景126

12.8 小结127

第二部分 TCP的其他应用

第13章 HTTP: 超文本传送协议129

13.1 概述129

13.2 HTTP和HTML概述130

13.3 HTTP132

 13.3.1 报文类型: 请求与响应132

 13.3.2 首部字段133

 13.3.3 响应代码133

 13.3.4 各种报文头举例134

 13.3.5 例子: 客户程序缓存135

 13.3.6 例子: 服务器重定向136

13.4 一个例子136

13.5 HTTP的统计资料138

13.6 性能问题139

13.7 小结141

第14章 在HTTP服务器上找到的分组142

14.1 概述142

14.2 多个HTTP服务器144

14.3 客户端SYN的到达间隔时间145

14.4 RTT的测量149

14.5 用listen设置入连接队列的容量150

14.6 客户端的SYN选项154

14.7 客户端的SYN重传156

14.8 域名157

14.9 超时的持续探测157

14.10 T/TCP路由表大小的模拟160

14.11 mbuf的交互162

14.12 TCP的PCB高速缓存和首部预测163

14.13 小结165

第15章 NNTP: 网络新闻传送协议166

15.1 概述166

15.2 NNTP167

15.3 一个简单的新闻客户端	170	17.13 pipe系统调用	202
15.4 一个复杂的新闻客户端	171	17.14 PRU_ACCEPT请求	203
15.5 NNTP的统计资料	172	17.15 PRU_DISCONNECT请求和 unp_disconnect函数	204
15.6 小结	173	17.16 PRU_SHUTDOWN请求和unp_shutdown 函数	205
第三部分 Unix域协议			
第16章 Unix域协议: 概述	175	17.17 PRU_ABORT请求和unp_drop函数	206
16.1 概述	175	17.18 其他各种请求	207
16.2 用途	176	17.19 小结	209
16.3 性能	177	第18章 Unix域协议: I/O和描述符的传递	210
16.4 编码举例	177	18.1 概述	210
16.5 小结	179	18.2 PRU_SEND和PRU_RCVD请求	210
第17章 Unix域协议: 实现	180	18.3 描述符的传递	214
17.1 概述	180	18.4 unp_internalize函数	218
17.2 代码介绍	180	18.5 unp_externalize函数	220
17.3 Unix domain和protosw结构	181	18.6 unp_discard函数	221
17.4 Unix域插口地址结构	182	18.7 unp_dispose函数	222
17.5 Unix域协议控制块	183	18.8 unp_scan函数	222
17.6 uipc_usrreq函数	185	18.9 unp_gc函数	223
17.7 PRU_ATTACH请求和unp_attach函数	186	18.10 unp_mark函数	230
17.8 PRU_DETACH请求和unp_detach函数	187	18.11 性能(再讨论)	231
17.9 PRU_BIND请求和unp_bind函数	189	18.12 小结	231
17.10 PRU_CONNECT请求和unp_connect 函数	191	附录A 测量网络时间	232
17.11 PRU_CONNECT2请求和unp_connect2 函数	195	附录B 编写T/TCP应用程序	242
17.12 socketpair系统调用	198	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作出应答。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" udpcli.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。

在头文件cliserv.h中，我们将SA定义为struct sockaddr*，即指向一般的插口地址结构的指针。每当有一个插口函数需要一个指向插口地址结构的指针时，该指针必须被置为指向一个一般性插口地址结构的指针。这是由于插口函数先于ANSI C标准出现，在80年代早期开发插口函数的时候，void*(空类型)指针类型尚不可用。问题是，“struct sockaddr*”总共有17个字符，这经常使这一行源代码超出屏幕(或书本页面)的右边界，因此我们将其缩写成为SA。这个缩写是从BSD内核源代码中借用过来的。

图1-2给出了在本章所有程序中都包含的头文件cliserv.h。

```

1 /* Common includes and defines for UDP, TCP, and T/TCP
2  * clients and servers */
3 #include <sys/types.h>
4 #include <sys/socket.h>
5 #include <netinet/in.h>
6 #include <arpa/inet.h>
7 #include <stdio.h>
8 #include <stdlib.h>
9 #include <string.h>
10 #include <unistd.h>
11 #define REQUEST 400 /* max size of request, in bytes */
12 #define REPLY 400 /* max size of reply, in bytes */
13 #define UDP_SERV_PORT 7777 /* UDP server's well-known port */
14 #define TCP_SERV_PORT 8888 /* TCP server's well-known port */
15 #define TTCP_SERV_PORT 9999 /* T/TCP server's well-known port */
16 /* Following shortens all the type casts of pointer arguments */
17 #define SA struct sockaddr *
18 void err_quit(const char *,...);
19 void err_sys(const char *,...);
20 int read_stream(int, char *, int);

```

cliserv.h

cliserv..

图1-2 本章各程序中均包含的头文件cliserv.h

图1-3给出了相应的UDP服务器程序。

```

1 #include "cliserv.h"
2 int
3 main()
4 {
5     /* simple UDP server */
6     struct sockaddr_in serv, cli;
7     char request[REQUEST], reply[REPLY];
8     int sockfd, n, cliilen;
9     if ((sockfd = socket(PF_INET, SOCK_DGRAM, 0)) < 0)
10        err_sys("socket error");
11    memset(&serv, 0, sizeof(serv));
12    serv.sin_family = AF_INET;
13    serv.sin_addr.s_addr = htonl(INADDR_ANY);
14    serv.sin_port = htons(UDP_SERV_PORT);

```

udpsero.

图1-3 与图1-1的UDP客户程序对应的UDP服务器程序

```

14     if (bind(sockfd, (SA) &serv, sizeof(serv)) < 0)
15         err_sys("bind error");

16     for (;;) {
17         clilen = sizeof(cli);
18         if ((n = recvfrom(sockfd, request, REQUEST, 0,
19                         (SA) &cli, &clilen)) < 0)
20             err_sys("recvfrom error");

21         /* process "n" bytes of request[] and create reply[] ... */

22         if (sendto(sockfd, reply, REPLY, 0,
23                  (SA) &cli, sizeof(cli)) != REPLY)
24             err_sys("sendto error");
25     }
26 }

```

udpserv.c

图1-3 (续)

5. 创建UDP插口和绑定本机地址

8-15 调用socket函数创建一个UDP插口，并在其Internet插口地址结构中填入服务器的本机地址。这里本机地址设置为通配符(INADDR_ANY)，这意味着服务器可以从任何一个本机接口接收数据报(假设服务器是多宿主的，即可以有多个网络接口)。端口号设为服务器的知名端口(UDP_SERV_PORT)，该常量也在前面讲过的头文件cliserv.h中定义。本机IP地址和知名端口用bind函数绑定到插口上。

6. 处理客户请求

16-25 接下来，服务器程序就进入一个无限循环：等待客户程序的请求到达(recvfrom)，处理该请求(我们只用一行注释来表示处理动作)，然后发出应答(sendto)。

这只是最简单的UDP客户-服务器应用。实际中常见的例子是域名服务系统(DNS)。DNS客户(称作解析器)通常是一般客户应用程序(例如，Telnet客户、FTP客户或WWW浏览器)的一个部分。解析器向DNS服务器发出一个UDP数据报，查询某一域名对应的IP地址。服务器发回的应答通常也是一个UDP数据报。

如果观察客户向服务器发送请求时双方交换的分组，我们就会得到图1-4这样的时间系列，页面上时间自上而下递增。服务器程序先启动，其行为过程给在图1-4的右半部，客户程序稍后启动。

我们分别来看客户和服务器程序中调用的函数及其相应内核执行的动作。在对socket函数的两次调用中，上下紧挨着的两个箭头表示内核执行请求的动作并立即返回。在调用sendto函数时，尽管内核也立即返回，但实际上已经发出了一个UDP数据报。为简明起见，我们假设客户程序的请求和服务器程序的应答所生成的IP数据报的长度都小于网络的最大传输单元(MTU)，IP数据报不必分段。

在这个图中，有两次调用recvfrom函数使进程睡眠，直到有数据报到达才被唤醒。我们把内核中相应的例程记为sleep和wakeup。

最后，我们还在图中标出了事务所耗费的时间。图1-4的左侧标示的是客户端测得的事务时间：从客户发出请求到收到服务器的应答所经历的时间。组成这段事务时间的数值标在图的右侧：RTT + SPT，其中RTT是网络往返时间，SPT是服务器处理客户请求的时间。UDP客户-服务器事务的最短时间就是RTT + SPT。