

Addison
Wesley

TURING

图灵原版计算机科学系列

UNIX

网络编程

卷2：进程间通信

英文版 · 第2版

UNIX Network Programming

Volume 2: Interprocess Communications, *Second Edition*

· [美] W. Richard Stevens · 著

 人民邮电出版社
POSTS & TELECOM PRESS

TURING

图灵原版计算机科学系列

UNIX

网络编程

卷2：进程间通信

英文版·第2版

UNIX Network Programming

Volume 2: Interprocess Communications, *Second Edition*

[美] W. Richard Stevens 著

人民邮电出版社
北京

图书在版编目 (CIP) 数据

UNIX网络编程=UNIX Network Programming:第2版. 第2卷. 进程间通信: 英文/ (美) 史蒂文斯 (Stevens, W.R.) 著. —北京: 人民邮电出版社, 2009.11

(图灵原版计算机科学系列)

ISBN 978-7-115-21511-6

I. ①U… II. ①史… III. ①UNIX操作系统—程序设计—英文IV. ①TP316.81

中国版本图书馆CIP数据核字 (2009) 第172154号

内 容 提 要

本书是一部UNIX网络编程的经典之作。进程间通信 (IPC) 几乎是所有Unix程序性能的关键, 理解IPC也是理解如何开发不同主机间网络应用程序的必要条件。本书从对Posix IPC和System V IPC的内部结构开始讨论, 全面深入地介绍了4种IPC形式: 消息传递 (管道、FIFO、消息队列)、同步 (互斥锁、条件变量、读写锁、文件与记录锁、信号量)、共享内存 (匿名共享内存、具名共享内存) 及远程过程调用 (Solaris 门、Sun RPC)。附录中给出了测量各种IPC形式性能的方法。

本书内容详尽且具权威性, 几乎每章都提供精选的习题, 并提供了部分习题的答案, 是网络研究和开发人员理想的参考书。

图灵原版计算机科学系列

UNIX网络编程 卷2: 进程间通信 (英文版·第2版)

-
- ◆ 著 [美] W. Richard Stevens
责任编辑 杨海玲
 - ◆ 人民邮电出版社出版发行 北京市崇文区夕照寺街14号
邮编 100061 电子函件 315@ptpress.com.cn
网址 <http://www.ptpress.com.cn>
北京隆昌伟业印刷有限公司印刷
 - ◆ 开本: 800×1000 1/16
印张: 36
字数: 692千字 2009年11月第1版
印数: 1-2500册 2009年11月北京第1次印刷
- 著作权合同登记号 图字: 01-2009-5714号
ISBN 978-7-115-21511-6
-

定价: 89.00元

读者服务热线: (010) 51095186 印装质量热线: (010) 67129223

反盗版热线: (010) 67171154

版 权 声 明

Original edition, entitled *UNIX Network Programming, Volume 2: Interprocess Communications, Second Edition*, 9780130810816 by W. Richard Stevens, published by Pearson Education, Inc., publishing as Prentice Hall PTR, Copyright © 1999 by Prentice Hall PTR.

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.

China edition published by PEARSON EDUCATION ASIA LTD. and POSTS & TELECOM PRESS Copyright © 2009.

This edition is manufactured in the People's Republic of China, and is authorized for sale only in the People's Republic of China excluding Hong Kong, Macao and Taiwan.

本书英文版由Pearson Education Asia Ltd.授权人民邮电出版社独家出版。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。

仅限于中华人民共和国境内（香港、澳门特别行政区和台湾地区除外）销售发行。

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

版权所有，侵权必究。

Preface

前言

概述

大多数重要的程序都涉及进程间通信 (Interprocess Communication, IPC)。这是受下述设计原则影响的自然结果：把应用程序设计为一组相互通信的小片断比将其设计为单个庞大的程序更好。从历史角度看，应用程序有如下几种构建方法。

(1) 用一个庞大的程序完成全部工作。程序的各部分可以实现为函数，函数之间通过参数、返回值和全局变量来交换信息。

(2) 使用多个程序，程序之间用某种形式的IPC进行通信。许多标准的Unix工具都是按这种风格设计的，它们使用shell管道 (IPC的一种形式) 在程序之间传递信息。

(3) 使用一个包含多个线程的程序，线程之间使用某种IPC。这里仍然使用术语IPC，尽管通信是在线程之间而不是在进程之间进行的。

还可以把后两种设计形式结合起来：用多个进程来实现，其中每个进程包含几个线程。在这种情况下，进程内部的线程之间可以通信，不同的进程之间也可以通信。

上面讲述了可以把完成给定任务所需的工作分到多个进程中，或许还可以进一步分到进程内的多个线程中。在包含多个处理器 (CPU) 的系统中，多个进程也许可以 (在不同的CPU上) 同时运行，或许给定进程内的多个线程也能同时运行。因此，把任务分到多个进程或线程中有望减少完成指定任务的时间。

本书详细描述了以下4种不同的IPC形式：

- (1) 消息传递 (管道、FIFO和消息队列)；
- (2) 同步 (互斥量、条件变量、读写锁、文件和记录锁、信号量)；
- (3) 共享内存 (匿名的和具名的)；
- (4) 远程过程调用 (Solaris门和Sun RPC)。

本书不讨论如何编写通过计算机网络通信的程序。这种通信通常涉及使用TCP/IP协议族的套接字API，相关主题在第1卷[Stevens 1998]中有详细讨论。

有人可能会提出质疑：不应该使用单主机或非网络IPC (本卷的主题)，所有程序都应该在网络上的多台主机上同时运行。但在日常实践中，单主机IPC往往比网络通信快得多，而且有时

还简单些。共享内存、同步等方法通常也只能用于单主机，跨网络时可能无法使用。经验和历史表明，非网络IPC（本卷）与跨网络IPC（第1卷）都是需要的。

本卷建立在第1卷和我写的另外4本书的基础上，这5本书在本书中简记如下：

- UNPv1: *UNIX Network Programming, Volume 1* [Stevens 1998]；
- APUE: *Advanced Programming in the UNIX Environment* [Stevens 1992]；
- TCPv1: *TCP/IP Illustrated, Volume 1* [Stevens 1994]；
- TCPv2: *TCP/IP Illustrated, Volume 2* [Wright and Stevens 1995]；
- TCPv3: *TCP/IP Illustrated, Volume 3* [Stevens 1996]。

在一本以“网络编程”为书名一部分的书中讨论IPC看似有点奇怪，但事实上IPC经常用于网络应用程序。我在《UNIX网络编程》1990年版的前言里就指出：“想知道如何为网络开发软件，必须先理解进程间通信（IPC）。”

与第1版的区别

本书完全重写并扩充了1990年版《UNIX网络编程》的第3章和第18章。字数统计表明，现在的内容是第1版的5倍。新版的主要改动归纳如下。

- 不仅讨论了“System V IPC”的三种形式（消息队列、信号量以及共享内存），还对实现了这些IPC的新的Posix函数进行了介绍。（1.7节将详细介绍Posix标准族。）我认为使用Posix IPC函数是大势所趋，因为它们比System V中的相应部分更具优势。
- 讨论了用于同步的Posix函数：互斥锁、条件变量以及读写锁。它们可用于线程或进程的同步，而且往往在访问共享内存时使用。
- 本卷假定使用Posix线程环境（称为“Pthreads”），许多示例都是用多线程而不是多进程构建的。
- 对管道、FIFO和记录锁的讨论侧重于从它们的Posix定义出发。
- 本卷不仅描述了IPC机制及其使用方法，还实现了Posix消息队列、读写锁与Posix信号量（都可以实现为用户库）。这些实现可以把多种不同的特性捆绑起来（例如，Posix信号量的一种实现用到了互斥量、条件变量和内存映射I/O），还强调了我们在应用程序中经常要处理的一些问题（如竞争状态、错误处理、内存泄漏和变长参数列表）。理解某种特性的实现通常有助于了解如何使用该特性。
- 对RPC的讨论侧重于Sun的RPC包。在此之前讲述了新的Solaris门API，它类似于RPC但用于单主机。这么一来我们就介绍了许多在调用其他进程中的过程时需要考虑的特性，而不关心网络方面的细节。

读者对象

本书既可以用作IPC的教程，也可以用作有经验的程序员的参考书。全书划分为以下4个主

要部分：

- 消息传递；
- 同步；
- 共享内存；
- 远程过程调用。

但许多读者可能只对特定的部分感兴趣。第2章总结了所有Posix IPC函数共有的特性，第3章归纳了所有System V IPC函数共有的特性，第12章介绍了Posix和System V的共享内存，但书中多数章节都可以独立于其他章节阅读。所有读者都应该阅读第1章，尤其是1.6节，该节介绍了一些贯穿全书的包装函数。讨论Posix IPC的各章与讨论System V IPC的各章彼此独立，有关管道、FIFO和记录锁的几章不属于上述两个阵营，关于RPC的两章也独立于其他IPC方法。

为了方便读者把本书作为参考书，本书提供了完整的全文索引，并在最后几页总结了每个函数和结构的详细描述在正文中的哪里可以找到。为了给不按顺序阅读本书的读者提供方便，我们在书中为各个主题提供了大量的交叉引用。

源代码与勘误

书中所有示例的源代码可以从作者主页（列在前言的最后）获得^①。学习本书讲述的IPC技术的最好方法就是下载这些程序，对其进行修改和改进。只有这样实际编写代码才能深入理解有关概念和方法。每章末尾提供了大量的习题，大部分在附录D中给出答案。

本书的最新勘误表也可以从作者主页获取。

致谢

尽管封面上只出现了作者一个人的名字，但一本高质量的书的出版需要许多人的共同努力。首先要感谢我的家人，他们在我写书的那段时间里承担了一切。再次感谢你们：Sally、Bill、Ellen和David。

感谢技术审稿人给出的宝贵的反馈意见（打印出来有135页）。他们发现了许多错误，指出了需要更多解释的地方，并对表达、用词和代码提出了许多修改建议，他们是Gavin Bowe、Allen Briggs、Dave Butenhof、Wan-Teh Chang、Chris Cleeland、Bob Friesenhahn、Andrew Gierth、Scott Johnson、Marty Leisner、Larry McVoy、Craig Metz、Bob Nelson、Steve Rago、Jim Reid、Swamy K. Sitarama、Jon C. Snader、Ian Lance Taylor、Rich Teer和Andy Tucker。

下列诸位通过电子邮件回答过我的问题，有人甚至回答过很多问题。澄清这些问题提高了本书的准确性并改进了语言表达，他们是David Bausum、Dave Butenhof、Bill Gallmeister、

^①书中所有示例的源代码也可以从图灵网站（www.turingbook.com）本书网页免费注册下载。——编者注

Mukesh Kacker、Brian Kernighan、Larry McVoy、Steve Rago、Keith Skowran、Bart Smaalders、Andy Tucker和John Wait。

特别感谢GSquared的Larry Rafsky提供了很多帮助。像以往一样，感谢国家光学天文台（NOAO）、Sidney Wolff、Richard Wolff和Steve Grandi，他们为我提供了网络与主机的访问权限。DEC公司的Jim Bound、Matt Thomas、Mary Clouter和Barb Glover提供了用于本书多数示例的Alpha系统。书中的一部分代码是在其他Unix系统上测试的：感谢Red Hat软件公司的Michael Johnson提供了最新版本的Red Hat Linux，感谢IBM奥斯汀实验室的Dave Marquardt和Jessie Haug提供了RS/6000系统以及最新版本的AIX的访问权限。

最后还要感谢Prentice Hall的优秀员工（本书的编辑Mary Franz，还有Noreen Regina、Sophie Papanikolaou和Patti Guerrieri）给予的帮助，尤其是在很紧的时间内完成一切所付出的努力。

版权说明

我制作了本书的最终电子版（PostScript格式），最后排版成现在的书。我用James Clark编写的优秀的groff包为本书排版，该软件包安装在一台运行Solaris 2.6的SparcStation工作站上。（认为troff已经过时的报导显然太夸张了。）我使用vi编辑器键入了所有的138 897个单词，用gpic程序绘制了72幅插图（其中用到了许多由Gary Wright编写的宏），用gtbl程序生成了35张表格，为全书添加了索引（用到了Jon Bentley与Brian Kernighan编写的一组awk脚本），并设计了最终的版式。我录入书中的8 046行C语言源代码，使用的是Dave Hanson的loom程序、GNU的indent程序和Gary Wright写的一些脚本。

欢迎读者以电子邮件的方式反馈意见、提出建议或订正错误。

W. Richard Stevens
1998年7月于亚利桑那州图森市
<http://www.kohala.com/~rstevens>

Contents

目 录

Part 1 Introduction	1	2.3 Creating and Opening IPC Channels	22
简介		创建与打开IPC通道	
Chapter 1 Introduction	3	2.4 IPC Permissions	25
简介		IPC权限	
1.1 Introduction	3	2.5 Summary	26
概述		小结	
1.2 Processes, Threads, and the Sharing of Information	5	Chapter 3 System V IPC	27
进程、线程与信息共享		3.1 Introduction	27
1.3 Persistence of IPC Objects	6	概述	
IPC对象的持续性		3.2 key_t Keys and ftok Function	28
1.4 Name Spaces	7	key_t键和ftok函数	
名字空间		3.3 ipc_perm Structure	30
1.5 Effect of fork, exec, and exit on IPC Objects	9	ipc_perm结构	
fork、exec和exit对IPC对象的影响		3.4 Creating and Opening IPC Channels	30
1.6 Error Handling: Wrapper Functions	11	创建与打开IPC通道	
错误处理：包装函数		3.5 IPC Permissions	32
1.7 Unix Standards	13	IPC权限	
Unix标准		3.6 Identifier Reuse	34
1.8 Road Map to IPC Examples in the Text	15	标识符重用	
本书中IPC示例的路线图		3.7 ipcs and ipcrm Programs	36
1.9 Summary	16	ipcs和ipcrm程序	
小结		3.8 Kernel Limits	36
Chapter 2 Posix IPC	19	内核限制	
2.1 Introduction	19	3.9 Summary	38
概述		小结	
2.2 IPC Names	19	Part 2 Message Passing	41
IPC名字		消息传递	
Chapter 2 Posix IPC	19	Chapter 4 Pipes and FIFOs	43
2.1 Introduction	19	管道和FIFO	
概述			
2.2 IPC Names	19		
IPC名字			

4.1 Introduction	43	5.7 Posix Realtime Signals	98
概述		Posix实时信号	
4.2 A Simple Client-Server Example	43	5.8 Implementation Using Memory-Mapped I/O	106
一个简单的客户-服务器示例		使用内存映射I/O实现	
4.3 Pipes	44	5.9 Summary	126
管道		小结	
4.4 Full-Duplex Pipes	50		
全双工管道		Chapter 6 System V Message Queues	129
4.5 popen and pclose Functions	52	System V消息队列	
popen和pclose函数		6.1 Introduction	129
4.6 FIFOs	54	概述	
4.7 Additional Properties of Pipes and FIFOs	58	6.2 msgget Function	130
管道和FIFO的额外属性		msgget函数	
4.8 One Server, Multiple Clients	60	6.3 msgsnd Function	131
单服务器, 多客户		msgsnd函数	
4.9 Iterative versus Concurrent Servers	66	6.4 msgrcv Function	132
迭代服务器与并发服务器的比较		msgrcv函数	
4.10 Streams and Messages	67	6.5 msgctl Function	134
流与消息		msgctl函数	
4.11 Pipe and FIFO Limits	72	6.6 Simple Programs	135
管道和FIFO限制		简单的程序	
4.12 Summary	73	6.7 Client-Server Example	140
小结		客户-服务器示例	
Chapter 5 Posix Message Queues	75	6.8 Multiplexing Messages	142
Posix消息队列		多路复用消息	
5.1 Introduction	75	6.9 Message Queues with select and poll	151
概述		消息队列上使用select和poll	
5.2 mq_open, mq_close, and mq_unlink Functions	76	6.10 Message Queue Limits	152
mq_open、mq_close和mq_unlink函数		消息队列限制	
5.3 mq_getattr and mq_setattr Functions	79	6.11 Summary	155
mq_getattr和mq_setattr函数		小结	
5.4 mq_send and mq_receive Functions	82	Part 3 Synchronization	157
mq_send和mq_receive函数		同步	
5.5 Message Queue Limits	86		
消息队列限制		Chapter 7 Mutexes and Condition Variables	159
5.6 mq_notify Function	87	互斥锁和条件变量	
mq_notify函数			

7.1 Introduction	159	9.1 Introduction	193
概述		概述	
7.2 Mutexes: Locking and Unlocking	159	9.2 Record Locking versus File Locking	197
互斥锁: 加锁与解锁		记录加锁与文件加锁	
7.3 Producer-Consumer Problem	161	9.3 Posix fcntl Record Locking	199
生产者-消费者问题		Posix fcntl记录加锁	
7.4 Locking versus Waiting	165	9.4 Advisory Locking	203
加锁与等待		劝告性加锁	
7.5 Condition Variables: Waiting and Signaling	167	9.5 Mandatory Locking	204
条件变量: 等待与信号发送		强制性加锁	
7.6 Condition Variables: Timed Waits and Broadcasts	171	9.6 Priorities of Readers and Writers	207
条件变量: 定时等待和广播		读出者和写入者的优先级	
7.7 Mutexes and Condition Variable Attributes	172	9.7 Starting Only One Copy of a Daemon	213
互斥锁和条件变量的属性		只启动守护进程的一个副本	
7.8 Summary	174	9.8 Lock Files	214
小结		锁文件	
Chapter 8 Read-Write Locks	177	9.9 NFS Locking	216
读写锁		NFS加锁	
8.1 Introduction	177	9.10 Summary	216
概述		小结	
8.2 Obtaining and Releasing Read-Write Locks	178	Chapter 10 Posix Semaphores	219
获取与释放读写锁		Posix信号量	
8.3 Read-Write Lock Attributes	179	10.1 Introduction	219
读写锁属性		概述	
8.4 Implementation Using Mutexes and Condition Variables	179	10.2 sem_open, sem_close, and sem_unlink Functions	225
使用互斥锁和条件变量实现		sem_open、sem_close和sem_unlink函数	
8.5 Thread Cancellation	187	10.3 sem_wait and sem_trywait Functions	226
线程取消		sem_wait和sem_trywait函数	
8.6 Summary	192	10.4 sem_post and sem_getvalue Functions	227
小结		sem_post和sem_getvalue函数	
Chapter 9 Record Locking	193	10.5 Simple Programs	228
记录加锁		简单的程序	
		10.6 Producer-Consumer Problem	233
		生产者-消费者问题	

10.7	File Locking	238	11.6	File Locking	294
	文件加锁			文件加锁	
10.8	sem_init and sem_destroy Functions	238	11.7	Semaphore Limits	296
	sem_init和sem_destroy函数			信号量限制	
10.9	Multiple Producers, One Consumer	242	11.8	Summary	300
	多生产者, 单消费者			小结	
10.10	Multiple Producers, Multiple Consumers	245	Part 4 Shared Memory	301	
	多生产者, 多消费者			共享内存	
10.11	Multiple Buffers	249	Chapter 12 Shared Memory		
	多缓冲区		Introduction	303	
10.12	Sharing Semaphores between Processes	256		共享内存简介	
	进程间共享信号量		12.1	Introduction	303
10.13	Semaphore Limits	257		概述	
	信号量限制		12.2	mmap, munmap, and msync Functions	307
10.14	Implementation Using FIFOs	257		mmap、munmap和msync函数	
	使用FIFO实现		12.3	Increment Counter in a Memory-Mapped File	311
10.15	Implementation Using Memory-Mapped I/O	262		内存映射文件中的计数器递加	
	使用内存映射I/O实现		12.4	4.BSD Anonymous Memory Mapping	315
10.16	Implementation Using System V Semaphores	271		4.BSD匿名内存映射	
	使用System V信号量实现		12.5	SVR4 /dev/zero Memory Mapping	316
10.17	Summary	278		SVR4 /dev/zero内存映射	
	小结		12.6	Referencing Memory-Mapped Objects	317
Chapter 11 System V Semaphores	281			引用内存映射的对象	
	System V信号量		12.7	Summary	322
11.1	Introduction	281		小结	
	概述		Chapter 13 Posix Shared Memory	325	
11.2	semget Function	282		Posix共享内存	
	semget函数		13.1	Introduction	325
11.3	semop Function	285		概述	
	semop函数		13.2	shm_open and shm_unlink Functions	326
11.4	semctl Function	287		shm_open和shm_unlink函数	
	semctl函数		13.3	ftruncate and fstat Functions	327
11.5	Simple Programs	289		ftruncate和fstat函数	
	简单的程序		13.4	Simple Programs	328
				简单的程序	

13.5	Incrementing a Shared Counter333 共享计数器递加	15.6	door_info Function365 door_info函数
13.6	Sending Messages to a Server336 向服务器发送消息	15.7	Examples366 示例
13.7	Summary342 小结	15.8	Descriptor Passing379 描述符传递
Chapter 14 System V Shared Memory ...343 System V共享内存		15.9	door_server_create Function384 door_server_create函数
14.1	Introduction343 概述	15.10	door_bind, door_unbind, and door_revoke Functions390 door_bind、door_unbind和 door_revoke函数
14.2	shmget Function343 shmget函数	15.11	Premature Termination of Client or Server390 客户或服务器的过早终止
14.3	shmat Function344 shmat函数	15.12	Summary397 小结
14.4	shmdt Function345 shmdt函数	Chapter 16 Sun RPC399	
14.5	shmctl Function345 shmctl函数	16.1	Introduction399 概述
14.6	Simple Programs346 简单的程序	16.2	Multithreading407 多线程技术
14.7	Shared Memory Limits349 共享内存限制	16.3	Server Binding411 服务器绑定
14.8	Summary351 小结	16.4	Authentication414 鉴别
Part 5 Remote Procedure Calls353 远程过程调用		16.5	Timeout and Retransmission417 超时和重传
Chapter 15 Doors355 门		16.6	Call Semantics422 调用语义
15.1	Introduction355 概述	16.7	Premature Termination of Client or Server424 客户或服务器的过早终止
15.2	door_call Function361 door_call函数	16.8	XDR: External Data Representation426 XDR: 外部数据表示
15.3	door_create Function363 door_create函数	16.9	RPC Packet Formats444 RPC分组格式
15.4	door_return Function364 door_return函数	16.10	Summary449 小结
15.5	door_cred Function365 door_cred函数		

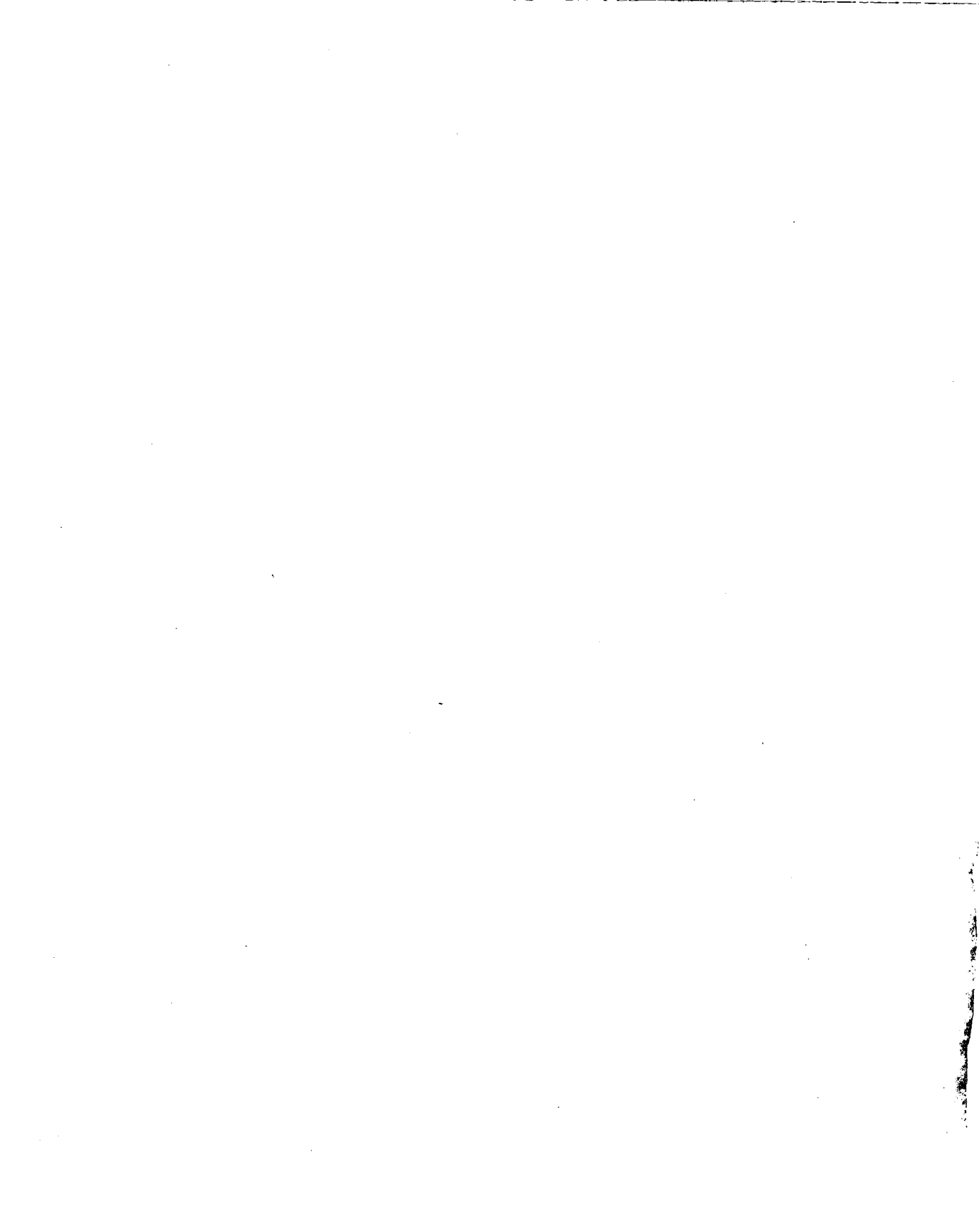
Epilogue	453	B.2 Basic Thread Functions: Creation and Termination	502
后记		基本线程函数：创建和终止	
Appendix A Performance Measurements	457	Appendix C Miscellaneous Source Code	505
性能测量		其他源代码	
A.1 Introduction	457	C.1 unipic.h Header	505
概述		unipic.h头文件	
A.2 Results	458	C.2 config.h Header	509
结果		config.h头文件	
A.3 Message Passing Bandwidth Programs	467	C.3 Standard Error Functions	510
消息传递带宽程序		标准错误处理函数	
A.4 Message Passing Latency Programs	480	Appendix D Solutions to Selected Exercises	515
消息传递延迟程序		精选习题答案	
A.5 Thread Synchronization Programs	486	Bibliography	535
线程同步程序		参考文献	
A.6 Process Synchronization Programs	497	Index	539
进程同步程序		索引	
Appendix B A Threads Primer	501		
线程入门			
B.1 Introduction	501		
概述			

Part 1

Introduction

简介





Introduction

简介

1.1 Introduction 概述

IPC stands for *interprocess communication*. Traditionally the term describes different ways of *message passing* between different processes that are running on some operating system. This text also describes numerous forms of *synchronization*, because newer forms of communication, such as shared memory, require some form of synchronization to operate.

In the evolution of the Unix operating system over the past 30 years, message passing has evolved through the following stages:

- *Pipes* (Chapter 4) were the first widely used form of IPC, available both within programs and from the shell. The problem with pipes is that they are usable only between processes that have a common ancestor (i.e., a parent-child relationship), but this was fixed with the introduction of *named pipes* or *FIFOs* (Chapter 4).
- *System V message queues* (Chapter 6) were added to System V kernels in the early 1980s. These can be used between related or unrelated processes on a given host. Although these are still referred to with the "System V" prefix, most versions of Unix today support them, regardless of whether their heritage is System V or not.

When describing Unix processes, the term *related* means the processes have some ancestor in common. This is another way of saying that these related processes were generated