



机密

Linux/UNIX

系统编程手册 (上册)

THE **LINUX**
PROGRAMMING
INTERFACE

A Linux and UNIX® System Programming Handbook

[德] Michael Kerrisk 著
孙剑 许从年 董健 孙余强 译



 人民邮电出版社
POSTS & TELECOM PRESS

异步社区会员 flyman150(2410757683@qq.com) 专享 尊重版权

Linux/UNIX 系统编程手册 (上册)

THE LINUX
PROGRAMMING
INTERFACE

A Linux and UNIX[®] System Programming Handbook

[德] Michael Kerrisk 著
孙剑 许从年 董健 孙余强 译

人民邮电出版社

北京

异步社区会员 flyman150(2410757683@qq.com) 专享 尊重版权

异步社区电子书

感谢您购买异步社区电子书！异步社区已上架电子书 500 余种，社区还会经常发布福利信息，对社区有贡献的读者赠送免费样书券、优惠码、积分等等，希望您在阅读过程中，把您的阅读体验传递给我们，让我们了解读者心声，有问题我们会及时修正。

社区网址：<http://www.epubit.com.cn/>

反馈邮箱：contact@epubit.com.cn

异步社区里有什么？

图书、电子书（半价电子书）、优秀作译者、访谈、技术会议播报、赠书活动、下载资源。

异步社区特色：

纸书、电子书同步上架、纸电捆绑超值优惠购买。

最新精品技术图书全网首发预售。

晒单有意外惊喜！

异步社区里可以做什么？

博客式写作发表文章，提交勘误赚取积分，积分兑换样书，写书评赢样书券等。

联系我们：

微博：

@ 人邮异步社区

@ 人民邮电出版社 - 信息技术分社

微信公众号：

人邮 IT 书坊

异步社区

QQ 群：368449889

图书在版编目 (C I P) 数据

Linux/UNIX系统编程手册 : 全2册 / (美) 凯利斯克 (Kerrisk, M.) 著 ; 孙剑等译. -- 北京 : 人民邮电出版社, 2014. 1
ISBN 978-7-115-32867-0

I. ①L… II. ①凯… ②孙… III. ①Linux操作系统—程序设计—手册②UNIX操作系统—程序设计—手册 IV. ①TP316.89-62②TP316.81-62

中国版本图书馆CIP数据核字(2013)第187440号

版权声明

Copyright © 2010 by Michael Kerrisk. Title of English-language original: The Linux Programming Interface, ISBN 978-1-59327-220-3, published by No Starch Press. Simplified Chinese-language edition copyright © 2014 by Posts and Telecom Press. All rights reserved.

本书中文简体字版由美国 **No Starch** 出版社授权人民邮电出版社出版。未经出版者书面许可, 对本书任何部分不得以任何方式复制或抄袭。

版权所有, 侵权必究。

-
- ◆ 著 [德] Michael Kerrisk
 - 译 孙 剑 许从年 董 健 孙余强 郭光伟
陈 舸
 - 责任编辑 傅道坤
 - 责任印制 程彦红 杨林杰
 - ◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路 11 号
邮编 100164 电子邮件 315@ptpress.com.cn
网址 <http://www.ptpress.com.cn>
北京艺辉印刷有限公司印刷
 - ◆ 开本: 787×1092 1/16
印张: 76.25
字数: 1 618 千字 2014 年 1 月第 1 版
印数: 1 - 3 500 册 2014 年 1 月北京第 1 次印刷

著作权合同登记号 图字: 01-2010-3829 号

定价: 158.00 元 (上、下册)

读者服务热线: (010)81055410 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京崇工商广字第 0021 号

内容提要

本书是介绍 Linux 与 UNIX 编程接口的权威著作。Linux 编程资深专家 Michael Kerrisk 在书中详细描述了 Linux/UNIX 系统编程所涉及的系统调用和库函数，并辅之以全面而清晰的代码示例。本书涵盖了逾 500 个系统调用及库函数，并给出逾 200 个程序示例，另含 88 张表格和 115 幅示意图。

本书总共分为 64 章，主要讲解了高效读写文件，对信号、时钟和定时器的运用，创建进程、执行程序，编写安全的应用程序，运用 POSIX 线程技术编写多线程程序，创建和使用共享库，运用管道、消息队列、共享内存和信号量技术来进行进程间通信，以及运用套接字 API 编写网络应用等内容。

本书在汇聚大批 Linux 专有特性（epoll、inotify、/proc）的同时，还特意强化了对 UNIX 标准（POSIX、SUS）的论述，彻底达到了“鱼与熊掌，二者得兼”的效果，这也堪称本书的最大亮点。

本书布局合理，论述清晰，说理透彻，尤其是作者对示例代码的构思巧妙，独具匠心，仔细研读定会受益良多。本书适合从事 Linux/UNIX 系统开发、运维工作的技术人员阅读，同时也可作为高校计算机专业学生的参考研习资料。

献 辞

谨将本书献给 Cecilia，你照亮了我的世界！

对本书的赞誉

编写 Linux 软件时如果只能选择一本参考书，则非本书莫属。

——MARTIN LANDERS, Google 公司软件工程师

本书描述精到，示例周详，涵盖了 Linux 底层 API 编程的详尽内容及个中细微之处——无论读者水平如何，都能从本书中受益。

——MEL GORMAN, *Understanding the Linux Virtual Memory Manager* 作者

Michael Kerrisk 的这本 Linux 编程巨著，不但论及 Linux 编程、其与各种标准之间的联系，而且还就作者所知，重点介绍了已获修正的 Linux 内核 bug 以及改进颇多的 Linux 手册页。凭此三点，足可让 Linux 编程更易上手。本书对各项主题的深入探讨使其成为必备的参考书籍——无论读者在 Linux 编程方面造诣如何。

——ANDREAS JAEGER, NOVELL 公司 OPENSUSE 项目经理

Michael 用他坚忍不拔的毅力为 Linux 程序员奉献了这本论述严谨、表述清晰、简洁的权威参考书。虽然本书针对 Linux 程序员而著，但对任何在 UNIX/POSIX 环境中编程的程序员来说都极具价值。

——DAVID BUTENHOF, *Programming with POSIX Threads* 作者、POSIX/UNIX 标准撰写者

本书在重点关注 Linux 系统的同时，对于 UNIX 系统和网络编程也阐述透彻，浅显易懂。无论是初涉 UNIX 编程的新丁，还是编程经验丰富的 UNIX 老手（想要了解大行其道的 GNU/Linux 系统有何新意），我都向他们力荐此书。

——FERNANDO GONT, 网络安全研究员、IETF 参与者、IETF RFC 作者

本书以百科全书般的叙述风格对 Linux 接口编程做了既深且广的覆盖，还提供了大量教科书风格的编程示例和练习。本书所包含的各项主题——从原理到可以实际运行的代码——都已描述清晰且易于理解。本书正是专业人士、学生以及教育工作者所期盼的 Linux/UNIX 参考书。

——ANTHONY ROBINS, 奥塔哥大学计算机科学副教授

无论从精确性、质量，还是详细程度来说，本书都令我印象深刻。身为 Linux 系统调用的行家，Michael Kerrisk 与我们分享了他对 Linux API 的认知和理解。

——CHRISTOPHE BLAESS, *Programmation système en C sous Linux* 作者

对于治学严谨的专业 Linux/UNIX 系统程序员而言，本书实为必备的参考书籍。本书涵盖了所有关键 API 的使用，同时兼顾 Linux 和 UNIX 系统接口，描述清晰，示例具体；除此之外，还强调了遵从诸如 SUS 和 POSIX 1003.1 等标准的重要性和益处。

——ANDREW JOSEY, The Open Group 标准部总监、POSIX 1003.1 工作组主席
由手册页的维护者亲自操刀，以系统程序员视角写出一本百科全书式的 Linux 系统编程巨著，还有比这更完美的吗？本书全面而又详实。我坚信本书将在我的书架上牢牢占据一席之地。

——BILL GALLMEISTER, *POSIX.4 Programmer's Guide: Programming for the Real World* 作者
本书是最新最全的 Linux/UNIX 系统编程参考书。无论读者是 Linux 系统编程新兵，还是关注 Linux 编程和程序移植性的 UNIX 系统编程老将，又或者只是在寻找一本 Linux 编程接口方面的优秀参考书的读者，Michael Kerrisk 的这本大作都笃定是其案头良伴。

——LOÏC DOMAIGNÉ, CORPULS.COM 首席软件架构师（嵌入式）

前 言

主题

本书将描述 Linux 编程接口：由 UNIX 操作系统的开源实现——Linux 所提供的系统调用、库函数以及其他底层接口。运行于 Linux 之上的每一个程序都会直接或间接地使用这些接口。这些接口允许应用程序去执行诸多任务：文件 I/O、创建/删除文件和目录、创建新进程、执行程序、设置定时器、在同一台计算机上发起进程或线程间通信，以及为联网计算机间的进程建立通信等等。有时，人们也将这一系列的底层接口称为系统编程接口。

尽管本书着眼于 Linux，但对于标准和可移植性问题也倍加关注。对于 Linux 所特有的技术细节，以及已由 POSIX 和 SUS 标准化的 UNIX 普遍特性，本书会在论述中清晰地加以区分。因此，本书也提供了对 UNIX/POSIX 编程接口的全面描述。对于那些在其他 UNIX 系统环境中编程，或者编写跨平台可移植应用的程序员来说，本书同样具有实用价值。

本书的读者

本书主要针对以下读者：

- 为 Linux 系统、其他 UNIX 系统，或兼容于 POSIX 的操作系统编写应用程序的程序员和软件设计人员。
- 在 Linux 和其他 UNIX 实现之间，以及 Linux 和其他操作系统之间进行应用程序移植的程序员。
- 教/学 Linux 和 UNIX 系统编程的高校师生。
- 意欲深入理解 Linux/UNIX 编程接口以及系统软件各模块实现细节的系统管理人员和高级用户（power users）。

作者假定读者之前有些许编程经验，但不必是在系统编程领域。此外，作者还假定读者具备阅读 C 语言源码的能力，并了解如何使用 shell 和 UNIX 或 Linux 的常用命令。对于不熟悉 UNIX 和 Linux 的读者来说，阅读第 2 章中面向程序员对 UNIX 和 Linux 系统基本概念所做的回顾会有所帮助。

提示: [Kernighan & Ritchie, 1988]是最具权威性的 C 语言参考书籍。[Harbison& Steele, 2002]一书对 C 语言的介绍则更为详细,并涵盖了由 C99 标准所带来的变化。[van der Linden, 1994]也是一本不错的 C 语言书籍,寓教于乐。[Peek et al., 2001]则对 UNIX 的使用做了简洁而完整的介绍。

贯穿本书,会以这种缩进小字体的文字形式用于旁注,其内容包括基本原理、实现细节、背景信息、史上轶闻以及与正文相关的其他辅助主题。

Linux 和 UNIX

其他 UNIX 实现的大多数特性同样见诸于 Linux,反之亦然。有鉴于此,本书本可只关注于标准 UNIX (即 POSIX) 的系统编程。编写可移植的应用程序固然是值得追求的目标,但描述 Linux 对标准 UNIX 编程接口的扩展也同样重要。Linux 的广受欢迎只是原因之一,而有时出于性能方面的考虑,或是需要访问标准 UNIX 编程接口所不支持的功能时,使用非标准扩展(正因如此,所有 UNIX 实现都提供有非标准扩展)就显得至为重要,此为原因之二。

综上所述,在构思本书时,作者不但力图使其对在各种 UNIX 实现中编程的程序员有所帮助,还全面介绍了 Linux 专有的编程特性,如下所示。

- `epoll`, 获取文件 I/O 事件通知的一种机制。
- `inotify`, 监控文件和目录变化的一种机制。
- `capabilities`, 为进程赋予超级用户的部分权限的一种机制。
- 扩展属性。
- `i-node` 标记。
- `clone()` 系统调用。
- `/proc` 文件系统。
- 在文件 I/O、信号、定时器、线程、共享库、进程间通信以及套接字方面, Linux 所专有的实现细节。

本书的用途和组织结构

本书主要有以下两方面的用途:

- 作为 Linux/UNIX 编程接口的入门教程,读者可循序阅读本书。后续各章内容均构建于之前诸章素材的基础之上,伴之以尽可能简短的前向引用。
- 作为 Linux/UNIX 编程接口的参考大全,读者可以根据书后的详细索引和频现于正文中的交叉引用,随机选择阅读主题。

本书各章可分为以下几个部分。

1. 背景知识及概念: UNIX、C 语言以及 Linux 的历史回顾,以及对 UNIX 标准的概述(第 1 章);以程序员为对象,对 Linux 和 UNIX 的概念进行介绍(第 2 章);Linux 和 UNIX 系统编程的基本概念(第 3 章)。
2. 系统编程接口的基本特性: 文件 I/O (第 4 章、第 5 章), 进程(第 6 章), 内存分配(第 7 章), 用户和组(第 8 章), 进程凭证(`process credential`)(第 9 章), 时间(第 10 章), 系统限制和选项(第 11 章), 以及获取系统和进程信息(第 12 章)。
3. 系统编程接口的高级特性: 文件 I/O 缓冲(第 13 章), 文件系统(第 14 章), 文件

属性（第 15 章），扩展属性（第 16 章），访问控制列表（第 17 章），目录和链接（第 18 章），监控文件事件（第 19 章），信号（signals）（第 20~22 章），以及定时器（第 23 章）。

4. 进程、程序及线程：进程的创建、终止，监控子进程，执行程序（第 24~28 章），以及 POSIX 线程（第 29~33 章）。
5. 进程及程序的高级主题：进程组、会话以及任务控制（第 34 章），进程优先级和进程调度（第 35 章），进程资源（第 36 章），守护进程（第 37 章），编写安全的特权程序（第 38 章），能力（capability）（第 39 章），登录记账（第 40 章），以及共享库（第 41 章和第 42 章）。
6. 进程间通信（IPC）：IPC 概览（第 43 章），管道和 FIFO（第 44 章），系统 V IPC 消息队列、信号量（semaphore）及共享内存（第 45~48 章），内存映射（第 49 章），虚拟内存操作（第 50 章），POSIX 消息队列、信号量及共享内存（第 51~54 章），以及文件锁定（第 55 章）。
7. 套接字和网络编程：使用套接字的 IPC 和网络编程（第 56~61 章）。
8. 高级 I/O 主题：终端（第 62 章），其他 I/O 模型（第 63 章），以及伪终端（第 64 章）。

程序示例

本书会以简短而完整的程序示例来描述大部分编程接口，以期读者通过命令行方便地体验这些程序，从而了解各种不同系统调用和库函数的运作方式。因此，本书包含了大量代码示例——约 15 000 行 C 语言代码和 shell 会话记录。

虽然阅读并执行上述示例代码是学习 Linux 编程接口的一个不错的起点，但巩固本书所述概念最为有效的方式还是动手编写代码——无论是修改示例程序以验证自己的编程思路，还是编写全新的程序。

书中所有源代码均可从本书网站上下载。网站所发布的源码中还包含了不少未见诸于本书的其他程序。这些程序的用途和详细信息在源码注释中均有描述。源码中还提供了 Makefile，用来构建相应的程序，附带的 README 文件则提供了相应程序的具体细节。

在 GNU Affero 通用公共许可证（Affero GPL）版本 3 条款的约束下，可自行重新发布和修改本书源码，源码中也提供了一份 GNU Affero GPL 版本 3 的文件副本。

习题

本书各章大都会在结尾处附有一组习题，其中一部分会利用书中的程序示例进行各种试验，另一些则与本章所讨论的概念相关，而其他习题则是引导读者亲自动手编程，意在巩固读者对所学内容的理解。附录 F 会有选择地给出部分习题的答案。

标准和可移植性

本书自始至终都对可移植性问题予以了特殊关注。对相关标准的引用在书中会反复出现，尤其是 POSIX.1-2001 和 SUSv3 的联合标准¹。此外，本书还包括了该标准最新版本（POSIX.1-2008 与 SUSv4 联合标准）的变更细节。（由于 SUSv3 较之于之前的版本做了大范围的修订，并且在写作本书之际，SUSv3 依然是影响最为广泛的 UNIX 标准，故而本书对标准的讨论一般都以

¹ 译者注：大致可视为一套标准的两种称谓。

SUSv3 为框架，并会注明其与 SUSv4 之间的差别。然而，对读者来说，除非另有说明，与 SUSv3 规范有关的论述同样适用于 SUSv4。）

对于那些尚未标准化的特性，本书会列出与其他 UNIX 实现间的差异范围。此外，本书也会强调那些独具 Linux 实现特色的主要特性，以及 Linux 和其他 UNIX 实现之间在系统调用和库函数方面的细微差别。任何特性，凡未注明为 Linux 所专有，读者通常可将其视为大部分或所有 UNIX 实现的标准特性。

书中的编程示例（除了注明为 Linux 所专有的特性）大多已在 Solaris、FreeBSD、Mac OS X、Tru64 UNIX 以及 HP-UX 中的所有或部分系统上进行了测试。为了改进针对其中某些系统的可移植性，本书的 Web 站点还为特定编程示例提供了其他版本，此类代码就不再于本书中列出。

Linux 内核和 C 语言函数库版本

本书主要着眼于 Linux 2.6.x，撰写本书之际，这一内核版本的使用也最为广泛。本书同样涵盖了 Linux 2.4 内核的详细信息，也会适时指明 Linux 2.4 和 2.6 内核间的特性差异。凡是见诸于 Linux 2.6.x 系列中的新特性，作者均会标出其（首度）出现的确切内核版本号（例如，2.6.34）。

至于 C 语言函数库，本书会重点关注 GNU C 语言库（glibc）版本 2。本书也会适时指出 glibc 2.x 版本之间所存在的差异。

本书付梓之际，Linux 内核版本 2.6.35 刚刚问世，不久又发布了 glibc 版本 2.12。本书目前的论述涵盖了以上两种软件的这两个版本。本书出版后，Linux 和 glibc 接口所发生的变化会在本书的 Web 站点上公布。

在其他语言中调用编程接口

虽然本书的程序示例都是以 C 语言编写而成的，但读者也能使用其他编程语言来调用本书所描述的编程接口。这些语言既包括编译语言，例如：C++、Pascal、Modula、Ada、FORTRAN 和 D 语言，也包括脚本语言，例如：Perl、Python 和 Ruby（如要使用 Java，则需另辟蹊径，可参阅[Rochkind, 2004]）。这需要运用不同的技术以获取必要的常量定义和函数申明（C++除外），而按 C 语言链接惯例所约定的方式来传递函数参数可能也需要额外的工作投入。虽然实现起来有差异，但基本原理却都相同，即便读者在使用其他语言进行编程，本书所含信息对他们也同样适用。

关于作者

本人于 1987 年开始使用 UNIX 和 C 语言。当时，作者连续几个礼拜都泡在一台 HP Bobcat 工作站旁，陪伴我的只有 Marc Rochkind 所著 *Advanced UNIX Programming*（第 1 版）一书，以及一本最终被翻得卷了边的 C shell 手册的印刷本。投入时间阅读文档（如果有的话），并编写一些小型的（规模可逐渐变大）测试程序进行试验，直至自己对软件的理解感到信心满满——这是作者当时所采用的编程学习方法，并一直沿用至今——作者也向任何试水新型软件技术的人们推荐这一做法。依作者拙见，从长远来看，这种自学方法能够大大节约时间。本书所载的许多编程示例正是在这一学习方法的激励之下设计而成。

作者的主要身份是软件工程师和设计师。然而，作者同样好为人师，并在学术或商业领域有过数年的教学经验。作者还开设过多门为期一周的 UNIX 系统编程课程，这一经验对本

书的写作也颇有影响。

作者使用 Linux 的时间大约只有与 UNIX 打交道的时一半长，在这段时间里，作者的兴也逐渐集中在了内核和用户空间的“分水岭”——Linux 编程接口上。这一兴趣也使作者投身于一系列紧密相联的活动中。作者会时不时地对 POSIX/SUS 标准提出自己的意见并提供 BUG 报告；对新加入 Linux 内核中的用户空间接口进行测试和设计评审（还能帮助发现并修复那些接口中的诸多代码和设计缺陷）；作者还经常在关于编程接口及其文档的主题会议上发言，并受邀多次出席 Linux 内核开发者年度峰会。将上述所有活动串接在一起的主线是作者对 Linux 领域最突出的贡献：作者为 Linux man-pages 项目 ([http:// www.kernel.org/doc/man-pages/](http://www.kernel.org/doc/man-pages/)) 所做的工作。

Linux 手册页中的第 2、3、4、5 以及 7 部分都属于 man-pages 项目。这几部分也是手册页中描述编程接口的内容，这些编程接口由 Linux 内核及 GNU C 语言库提供，本书所要介绍的正是这方面的内容。作者参与 man-pages 项目已逾十载。自 2004 年起，作者成为了该项目的维护人，所承担的任务大致包括：撰写文档、阅读内核和 C 语言库源码，以及通过编程来验证文档细节（通过为接口撰写文档来发现相关接口中的 BUG，效果颇为不俗）。此外，作者对 man-pages 项目的贡献也最多——在约 900 页的手册页中，作者独自编写了其中的 140 页，并与他人合著了另外的 125 页。因此，在购买本书之前，读者想必已经阅读过本人的工作成果了。作者希望这些成果能对读者有所帮助，希望本书更是如此。

致谢

没有一千人等的支持，本书的质量绝不会如此之高。我要向他们致以最诚挚的谢意。

来自世界各地的多位技术审稿人都参与了本书初稿的阅读，找出错误，指出含糊不清的解释，对措辞、插图以及习题提出建议，测试程序，发现不为作者所知的 Linux 和其他 UNIX 实现间的行为差异，并不时为作者打气助威。在本书中，作者将许多审稿人无私奉献的真知灼见一并收纳，实则作者的知识并非如此渊博。当然，书中的任何错误都是作者一人之过。

无论以下技术审稿人（按姓氏字母排序）对本书手稿的审校巨细与否，篇幅多少，作者都要向他们致以由衷的谢意。

- **Christophe Blaesss** 是一名软件咨询工程师和培训专家，专长是 Linux 在工业（实时和嵌入）方面的应用。Christophe 是 *Programmation système en C sous Linux* 一书的作者，这本法文杰作涵盖了与本书相同的多项主题。他不吝审阅了本书的众多章节。
- **David Butenhof** (HP 公司) 是原 POSIX 线程工作组和 SUS 线程扩展工作组的成员之一，也是 *Programming with POSIX Threads* 一书的作者。他曾为开放软件基金会编写了最初的 DCE 线程参考实现，并曾担任 OpenVMS 和 Digital UNIX 线程实现的首席架构师。David 审校了本书与线程相关的章节，提出了诸多改进意见，还耐心地纠正了几处作者对 POSIX 线程 API 的理解错误。
- **Geoff Clare** 目前在为 The Open Group 开发 UNIX 一致性测试包，从事 UNIX 标准化工作已逾 20 年，是 Austin Group 的 6 位关键参与者之一，该组的宗旨是开发出构成 POSIX.1 和 Single UNIX Specification 基础卷的共同标准。Geoff 仔细审校了手稿中与 UNIX 编程接口相关的内容，耐心细致地提出了众多修正和改进意见，发现了诸多潜藏于手稿中的错误，在突出标准对可移植编程的重要性方面助益良多。
- **Loïc Dornaigné** (当时供职于德国空中交通管制中心[German Air Traffic Control]) 是一名系统软件工程师，主要从事分布式、多并发、容错型嵌入式系统的设计和开发，此

类系统对实时性有着严苛的要求。他针对 SUSv3 中与线程规范有关的内容发表过评论和建议，在多个网上技术论坛中古道热肠，诲人不倦，无私地分享自己的编程心得。Loïc 细致地审校了本书与线程相关各章节，以及多处其他内容。除了编写若干精巧的程序来验证 Linux 线程实现的细节以外，他还倾注了巨大的热情并鼓励作者，提出了许多建议以改进本书整体的表现形式。

- **Gert Döring** mgetty 和 sendfax 程序的开发者，这一“双子星座”也是 Linux/UNIX 系统上使用最为广泛的开源传真软件包。最近，他主要忙于搭建并维护基于 IPv4 和 IPv6 的大型网络，其肩负的主要任务是：与全欧洲的同事一起定义有效的网络策略，以确保 Internet 基础设施顺畅运行。Gert 审校了本书与终端、登录记账、进程组、会话以及任务控制相关各章节，并反馈了大量的有用信息。
- **Wolfram Gloger** 是一名 IT 顾问，过去 15 年，他参与过许多自由和开源软件项目 (Free and Open Source Software, FOSS)。除此之外，Wolfram 还是 GNU C 语言库中 malloc 软件包的实现者。目前，他主要从事于 Web 服务的开发，尤其专注于网上教学，当然，在内核和系统库方面，他仍会偶露峥嵘。Wolfram 审校了本书诸多章节，尤其是侧重于内存方面的内容。
- **Fernando Gont** 是阿根廷国家科技大学 (Universidad Tecnológica Nacional, Argentina) 电子信息中心 (Centro de Estudios de Informática, CEDI) 成员。Internet 工程技术 (Internet engineering) 是其兴趣所在，在 Internet 工程任务组 (IETF) 也能见到其活跃的身影，他还是多个 RFC 文档的作者。此外，Fernando 还为英国 CPNI (国家基础设施保护机构) 中心效力，以提供对通信协议安全方面的评估，而首个完整的 TCP 和 IP 协议安全评估报告也正是由他提出的。Fernando 仔细审校了本书涉及网络编程的相关章节，不厌其烦地向作者解释了 TCP/IP 协议的诸多细节，并对相关内容提出了不少改进意见。
- **Andreas Grünbacher** (SUSE 实验室) 是位内核高手，还是 Linux 扩展属性和 POSIX 访问控制列表的实现者。Andreas 除了仔细审校了本书多章内容以外，还对作者勉励有加，有时，他的只言片语便极有可能改变这部书的整体结构。
- **Christoph Hellwig** 是 Linux 存储和文件系统咨询师，也是内核方面公认的行家，参与过 Linux 内核中多个部分的开发工作。在忙于编写及审查 Linux 内核补丁代码之余，Christoph 抽空审校了本书若干章节，并提出了诸多有益的改进和修正意见。
- **Andreas Jaeger** 曾领导过 Linux 向 x86-64 架构的移植开发工作。身为 GNU C 语言库的开发者，他不但将该库移植到了 x86-64 平台，而且还促成该库符合多个领域的标准，尤其是在数学库方面。他当前效力于 Novell 公司，是 openSUSE 的程序经理。Andreas 审校的章节之多，超乎作者预期。在本书的写作过程中，他除了提出诸多改进意见之外，也给予作者热情的鼓励。
- **Rick Jones** 绰号“Mr. Netperf” (HP 公司联网系统的性能偏执狂)，对本书的网络编程相关章节提出了宝贵意见。
- **Andi Kleen** (当时效力于 SUSE 实验室) 长期以来，一直是内核方面公认的行家里手，对 Linux 内核诸多不同领域贡献颇多，包括：网络、错误处理以及底层架构代码等方面。Andi 对网络编程相关内容做了全面审校，使作者在 Linux TCP/IP 实现方面大开眼界，此外，他还提供了许多建议，用以改善本书主题的展现形式。
- **Martin Landers** (Google) 在我有幸与他共事时，他还是一名学生。其后，他在短期内便集诸多技能于一身，且形象百变——软件架构师、IT 培训师以及职业黑客。劳

Martin 大驾审校本书，实为作者之幸。他对本书的评论和更正往往一针见血，对多章内容质量的改进功莫大焉。

- **Jamie Lokier** 是公认的内核高手，投身于 Linux 开发已达 15 年之久。如今，他自封为“专家，长于解决潜伏于 Linux 系统中的疑难杂症”。Jamie 极其全面地审校了本书涉及内存映射、POSIX 共享内存以及虚拟内存操作等方面的章节。他的审校工作不但纠正了作者对相关主题细节方面的许多误解，相应各章的结构也得以大为改观。
- **Barry Margolin** 在其 25 年职业生涯中，从事过系统程序员、系统管理员以及技术支持工程师。当前，他作为一名系统性能工程师，供职于 Akamai 技术公司。在各种讨论 UNIX 和 Internet 技术主题的网络论坛上，他频频现身，威名素著。他还是多本相关技术主题书籍的技术审稿人。Barry 审校了本书若干章节，并提出了诸多改进意见。
- **Paul Pluzhnikov** (Google) 之前曾是 Insure++ 内存调试工具的技术带头人和主要开发者。有时，他也会以 GDB 黑客的身份现身，在网上论坛里积极地回复有关调试、内存分配、共享库以及运行时环境方面的问题。Paul 审校了本书多章内容，提出了许多宝贵意见。
- **John Reiser** (与 Tom London) 实现了 UNIX 向 32 位架构移植的早期版本之一 VAX-11/780。他还是 mmap() 系统调用的编写者。John 审校了本书多章内容，自然也包含 mmap() 所在的章节。他所提供的大量历史洞见，及其对技术透彻地阐释为本书增色不少。
- **Anthony Robins** (新西兰 Otago 大学计算机科学副教授) 笔者 30 年的密友，本书某些章节的第一个读者，也是最早提出宝贵意见的技术审校者，在本书的写作过程中，一直给予作者以激励。
- **Michael Schröder** (Novell) GNU screen 程序的主要开发者之一，这项编程工作已令 Michael 在终端驱动程序的实现方面达到了“巨细靡遗，了如指掌”的境界。Michael 除了审校本书与终端和伪终端相关的章节以外，还针对进程组、会话以及任务控制等章节反馈了极为有益的意见。
- **Manfred Spraul** 曾从事过 Linux 内核中 (包括但不限于) IPC 的开发工作，不吝审校了本书与 IPC 相关的若干章节，并提出了许多改进意见。
- **Tom Swigg**, 作者在 DEC 从事培训工作时曾与他共事，作为本书最早的技术审校者之一，他对许多章节都反馈了极其重要的意见。Tom 从事软件工程师和 IT 培训师的工作已逾 25 年，目前就职于伦敦南岸大学 (London South Bank University)，在 VMware 环境下从事 Linux 编程和技术支持工作。
- **Jens Thoms Törning** 继承了物理学家改学编程的优良传统，大批开源的设备驱动程序和其他软件都出自他手。Jens 审校的章节之多，在技术方面跨度之大，着实令人瞠目，他对各章内容的改进都提出了独特而又弥足珍贵的见解。

还有许多其他技术审稿人也审校了本书的不同内容，并提出了诸多宝贵意见。在此，作者向以下一干技术审稿人表示感谢 (以姓氏字母顺序排列): George Anzinger (MontaVista Software)、Stefan Becher、Krzysztof Benedyczak、Daniel Brahneborg、Andries Brouwer、Annabel Church、Dragan Cvetkovic、Floyd L. Davidson、Stuart Davidson (Hewlett-Packard Consulting)、Kasper Dupont、Peter Fellingner (jambit GmbH)、Mel Gorman (IBM)、Niels Göllesch、Claus Gratzl、Serge Hallyn (IBM)、Markus Hartinger (jambit GmbH)、Richard Henderson (Red Hat)、Andrew Josey (The Open Group)、Dan Kegel (Google)、Davide Libenzi、Robert Love (Google)、

H.J. Lu (Intel Corporation)、Paul Marshall、Chris Mason、Michael Matz (SUSE)、Trond Myklebust、James Peach、Mark Phillips (Automated Test Systems)、Nick Piggin (Novell SUSE 实验室)、Kay Johannes Potthoff、Florian Rampp、Stephen Rothwell (IBM Linux 技术中心)、Markus Schwaiger、Stephen Tweedie (Red Hat)、Britta Vargas、Chris Wright、Michal Wronski 以及 Umberto Zamuner。

除了技术审稿人之外，作者还得到了各界人士及组织在其他方面的帮助。

我要感谢以下人等为我解答技术难题，他们是：Jan Kara、Dave Klekamp 和 Jon Snader。我要感谢 Claus Gatzl 和 Paul Marshall 在系统管理方面对我的帮助。

我要感谢 Linux 基金会 (LF)。2008 年间，LF 资助我作为一名全职研究人员参与 man-pages 项目，并从事 Linux 编程接口的测试和设计评审工作。虽然 LF 全职研究员的身份不能为本书的写作提供直接的资金支持，但却使作者得以养家糊口，这一助力本意在于令我全身心投入对 Linux 编程接口的测试以及对文档的编纂工作，却也惠及作者的“私活”。抛开公事不谈，我要感谢 Jim Zemlin——我在 LF 的“接口”人，还要感谢 LF 技术咨询委员会的一干专家，感谢他们对我的聘任。

感谢 Alejandro Forero Cuervo 对本书书名的建议！

25 年前，在我为第一个学位拼搏之时，Robert Biddle 激起了我对 UNIX、C 以及 Ratfor 的兴趣，谢谢你，老兄。虽然以下诸君与本书并无直接干系，但当我在新西兰坎特伯雷大学攻读第二学位时，他们就鼓励我在写作道路上坚持下去，在此，我要向他们表示感谢，他们是 Michael Howard、Jonathan Mane-Wheoki、Ken Strongman、Garth Fletcher、Jim Pollard，以及 Brian Haig。

由 Richard Stevens 所著的几部关于 UNIX 编程和 TCP/IP 方面的杰作，数年来一直被我辈程序员奉为圭臬，只可惜先贤已逝。凡是读过上述书籍的读者势必会注意到，本书与 Richard Stevens 的那几本巨著看起来有些相似。这并非偶然。在构思本书时，作者曾从较为宏观的角度就书籍设计反复斟酌，可最终发现 Richard Stevens 所采用的方法才是正解，正因如此，本书采用了与其相同的展示方式。

感谢下列人士和组织为我提供 UNIX 系统，使我得以运行测试程序，并验证其他 UNIX 实现的细节，感谢 Anthony Robins 和 Cathy Chandra 在新西兰 Otago 大学所提供的多种 UNIX 测试系统，感谢 Martin Landers、Ralf Ebner 和 Klaus Tilk 在德国慕尼黑技术大学 (Technische Universität) 所提供的多种 UNIX 测试系统，感谢 HP 公司在 Internet 上免费开放他们的 testdrive 系统，感谢 Paul de Weerd 使我得以访问 OpenBSD 系统。

要衷心感谢两家慕尼黑公司及其老板，这两家公司除了为我提供了工作机会（还是弹性工作制）和热情的同事，还格外开恩，允许我在写作本书时使用他们的办公室。感谢 exolution 有限公司的 Thomas Kahabka 和 Thomas Gmelch，特别要感谢 jambit 有限公司的 Peter Fellinger 和 Markus Hartinger。

感谢下列人士对我提供的各种帮助，他们是 Dan Randow、Karen Korrel、Claudio Scalmazzi、Michael Schüpbach 和 Liz Wright。感谢 Rob Suisted 和 Lynley Cook 为封面和封底所提供的照片。

感谢下列人士以不同方式给作者以鼓励和支持，他们是 Deborah Church、Doris Church 和 Annie Currie。

感谢 No Starch 出版社大队人马为这一庞大创作项目所提供的各种帮助。Bill Pollock 从项目之初就一直秉持直言不讳的风格，始终对本书的完成充满信心，并耐心地关注着项目的进展，我要对他表示感谢。感谢本书最初的责任编辑 Megan Dunchak。感谢本书的文字编辑

Marilyn Smith, 无论我如何殚精竭虑以求文字的清晰与一致, 此君总能从鸡蛋里挑出骨头。本书的版面和设计由 Riley Hoffman 全面负责, 在“上了同一条船”后又挑起了制作编辑的重担。Riley 总是不厌其烦地满足我的请求, 以求本书的排版无误——最终结果堪称完美。谢谢你。

现在, 我才体味出下面这句老话的真正含义: 一人写作, 全家受累。感谢 Britta 和 Cecilia 对我的支持, 感谢你们能容忍我因写作本书而长时间地不着家。

许可

承蒙 IEEE (美国电气电子工程师学会) 和 The Open Group 惠允, 本书得以引用 IEEE Std 1003.1, 2004 版以及 The Open Group 基础规范第 6 号 (Issue 6) 中 POSIX (可移植性操作系统接口) ——信息技术标准的部分文字。可通过 <http://www.unix.org/version3/online.html> 在线查阅规范的完整版本。

Web 站点和程序示例的源码

读者可在 <http://man7.org/tlpi> 上找到更多有关本书的信息, 包括本书的勘误表和程序示例源码。

反馈

欢迎读者提供 BUG 报告、对代码的改进建议, 以及为进一步提高代码可移植性而提出的修订意见。同样欢迎读者提供针对本书内容的缺陷报告和改进叙述方式的一般性建议。当前的勘误列表可参见 <http://man7.org/tlpi/errata/>。由于 Linux 编程接口变化无常、且变更有时极为频繁, 仅凭作者一己之力很难“与时俱进”, 因此读者就全新或已变更的 Linux 编程接口特性所提供的反馈信息, 作者也将乐于收到, 并会纳入本书的下一版中。

Michael Timothy Kerrisk

于德国慕尼黑和新西兰克赖斯特彻奇

2010 年 8 月

mtk@man7.org

目 录

第 1 章 历史和标准	1
1.1 UNIX 和 C 语言简史	1
1.2 Linux 简史	4
1.2.1 GNU 项目	4
1.2.2 Linux 内核	5
1.3 标准化	8
1.3.1 C 编程语言	8
1.3.2 首个 POSIX 标准	9
1.3.3 X/Open 公司和 The Open Group	10
1.3.4 SUSv3 和 POSIX.1-2001	10
1.3.5 SUSv4 和 POSIX.1-2008	12
1.3.6 UNIX 标准时间表	12
1.3.7 实现标准	14
1.3.8 Linux、标准、Linux 标准规范 (Linux Standard Base)	14
1.4 总结	15
第 2 章 基本概念	17
2.1 操作系统的核心——内核	17
2.2 shell	19
2.3 用户和组	20
2.4 单根目录层级、目录、链接及文件	21
2.5 文件 I/O 模型	23
2.6 程序	24
2.7 进程	25
2.8 内存映射	27
2.9 静态库和共享库	28
2.10 进程间通信及同步	28

2.11	信号	29
2.12	线程	30
2.13	进程组和 shell 任务控制	30
2.14	会话、控制终端和控制进程	30
2.15	伪终端	31
2.16	日期和时间	31
2.17	客户端服务器架构	32
2.18	实时性	32
2.19	/proc 文件系统	33
2.20	总结	33
第 3 章	系统编程概念	34
3.1	系统调用	34
3.2	库函数	36
3.3	标准 C 语言函数库；GNU C 语言函数库 (glibc)	37
3.4	处理来自系统调用和库函数的错误	38
3.5	关于本书示例程序的注意事项	40
3.5.1	命令行选项及参数	40
3.5.2	常用的函数及头文件	40
3.6	可移植性问题	49
3.6.1	特性测试宏	49
3.6.2	系统数据类型	51
3.6.3	其他的可移植性问题	53
3.7	总结	54
3.8	练习	55
第 4 章	文件 I/O：通用的 I/O 模型	56
4.1	概述	56
4.2	通用 I/O	58
4.3	打开一个文件：open()	58
4.3.1	open()调用中的 flags 参数	60
4.3.2	open()函数的错误	63
4.3.3	creat()系统调用	64
4.4	读取文件内容：read()	64
4.5	数据写入文件：write()	65
4.6	关闭文件：close()	66
4.7	改变文件偏移量：lseek()	66
4.8	通用 I/O 模型以外的操作：ioctl()	70
4.9	总结	71
4.10	练习	71

第 5 章 深入探究文件 I/O	72
5.1 原子操作和竞争条件	72
5.2 文件控制操作: <code>fcntl()</code>	75
5.3 打开文件的状态标志	75
5.4 文件描述符和打开文件之间的关系	76
5.5 复制文件描述符	78
5.6 在文件特定偏移量处的 I/O: <code>pread()</code> 和 <code>pwrite()</code>	80
5.7 分散输入和集中输出 (Scatter-Gather I/O): <code>readv()</code> 和 <code>writev()</code>	81
5.8 截断文件: <code>truncate()</code> 和 <code>ftruncate()</code> 系统调用	84
5.9 非阻塞 I/O	84
5.10 大文件 I/O	85
5.11 <code>/dev/fd</code> 目录	88
5.12 创建临时文件	88
5.13 总结	90
5.14 练习	90
第 6 章 进程	92
6.1 进程和程序	92
6.2 进程号和父进程号	93
6.3 进程内存布局	94
6.4 虚拟内存管理	97
6.5 栈和栈帧	99
6.6 命令行参数 (<code>argc, argv</code>)	99
6.7 环境列表	101
6.8 执行非局部跳转: <code>setjmp()</code> 和 <code>longjmp()</code>	106
6.9 总结	111
6.9 练习	112
第 7 章 内存分配	113
7.1 在堆上分配内存	113
7.1.1 调整 program break: <code>brk()</code> 和 <code>sbrk()</code>	113
7.1.2 在堆上分配内存: <code>malloc()</code> 和 <code>free()</code>	114
7.1.3 <code>malloc()</code> 和 <code>free()</code> 的实现	117
7.1.4 在堆上分配内存的其他方法	120
7.2 在堆栈上分配内存: <code>alloca()</code>	122
7.3 总结	123
7.4 练习	123
第 8 章 用户和组	124
8.1 密码文件: <code>/etc/passwd</code>	124

8.2	shadow 密码文件: /etc/shadow	125
8.3	组文件: /etc/group	126
8.4	获取用户和组的信息	127
8.5	密码加密和用户认证	132
8.6	总结	135
8.7	练习	135
第 9 章	进程凭证	136
9.1	实际用户 ID 和实际组 ID	136
9.2	有效用户 ID 和有效组 ID	136
9.3	Set-User-ID 和 Set-Group-ID 程序	137
9.4	保存 set-user-ID 和保存 set-group-ID	138
9.5	文件系统用户 ID 和组 ID	139
9.6	辅助组 ID	140
9.7	获取和修改进程凭证	140
9.7.1	获取和修改实际、有效和保存设置标识	140
9.7.2	获取和修改文件系统 ID	145
9.7.3	获取和修改辅助组 ID	145
9.7.4	修改进程凭证的系统调用总结	146
9.7.5	示例: 显示进程凭证	148
9.8	总结	149
9.9	习题	150
第 10 章	时间	151
10.1	日历时间 (Calendar Time)	151
10.2	时间转换函数	153
10.2.1	将 time_t 转换为可打印格式	153
10.2.2	time_t 和分解时间之间的转换	154
10.2.3	分解时间和打印格式之间的转换	155
10.3	时区	161
10.4	地区 (Locale)	163
10.5	更新系统时钟	167
10.6	软件时钟 (jiffies)	168
10.7	进程时间	168
10.8	总结	171
10.9	练习	172
第 11 章	系统限制和选项	173
11.1	系统限制	174
11.2	在运行时获取系统限制 (和选项)	176
11.3	运行时获取与文件相关的限制 (和选项)	178

11.4	不确定的限制	179
11.5	系统选项	180
11.6	总结	181
11.7	练习	182
第 12 章	系统和进程信息	183
12.1	/proc 文件系统	183
12.1.1	获取与进程有关的信息: /proc/PID	183
12.1.2	/proc 目录下的系统信息	185
12.1.3	访问/proc 文件	186
12.2	系统标识: uname()	188
12.3	总结	190
12.4	练习	190
第 13 章	文件 I/O 缓冲	191
13.1	文件 I/O 的内核缓冲: 缓冲区高速缓存	191
13.2	stdio 库的缓冲	194
13.3	控制文件 I/O 的内核缓冲	196
13.4	I/O 缓冲小结	200
13.5	就 I/O 模式向内核提出建议	201
13.6	绕过缓冲区高速缓存: 直接 I/O	202
13.7	混合使用库函数和系统调用进行文件 I/O	204
13.8	总结	205
13.9	练习	205
第 14 章	系统编程概念	207
14.1	设备专用文件 (设备文件)	207
14.2	磁盘和分区	208
14.3	文件系统	209
14.4	i 节点	211
14.5	虚拟文件系统 (VFS)	213
14.6	日志文件系统	214
14.7	单根目录层级和挂载点	215
14.8	文件系统的挂载和卸载	216
14.8.1	挂载文件系统: mount()	217
14.8.2	卸载文件系统: umount()和 umount2()	222
14.9	高级挂载特性	223
14.9.1	在多个挂载点挂载文件系统	224
14.9.2	多次挂载同一挂载点	224
14.9.3	基于每次挂载的挂载标志	225
14.9.4	绑定挂载	225

14.9.5 递归绑定挂载	226
14.10 虚拟内存文件系统: tmpfs	227
14.11 获得与文件系统有关的信息: statvfs()	228
14.12 总结	229
14.13 练习	230
第 15 章 文件属性	231
15.1 获取文件信息: stat()	231
15.2 文件时间戳	236
15.2.1 使用 utime()和 utimes()来改变文件时间戳	238
15.2.2 使用 utimensat()和 futimens()改变文件时间戳	239
15.3 文件属主	241
15.3.1 新建文件的属主	241
15.3.2 改变文件属主: chown()、fchown()和 lchown()	241
15.4 文件权限	244
15.4.1 普通文件的权限	244
15.4.2 目录权限	246
15.4.3 权限检查算法	246
15.4.4 检查对文件的访问权限: access()	248
15.4.5 Set-User-ID、Set-Group-ID 和 Sticky 位	249
15.4.6 进程的文件模式创建掩码: umask()	249
15.4.7 更改文件权限: chmod()和 fchmod()	251
15.5 I 节点标志 (ext2 扩展文件属性)	252
15.6 总结	256
15.7 练习	256
第 16 章 扩展属性	258
16.1 概述	258
16.2 扩展属性的实现细节	260
16.3 操控扩展属性的系统调用	260
16.4 总结	264
16.5 练习	264
第 17 章 访问控制列表	265
17.1 概述	265
17.2 ACL 权限检查算法	267
17.3 ACL 的长、短文本格式	268
17.4 ACL_mask 型 ACE 和 ACL 组分类	269
17.5 getfacl 和 setfacl 命令	270
17.6 默认 ACL 与文件创建	271

17.7	ACL 在实现方面的限制	272
17.8	ACL API	273
17.9	总结	280
17.10	练习	280
第 18 章	目录与链接	281
18.1	目录和（硬）链接	281
18.2	符号（软）链接	283
18.3	创建和移除（硬）链接：link()和 unlink()	286
18.4	更改文件名：rename()	289
18.5	使用符号链接：symlink()和 readlink()	290
18.6	创建和移除目录：mkdir()和 rmdir()	291
18.7	移除一个文件或目录：remove()	292
18.8	读目录：opendir()和 readdir()	292
18.9	文件树遍历：nftw()	297
18.10	进程的当前工作目录	301
18.11	针对目录文件描述符的相关操作	303
18.12	改变进程的根目录：chroot()	304
18.13	解析路径名：realpath()	306
18.14	解析路径名字符串：dirname()和 basename()	307
18.15	总结	309
18.16	练习	309
第 19 章	监控文件事件	311
19.1	概述	311
19.2	inotify API	312
19.3	inotify 事件	313
19.4	读取 inotify 事件	315
19.5	队列限制和/proc 文件	319
19.6	监控文件的旧有系统：dnotify	320
19.7	总结	320
19.8	练习	320
第 20 章	信号：基本概念	321
20.1	概念和概述	321
20.2	信号类型和默认行为	323
20.3	改变信号处置：signal()	329
20.4	信号处理器简介	330
20.5	发送信号：kill()	333
20.6	检查进程的存在	334
20.7	发送信号的其他方式：raise()和 killpg()	335

20.8	显示信号描述	336
20.9	信号集	337
20.10	信号掩码（阻塞信号传递）	339
20.11	处于等待状态的信号	341
20.12	不对信号进行排队处理	341
20.13	改变信号处置：sigaction ()	345
20.14	等待信号：pause()	346
20.15	总结	347
20.16	练习	347
第 21 章	信号：信号处理器函数	348
21.1	设计信号处理器函数	348
21.1.1	再论信号的非队列化处理	348
21.1.2	可重入函数和异步信号安全函数	349
21.1.3	全局变量和 sig_atomic_t 数据类型	353
21.2	终止信号处理器函数的其他方法	354
21.2.1	在信号处理器函数中执行非本地跳转	354
21.2.2	异常终止进程：abort()	358
21.3	在备选栈中处理信号：sigaltstack()	358
21.4	SA_SIGINFO 标志	361
21.5	系统调用的中断和重启	366
21.6	总结	368
21.7	练习	369
第 22 章	信号：高级特性	370
22.1	核心转储文件	370
22.2	传递、处置及处理的特殊情况	372
22.3	可中断和不可中断的进程睡眠状态	373
22.4	硬件产生的信号	374
22.5	信号的同步生成和异步生成	374
22.6	信号传递的时机与顺序	375
22.7	signal()的实现及可移植性	376
22.8	实时信号	378
22.8.1	发送实时信号	379
22.8.2	处理实时信号	380
22.9	使用掩码来等待信号：sigsuspend()	384
22.10	以同步方式等待信号	387
22.11	通过文件描述符来获取信号	390
22.12	利用信号进行进程间通信	393
22.13	早期的信号 API（System V 和 BSD）	393
22.14	总结	395

22.15 练习	396
第 23 章 定时器与休眠	397
23.1 间隔定时器	397
23.2 定时器的调度及精度	402
23.3 为阻塞操作设置超时	402
23.4 暂停运行（休眠）一段固定时间	404
23.4.1 低分辨率休眠: <code>sleep()</code>	404
23.4.2 高分辨率休眠: <code>nanosleep()</code>	404
23.5 POSIX 时钟	407
23.5.1 获取时钟的值: <code>clock_gettime()</code>	407
23.5.2 设置时钟的值: <code>clock_settime()</code>	408
23.5.3 获取特定进程或线程的时钟 ID	408
23.5.4 高分辨率休眠的改进版: <code>clock_nanosleep()</code>	409
23.6 POSIX 间隔式定时器	410
23.6.1 创建定时器: <code>timer_create()</code>	410
23.6.2 配备和解除定时器: <code>timer_settime()</code>	412
23.6.3 获取定时器的当前值: <code>timer_gettime()</code>	413
23.6.4 删除定时器: <code>timer_delete()</code>	413
23.6.5 通过信号发出通知	414
23.6.6 定时器溢出	417
23.6.7 通过线程来通知	417
23.7 利用文件描述符进行通知的定时器: <code>timerfd</code> API	420
23.8 总结	423
23.9 练习	424
第 24 章 进程的创建	425
24.1 <code>fork()</code> 、 <code>exit()</code> 、 <code>wait()</code> 以及 <code>execve()</code> 的简介	425
24.2 创建新进程: <code>fork()</code>	427
24.2.1 父、子进程间的文件共享	428
24.2.2 <code>fork()</code> 的内存语义	430
24.3 系统调用 <code>vfork()</code>	433
24.4 <code>fork()</code> 之后的竞争条件（Race Condition）	434
24.5 同步信号以规避竞争条件	436
24.6 总结	438
24.7 练习	439
第 25 章 进程的终止	440
25.1 进程的终止: <code>_exit()</code> 和 <code>exit()</code>	440
25.2 进程终止的细节	441

25.3	退出处理程序	442
25.4	fork()、stdio 缓冲区以及_exit()之间的交互	445
25.5	总结	446
25.6	练习	446
第 26 章	监控子进程	447
26.1	等待子进程	447
26.1.1	系统调用 wait()	447
26.1.2	系统调用 waitpid()	449
26.1.3	等待状态值	450
26.1.4	从信号处理程序中终止进程	454
26.1.5	系统调用 waitid()	455
26.1.6	系统调用 wait3()和 wait4()	456
26.2	孤儿进程与僵尸进程	457
26.3	SIGCHLD 信号	459
26.3.1	为 SIGCHLD 建立信号处理程序	459
26.3.2	向已停止的子进程发送 SIGCHLD 信号	462
26.3.3	忽略终止的子进程	462
26.4	总结	464
26.5	练习	464
第 27 章	程序的执行	465
27.1	执行新程序: execve()	465
27.2	exec()库函数	468
27.2.1	环境变量 PATH	469
27.2.2	将程序参数指定为列表	470
27.2.3	将调用者的环境传递给新程序	471
27.2.4	执行由文件描述符指代的程序: fexecve()	471
27.3	解释器脚本	472
27.4	文件描述符与 exec()	474
27.5	信号与 exec()	477
27.6	执行 shell 命令: system()	477
27.7	system()的实现	480
27.8	总结	485
27.9	练习	485
第 28 章	详述进程创建和程序执行	487
28.1	进程记账	487
28.2	系统调用 clone()	493
28.2.1	clone()的 flags 参数	497

28.2.2	因克隆生成的子进程而对 waitpid()进行的扩展	503
28.3	进程的创建速度	503
28.4	exec()和 fork()对进程属性的影响	505
28.5	总结	508
28.6	练习	508
第 29 章	线程：介绍	509
29.1	概述	509
29.2	Pthreads API 的详细背景	511
29.3	创建线程	513
29.4	终止线程	514
29.5	线程 ID (Thread ID)	514
29.6	连接 (joining) 已终止的线程	515
29.7	线程的分离	517
29.8	线程属性	518
29.9	线程 VS 进程	518
29.10	总结	519
29.11	练习	519
第 30 章	线程：线程同步	521
30.1	保护对共享变量的访问：互斥量	521
30.1.1	静态分配的互斥量	524
30.1.2	加锁和解锁互斥量	524
30.1.3	互斥量的性能	526
30.1.4	互斥量的死锁	527
30.1.5	动态初始化互斥量	527
30.1.6	互斥量的属性	528
30.1.7	互斥量类型	528
30.2	通知状态的变化：条件变量 (Condition Variable)	529
30.2.1	由静态分配的条件变量	530
30.2.2	通知和等待条件变量	531
30.2.3	测试条件变量的判断条件 (predicate)	534
30.2.4	示例程序：连接任意已终止线程	534
30.2.5	经由动态分配的条件变量	537
30.3	总结	538
30.4	练习	538
第 31 章	线程：线程安全和每线程存储	539
31.1	线程安全 (再论可重入性)	539
31.2	一次性初始化	541

31.3	线程特有数据	542
31.3.1	库函数视角下的线程特有数据	542
31.3.2	线程特有数据 API 概述	543
31.3.3	线程特有数据 API 详述	543
31.3.4	使用线程特有数据 API	545
31.3.5	线程特有数据的实现限制	549
31.4	线程局部存储	549
31.5	总结	550
31.6	练习	551
第 32 章	线程：线程取消	552
32.1	取消一个线程	552
32.2	取消状态及类型	552
32.3	取消点	553
32.4	线程可取消性的检测	556
32.5	清理函数（cleanup handler）	556
32.6	异步取消	559
32.7	总结	560
第 33 章	线程：更多细节	561
33.1	线程栈	561
33.2	线程和信号	562
33.2.1	UNIX 信号模型如何映射到线程中	562
33.2.2	操作线程信号掩码	563
33.2.3	向线程发送信号	563
33.2.4	妥善处理异步信号	564
33.3	线程和进程控制	564
33.4	线程实现模型	566
33.5	Linux POSIX 线程的实现	567
33.5.1	LinuxThreads	567
33.5.2	NPTL	569
33.5.3	哪一种线程实现	570
33.6	Pthread API 的高级特性	572
33.7	总结	572
33.8	练习	572

第 1 章

历史和标准

Linux 是 UNIX 操作系统家族中的一员。就计算机的发展而言，UNIX 历史悠久。本章的第一部分会简要介绍 UNIX 的历史——以对 UNIX 系统和 C 编程语言起源的回顾拉开序幕，接着会述及成就今日 Linux 系统的两大关键因素：GNU 项目和 Linux 内核的开发。

UNIX 系统最引人关注的特征之一，是其开发不受控于某一厂商或组织。相反，许多团体——既有商业团体，也有非商业团体——都曾为 UNIX 的演进做出过贡献。这一渊源使 UNIX 集多种开创性的特性于一身，但同时也带来了负面影响——随着时间的推移，UNIX 的实现渐趋分裂。因此，要编写出能够运行于所有 UNIX 实现之上的应用程序愈发困难。这又导致了人们对 UNIX 实现的标准化呼声越来越高，本章的第二部分将讨论这一问题。

对 UNIX 的定义通常有两种。其一是指通过 SUS 所规范的官方一致性测试，且由 OPEN GROUP（UNIX 商标的持有者）正式授权冠以“UNIX”的操作系统。在写作本书之际，尚无开源的 UNIX 实现（比如，Linux 和 FreeBSD）获得了“UNIX”冠名。

在第二种定义中，UNIX 是指那种运作方式类似于经典 UNIX 系统（比如，最初的 Bell 实验室 UNIX 系统，及其后来的主要分支 System V 和 BSD）的操作系统。根据这一定义，一般将 Linux 视为 UNIX 系统（如同现代 BSD 系统一样）。尽管本书会密切关注 SUS，但也会遵循对 UNIX 的第二种定义，因此诸如“Linux，像其他 UNIX 实现一样……”这样的说法，会在书中频繁出现。

1.1 UNIX 和 C 语言简史

1969 年，在 AT&T 电话公司下辖的 bell 实验室中，Ken Thompson 开发出了首个 UNIX 实现。该实现是使用 Digital PDP-7 小型机的汇编语言开发而成的。其名称 UNIX 是“MULTICS（多信息及计算服务，Multiplexed Information and Computing Service）”一词的双关语，而 MULTICS 之名则出自一个早期的操作系统开发项目，该项目由 AT&T、MIT（麻省理工学院）以及通用电器公司联合开发。（因为未能开发出一款经济实用的操作系统，该项目首战失利。沮丧之余，AT&T 随即退出这一项目中。）Thompson 设计新操作系统的某些灵感正源于

MULTICS，其中包括：树形结构的文件系统、设立单独的程序用于解释命令（shell），以及将文件作为无结构字节流看待的概念。

1970年，AT&T的工程师们又在刚购进的 Digital PDP-11 小型机上，以汇编语言重写了 UNIX，当时，Digital PDP-11 算得上是最新颖、功能也最为强劲的计算机了。从大多数 UNIX 实现（包括 Linux）沿用至今的各种名称上，仍能发现这一 PDP-11 实现所残留的历史遗迹。

未过多久，Dennis Ritchie（Thompson 在 bell 实验室的同事，UNIX 开发的早期合作者）设计并实现出了 C 编程语言。这里有一个演变过程：C 语言传承自早期的解释型语言——B 语言；B 语言最初由 Thompson 实现，但其所包含的许多理念却来自于更早期的编程语言——BCPL。到了 1973 年，C 语言步入了成熟期，人们能够使用这一新语言重写几乎整个 UNIX 内核。UNIX 因此也一变而为最早以高级语言开发而成的操作系统之一，这也促成了 UNIX 系统后续向其他硬件架构的移植。

从 C 语言的起源不难看出为什么 C 语言及其“后裔”C++是当今使用最为广泛的系统编程语言。早期流行的编程语言其设计初衷并不在于此，例如：FORTRAN 语言意在帮助工程师和科研工作者们进行数学计算，COBOL 语言则是在商业系统中用来处理面向记录的数据流。C 语言的出现，填补了当时系统编程方面的语言空白。与 FORTRAN 和 COBOL 不同（这两种编程语言均由大型组织设计开发），C 语言的设计理念和设计需求出自于几位程序员的构思，他们的目标很单纯：为实现 UNIX 内核及其相关软件而开发一种高层语言。像 UNIX 操作系统本身一样，C 语言由专业程序员设计而为己用。其最终结果堪称完美：C 语言的设计前后连贯，且支持模块化设计，成为短小精干、高效实用、功能强大的编程语言。

UNIX 的第一版到第六版

1969~1979 年间，UNIX 历经了多次发布，也称为版本（edition）。实质上，这些发布是 AT&T 对 UNIX 进行演进开发时的一系列版本快照。[Salus, 1994]记录了 UNIX 前六版的发布日期如下。

- 1971 年 11 月发布的第一版：当时，UNIX 还运行在 PDP-11 上，但已附带了 FORTRAN 编译器，许多被沿用至今的程序都已有了雏形，这包括：ar、cat、chmod、chown、cp、dc、ed、find、ln、ls、mail、mkdir、mv、rm、sh、su 以及 who。
- 1972 年 6 月发布的第二版：当时，AT&T 内有 10 台计算机安装了 UNIX。
- 1973 年 2 月发布的第三版：该版本包括了 C 编译器，以及管道的首个实现。
- 1973 年 11 月发布的第四版：这也是几乎完全以 C 语言重写的首个 UNIX 版本。
- 1974 年 6 月发布的第五版：当时，UNIX 的装机数已经超过了 50 台。
- 1975 年 3 月发布的第六版：这也是在 AT&T 之外广泛使用的首个 UNIX 版本。

在此期间，UNIX 使用范围从 AT&T 自内而外逐步扩展，声名也随之远播。读者甚众的《ACM 通信》杂志刊载了一篇关于 UNIX 的论文（[Ritchie & Thompson, 1974]），这对 UNIX 知名度的提升功莫大焉。

当时，在美国政府的授权下，AT&T 垄断着全美电信市场。AT&T 与美国政府达成的协议条款禁止 AT&T 涉足软件销售行业——这意味着，AT&T 不能将 UNIX 作为产品销售。相反，从 1974 年的 UNIX 第五版开始，AT&T 准许高校在支付象征性的发布费用后使用 UNIX 系统——这一现象尤以第六版为烈。UNIX 系统的高校发布版包括了相关文档及内核源码（当时，内核源码约为 10 000 行左右）。

AT&T 对高校发布的 UNIX 极大促进了这一操作系统的普及和使用。时至 1977 年，UNIX 已经在约 500 个站点中运行，其中包括了全美及其他国家的 125 所大学。当时的商业操作系

统非常昂贵，UNIX 则为高校提供了一种交互式多用户操作系统，可谓物美价廉。此外，各校的计算机系还藉此获得了“鲜活”的操作系统源码，可以对源码进行修改，还可供学生们学习、实验之用。一些以 UNIX 知识为武装的学生后来成为 UNIX “传教士”。另外一些学生则组建或加盟了大量新兴公司，其业务主要是销售廉价的计算机工作站，而运行于其上的正是易于移植的 UNIX 操作系统。

BSD 和 System V 的诞生

发布于 1979 年 1 月的 UNIX 第七版改善了系统的可靠性，配备了增强型的文件系统。该版本还附带了不少新的工具软件，其中包括：awk、make、sed、tar、uucp、Bourne shell 以及 FORTRAN 77 编译器。第七版 UNIX 发布的重要意义还在于，从该版本起，UNIX 分裂为了两大分支：BSD 和 System V。接下来会简要描述二者的由来。

受母校加州大学伯克利分校之邀，Thompson 于 1975/1976 学年曾担任该校的客座教授。在此期间，他与研究生们一起为 UNIX 开发了许多新特性。（他的学生之一，Bill Joy，后来与人共同组建了 SUN 微系统公司——一家最早涉足 UNIX 工作站市场的公司。）光阴荏苒，许多 UNIX 的新工具和新特性又陆续在伯克利分校问世，这包括：C shell、vi 编辑器、一种改进型的文件系统（伯克利快速文件系统）、sendmail、Pascal 语言编译器，以及用于新型 Digital VAX 架构的虚拟内存管理机制。

这一命名为 BSD（伯克利软件发布，Berkeley Software Distribution）的 UNIX 版本（包括源码在内）分发颇广。1979 年 12 月，诞生了首个完整的 UNIX 发布版 3BSD。（之前发布的 Berkeley-BSD 和 2BSD 并非完整的 UNIX 发布版，仅含由伯克利分校开发的新工具。）

1983 年，加州大学伯克利分校的计算机系统研究组（Computer Systems Research Group）发布了 4.2BSD。该版本的发布意义深远，因为其包含了完整的 TCP/IP 实现，其中包括套接字应用编程接口（API）以及各种网络工具。4.2BSD 及其前身 4.1BSD 在世界上多所大学开始广为流传。以这两者为基础，还形成了 SunOS 操作系统（首发于 1983 年）——这一由 SUN 公司销售的 UNIX 变种。其他重要的 BSD 版本还有发布于 1986 年的 4.3BSD，以及发布于 1993 年的最终版本 4.4BSD。

首批 UNIX 向非 PDP-11 硬件机型的移植发生在 1977 年和 1978 年，当时，Dennis Ritchie 和 Steve Johnson 将 UNIX 移植到了 Interdata 8/32 上，与此同时，澳大利亚 Wollongong 大学的 Richard Miller 也将其移植到了 Interdata 7/32 上。伯克利分校针对 Digital Vax 架构的移植——也称为 32V，则基于 John Reiser 和 Tom Lodon 较早前（1978 年）的工作成果。该移植本质上与 PDP-11 上的 UNIX 第七版相同，只是支持的地址空间更大、数据类型更宽罢了。

与此同时，美国的反托拉斯法案强制对 AT&T 进行拆分（于 20 世纪 70 年代中期开始立案，到 1982 年 AT&T 正式解体）。随着其在电话系统市场垄断地位的丧失，AT&T 也因而获准销售 UNIX。这也催生了 1981 年 System III（3）的发布。System III 由 AT&T 所属的 UNIX 支撑团队（UNIX Support Group, USG）研发，该团队雇佣了数以百计的研发人员来从事 UNIX 系统的增强以及应用开发（尤其针对文档预备软件包和软件开发工具）。1983 年，System V 的首个发布版又接踵而至，在经过一系列发布后，USG 最终于 1989 年推出了 System V Release 4（SVR4），此时的 System V 纳入了 BSD 的诸多特性，包含联网能力。AT&T 将 System V 授权给不同厂商，这些厂商又将其作为自身 UNIX 实现的基础。

因此，除了遍布于学术界的各种 BSD 发布版外，到 20 世纪 80 年代末，商业性质的 UNIX 实现在各种硬件架构上都有了广泛应用。这包括：SUN 公司的 SunOS，以及后来的 Solaris；Digital 公司的 Ultrix 和 OSF/1（在历经一系列更名和收购后，现称为 HP Tru64 UNIX）；IBM

公司的 AIX；HP 公司的 HP-UX；NeXT 公司的 NeXTStep；在 Apple Macintosh 机上的 A/UX；以及 Microsoft 和 SCO 公司联合为 Intel x86-32 架构开发的 XENIX。（贯穿本书，我们将 x86-32 架构上的 Linux 实现称为 Linux/x86-32。）这一局面与当时典型的专有硬件搭配专有操作系统的模式形成了鲜明对照，那时，每个厂商只生产一种或至多几种专有的计算机芯片架构，然后再销售运行于该硬件架构之上的专有操作系统。大多数厂商系统的这种专有性，意味着消费者只能在一棵树上“吊死”。转换到另一专有操作系统和硬件平台，其代价十分高昂，不但需要移植现有应用，还需要对操作人员进行重新培训。从商业角度来看，考虑到上述因素，加之各厂商纷纷推出了廉价的单用户 UNIX 工作站，具备可移植性的 UNIX 系统魅力逐渐开始“凸显”。

1.2 Linux 简史

术语 Linux 通常用来指代完整的类 UNIX (UNIX-like) 操作系统，Linux 内核只是其中的一部分。这么定义多少有些措辞不当，因为一般商业 Linux 发布版中所含的诸多关键组件实际上发源于另一项目，早在 Linux 问世前几年就已经启动了。

1.2.1 GNU 项目

1984 年，Richard Stallman 之前一直供职于 MIT 的一位天赋异禀的程序员，开始着手创建一个“自由的 (free)” UNIX 实现。Stallman 的观点属于道德层面，而对“free”一词的定义则属于法律范畴而非经济范畴（请参见 <http://www.gnu.org/philosophy/free-sw.html>）。然而，Stallman 所描述的这一法律意义上的“自由 (freedom)”却蕴含着言外之意：应可免费或以低价获得诸如操作系统之类的软件。

对于那些在专有操作系统上加强限制条款的计算机厂商来说，Stallman 的这一举动无疑妨害了他们。所谓的限制条款是指：在一般情况下，计算机软件的消费者不但无权阅读自己所购软件的源码，而且还不能复制、更改及重新发行所购软件。Stallman 指出，在这种体制之下，只会造成程序员之间勾心斗角、敝帚自珍的局面，无法实现工作协同和成果共享。

与之针锋相对，为开发出一套完整而又可自由获取，包含内核以及所有相关软件包的类 UNIX 系统，Stallman 发起了 GNU 项目（“GNU's not UNIX”的递归缩写形式），并积极邀请有志之士加盟。1985 年，Stallman 创立了非盈利机构——自由软件基金会 (FSF)，以支持 GNU 项目和广义意义上的自由软件开发。

GNU 项目启动之时，BSD 还不具备 Stallman 所指的那种“free”属性。使用 BSD 不但仍需获得 AT&T 的许可，而且用户不得随意修改并重新发布 BSD 中 AT&T 拥有产权的代码部分。

GNU 项目的重要成果之一是制定了 GNU GPL (通用公共许可协议)，这也是 Stallman 倡导的自由 (free) 软件概念在法律上的体现。Linux 发布版中的大多数软件，包括 Linux 内核，都是以 GPL 或与之类似的许可协议发布的。以 GPL 许可协议发布的软件不但必须开放源码，而且应能在 GPL 条款的约束下自由对其进行重新发布。可以不受限制的修改以 GPL 许可协议发布的软件，但任何经修改后发布的软件仍需遵守 GPL 条款。若经过修改的软件以二进制（可执行）形式发布，那么软件的修改者必需满足软件使用者的以下要求：以不高于发行成本的价格，获得修改后的软件源码。GPL 的第一版发布于 1989 年。当前的许可协议版本为 2007 年发布的第三版。此许可协议的第二版于 1991 年发布，至今仍在广泛使用，Linux 内核就是以该版许可协议发布的。（对各种自由软件许可协议的讨论可见诸于[St. Laurent, 2004]和[Rosen, 2005]。）

最初，GNU 项目未能开发出能够有效运作的 UNIX 内核，但却开发了大量其他程序。由于这些程序全都针对类 UNIX 系统而设计，因此（理论上）均有可能在现有的 UNIX 实现上运行（实际情况也的确如此），更有甚者，有时还被移植到了其他操作系统上。Emacs 文本编辑器、GCC（原名为 GNU C 编译器，现更名为 GNU 编译器集合，集 C、C++，以及其他编程语言的编译器于一身）、bash shell 以及 glibc（GNU C 语言库）便是 GNU 项目结出的硕果。

到了 20 世纪 90 年代早期，GNU 项目已经开发出了一套几乎完整的操作系统，除了还缺少其中最重要的一环：能够有效运转的 UNIX 内核。于是，GNU 项目以 Mach 微内核为基础，发起了一项雄心勃勃的内核设计计划，史称 GNU/HURD 计划。然而，时至今日，HURD 的发布还遥遥无期（写作本书之际，HURD 的研发尚在进行中，该内核目前只能运行于 x86-32 架构之上）。

在构成通常所说的“Linux 系统”的程序代码中，由于有相当一部分都源自 GNU 项目，因此 Stallman 更愿意用“GNU/Linux”一词来称呼整个系统。这一称谓问题（Linux Vs. GNU/Linux）也在自由软件社区中引发了一些口舌之争。因为本书主要关注 Linux 内核的 API，故而通常会采用术语“Linux”。

万事具备，独缺内核。只要再拥有一个能够有效运作的内核，就能使 GNU 项目开发出的 UNIX 系统“功德圆满”。

1.2.2 Linux 内核

1991 年，Linus Torvalds，一位芬兰赫尔辛基大学的学生，在外界的激励下为自己的 Intel 80386 PC 开发了一个操作系统。在一门学习课程中，Torvalds 开始接触 Minix——由荷兰大学教授 Andrew Tanenbaum 于 20 世纪 80 年代中期开发的一款小型、类 UNIX 的操作系统内核。Tanenbaum 将 Minix 连同源码完全开放，作为大学操作系统设计课程的教学工具。人们可以在 386 系统上构建并运行 Minix 内核。当然，正因为其主要用于教学，Minix 在设计上几乎独立于硬件架构，故而也未对 386 处理器的能力充分加以利用。

因此，为了开发出一个高效而又功能齐备的 UNIX 内核，Torvalds 开始“自力更生”。数月之后，Torvalds 开发出一个内核“雏形”，可以编译并运行各种 GNU 程序。随之，于 1991 年 10 月 5 日，为求得其他程序员的帮助，Torvalds 在 Usenet 新闻组 comp.os.minix 上就其内核 0.02 版发表了如下申明，如今已被广为引用。

还在念叨 minix1.1 的好日子——人人都能给自个儿写设备驱动，不用看别人的脸色？手头没有称心的项目？是不是特想有一个操作系统，能依着自个的想法来回折腾，还能长见识？瞧瞧 minix 上面跑的那些玩意吧，是不是挺没劲？不想再为调个酷毙了的程序，一宿一宿熬个没完？真这样，那我可找对人了。一个月前在帖子里就提过，我正在写一个操作系统，在 AT-386 上面跑，免费的，挺像 minix。现在总算到了这份上，凑合能用（当然，这得看您想干吗）。现在，我愿意公布系统的源码，请大家多瞅瞅，多用用。系统版本只是 0.02，不过 bash、gcc、gnu-make、gnu-sed 还有 compress 等等倒是都跑通了。

为了传承 UNIX 历史悠久的光荣传统，在为 UNIX 系统克隆命名时，总以字母“X”结尾，故而，人们最终将这一内核命名为 Linux。最初，Linux 的使用许可协议要严格得多，但 Torvalds 很快就将其归于 GNU GPL 阵营。

Torvalds 做到了一呼百应。其他程序员与 Torvalds 一起加入到 Linux 的开发行列，添加了

很多新特性，诸如：改进型的文件系统、对网络的支持、设备驱动程序以及对多处理器的支持等。到了 1994 年 3 月，开发者们发布了 Linux 1.0 版本。随之，Linux 1.2 发布于 1995 年 3 月，Linux 2.0 发布于 1996 年 6 月，Linux 2.2 发布于 1999 年 1 月，Linux 2.4 发布于 2001 年 1 月。对内核 2.5 版本的开发始于 2001 年 11 月，并最终于 2003 年 12 月发布了 Linux 内核 2.6。

题外话：BSD

值得一提的是，20 世纪 90 年代初，另一种可以免费获得的 UNIX 也能在 x86-32 硬件架构上运行。Bill 和 Lynne Jolitz 将业已成熟的 BSD 系统移植到 32 位的 x86 cpu 上，命名为 386/bsd。这项移植工作基于 BSD Net/2（发布于 1991 年 6 月），即 4.3BSD 源码的版本之一，该版本中残存的所有 AT&T 专有源码要么被全部替换，要么予以删除——主要针对 6 个无法轻易更换的源码文件而言。Jolitizes 夫妇将 Net/2 代码移植到了 x86-32 硬件架构，重写了缺失的源码，并于 1992 年 2 月发布了 386/BSD 的首个版本（0.0 版本）。

在初战告捷后，对 386/BSD 的开发工作便出于各种原因而停滞不前。面对日渐积压的大量补丁程序，另外两组开发团队相机而动，基于 386/BSD 分别创建了自己的版本：NetBSD 和 FreeBSD。前者侧重于对大量硬件平台的可移植性；后者则主要关注性能，并成为如今应用最为广泛的 BSD。1993 年 4 月，NetBSD 首版（版本号为 0.8）发布。FreeBSD 的首个 CD-ROM 版本（版本号为 1.0）则发布于 1993 年 12 月。1996 年，OpenBSD 在从 NetBSD 项目分离出去之后，也发布了最初版本（版本号 2.0）。相比较而言，OpenBSD 偏重于安全性。2003 年中段，在与 FreeBSD 4.x 分道扬镳之后，一款新型 BSD——DragonFly BSD 又浮出水面。DragonFly BSD 采用的设计方法与 FreeBSD 5.x 有所不同，能够支持对称多处理器（SMP）架构。

若是不提及 20 世纪 90 年代初 UNIX System Laboratories（USL，派生自 AT&T 的子公司，专门从事 UNIX 的开发和销售）和 Berkeley 之间的那场官司，那么对 BSD 的介绍恐怕就算不得完整。1992 年初，Berkeley Software Design, Incorporated 公司（BSDi，如今隶属于 Wind River 公司）开始发行受商业支持的 BSD UNIX——BSD/OS——以 Net/2 发布版以及 Jolitze 夫妇所开发的 386/BSD 特性为基础。BSDi 的发布版包含二进制和源代码，售价 995 美元，此外，BSDi 还建议潜在客户使用其电话号码 1-800-ITS-UNIX。

1992 年 4 月，USL 对 BSDi 发起诉讼，诉状称 BSDi 售出产品中含有 USL 专有源码及商业机密，要求其停止销售。此外，诉状还指称 BSDi 的电话号码容易误导消费者，要求 BSDi 停止使用。这场诉讼愈演愈烈，最终还加入了对加州大学的索赔请求。法院最终驳回了 USL 几乎所有的诉讼请求，仅对其中的两项请求予以支持。随后，加州大学又针对 USL 发起发诉，诉称：USL 没有为 System V 中使用的 BSD 代码支付费用。

这场诉讼悬而未决之际，USL 已被 Novell 收购，Novell 时任 CEO——Ray Noorda 公开声称：较之于法庭辩论，自己的公司更愿意参与市场竞争。双方最终于 1994 年 1 月达成庭外和解。在删除 Net/2 release 源码 18000 个文件中的 3 个文件，对若干其他文件做出细微改动，并为其他大约 70 个文件添加 USL 版权注意事项后，加州大学仍可继续自由发布 BSD。1994 年 6 月，经过修改的系统以 4.4BSD-Lite 之名发布（1995 年 6 月，加州大学发布了最后一版 4.4BSD-Lite，版本号为 Release 2）。此时，根据和解条款，BSDi、FreeBSD 以及 NetBSD 纷纷以经过修改的 4.4BSD-Lite 源码替换了各自的 Net/2 基础源码。据[McKusick et al., 1996]一书记述，尽管这在一定程度上延误了 BSD 衍生系统的开发，但也有其积极意义。加州大学计算机研究组（Computer Systems Research Group）自 Net/2 发布后 3 年的开发成果，被重新同步到上述系统中。

Linux 内核版本号

与大多自由软件项目一样，Linux 也遵循及早、经常的发布模式，因而对内核的修订会频繁出现（有时甚至是每天都有）。随着 Linux 用户群的激增，对这一发布模式有所调整，意在降低对现有用户的干扰。具体来说，在 Linux 1.0 版本之后，内核开发者针对每次发布所采用的内核版本编号方案为 x.y.z。x 表示主版本号，y 为附属于主版本号的次版本号，z 是从属于次版本号的修订版本号（细微的改进和 BUG 修复）。

采用这一发布模式，内核的两个版本会一直处于开发之中。一个是用于生产系统的稳定（stable）分支，其次版本号为偶数；另一个是经常变动的开发（development）分支，其次版本号为奇数（当前稳定版次版本号+1）。指导思想是（在实践中并未严格执行）应将所有新特性添加到内核当前的开发分支系列中，而对内核稳定分支系列的修订应严格限定为细微的改进及 bug 修复。当开发者认为当前的开发分支已宜于发布时，会将该开发分支转换成新的稳定分支，并为其分配一个偶数的次版本号。例如，内核开发分支 2.3.z 会“进化”为内核稳定分支 2.4。

随着 2.6 内核的发布，内核开发模式再次发生改变。稳定内核版本之间发布间隔过长，因而导致诸多问题和不便，这是内核开发模型改变的主要原因（从 Linux 2.4.0 到 2.6.0 的发布历时近 3 年）。虽然还会就该模型的微调定期开展讨论，但基本细节已经确定如下。

- 不再有稳定内核和开发内核的概念。每个新的 2.6.z 发布版都可以包含新特性，其生命周期始于对新特性的追加，然后历经一系列候选发布版本让新特性稳定下来。当开发者认为某个候选版本足够稳定时，便可将其作为内核 2.6.z 发布。一般情况下，发布周期约为 3 个月。
- 有时，也可能需要为某个稳定的 2.6.z 发布版打上些小补丁程序，以修复 bug 或安全问题。如果这样的修复工作具有足够高的优先级，并且补丁程序的正确性也“毋庸置疑”，那么无需等待下一个 2.6.z 发布版，可以直接应用补丁创建一个版本号形如 2.6.z.r 的发布版本，其中，r 作为该 2.6.z 内核版本的次修订版序号。
- 额外责任将转嫁给 Linux 发行厂商，由他们来确保随 Linux 发行版一同发行内核的稳定性。

本书后续各章有时会提及 API 发生特定变化（比如，新增了系统调用或者系统调用发生变化时）的相应内核版本。在 2.6.z 系列之前，虽然大多数内核变化都见于具有奇数版本号的开发分支，但本书通常所指的是那些变化初次出现的稳定内核版本，这是因为大多数应用开发者一般都会使用稳定版的内核，而非开发版本。很多情况下，手册页会注明某一具体特性出现或发生变化时开发版内核的确切版本号。

对 2.6.z 系列内核所发生的改变，本书会注明确切的内核版本号。当书中言及 2.6 版本内核的新特性，且版本号又不带“z”这一修订版本号时，意指该特性是在 2.5 开发版内核中实现，并首度出现于稳定内核版本 2.6.0。

写作本书之际，Linux 内核 2.4 的稳定版尚处于维护期，维护者们仍在将关键性的补丁和缺陷修正合并起来，定期发布新的修订版。这使得已安装系统能继续使用 2.4 内核，而不必非要升级到新的内核系列（有时候，升级起来并不轻松）。

向其他硬件架构的移植

Linux 开发之初，主要目标是针对 Intel 80386 的高效系统实现，而非向其他处理器架构迁移的可移植性。然而，随着 Linux 的日益普及，针对其他处理器架构的移植版本开始出现，首

先就是向 Digital Alpha 芯片的移植。Linux 所支持的硬件架构队伍在持续壮大，其中包括：x86-64、Motorola/IBM PowerPC 和 PowerPC64、Sun SPARC 和 SPARC64 (UltraSPARC)、MIPS、ARM (Acorn)、IBM zSeries (formerly System/390)、Intel IA-64 (Itanium, 请参阅[Mosberger & Eranian, 2002])、Hitachi SuperH、HP PA-RISC, 以及 Motorola 68000。

Linux 发行版

准确说来，术语 Linux 只是指由 Linus Torvalds 和其他人所开发出的内核。可是，也常使用该术语来指代内核外加一大堆其他软件（工具和库）所构成的完整操作系统。Linux 草创之际，需要用户自行组装上述所有软件，创建文件系统，在文件系统上正确地安置并配置所有软件。用户不但要具备专业知识，还需为此耗费大量时间。如此一来，这便为 Linux 发行商们开启了市场，他们创建软件包（发行版），来自动完成大部分安装过程，其中包括了建立文件系统以及安装内核和其他所需软件等。

Linux 的发行版最早出现于 1992 年，包括 MCC Interim Linux（英国，曼彻斯特计算机中心）、TAMU（德克萨斯 A&M 大学）以及 SLS (SoftLanding Linux System)。至今健在的商业发行版 Slackware 诞生于 1993 年。几乎与此同时，也诞生了非商业的 Debian 发行版，SUSE 和 Red Hat 紧随其后。时下最流行的 Ubuntu 发行版问世于 2004 年。如今，对于那些在自由软件项目中表现活跃的程序员，许多 Linux 发行公司也会加以雇佣。

1.3 标准化

20 世纪 80 年代末，可用的 UNIX 实现层出不穷，由此也带来了种种弊端。有些 UNIX 实现基于 BSD，而另一些则基于 System V，还有一些则是对两大“流派”“兼容并蓄”。更有甚者，每个厂商都在自己的 UNIX 实现中添加了额外特性。其结果是将软件及技术人员在不同 UNIX 实现间转移就变得异常困难。这一形式有力地推动了 C 语言和 UNIX 系统的标准化进程，使得应用程序能够在不同操作系统间很方便地进行移植。接下来，将介绍由此而产生的各种标准。

1.3.1 C 编程语言

20 世纪 80 年代初，C 语言问世已达 10 年之久，在大量 UNIX 系统以及其他操作系统上都有实现。各种 C 语言的实现之间存在着细微差别，这部分是由于在当时 C 语言事实上的标准——Kernighan 和 Ritchie 于 1978 年所著的 *The C Programming Language* 一书中（有时，人们将书中所记载的老式 C 语言语法称为传统 C 或 K&R C），并未就 C 语言在某些方面的运作方式进行细化。此外，借鉴于 1985 年面世的 C++ 语言，在不破坏现有程序的前提下，C 语言得以进一步丰富和完善，其中最知名的莫过于函数原型、结构赋值、类型限定符（const and volatile）、枚举类型以及 void 关键字。

上述因素形成了 C 语言标准化进程的强力推手，ANSI（美国国家标准委员会）C 语言标准（X3.159-1989）最终于 1989 年获批，随之于 1990 年被 ISO（国际标准化组织）所采纳（ISO/IEC 9899:1990）。这份标准在定义 C 语言语法和语义的同时，还对标准 C 语言库操作进行了描述，这包括 stdio 函数、字符串处理函数、数学函数、各种头文件等等。通常将 C 语言的这一版本称为 C89 或者（不太常见的）ISO C90，Kernighan 和 Ritchie 所著的 *The C Programming Language* 第 2 版（1988）对其有完整描述。

1999 年，ISO 又正式批准了对 C 语言标准的修订版（ISO/IEC 9899:1999，请见 <http://www>。

open-std.org/jtc1/sc22/wg14/www/standards)。通常将这一标准称为 C99，其中包括了对 C 语言及其标准库的一系列修改，诸如，增加了 long long 和布尔数据类型、C++风格的注释 (//)、受限指针以及可变量数组等。写作本书之际，对 C 语言标准的进一步修订（非正式命名为 C1X）仍在进行之中，预计将于 2011 年正式获批。

C 语言标准独立于任何操作系统，换言之，C 语言并不依附于 UNIX 系统。这也意味着仅仅利用标准 C 语言库编写而成的 C 语言程序可以移植到支持 C 语言实现的任何计算机或操作系统上。

回顾历史，过去的 ANSI C 通常指 C89，时至今日，这一用法还时有所见。GCC 就是一例，其限定符-ansi 意指“支持所有 ISO C90 程序”。然而，本书会避免这种用法，因为如今该术语的含义有些含糊不清。自从 ANSI 委员会批准了 C99 修订版之后，确切说来，现在的 ANSI C 应该是 C99。

1.3.2 首个 POSIX 标准

术语“POSIX（可移植操作系统 Portable Operating System Interface 的缩写）”是指在 IEEE（电器及电子工程师协会），确切地说是其下属的可移植应用标准委员会（PASC, <http://www.pasc.org/>）赞助下所开发的一系列标准。PASC 标准的目标是提升应用程序在源码级别的可移植性。

POSIX 之名来自于 Richard Stallman 的建议。最后一个字母之所以是“X”是因为大多数 UNIX 变体之名总以“X”结尾。该标准特别注明，POSIX 应发音为“pahz-icks”，类似于“positive”。

本书会关注名为 POSIX.1 的第一个 POSIX 标准，以及后续的 POSIX.2 标准。

POSIX.1 和 POSIX.2

POSIX.1 于 1989 年成为 IEEE 标准，并在稍作修订后于 1990 年被正式采纳为 ISO 标准（ISO/IEC 9945-1:1990）。无法在线获得这一 POSIX 标准，但能从 IEEE（<http://www.ieee.org/>）购得。

POSIX.1 一开始是基于一个更早期的（1984 年）非官方标准，由名为/usr/group 的 UNIX 厂商协会制定。

符合 POSIX.1 标准的操作系统应向程序提供调用各项服务的 API，POSIX.1 文档对此作了规范。凡是提供了上述 API 的操作系统都可被认定为符合 POSIX.1 标准。

POSIX.1 基于 UNIX 系统调用和 C 语言库函数，但无需与任何特殊实现相关。这意味着任何操作系统都可以实现该接口，而不一定要是 UNIX 操作系统。实际上，在不对底层操作系统大加改动的同时，一些厂商通过添加 API 已经使自己的专有操作系统符合了 POSIX.1 标准。

对原有 POSIX.1 标准的若干扩展也同样重要。正式获批于 1993 年的 IEEE POSIX 1003.1b（POSIX.1b，原名 POSIX.4 或 POSIX 1003.4）包含了一系列对基本 POSIX 标准的实时性扩展。正式获批于 1995 年的 IEEE POSIX 1003.1c（POSIX.1c）对 POSIX 线程作了定义。1996 年，一个经过修订 POSIX.1 版本诞生，在核心内容保持不变的同时，并入了实时性和线程扩展。IEEE POSIX 1003.1g（POSIX.1g）定义了包括套接字在内的网络 API。分别获批于 1999 年和 2000 年的 IEEE POSIX 1003.1d（POSIX.1d）和 POSIX.1j 在 POSIX 基本标准的基础上，定义了附加的实时性扩展。

POSIX.1b 实时性扩展包括文件同步、异步 I/O、进程调度、高精度时钟和定时器、采用信号量、共享内存，以及消息队列的进程间通信。这 3 种进程间通信方法的称谓前通常冠以 POSIX，以示其有别于与之类似而又较为古老的 System V 信号、共享内存以及消息队列。

POSIX.2 (1992, ISO/IEC 9945-2:1993) 这一与 POSIX.1 相关的标准，对 shell 和包括 C 编译器命令行接口在内的各种 UNIX 工具进行了标准化。

F151-1 和 FIPS 151-2

FIPS 是 Federal Information Processing Standard (联邦信息处理标准) 的缩写，这套标准由美国政府为规范其对计算机系统的采购而制定。FIPS 151-1 于 1989 年发布。这份标准基于 1988 年的 IEEE POSIX.1 标准和 ANSI C 语言标准草案。FIPS 151-1 和 POSIX.1 (1988) 之间的主要差别在于：某些对后者来说是可选的特性，对于前者来说是必须的。由于美国政府是计算机系统的“大买家”，大多数计算机厂商都会确保其 UNIX 系统符合 FIPS 151-1 版本的 POSIX.1 规范。

FIPS 151-2 与 POSIX.1 的 1990 ISO 版保持一致，但在其他方面则保持不变。2000 年 2 月，已然过时的 FIPS 151-2 标准被废止。

1.3.3 X/Open 公司和 The Open Group

X/Open 公司是由多家国际计算机厂商所组成的联盟，致力于采纳和改进现有标准，以制定出一套全面而又一致的开放系统标准。该公司编纂的《X/Open 可移植性指南》是一套基于 POSIX 标准的可移植性指导丛书。这份指南的首个重要版本是 1989 年发布的第三号 (XPG3)，XPG4 随之于 1992 年发布。1994 年，X/Open 又对 XPG4 做了修订，从而诞生了 XPG4 版本 2，其中吸收了 1.3.7 节所述 AT&T System V 接口定义第三号中的重要内容。也将这一修订版称为 Spec 1170，而 1170 是指标准中所定义的接口（函数、头文件及命令）数量。

1993 年初，Novell 从 AT&T 收购了 UNIX 系统的相关业务，又在稍后放弃了这项业务，并将 UNIX 商标权转让给了 X/Open。（这一转让计划公布于 1993 年，但法律限制将这一转让推迟到 1994 年初。）随后，X/Open 又将 XPG4 版本 2 “重新包装”为 SUS (Single UNIX Specification)（有时，也叫 SUSv1）或称之为 UNIX95。其内容包括：XPG4 版本 2，X/Open Curses 规范第 4 号版本 2，以及 X/Open 联网服务 (XNS) 规范第 4 号。SUS 版本 2 (SUSv2, <http://www.unix.org/version2/online.html>) 发布于 1997 年，人们也将经过该规范认证的 UNIX 实现称为 UNIX 98。（有时，该规范也被称之为 XPG5。）

1996 年，X/Open 与开放软件基金会 (OSF) 合并，成立 The Open Group。如今，几乎每家与 UNIX 系统有关的公司或组织都是 The Open Group 的会员，该组织持续着对 API 标准的开发。

OSF 是 20 世纪 80 年代末 UNIX 纷争期间成立的两家厂商联盟之一。OSF 的主要成员包括 Digital、IBM、HP、Apollo、Bull、Nixdorf 和 Siemens。OSF 成立的主要目的是为了应对由 AT&T (UNIX 的发明者) 和 SUN 公司 (UNIX 工作站市场的领跑者) 结盟所带来的威胁。随之，AT&T、SUN 和其他公司结成了与 OSF 对抗的 UNIX International 联盟。

1.3.4 SUSv3 和 POSIX.1-2001

始于 1999 年，出于修订并加强 POSIX 标准和 SUS 规范的目的，IEEE、Open 集团以及 ISO/IEC 联合技术委员会共同成立了奥斯丁公共标准修订工作组 (CSRG, <http://www.opengroup.org/>)

austin/)。(该工作组的首次会议于 1998 年 9 月在德州奥斯丁召开,这也是奥斯丁工作组名称的由来。)2001 年 12 月,该工作组正式批准了 POSIX 1003.1-2001,有时简称为 POSIX.1-2001 (随后,又获批为 ISO 标准:ISO/IEC 9945:2002)。

POSIX 1003.1-2001 取代了 SUSv2、POSIX.1、POSIX.2 以及大批的早期 POSIX 标准。有时,人们也将该标准称为 Single Unix Specification 版本 3,本书在后续内容中将称其为 SUSv3。

SUSv3 基本规范约有 3700 页,分为以下 4 部分。

- 基本定义 (XBD),包含了定义、术语、概念以及对头文件内容的规范。总计提供了 84 个头文件的规范。
- 系统接口 (XSH),首先介绍了各种有用的背景信息。主要内容包含对各种函数(在特定的 UNIX 实现中,这些函数要么是作为系统调用,要么是作为库函数来实现的)的定义。总计包括了 1123 个系统接口。
- Shell 和实用工具 (XCU),明确定义了 shell 和各种 UNIX 命令的行为。总共定义了 160 个实用工具的行为。
- 基本原理 (XRAT),包括了与前三部分有关的描述性文字和原理说明。

此外,SUSv3 还包含了 X/Open CURSES 第 4 号版本 2 (XCURSES)规范,该规范针对 curses 屏幕处理 API 定义了 372 个函数和 3 个头文件。

在 SUSv3 中共计定义了 1742 个接口。与之形成鲜明对照的是,POSIX.1-1990 (连同 FIPS 151-2)定义了 199 个接口,POSIX.2-1992 定义了 130 个实用工具。

SUSv3 规范可在线获得,网址是 <http://www.unix.org/version3/online.html>。通过 SUSv3 认证的 UNIX 实现可被称为 UNIX 03。

自 SUSv3 获批以来,人们针对规范文本中所发现的问题进行了多次小规模修复和改进。因此而诞生的 1 号技术勘误表并入了 2003 年发布的 SUSv3 修订版,而 2 号技术勘误表的改进成果则并入了 SUSv3 2004 修订版。

符合 POSIX、XSI 规范和 XSI 扩展

回顾历史,SUS (和 XPG)标准顺应了相应 POSIX 标准,并被组织为 POSIX 的功能超集。除了对许多额外接口作出规范外,SUS 标准还将诸多被 POSIX 视为可选的接口和行为规范作为必备项。

对于身兼 IEEE 标准和 OPEN 集团技术标准的 POSIX 1003.1-2001 来说(如前所述,POSIX 1003.1-2001 是由早期的 POSIX 和 SUS 标准合并而成),上述区别的存在方式更显微妙。该文档定义了对规范的两级符合度。

- POSIX 规范符合度,就符合该规范的 UNIX 实现所必须提供的接口定义了基线。规范允许符合度达标的 UNIX 实现提供其他可选接口。
- XSI (X/Open 系统接口[X/Open System Interface])规范符合度,对 UNIX 实现来说,要想完全符合 XSI 规范,除了必须满足 POSIX 规范的所有规定之外,还要提供若干 POSIX 规范中的可选接口和行为。只有这一规范符合度达标,才能从 OPEN GROUP 获得 UNIX03 称号。

人们将 XSI 规范符合度达标所需的额外接口和行为统称为 XSI 扩展。这些扩展支持以下特性:线程、mmap()和 munmap()、dlopen API、资源限制、伪终端、System V IPC、syslog API、poll()以及登录记账。

后续各章所言及的“符合 SUSv3 规范”是指“符合 XSI 规范”。

由于 POSIX 和 SUSv3 目前由同一份文档描述，故而在文档的正文中，对于满足 SUSv3 符合度所需的额外接口和强制选项都以阴影和边注形式加以标明。

未定义和未明确定义的接口

有时，我们会称某些接口在 SUSv3 中“未定义”或“未明确定义”。

未定义的接口是指尽管偶尔会在背景和原理描述中提及，却根本未经正式标准定义过的接口。

未明确定义的接口是指标准虽然包括了该接口，但却未对其重要细节进行规范。（通常是由于现有接口的实现差异导致标准委员会成员无法达成一致意见。）

“未定义”或“未明确定义”的接口一经使用，在不同 UNIX 实现间移植应用就很难得到保证。尽管如此，少数此类接口在不同实现下的表现又相当一致。针对这些接口，本书通常会在提及是时一一指出。

LEGACY（传统）特性

书中有时会指出 SUSv3 将某个特定特性标记为 LEGACY。这一术语意味着保留此特性在与老应用程序保持兼容，而在新应用程序中应避免使用。这也是标准对此特性的限制所在。大多数情况下，都能找到与 LEGACY 特性等效的其他 API。

1.3.5 SUSv4 和 POSIX.1-2008

2008 年，奥斯丁工作组完成了对已合并的 POSIX 和 SUS 规范的修订工作。较之于之先前版本，该标准包含了基本规范以及 XSI 扩展。人们将这一修订版本称为 SUSv4。

与 SUSv3 的变化相比，SUSv4 的变化范围不算太大。最显著的变化如下所示。

- SUSv4 为一系列函数添加了新规范。本书将会介绍以下新标准中定义的如下函数：`dirfd()`、`fdopendir()`、`fexecve()`、`futimens()`、`mkdtemp()`、`psignal()`、`strsignal()`以及 `utimensat()`。另一组与文件相关的函数，例如：`openat()`，参见 18.11 节，和现有函数，例如：`open()`，功能相同，其区别在于前者对相对路径的解释是相对于打开文件描述符的所属目录而言，而非相对于进程的当前工作目录。
- 某些在 SUSv3 中被定义为可选的函数在 SUSv4 中成为基本标准的必备部分。例如，某些原本在 SUSv3 中属于 XSI 扩展的函数，在 SUSv4 中转而隶属于基本标准。在 SUSv4 中转变为必备的函数中包括了 `dlopen` API (42.1 节)、实时信号 API (22.8 节)、POSIX 信号量 API (53 章) 以及 POSIX 定时器 API (23.6 节)。
- SUSv4 废止了 SUSv3 中的某些函数，这包括 `asctime()`、`ctime()`、`ftw()`、`gettimeofday()`、`getitimer()`、`setitimer()`以及 `siginterrupt()`。
- SUSv4 删除了在 SUSv3 中被标记为作废的一些函数，这包括 `gethostbyname()`、`gethostbyaddr()`以及 `vfork()`。
- SUSv4 对 SUSv3 现有规范的各方面细节进行了修改。例如，对于应满足异步信号安全 (async-signal-safe) 的函数列表，二者内容就有所不同 (见表 21-1)。

本书后文会就所论及的相关主题指出其在 SUSv4 中的变化。

1.3.6 UNIX 标准时间表

图 1-1 总结了上述各节所述及各种标准之间的关系，并按时间顺序对标准进行了排列。图

中的实线表示标准间的直接过渡，虚线则表示标准间有一定的瓜葛，这无非有两种情况：其一，一个标准被并入了另一标准；其二，一个标准依附于另一个标准。

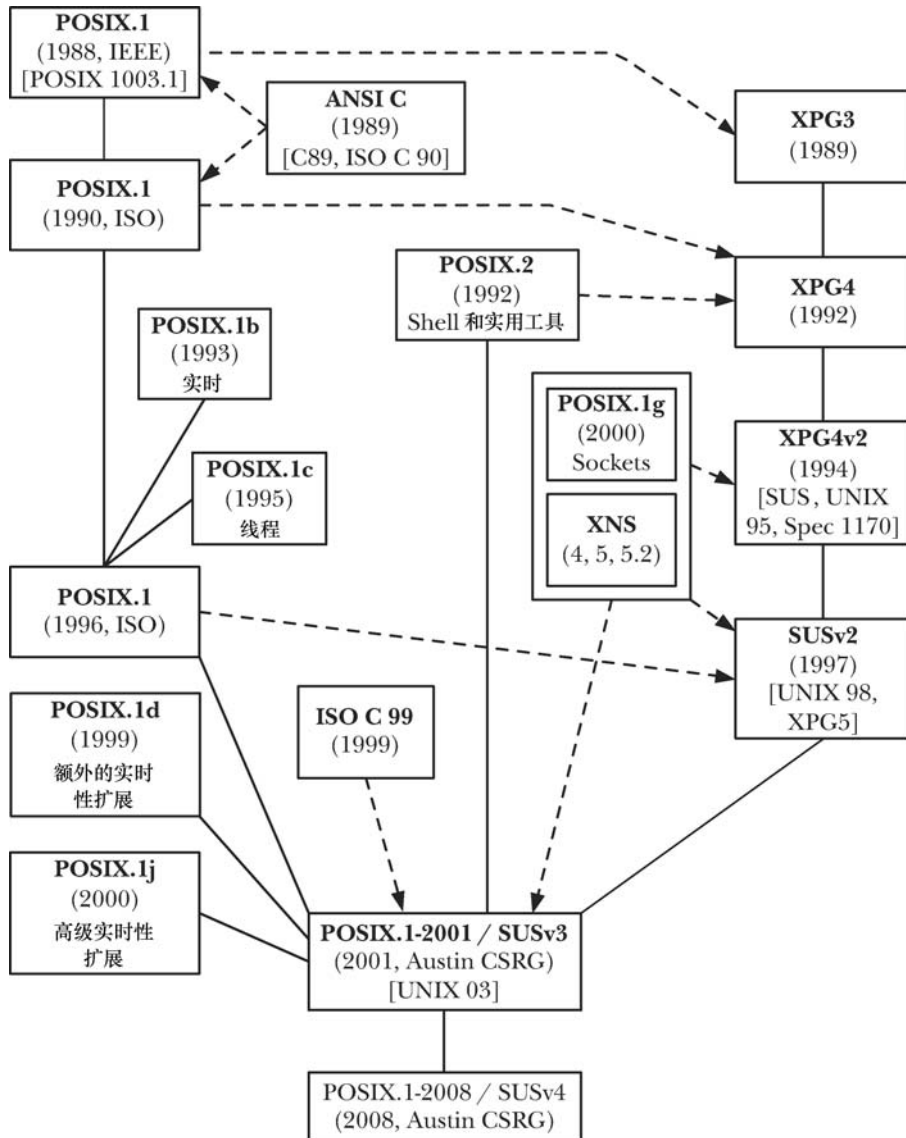


图 1-1: 各种 UNIX 和 C 标准之间的关系图

在网络标准方面，情况稍微有些复杂。该领域的标准化工作始于 20 世纪 80 年代末期，成立了 POSIX 1003.12 委员会，对套接字 API、XTI (X/Open 传输接口) API (另一套基于 System V 传输层接口的网络编程 API) 以及各种相关的 API 进行规范。该标准的酝酿历时数年，并于 2000 年获得了批准。其间，POSIX 1003.12 被更名为 POSIX 1003.1g。

在开发 POSIX 1003.1g 的同时，X/Open 也在开发自己的 X/Open 网络规范 (XNS)。该规范的第一版 XNS 第 4 号隶属于 SUS 首版。其后继版本为 XNS 第 5 号，隶属于 SUSv2。XNS 第 5 号与当时的 POSIX.1g 草案 (6.6) 基本相同。紧随其后的 XNS 第 5.2 号与 XNS 第 5 号以及获批为标准的 POSIX.1g 有所不同，将 XTI API 标记为作废，并纳入了于 20 世纪 90 年代中

期开发出的 IPv6。XNS 第 5.2 号构成了 SUSv3 中网络编程相关内容的基础，如今已被取代。出于类似原因，POSIX.1g 在获批后不久也退出了历史舞台。

1.3.7 实现标准

除了由独立或多边组织所制定的标准以外，有时，人们也会提到由 4.4BSD (BSD 的最终版) 和 SVR4 (AT&T 的 System V Release 4) 所定义的两类实现标准。后者随 AT&T 所发布的 SVID (System V 定义) 而正式出台。1989 年，AT&T 发布了 SVID 第 3 号，定义了自称为 System V Release 4 的 UNIX 实现所必须提供的接口。(从 [http://www.sco.com/ developers/devspecs/](http://www.sco.com/developers/devspecs/) 可以下载到 SVID。)

在 BSD 和 SVR4 之间，某些系统调用和库函数的行为各不相同，因此，许多 UNIX 实现都提供了兼容函数库和条件编译工具，可仿效并非特定 UNIX 实现“本色”的任意一种 UNIX 特性 (请参见 3.6.1 节)。这减轻了从另一 UNIX 实现移植应用程序的负担。

1.3.8 Linux、标准、Linux 标准规范 (Linux Standard Base)

遵守各种 UNIX 标准，尤其是符合 POSIX 和 SUS 规范，是 Linux (即内核、glibc 以及工具) 开发的总体目标。可是，在写作本书之际，尚无 Linux 发行版被 The Open group 授予“UNIX”商标。造成这一问题的主要原因不外乎是时间和费用。为了获得这一冠名，每个厂商的发行版都要经受规范符合度检查，每当有新的发行版诞生，还需重复执行上述检查。不过，正是由于 Linux 实际上几近于符合各种 UNIX 标准，才令其在 UNIX 市场上如此成功。

对于大多数商业 UNIX 实现来说，都是由同一家公司来开发和发布操作系统的。Linux 则有所不同，其实现与发行是分开的，多家组织——无论是商业性质还是非商业性质——都握有 Linux 的发行权。

Linus Torvalds 并不参与或支持任一特定 Linux 发行版的发行。然而，就参与 Linux 开发的其他人而言，情况更为复杂。许多从事 Linux 内核及其他自由软件项目开发的人员要么受雇于各家 Linux 发行商，要么就职于对 Linux 抱有浓厚兴趣的某些公司 (诸如 IBM 和 HP)。这些公司允许其程序员为特定 Linux 项目的开发投入一定的工作时间，这虽然对 Linux 的发展方向有所影响，但还没有哪家公司能够真正左右 Linux 的开发。更何况，很多参与 Linux 和 GUN 项目的其他开发者都是义工。

写作本书之际，Torvalds 受雇成为 Linux 基金会会员 (<http://www.linux-foundation.org>，之前的开源码发展实验室 OSDL)，该基金会是一家由多个商业和非商业组织组成的非赢利性联盟，旨在推动 Linux 的成长。

由于 Linux 的发行商众多，并且内核的开发者又无法控制 Linux 发布版的内容，因此还没有诞生“标准”的商业 Linux。一般情况下，每家 Linux 发行商所提供的内核都是基于某特定时间点发布的主要内核 (比如 Torvalds) 版本的快照，最多不过针对其打上几个补丁。发行商普遍认为，这些补丁所提供的特性可以在一定程度上迎合商业需求，从而能够提高市场竞争力。在某些情况下，主要内核版本稍后会打上这些补丁。实际上，某些新内核特性最初正是由某个 Linux 发行商开发而成，最终被纳入主要内核版本之前，这些新特性早已随着发行商的 Linux 发布版销售了。例如，被正式纳入主线 2.4 内核版本之前，版本 3 的 Reiserfs 日志文件服务器已经随着某些 Linux 发布版销售很长时间了。

上面的论述所要说明的就是由不同 Linux 发行公司提供的系统（往往）存在（细微的）差别。这使人在一定程度上不禁想起在 UNIX 发展之初，其实现方面所存在的各种差异。为了保证不同 Linux 发布版之间的兼容性，LSB 付出了不懈的努力。为了达成上述愿望，LSB (<http://www.linux-foundation.org/en/LSB>) 开发并推广了一套 Linux 系统标准，其主要目的是用来确保让二进制应用程序（即编译过的程序）能够在任何符合 LSB 规范的系统上运行。

由 LSB 所推广的二进制可移植性与 POSIX 所推广的源码可移植性可谓“一时瑜亮”。源码可移植性是指以 C 语言编写的程序可在任何符合 POSIX 规范的系统上编译并运行。而二进制可移植性则要苛刻得多，通常，只要硬件平台不一，便无法实现。二进制可移植性允许我们在某特定平台上将程序一次编译“成型”，然后，便可在任何符合 LSB 标准的 Linux 实现上运行该编译好的程序，当然，符合 LSB 标准的 Linux 实现必须运行在相同的硬件平台之上。对于在 Linux 上开发应用程序的独立软件开发商来说，二进制可移植性是其生存的基本前提。

1.4 总结

1969 年，贝尔实验室（AT&T 的一个部门）的 Ken Thompson 在 Digital PDP-7 小型机上首次实现了 UNIX 系统。对该操作系统而言，无论是理念还是其双关语的称谓都来源于早期的 MULTICS 系统。时至 1973 年，UNIX 已经被移植到了 PDP-11 小型机上，并以 C 语言对其进行了重写，C 编程语言是由贝尔实验室的 Dennis Ritchie 设计并实现的。因为法律禁止 AT&T 销售 UNIX，于是，在象征性地收取了一定的费用之后，AT&T 索性将 UNIX 系统散布进了大学。这其中便包括了源码，因为这一廉价操作系统的代码可供大学计算机系的师生研究和修改，故而这一操作系统在校园内广受欢迎。

在 UNIX 系统的开发方面，加州大学伯克利分校扮演了“关键先生”。在该校，Ken Thompson 及一干研究生又对这一操作系统进行了“精雕细琢”。到了 1979 年，这所大学发布了属于自己的 UNIX 发布版——BSD。这一发布版在学术界广为流传，并在日后成为某些商业 UNIX 实现的基石。

在此期间，随着 AT&T 不再对电信市场形成垄断，该公司被获准销售 UNIX。这也就催生出了另一种 UNIX 的变种——System V，日后，它也成为了某些商业 UNIX 实现的基石。

有两股不同的潮流引领着（GNU/）Linux 的开发。其中之一便是由 Richard Stallman 所创的 GNU 项目。20 世纪 80 年代末，GNU 项目已经开发出了一套几乎完备且可以自由分发的 UNIX 实现，但独缺一颗能够有效运作的内核。1991 年，Linus Torvalds 被 Minix 内核（由 Andrew Tanenbaum 编写）“灵魂附体”，于是便开发出了一颗能够在 Intel x86-32 架构上正常运作的内核。应 Torvalds 之邀，许多其他程序员也加入到了改进内核的行列中。随着时光的流逝，在一干程序员的不懈努力下，Linux 逐渐发展壮大，并被移植到了多种硬件架构之上。

20 世纪 80 年代末，UNIX 和 C 语言的实现“百花齐放”，所引发的可移植性问题迫使人们开展针对以上两者的标准化工作。1989 年，对 C 语言的标准化工作完成（C89 颁布），在 1999 年，对 C89 这一标准进行了修订（C99 颁布）。在操作系统接口方面，对其标准化的“第一次吃螃蟹”便催生出了 POSIX.1，1988 年和 1990 年，IEEE 和 ISO 先后将 POSIX.1 采纳为标准。20 世纪 90 年代，人们又开始酝酿一个囊括各版 SUS 在内的更为详尽的标准。2001 年，合二为一的 POSIX 1003.1-2001 和 SUSv3 标准颁布。该标准合并并扩展了先前的 POSIX 标准和

各版 SUS。2008 年，人们完成了对该标准的修订（改动幅度不算太大）工作，于是，合二为一的 POSIX 1003.1-2008 和 SUSv4 标准浮出水面。

与大多数商业 UNIX 实现不同，Linux 的开发与发行可谓“风马牛不相及”。因此，并无单一的“官方”Linux 发布版。各家 Linux 发行商所提供的只是当前稳定内核的快照，最多针对其打几个补丁。LSB 开发并推广了一套 Linux 系统标准，其主要目的是用来保证二进制应用程序（即编译过的过程）在不同 Linux 发布版之间的兼容性，以便编译过的应用程序能够运行在任何符合 LSB 规范的操作系统上，但前提是操作系统所运行的硬件平台必须相同。

进阶阅读

欲知更多有关 UNIX 历史及标准的信息，请参阅[Ritchie,1984]、[McKusick et al.,1996]、[McKusick & Neville-Neil, 2005]、[Libes & Ressler,1989]、[Garfinkel et al., 2003]、[Stevens & Rago, 2005]、[Stevens, 1999]、[Quartermann & Wilhelm, 1993]、[Goodheart & Cox, 1994]以及 [McKusick, 1999]。

[Salus, 1994]是一本详尽的 UNIX 编年史，本章开篇的许多内容均取自该书。[Salus, 2008]回顾了 Linux 和其他自由软件项目的简史。此外，与 UNIX 历史相关的许多细节都可以在 Ronda Hauben 所著的在线书籍 *History of UNIX* 中找到。在 <http://www.dei.isep.ipp.pt/~acc/docs/unix.html> 上，可下载到该书。在 <http://www.levenez.com/unix/> 上，刊载了一张非常详尽的、显示了各种 UNIX 实现版本变迁的时间表。

[Josey, 2004]概括了 UNIX 系统和 SUSv3 发展的历史，在指导读者如何使用 SUSv3 规范的同时，还提供了 SUSv3 所含接口的汇总表，除此之外，还给出了从 SUSv2 和 C89 升级到 SUSv3 和 C99 的迁移指南。

除了提供软件和文档之外，GNU Web 站点 (<http://www.gnu.org/>) 还刊载了许多与自由软件项目有关的哲学性文章。[Williams, 2002]是一本 Richard Stallman 的个人传记。

在[Torvalds & Diamond, 2001]中，Torvalds 提供了自己用作 Linux 开发的私人账户。

第 2 章

基本概念

本章旨在向 Linux 和 UNIX “生手” 们介绍一系列与 Linux 系统编程有关的概念。

2.1 操作系统的核心——内核

术语“操作系统”通常包含两种不同含义。

1. 指完整的软件包，这包括用来管理计算机资源的核心层软件，以及附带的所有标准软件工具，诸如命令行解释器、图形用户界面、文件操作工具和文本编辑器等。
2. 在更狭义的范围内，是指管理和分配计算机资源（即 CPU、RAM 和设备）的核心层软件。

术语“内核”通常是第二种含义，本书中的“操作系统”一词也是这层意思。

虽然在没有内核的情况下，计算机也能运行程序，但有了内核会极大简化其他程序的编写和使用，令程序员“功力”大进、游刃有余。这要归功于内核为管理计算机的有限资源所提供的软件层。

一般情况下，Linux 内核可执行文件采用 `/boot/vmlinuz` 或与之类似的路径名。而文件名的来历也颇有渊源。早期的 UNIX 实现称其内核为 UNIX。在后续实现了虚拟内存机制的 UNIX 系统中，其内核名称变更为 `vmunix`。对 Linux 来说，文件名称中的系统名需要调整，而以“z”替换“linux”末尾的“x”，意在表明内核是经过压缩的可执行文件。

内核的职责

内核所能执行的主要任务如下所示。

- 进程调度：计算机内均配备有一个或多个 CPU（中央处理单元），以执行程序指令。与其他 UNIX 系统一样，Linux 属于抢占式多任务操作系统。“多任务”意指多个进程（即运行中的程序）可同时驻留于内存，且每个进程都能获得对 CPU 的使用权。“抢占”则是指一组规则。这组规则控制着哪些进程获得对 CPU 的使用，以及每个进程能使用多长时间，这两者都由内核进程调度程序（而非进程本身）决定。

- 内存管理：以一二十年前的标准来看，如今计算机的内存容量可谓相当可观，但软件的规模也保持了相应地增长，故而物理内存（RAM）仍然属于有限资源，内核必须以公平、高效地方式在进程间共享这一资源。与大多数现代操作系统一样，Linux 也采用了虚拟内存管理机制（6.4 节），这项技术主要具有以下两方面的优势。
 - 进程与进程之间、进程与内核之间彼此隔离，因此一个进程无法读取或修改内核或其他进程的内存内容。
 - 只需将进程的一部分保持在内存中，这不但降低了每个进程对内存的需求量，而且还能在 RAM 中同时加载更多的进程。这也大幅提升了如下事件的发生概率，在任一时刻，CPU 都有至少一个进程可以执行，从而使得对 CPU 资源的利用更加充分。
- 提供了文件系统：内核在磁盘之上提供有文件系统，允许对文件执行创建、获取、更新以及删除等操作。
- 创建和终止进程：内核可将新程序载入内存，为其提供运行所需的资源（比如，CPU、内存以及对文件的访问等）。这样一个运行中的程序我们称之为“进程”。一旦进程执行完毕，内核还要确保释放其占用资源，以供后续程序重新使用。
- 对设备的访问：计算机外接设备（鼠标、键盘、磁盘和磁带驱动器等）可实现计算机与外部世界的通信，这一通信机制包括输入、输出或是两者兼而有之。内核既为程序访问设备提供了简化版的标准接口，同时还要仲裁多个进程对每一个设备的访问。
- 联网：内核以用户进程的名义收发网络消息（数据包）。该任务包括将网络数据包路由至目标系统。
- 提供系统调用应用编程接口（API）：进程可利用内核入口点（也称为系统调用）请求内核去执行各种任务。Linux 系统调用 API 是本书的主题。3.1 节会详细描述进程在执行系统调用时所经历的步骤。

除了上述特性外，一般而言，诸如 Linux 之类的多用户操作系统会为每个用户营造一种抽象：虚拟私有计算机（virtual private computer）。这就是说，每个用户都可以登录进入系统，独立操作，而与其他用户大致无干。例如，每个用户都有属于自己的磁盘存储空间（主目录）。再者，用户能够运行程序，而每一程序都能从 CPU 资源中“分得一杯羹”，运转于自有的虚拟地址空间中。而且这些程序还能独立访问设备，并通过网络传递信息。内核负责解决（多进程）访问硬件资源时可能引发的冲突，用户和进程对此则往往一无所知。

内核态和用户态

现代处理器架构一般允许 CPU 至少在两种不同状态下运行，即：用户态和核心态（有时也称之为监管态 supervisor mode）。执行硬件指令可使 CPU 在两种状态间来回切换。与之对应，可将虚拟内存区域划分（标记）为用户空间部分或内核空间部分。在用户态下运行时，CPU 只能访问被标记为用户空间的内存，试图访问属于内核空间的内存会引发硬件异常。当运行于核心态时，CPU 既能访问用户空间内存，也能访问内核空间内存。

仅当处理器在核心态运行时，才能执行某些特定操作。这样的例子包括：执行宕机（halt）指令去关闭系统，访问内存管理硬件，以及设备 I/O 操作的初始化等。实现者们利用这一硬件设计，将操作系统置于内核空间。这确保了用户进程既不能访问内核指令和数据结构，也无法执行不利于系统运行的操作。

以进程及内核视角检视系统

在完成诸多日常编程任务时，程序员们习惯于以面向进程（process-oriented）的思维方式来考虑编程问题。然而，在研究本书后续所涵盖的各种主题时，读者有必要转换视角，站在内核的角度上来看问题。为突显二者间的差异，本书接下来会分别从进程和内核视角来检视系统。

一个运行系统通常会有多个进程并行其中。对进程来说，许多事件的发生都无法预期。执行中的进程不清楚自己对 CPU 的占用何时“到期”，系统随之又会调度哪个进程来使用 CPU（以及以何种顺序来调度），也不知道自己何时会再次获得对 CPU 的使用。信号的传递和进程间通信事件的触发由内核统一协调，对进程而言，随时可能发生。诸如此类，进程都一无所知。进程不清楚自己在 RAM 中的位置。或者换种更通用的说法，进程内存空间的某块特定部分如今到底是驻留在内存中还是被保存在交换空间（磁盘空间中的保留区域，作为计算机 RAM 的补充）里，进程本身并不知晓。与之类似，进程也闹不清自己所访问的文件“居于”磁盘驱动器的何处，只是通过名称来引用文件而已。进程的运作方式堪称“与世隔绝”——进程间彼此不能直接通信。进程本身无法创建出新进程，哪怕“自行了断”都不行。最后还有一点，进程也不能与计算机外接的输入输出设备直接通信。

相形之下，内核则是运行系统的中枢所在，对于系统的一切无所不知、无所不能，为系统上所有进程的运行提供便利。由哪个进程来接掌对 CPU 的使用，何时“接任”，“任期”多久，都由内核说了算。在内核维护的数据结构中，包含了与所有正在运行的进程有关的信息。随着进程的创建、状态发生变化或者终结，内核会及时更新这些数据结构。内核所维护的底层数据结构可将程序使用的文件名转换为磁盘的物理位置。此外，每个进程的虚拟内存与计算机物理内存及磁盘交换区之间的映射关系，也在内核维护的数据结构之列。进程间的所有通信都要通过内核提供的通信机制来完成。响应进程发出的请求，内核会创建新的进程，终结现有进程。最后，由内核（特别是设备驱动程序）来执行与输入/输出设备之间的所有直接通信，按需与用户进程交互信息。

本书后续内容中会出现如下措辞，例如：“某进程可创建另一个进程”、“某进程可创建管道”、“某进程可将数据写入文件”，以及“调用 `exit()` 以终止某进程”。请务必牢记，以上所有动作都是由内核来居中“调停”，上面的说法不过是“某进程可以请求内核创建另一个进程”的缩略语，以此类推。

进阶阅读

涵盖操作系统概念和设计，尤其是对 UNIX 操作系统加以重点关注的现代教科书包括：[Tanenbaum, 2007]、[Tanenbaum & Woodhull, 2006] 以及 [Vahalia, 1996]，最后一本包含了与虚拟内存架构有关的详细内容。[Goodheart & Cox, 1994] 详细介绍了 System V Release 4。[Maxwell, 1999] 则是有选择性地针对 Linux 2.2.5 的部分内核源码进行了注释。[Lions, 1996] 对第六版 UNIX 源码进行了详尽阐释，并一直是研究 UNIX 操作系统内幕的入门级经典。[Bovet & Cesati, 2005] 描述了 Linux 2.6 内核的实现。

2.2 shell

shell 是一种具有特殊用途的程序，主要用于读取用户输入的命令，并执行相应的程序以响应命令。有时，人们也称之为命令解释器。

术语登录 shell（login shell）是指用户刚登录系统时，由系统创建，用以运行 shell 的进程。

尽管某些操作系统将命令解释器集成于内核中，而对 UNIX 系统而言，shell 只是一个用户进程。shell 的种类繁多，登入同一台计算机的不同用户同时可使用不同的 shell（就单个用户来说，情况也一样）。纵观 UNIX 历史，出现过以下几种重要的 shell。

- **Bourne shell (sh)**: 这款由 Steve Bourne 编写的 shell 历史最为悠久，且应用广泛，曾是第七版 UNIX 的标配 shell。Bourne shell 包含了在其他 shell 中常见的许多特性，I/O 重定向、管道、文件名生成（通配符）、变量、环境变量处理、命令替换、后台命令执行以及函数。对于所有问世于第七版 UNIX 之后的实现而言，除了可能提供有其他 shell 之外，都附带了 Bourne shell。
- **C shell (csh)**: 由 Bill Joy 于加州大学伯克利分校编写而成。其命名则源于该脚本语言的流控制语法与 C 语言有着许多相似之处。C shell 当时提供了若干极为实用的交互式特性，并不为 Bourne shell 所支持，这其中包括命令的历史记录、命令行编辑功能、任务控制和别名等。C shell 与 Bourne shell 并不兼容。尽管 C shell 曾是 BSD 系统标配的交互式 shell，但一般情况下，人们还是喜欢针对 Bourne shell 编写 shell 脚本（稍后介绍），以便其能够在所有 UNIX 实现上移植。
- **Korn shell (ksh)**: AT&T 贝尔实验室的 David Korn 编写了这款 shell，作为 Bourne shell 的“继任者”。在保持与 Bourne shell 兼容的同时，Korn shell 还吸收了那些与 C shell 相类似的交互式特性。
- **Bourne again shell (bash)**: 这款 shell 是 GNU 项目对 Bourne shell 的重新实现。Bash 提供了与 C shell 和 Korn shell 所类似的交互式特性。Brian Fox 和 Chet Ramey 是 bash 的主要作者。bash 或许是 Linux 上应用最为广泛的 shell 了。在 Linux 上，Bourne shell (sh) 其实正是由 bash 仿真提供的。

POSIX.2-1992 基于当时的 Korn shell 版本定义了一个 shell 标准。如今，Korn shell 和 bash 都符合 POSIX 规范，但两者都提供了大量对标准的扩展，其扩展之间存在许多差异。

设计 shell 的目的不仅仅是用于人机交互，对 shell 脚本（包含 shell 命令的文本文件）进行解释也是其用途之一。为实现这一目的，每款 shell 都内置有许多通常与编程语言相关的功能，其中包括变量、循环和条件语句、I/O 命令以及函数等。

尽管在语法方面有所差异，每款 shell 执行的任务都大致相同。除非指明是某款特定 shell 的操作，否则书中的“shell”都会按所描述的方式运作。本书绝大多数需要用到 shell 的示例都会使用 bash，若无其他说明，读者可假定这些示例也能以相同方式在其他类 Bourne 的 shell 上运行。

2.3 用户和组

系统会对每个用户的身份做唯一标识，用户可隶属于多个组。

用户

系统的每个用户都拥有唯一的登录名（用户名）和与之相对应的整数型用户 ID（UID）。系统密码文件/etc/passwd 为每个用户都定义有一行记录，除了上述两项信息外，该记录还包含如下信息。

- **组 ID**: 用户所属第一个组的整数型组 ID。
- **主目录**: 用户登录后所居于的初始目录。

- 登录 shell: 执行以解释用户命令的程序名称。

该记录还能以加密形式保存用户密码。然而，出于安全考虑，用户密码往往存储于单独的 shadow 密码文件中，仅供特权用户阅读。

组

出于管理目的，尤其是为了控制对文件和其他资源的访问，将多个用户分组是非常实用的做法。例如，某项目的开发团队人员需要共享同一组文件，就可以将他们编为同一组的成员。在早期的 UNIX 实现中，一个用户只能隶属于一个组。BSD 率先允许一个用户同时属于多个组，这一理念后来被其他 UNIX 实现纷纷效仿，并最终成为 POSIX.1-1990 标准。每个用户组都对应着系统组文件/etc/group 中的一行记录，该记录包含如下信息。

- 组名: (唯一的) 组名称。
- 组 ID (GID): 与组相关的整数型 ID。
- 用户列表: 隶属于该组的用户登录名列表 (通过密码文件记录的 group ID 字段未能标识出的该组其他成员，也在此列)，以逗号分隔。

超级用户

超级用户在系统中享有特权。超级用户账号的用户 ID 为 0，通常登录名为 root。在一般的 UNIX 系统上，超级用户凌驾于系统的权限检查之上。因此，无论对文件施以何种访问权限限制，超级用户都可以访问系统中的任何文件，也能发送信号干预系统运行的所有用户进程。系统管理员可以使用超级用户账号来执行各种系统管理任务。

2.4 单根目录层级、目录、链接及文件

内核维护着一套单根目录结构，以放置系统的所有文件。(这与微软 Windows 之类的操作系统形成了鲜明对照，Windows 系统的每个磁盘设备都有各自的目录层级。) 这一目录层级的根基就是名为“/”的根目录。所有的文件和目录都是根目录的“子孙”。图 1-2 所示为这种文件层级结构的示例。

文件类型

在文件系统内，会对文件类型进行标记，以表明其种类。其中一种用来表示普通数据文件，人们常称之为“普通文件”或“纯文本文件”，以示与其他种类的文件有所区别。其他文件类型包括设备、管道、套接字、目录以及符号链接。

术语“文件”常用来指代任意类型的文件，不仅仅指普通文件。

路径和链接

目录是一种特殊类型的文件，内容采用表格形式，数据项包括文件名以及对相应文件的引用。这一“文件名+引用”的组合被称为链接。每个文件都可以有多条链接，因而也可以有多个名称，在相同或不同的目录中出现。

目录可包含指向文件或其他目录的链接。路径间的链接建立起如图 2-1 所示的目录层级。

每个目录至少包含两条记录: .和..，前者是指向目录自身的链接，后者是指向其上级目录——父目录的链接。除根目录外，每个目录都有父目录。对于根目录而言，..是指向根目录自身的

链接（因此，`./`等于`/`）。

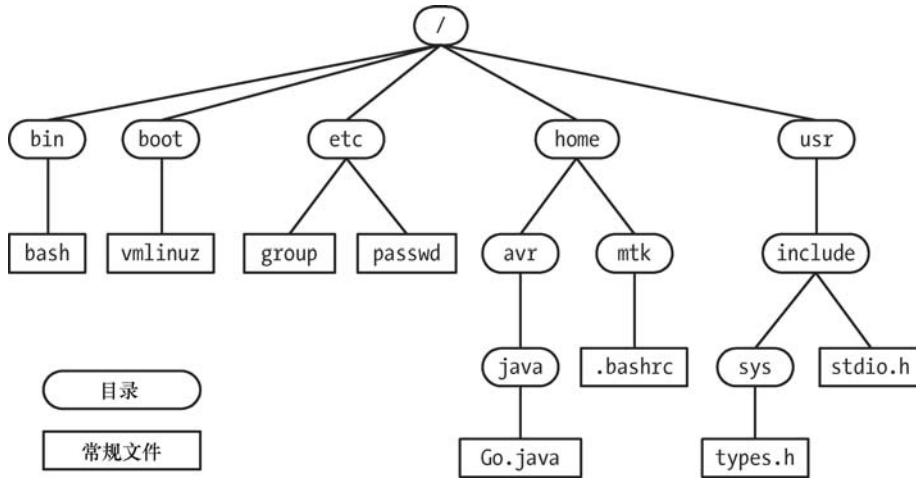


图 2-1: Linux 单根目录层级的一部分

符号链接

类似于普通链接，符号链接给文件起了一个“别号（alternative name）”。在目录列表中，普通链接是内容为“文件名+指针”的一条记录，而符号链接则是经过特殊标记的文件，内容包含了另一文件的名称。（换言之，一个符号链接对应着目录中内容为“文件名+指针”的一条记录，指针指向的文件内容¹为另一个文件名的字符串。）所谓“另一文件”通常被称为符号链接的目标，人们一般会说符号链接“指向”或“引用”目标文件。在多数情况下，只要系统调用用到了路径名，内核会自动解除（换言之，按照）该路径名中符号链接的引用，以符号链接所指向的文件名来替换符号链接。若符号链接的目标文件自身也是一个符号链接，那么上述过程会以递归方式重复下去。（为了应对可能出现的循环引用，内核对解除引用的次数作了限制。）如果符号链接指向的文件并不存在，那么可将该链接视为空链接（dangling link）。

通常，人们会分别使用硬链接（hard link）或软链接（soft link）这样的术语来指代正常链接和符号链接。之所以存在这两种不同类型的链接，将在第 18 章做出解释。

文件名

在大多数 Linux 文件系统上，文件名最长可达 255 个字符。文件名可以包含除“/”和空字符（\0）外的所有字符。但是，只建议使用字母、数字、点（“.”）、下划线（“_”）以及连字符（“-”）。SUSv3 将这 65 个字符的集合[-_a-zA-Z0-9]称为可移植文件名字符集（portable filename character set）。

对于可移植文件名字符集以外的字符，由于其可能会在 shell、正则表达式或其他场景中具有特殊含义，故而应避免在文件名中使用。如在上述环境中出现了包含特殊含义字符的文件名，则需要进行转义，即对此类字符进行特殊标记（一般会在特殊字符前插入一个“\”），以指明不应以特殊含义对其进行解释。若场景不支持转义机制，则不能使用此类文件名。

此外，还应避免以连字符（“-”）作为文件名的起始字符，因为一旦在 shell 命令中使用这种文件名，会被误认为命令行选项开关。

¹ 译者注：及该文件所属数据块存储的内容。

路径名

路径名是由一系列文件名组成的字符串，彼此以“/”分隔，首字符可以为“/”（非强制）¹。除却最后一个文件名外，该系列文件名均为目录名称（或为指向目录的符号链接）。路径名的尾部²可标识任意类型的文件，包括目录在内。有时将该字符串中最后一个“/”字符之前的部分称为路径名的目录部分，将其之后的部分称为路径名的文件部分或基础部分。

路径名应按从左至右的顺序阅读，路径名中每个文件名之前的部分，即为该文件所处目录。可在路径名中任意位置后引入字符串“.”³，用以指代路径名中当前位置的父目录。

路径名描述了单根目录层级下的文件位置，又可分为绝对路径名和相对路径名：

- **绝对路径名**以“/”开始，指明文件相对于根目录的位置。图 2-1 中的/home/mtk/.bashrc、/usr/include 以及/（根路径的路径名）都是绝对路径名的例子。
- **相对路径名**定义了相对于进程当前工作目录（见下文）的文件位置，与绝对路径名相比，相对路径名缺少了起始的“/”。如图 2-1 所示，在目录usr下，可使用相对路径名include/sys/types.h来引用文件types.h，在目录avr下，可使用相对路径名../mtk/.bashrc来访问文件.bashrc。

当前工作目录

每个进程都有一个当前工作目录（有时简称为进程工作目录或当前目录）。这就是单根目录层级下进程的“当前位置”，也是进程解释相对路径名的参照点。

进程的当前工作目录继承自其父进程。对登录shell来说，其初始当前工作目录，是依据密码文件中该用户记录的主目录字段来设置。可使用cd命令来改变shell的当前工作目录。

文件的所有权和权限

每个文件都有一个与之相关的用户ID和组ID，分别定义文件的属主和属组。系统根据文件的所有权来判定用户对文件的访问权限。

为了访问文件，系统把用户分为3类：文件的属主（有时，也称为文件的用户）、与文件组（group）ID相匹配的属组成员用户以及其他用户。可为以上3类用户分别设置3种权限（共计9种权限位）：只允许查看文件内容的读权限；允许修改文件内容的写权限；允许执行文件的执行权限。这里的文件要么指程序，要么是交由某种解释程序（通常指shell的一种，但也有例外）处理的脚本。

也可针对目录进行上述权限设置，但意义稍有不同。读权限允许列出目录内容（即该目录下的文件名），写权限允许对目录内容进行更改（比如，添加、修改或删除文件名），执行（有时也称为搜索）权限允许对目录中的文件进行访问（但需受文件自身访问权限的约束）。

2.5 文件 I/O 模型

UNIX 系统 I/O 模型最为显著的特性之一是其 I/O 通用性概念。也就是说，同一套系统调用（open()、read()、write()、close()等）所执行的 I/O 操作，可施之于所有文件类型，包括设

1 译者注：此处“文件”的含义中包括了目录——如前文所述：“目录是一种特殊类型的文件”。

2 译者注：即最后一个文件名。

3 译者注：需以“/”分隔。

备文件在内。（应用程序发起的 I/O 请求，内核会将其转化为相应的文件系统操作，或者设备驱动程序操作，以此来执行针对目标文件或设备的 I/O 操作。）因此，采用这些系统调用的程序能够处理任何类型的文件。

就本质而言，内核只提供一种文件类型：字节流序列，在处理磁盘文件、磁盘或磁带设备时，可通过 `lseek()` 系统调用来随机访问。

许多应用程序和函数库都将新行符（十进制 ASCII 码为 10，有时亦称其为换行）视为文本中一行的结束和另一行的开始。UNIX 系统没有文件结束符的概念，读取文件时如无数据返回，便会认定抵达文件末尾。

文件描述符

I/O 系统调用使用文件描述符——（往往是数值很小的）非负整数——来指代打开的文件。获取文件描述符的常用手法是调用 `open()`，在参数中指定 I/O 操作目标文件的路径名。

通常，由 shell 启动的进程会继承 3 个已打开的文件描述符：描述符 0 为标准输入，指代为进程提供输入的文件；描述符 1 为标准输出，指代供进程写入输出的文件；描述符 2 为标准错误，指代供进程写入错误消息或异常通告的文件。在交互式 shell 或程序中，上述三者一般都指向终端。在 `stdio` 函数库中，这几种描述符分别与文件流 `stdin`、`stdout` 和 `stderr` 相对应。

stdio 函数库

C 编程语言在执行文件 I/O 操作时，往往会调用 C 语言标准库的 I/O 函数。也将这样一组 I/O 函数称为 `stdio` 函数库，其中包括 `fopen()`、`fclose()`、`scanf()`、`printf()`、`fgets()`、`fputs()` 等。`stdio` 函数位于 I/O 系统调用层（`open()`、`close()`、`read()`、`write()` 等）之上。

本书假定读者已经了解了 C 语言的标准 I/O (`stdio`) 函数，因此也不会介绍这方面的内容。更多与 `stdio` 函数库有关的信息请参考 [Kernighan & Ritchie, 1988]、[Harbison & Steele, 2002]、[Plauger, 1992] 和 [Stevens & Rago, 2005]。

2.6 程序

程序通常以两种面目示人。其一为源码形式，由使用编程语言（比如，C 语言）写成的一系列语句组成，是人类可以阅读的文本文件。要想执行程序，则需将源码转换为第二种形式——计算机可以理解的二进制机器语言指令。（这与脚本形成了鲜明对照，脚本是包含命令的文本文件，可以由 shell 或其他命令解释器之类的程序直接处理。）一般认为，术语“程序”的上述两种含义几近相同，因为经过编译和链接处理，会将源码转换为语义相同的二进制机器码。

过滤器

从 `stdin` 读取输入，加以转换，再将转换后的数据输出到 `stdout`，常常将拥有上述行为的程序称为过滤器，`cat`、`grep`、`tr`、`sort`、`wc`、`sed`、`awk` 均在其列。

命令行参数

C 语言程序可以访问命令行参数，即程序运行时在命令行中输入的内容。要访问命令行参数，程序的 `main()` 函数需做如下声明：

```
int main(int argc, char *argv[])
```

`argc` 变量包含命令行参数的总个数，`argv` 指针数组的成员指针则逐一指向每个命令行参数字符串。首个字符串 `argv[0]`，标识程序名本身。

2.7 进程

简而言之，进程是正在执行的程序实例。执行程序时，内核会将程序代码载入虚拟内存，为程序变量分配空间，建立内核记账（bookkeeping）数据结构，以记录与进程有关的各种信息（比如，进程 ID、用户 ID、组 ID 以及终止状态等）。

在内核看来，进程是一个个实体，内核必须在它们之间共享各种计算机资源。对于像内存这样的受限资源来说，内核一开始会为进程分配一定数量的资源，并在进程的生命周期内，统筹该进程和整个系统对资源的需求，对这一分配进行调整。程序终止时，内核会释放所有此类资源，供其他进程重新使用。其他资源（如 CPU、网络带宽等）都属于可再生资源，但必须在所有进程间平等共享。

进程的内存布局

逻辑上将一个进程划分为以下几部分（也称为段）。

- 文本：程序的指令。
- 数据：程序使用的静态变量。
- 堆：程序可从该区域动态分配额外内存。
- 栈：随函数调用、返回而增减的一片内存，用于为局部变量和函数调用链接信息分配存储空间。

创建进程和执行程序

进程可使用系统调用 `fork()` 来创建一个新进程。调用 `fork()` 的进程被称为父进程，新创建的进程则被称为子进程。内核通过对父进程的复制来创建子进程。子进程从父进程处继承数据段、栈段以及堆段的副本后，可以修改这些内容，不会影响父进程的“原版”内容。（在内存中被标记为只读的程序文本段则由父、子进程共享。）

然后，子进程要么去执行与父进程共享代码段中的另一组不同函数，或者，更为常见的情况是使用系统调用 `execve()` 去加载并执行一个全新程序。`execve()` 会销毁现有的文本段、数据段、栈段及堆段，并根据新程序的代码，创建新段来替换它们。

以 `execve()` 为基础，C 语言库还提供了几个相关函数，接口虽然略有不同，但功能全都相同。以上所有库函数的名称均以字符串“`exec`”打头，在函数间差异无关宏旨的场合，本书会用符号 `exec()` 作为这些库函数的统称。不过，请读者牢记，实际上根本不存在名为 `exec()` 的库函数。

一般情况下，书中会使用“执行”一词来指代 `execve()` 及其衍生函数所实施的操作。

进程 ID 和父进程 ID

每一进程都有一个唯一的整数型进程标识符（PID）。此外，每一进程还具有一个父进程标识符（PPID）属性，用以标识请求内核创建自己的进程。

进程终止和终止状态

可使用以下两种方式之一来终止一个进程：其一，进程可使用 `_exit()` 系统调用（或相关的

`exit()`库函数), 请求退出; 其二, 向进程传递信号, 将其“杀死”。无论以何种方式退出, 进程都会生成“终止状态”, 一个非负小整数, 可供父进程的 `wait()`系统调用检测。在调用`_exit()`的情况下, 进程会指明自己的终止状态。若由信号来“杀死”进程, 则会根据导致进程“死亡”的信号类型来设置进程的终止状态。(有时会将传递给`_exit()`的参数称为进程的“退出状态”, 以示与终止状态有所不同, 后者要么指传递给`_exit()`的参数值, 要么表示“杀死”进程的信号。)

根据惯例, 终止状态为 0 表示进程“功成身退”, 非 0 则表示有错误发生。大多数 shell 会将前一执行程序的终止状态保存于 shell 变量`?`中。

进程的用户和组标识符 (凭证)

每个进程都有一组与之相关的用户 ID (UID)和组 ID (GID), 如下所示。

- **真实用户 ID 和组 ID:** 用来标识进程所属的用户和组。新进程从其父进程处继承这些 ID。登录 shell 则会从系统密码文件的相应字段中获取其真实用户 ID 和组 ID。
- **有效用户 ID 和组 ID:** 进程在访问受保护资源 (比如, 文件和进程间通信对象) 时, 会使用这两个 ID (并结合下述的补充组 ID) 来确定访问权限。一般情况下, 进程的有效 ID 与相应的真实 ID 值相同。正如即将讨论的那样, 改变进程的有效 ID 实为一种机制, 可使进程具有其他用户或组的权限。
- **补充组 ID:** 用来标识进程所属的额外组。新进程从其父进程处继承补充组 ID。登录 shell 则从系统组文件中获取其补充组 ID。

特权进程

在 UNIX 系统上, 就传统意义而言, 特权进程是指有效用户 ID 为 0 (超级用户) 的进程。通常由内核所施加的权限限制对此类进程无效。与之相反, 术语“无特权”(或非特权)进程是指由其他用户运行的进程。此类进程的有效用户 ID 为非 0 值, 且必须遵守由内核所强加的权限规则。

由某一特权进程创建的进程, 也可以是特权进程。例如, 一个由 `root` (超级用户) 发起的登录 shell。成为特权进程的另一方法是利用 `set-user-ID` 机制, 该机制允许某进程的有效用户 ID 等同于该进程所执行程序文件的用户 ID。

能力 (Capabilities)

始于内核 2.2, Linux 把传统上赋予超级用户的权限划分为一组相互独立的单元 (称之为“能力”)。每次特权操作都与特定的能力相关, 仅当进程具有特定能力时, 才能执行相应操作。传统意义上的超级用户进程 (有效用户 ID 为 0) 则相应开启了所有能力。

赋予某进程部分能力, 使得其既能够执行某些特权级操作, 又防止其执行其他特权级操作。

本书第 39 章会对能力做深入讨论。在本书后文中, 当谈及只能由特权进程执行的特殊操作时, 一般都会在括号中标明其具体能力。能力的命名以 `CAP_` 为前缀, 例如, `CAP_KILL`。

init 进程

系统引导时, 内核会创建一个名为 `init` 的特殊进程, 即“所有进程之父”, 该进程的相应程序文件为 `/sbin/init`。系统的所有进程不是由 `init` (使用 `fork()`) “亲自”创建, 就是由其后代进程创建。`init` 进程的进程号总为 1, 且总是以超级用户权限运行。谁 (哪怕是超级用户) 都不能“杀死” `init` 进程, 只有关闭系统才能终止该进程。`init` 的主要任务是创建并监控系统运行所需的一系列进程。(手册页 `init(8)` 中包含了 `init` 进程的详细信息。)

守护进程

守护进程指的是具有特殊用途的进程，系统创建和处理此类进程的方式与其他进程相同，但以下特征是其所独有的：

- “长生不老”。守护进程通常在系统引导时启动，直至系统关闭前，会一直“健在”。
- 守护进程在后台运行，且无控制终端供其读取或写入数据。

守护进程中的例子有 `syslogd`（在系统日志中记录消息）和 `httpd`（利用 HTTP 分发 Web 页面）。

环境列表

每个进程都有一份环境列表，即在进程用户空间内存中维护的一组环境变量。这份列表的每一元素都由一个名称及其相关值组成。由 `fork()` 创建的新进程，会继承父进程的环境副本。这也为父子进程间通信提供了一种机制。当进程调用 `exec()` 替换当前正在运行的程序时，新程序要么继承老程序的环境，要么在 `exec()` 调用的参数中指定新环境并加以接收。

在绝大多数 shell 中，可使用 `export` 命令来创建环境变量（C shell 使用 `setenv` 命令），如下所示：
`$ export MYVAR='Hello world'`

本书在展示交互式输入、输出的 shell 会话日志时，总是以黑体字来呈现输入文本。有时也会在日志中以斜体字形式加注，以解释输入的命令和产生的输出。

C 语言程序可使用外部变量（`char **environ`）来访问环境，而库函数也允许进程去获取或修改自己环境中的值。

环境变量的用途多种多样。例如，shell 定义并使用了一系列变量，供 shell 执行的脚本和程序访问。其中包括：变量 `HOME`（明确定义了用户登录目录的路径名）、变量 `PATH`（指明了用户输入命令后，shell 查找与之相应程序时所搜索的目录列表）。

资源限制

每个进程都会消耗诸如打开文件、内存以及 CPU 时间之类的资源。使用系统调用 `setrlimit()`，进程可为自己消耗的各类资源设定一个上限。此类资源限制的每一项均有两个相关值：软限制（`soft limit`）限制了进程可以消耗的资源总量，硬限制（`hard limit`）软限制的调整上限。非特权进程在针对特定资源调整软限制值时，可将其设置为 0 到相应硬限制值之间的任意值，但硬限制值则只能调低，不能调高。

由 `fork()` 创建的新进程，会继承其父进程对资源限制的设置。

使用 `ulimit` 命令（在 C shell 中为 `limit`）可调整 shell 的资源限制。shell 为执行命令所创建的子进程会继承上述资源设置。

2.8 内存映射

调用系统函数 `mmap()` 的进程，会在其虚拟地址空间中创建一个新的内存映射。

映射分为两类。

- 文件映射：将文件的部分区域映射入调用进程的虚拟内存。映射一旦完成，对文件映射内容的访问则转化为对相应内存区域的字节操作。映射页面会按需自动从文件中加载。
- 相映成趣的是并无文件与之相对应的匿名映射，其映射页面的内容会被初始化为 0。

由某一进程所映射的内存可以与其他进程的映射共享。达成共享的方式有二：其一是两个进程都针对某一文件的相同部分加以映射，其二是由 `fork()` 创建的子进程自父进程处继承映射。当两个或多个进程共享的页面相同时，进程之一对页面内容的改动是否对其他进程所见呢？这取决于创建映射时所传入的标志参数。若传入标志为私有，则某进程对映射内容的修改对于其他进程是不可见的，而且这些改动也不会真地落实到文件上；若传入标志为共享，对映射内容的修改就会为其他进程所见，并且这些修改也会造成对文件的改动。内存映射用途很多，其中包括：以可执行文件的相应段来初始化进程的文本段、内存（内容填充为 0）分配、文件 I/O（即映射内存 I/O）以及进程间通信（通过共享映射）。

2.9 静态库和共享库

所谓目标库是这样一种文件：将（通常是逻辑相关的）一组函数代码加以编译，并置于一个文件中，供其他应用程序调用。这一做法有利于程序的开发和维护。现代 UNIX 系统提供两种类型的对象库：静态库和共享库。

静态库

静态库（有时，也称之为档案文件[archives]）是早期 UNIX 系统中唯一的一种目标库。本质上说来，静态库是对已编译目标模块的一种结构化整合。要使用静态库中的函数，需要在创建程序的链接命令中指定相应的库。主程序会对静态库中隶属于各目标模块的不同函数加以引用。链接器在解析了引用情况后，会从库中抽取所需目标模块的副本，将其复制到最终的可执行文件中，这就是所谓静态链接。对于所需库内的各目标模块，采用静态链接方式生成的程序都存有一份副本。这会引起诸多不便。其一，在不同的可执行文件中，可能都存有相同目标代码的副本，这是对磁盘空间的浪费。同理，调用同一库函数的程序，若均以静态链接方式生成，且又于同时加以执行，这会造成内存浪费，因为每个程序所调用的函数都各有一份副本驻留在内存中，此其二。此外，如果对库函数进行了修改，需要重新加以编译、生成新的静态库，而所有需要调用该函数“更新版”的应用，都必须与新生成的静态库重新链接。

共享库

设计共享库的目的是为了解决静态库所存在的问题。

如果将程序链接到共享库，那么链接器就不会把库中的目标模块复制到可执行文件中，而是在可执行文件中写入一条记录，以表明可执行文件在运行时需要使用该共享库。一旦在运行时将可执行文件载入内存，一款名为“动态链接器”的程序会确保将可执行文件所需的动态库找到，并载入内存，随后实施运行时链接，解析可执行文件中的函数调用，将其与共享库中相应的函数定义关联起来。在运行时，共享库代码在内存中只需保留一份，且可供所有运行中的程序使用。

经过编译处理的函数仅在共享库内保存一份，从而节约了磁盘空间。另外，这一设计还能确保各类程序及时使用到函数的最新版本，功莫大焉，只需将带有函数新定义体的共享库重新加以编译即可，程序会在下次执行时自动使用新函数。

2.10 进程间通信及同步

Linux 系统上运行有多个进程，其中许多都是独立运行。然而，有些进程必须相互合作以

达成预期目的，因此彼此间需要通信和同步机制。

读写磁盘文件中的信息是进程间通信的方法之一。可是，对许多程序来说，这种方法既慢又缺乏灵活性。因此，像所有现代 UNIX 实现那样，Linux 也提供了丰富的进程间通信（IPC）机制，如下所示。

- 信号（signal），用来表示事件的发生。
- 管道（亦即 shell 用户所熟悉的“|”操作符）和 FIFO，用于在进程间传递数据。
- 套接字，供同一台主机或是联网的不同主机上所运行的进程之间传递数据。
- 文件锁定，为防止其他进程读取或更新文件内容，允许某进程对文件的部分区域加以锁定。
- 消息队列，用于在进程间交换消息（数据包）。
- 信号量（semaphore），用来同步进程动作。
- 共享内存，允许两个及两个以上进程共享一块内存。当某进程改变了共享内存的内容时，其他所有进程会立即了解到这一变化。

UNIX 系统的 IPC 机制种类如此繁多，有些功能还互有重叠，部分原因是由于各种 IPC 机制是在不同的 UNIX 实现上演变而来的，需要遵循的标准也各不相同。例如，就本质而言，FIFO 和 UNIX 套接字功能相同，允许同一系统上并无关联的进程彼此交换数据。二者之所以并存于现代 UNIX 系统之中，是由于 FIFO 来自 System V，而套接字则源于 BSD。

2.11 信号

尽管上一节将信号视为 IPC 的方法之一，但其在其他方面的广泛应用则更为普遍，因此值得深入讨论。

人们往往将信号称为“软件中断”。进程收到信号，就意味着某一事件或异常情况的发生。信号的类型很多，每一种分别标识不同的事件或情况。采用不同的整数来标识各种信号类型，并以 SIGxxxx 形式的符号名加以定义。

内核、其他进程（只要具有相应的权限）或进程自身均可向进程发送信号。例如，发生下列情况之一时，内核可向进程发送信号。

- 用户键入中断字符（通常为 Control-C）。
- 进程的子进程之一已经终止。
- 由进程设定的定时器（告警时钟）已经到期。
- 进程尝试访问无效的内存地址。

在 shell 中，可使用 kill 命令向进程发送信号。在程序内部，系统调用 kill() 可提供相同的功能。

收到信号时，进程会根据信号采取如下动作之一。

- 忽略信号。
- 被信号“杀死”。
- 先挂起，之后再被专用信号唤醒。

就大多数信号类型而言，程序可选择不采取默认的信号动作，而是忽略信号（当信号的默认处理行为并非忽略此信号时，会派上用场）或者建立自己的信号处理器。信号处理器是由程序员定义的函数，会在进程收到信号时自动调用，根据信号的产生条件执行相应动作。

信号从产生直至送达进程期间，一直处于挂起状态。通常，系统会在接收进程下次获得

调度时，将处于挂起状态的信号同时送达。如果接收进程正在运行，则会立即将信号送达。然而，程序可以将信号纳入所谓“信号屏蔽”¹以求阻塞该信号。如果产生的信号处于“信号屏蔽”之列，那么此信号将一直保持挂起状态，直至解除对该信号的阻塞。（亦即从信号屏蔽中移除。）

2.12 线程

在现代 UNIX 实现中，每个进程都可执行多个线程。可将线程想象为共享同一虚拟内存及若干其他属性的进程。每个线程都会执行相同的程序代码，共享同一数据区域和堆。可是，每个线程都拥有属于自己的栈，用来装载本地变量和函数调用链接信息。

线程之间可通过共享的全局变量进行通信。借助于线程 API 所提供的条件变量和互斥机制，进程所属的线程之间得以相互通信并同步行为——尤其是在对共享变量的使用方面。此外，利用 2.10 节所述的 IPC 和同步机制，线程间也能彼此通信。

线程的主要优点在于协同线程之间的数据共享（通过全局变量）更为容易，而且就某些算法而论，以多线程来实现比之以多进程实现要更加自然。再者，显而易见，多线程应用能从多处理器硬件的并行处理中获益匪浅。

2.13 进程组和 shell 任务控制

shell 执行的每个程序都会在一个新进程内发起。比如，shell 创建了 3 个进程来执行以下管道命令（在当前的工作目录下，根据文件大小对文件进行排序并显示）：

```
$ ls -l | sort -k5n | less
```

除 Bourne shell 以外，几乎所有的主流 shell 都提供了一种交互式特性，名为任务控制。该特性允许用户同时执行并操纵多条命令或管道。在支持任务控制的 shell 中，会将管道内的所有进程置于一个新进程组或任务中。（如果情况很简单，shell 命令行只包含一条命令，那么就会创建一个只包含单个进程的新进程组。）进程组中的每个进程都具有相同的进程组标识符（以整数形式），其实就是进程组中某个进程（也称为进程组组长 *process group leader*）的进程 ID。

内核可对进程组中的所有成员执行各种动作，尤其是信号的传递。如下节所述，支持任务控制的 shell 会利用这一特性，以挂起或恢复执行管道中的所有进程。

2.14 会话、控制终端和控制进程

会话指的是一组进程组（任务）。会话中的所有进程都具有相同的会话标识符。会话首进程（*session leader*）是指创建会话的进程，其进程 ID 会成为会话 ID。

使用会话最多的是支持任务控制的 shell，由 shell 创建的所有进程组与 shell 自身隶属于同一会话，shell 是此会话的会话首进程。

通常，会话都会与某个控制终端相关。控制终端建立于会话首进程初次打开终端设备之时。对于由交互式 shell 所创建的会话，这恰恰是用户的登录终端。一个终端至多只能成为一个会话的控制终端。

¹ 译者注：即一组进程希望阻塞的信号。

打开控制终端会致使会话首进程成为终端的控制进程。一旦断开了与终端的连接（比如，关闭了终端窗口），控制进程将会收到 `SIGHUP` 信号。

在任一时点，会话中总有一个前台进程组（前台任务），可以从终端中读取输入，向终端发送输出。如果用户在控制终端中输入了“中断”（通常是 `Control-C`）或“挂起”字符（通常是 `Control-Z`），那么终端驱动程序会发送信号以终止或挂起（亦即停止）前台进程组。一个会话可以拥有任意数量的后台进程组（后台任务），由以“&”字符结尾的行命令来创建。

支持任务控制的 `shell` 提供如下命令：列出所有任务，向任务发送信号，以及在前后台任务之间来回切换。

2.15 伪终端

伪终端是一对相互连接的虚拟设备，也称为主从设备。在这对设备之间，设有一条 `IPC` 信道，可供数据进行双向传递。

从设备（`slave device`）所提供的接口，其行为方式与终端相类似，基于这一特点，可以将某个为终端编写的程序与从设备连接起来，然后，再利用连接到主设备的另一程序来驱动这一“面向终端”的程序，这是伪终端的一个关键用途。由“驱动程序”¹所产生的输出，在经由终端驱动程序的常规输入处理（例如，默认情况下，会把回车符映射为换行符）后，会作为输入传递给与从设备相连的面向终端的程序。而由面向终端的程序向从设备写入的任何数据又作为“驱动程序”的输入来传递（在执行完所有常规的终端输入处理后）。换句话说，“驱动程序”所履行的功能，在效果上等同于用户通常在传统终端上所执行的操作。

伪终端广泛应用于各种应用领域，最知名的要数 `telnet` 和 `ssh` 之类提供网络登录服务的应用，以及 `X Window` 系统所提供的终端窗口实现。

2.16 日期和时间

进程涉及两种类型的时间。

- 真实时间：指的是在进程的生命期内（所经历的时间或时钟时间），以某个标准时间点（日历时间）或固定时间点（通常是进程的启动时间）为起点测量得出的时间。在 `UNIX` 系统上，日历时间是以国际协调时间（简称 `UTC`）1970 年 1 月 1 日凌晨为起始点，按秒测量得出的时间，再进行时区调整（定义时区的基准点为穿过英格兰格林威治的经线）²。这一日期与 `UNIX` 系统的生日很接近，也被称为纪元（`Epoch`）。
- 进程时间：亦称为 `CPU` 时间，指的是进程自启动起来，所占用的 `CPU` 时间总量。可进一步将 `CPU` 时间划分为系统 `CPU` 时间和用户 `CPU` 时间。前者是指在内核模式中，执行代码所花费的时间（比如，执行系统调用，或代表进程执行其他的内核服务）。后者是指在用户模式中，执行代码所花费的时间（比如，执行常规的程序代码）。

`time` 命令会显示出真实时间、系统 `CPU` 时间，以及为执行管道中的多个进程而花费的用户 `CPU` 时间。

¹ 译者注：此处专指与主设备相连的程序，而非设备驱动程序之类的含义。

² 译者注：即本初子午线。

2.17 客户端/服务器架构

本书有多处论及客户端/服务器应用程序的设计和实现。

客户端/服务器应用由两个组件进程组成。

- 客户端：向服务器发送请求消息，请求服务器执行某些服务。
- 服务器：分析客户端的请求，执行相应的动作，然后，向客户端回发响应消息。

有时，服务器与客户端之间可能需要就一次服务而进行多次交互。

客户端应用通常与用户打交道，而服务器应用则提供对某些共享资源的访问。一般说来，都是众多客户端进程与为数不多的一个或几个服务器端进程进行通信。

客户端和服务器既可以驻留于同一台计算机上，也可以位于联网的不同计算机上。客户端和服务器使用 2.10 节所讨论的 IPC 机制来实现彼此通信。

服务器可以提供各种服务，如下所示。

- 提供对数据库或其他共享信息资源的访问。
- 提供对远程文件的跨网访问。
- 对某些商业逻辑进行封装。
- 提供对共享硬件资源的访问（比如，打印机）。
- 提供 WWW 服务。

将某项服务封装于单独的服务器应用中，这一做法原因很多，举例如下。

- 效率：较之于在本地的每台计算上提供相同资源，在服务器应用管理之下提供资源的一份实例，则要节约许多。
- 控制、协调和安全：由于资源（尤其是信息资源）的统一存放，服务器既可以协调对资源的访问（例如，两个客户端不能同时更新同一信息），还可以保护资源安全，令其只对特定客户端开放。
- 在异构环境中运行：在网络中，客户端和服务器应用所运行的硬件平台和操作系统可以不同。

2.18 实时性

实时性应用程序是指那些需要对输入做出及时响应的程序。此类输入往往来自于外接的传感器或某些专门的输入设备，而输出则会去控制外接硬件。具有实时性需求的应用程序示例包括自动化装配流水线、银行 ATM 机，以及飞机导航系统等。

虽然许多实时性应用程序都要求对输入做出快速响应，但决定性因素却在于要在事件触发后的一定时限内，保证响应的交付。

要提供实时响应，特别是在短时间内加以响应，就需要底层操作系统的支持。由于实时响应的需求与多用户分时操作系统的需求存在冲突，大多数操作系统“天生”并不提供这样的支持。虽然已经设计出不少实时性的 UNIX 变体，但传统的 UNIX 实现都不是实时操作系统。Linux 的实时性变体也早已诞生，而近期的 Linux 内核正转向对实时性应用原生而全面的支持。

为支持实时性应用，POSIX.1b 定义了多个 POSIX.1 扩展，其中包括异步 I/O、共享内存、内存映射文件、内存锁定、实时性时钟和定时器、备选调度策略、实时性信号、消息队列，

以及信号量等。虽然这些扩展还不具备严格意义上的“实时性”，但当今的大多数 UNIX 实现都支持上面提到的全部或部分扩展（本书将讲解 Linux 所支持的 POSIX.1b 特性）。

本书会以术语“真实时间（real time）”来指代日历时间或经历时间的概念，而术语“实时性（realtime）”则是指操作系统或应用程序具备本节所述的响应能力。

2.19 /proc 文件系统

类似于其他的几种 UNIX 实现，Linux 也提供了 /proc 文件系统，由一组目录和文件组成，装配（mount）于 /proc 目录下。

/proc 文件系统是一种虚拟文件系统，以文件系统目录和文件形式，提供一个指向内核数据结构接口。这为查看和改变各种系统属性开启了方便之门。此外，还能通过一组以 /proc/PID 形式命名的目录（PID 即进程 ID）查看系统中运行各进程的相关信息。

通常，/proc 目录下的文件内容都采取人类可读的文本形式，shell 脚本也能对其进行解析。程序可以打开、读取和写入 /proc 目录下的既定文件。大多数情况下，只有特权级进程才能修改 /proc 目录下的文件内容。

本书在讲解各种 Linux 编程接口的同时，也会对相关的 /proc 文件进行介绍。12.1 节将就该文件系统的总体信息做进一步介绍。尚无任何标准对 /proc 文件系统进行过规范，书中与该文件系统相关的细节均为 Linux 专有。

2.20 总结

本章纵览了一系列与 Linux 系统编程相关的基本概念。对于 Linux 或 UNIX “生手”而言，理解这些基本概念将为学习系统编程提供足够的背景知识。

第 3 章

系统编程概念

了解本章所涵盖的各个主题是掌握系统编程的先决条件。本章首先会对系统调用加以介绍，并详述系统调用执行期间所发生的每个步骤。接下来，会讨论库函数及其与系统调用之间的差别，并结合这一差异，对（GNU）C 语言函数库进行描述。

无论何时，只要执行了系统调用或者库函数，检查调用的返回状态以确定调用是否成功，这是一条编程铁律。本章不但讲述了如何执行上述检查，还会介绍一组函数，本书刊载的编程示例大多都通过调用它们来诊断系统调用或库函数的错误。

本章的最后会探讨与可移植性编程相关的各种问题，尤其会关注对特性测试宏以及 SUSv3 中定义的标准系统数据类型的运用。

3.1 系统调用

系统调用是受控的内核入口，借助于这一机制，进程可以请求内核以自己的名义去执行某些动作。以应用程序编程接口（API）的形式，内核提供有一系列服务供程序访问。这包括创建新进程、执行 I/O，以及为进程间通信创建管道等。（手册页 `syscalls(2)` 列出了 Linux 系统调用。）

在深入系统调用的运作方式之前，务必关注以下几点。

- 系统调用将处理器从用户态切换到核心态，以便 CPU 访问受到保护的内核内存。
- 系统调用的组成是固定的，每个系统调用都由一个唯一的数字来标识。（程序通过名称来标识系统调用，对这一编号方案往往一无所知。）
- 每个系统调用可辅之以一套参数，对用户空间（亦即进程的虚拟地址空间）与内核空间之间（相互）传递的信息加以规范。

从编程角度来看，系统调用与 C 语言函数的调用很相似。然而，在执行系统调用时，其幕后会历经诸多步骤。为说明这点，下面以一个具体的硬件平台——x86-32 为例，按事件发生的顺序对这些步骤加以分析。

1. 应用程序通过调用 C 语言函数库中的外壳（wrapper）函数，来发起系统调用。
2. 对系统调用中断处理例程（稍后介绍）来说，外壳函数必须保证所有的系统调用参数可用。通过堆栈，这些参数传入外壳函数，但内核却希望将这些参数置入特定寄存器。因此，外

壳函数会将上述参数复制到寄存器。

3. 由于所有系统调用进入内核的方式相同，内核需要设法区分每个系统调用。为此，外壳函数会将系统调用编号复制到一个特殊的 CPU 寄存器（%eax）中。
4. 外壳函数执行一条中断机器指令（int 0x80），引发处理器从用户态切换到核心态，并执行系统中断 0x80（十进制数 128）的中断矢量所指向的代码。

较新的 x86-32 硬件平台实现了 `sysenter` 指令，较之传统的 `int 0x80` 中断指令，`sysenter` 指令进入内核的速度更快。2.6 内核及 `glibc 2.3.2` 以后的版本都支持 `sysenter` 指令。

5. 为响应中断 0x80，内核会调用 `system_call()` 例程（位于汇编文件 `arch/i386/entry.S` 中）来处理这次中断，具体如下。
 - a) 在内核栈中保存寄存器值（参见 6.5 节）。
 - b) 审核系统调用编号的有效性。
 - c) 以系统调用编号对存放所有调用服务例程的列表（内核变量 `sys_call_table`）进行索引，发现并调用相应的系统调用服务例程。若系统调用服务例程带有参数，那么将首先检查参数的有效性。例如，会检查地址指向用户空间的内存位置是否有效。随后，该服务例程会执行必要的任务，这可能涉及对特定参数中指定地址处的值进行修改，以及在用户内存和内核内存间传递数据（比如，在 I/O 操作中）。最后，该服务例程会将结果状态返回给 `system_call()` 例程。
 - d) 从内核栈中恢复各寄存器值，并将系统调用返回值置于栈中。
 - e) 返回至外壳函数，同时将处理器切换回用户态。
6. 若系统调用服务例程的返回值表明调用有误，外壳函数会使用该值来设置全局变量 `errno`（参见 3.4 节）。然后，外壳函数会返回到调用程序，并同时返回一个整型值，以表明系统调用是否成功。

在 Linux 上，系统调用服务例程遵循的惯例是调用成功则返回非负值。发生错误时，例程会对相应 `errno` 常量取反，返回一负值。C 语言函数库的外壳函数随即对其再次取反（负负得正），将结果拷贝至 `errno`，同时以 -1 作为外壳函数的返回值返回，向调用程序表明有错误发生。

上述惯例所依赖的前提条件是系统调用服务例程，若调用成功则不会返回负值。可是，对于少数例程来说，这一前提并不成立。一般情况下，这也不会有问题，因为取反的 `errno` 值范围不会与调用成功返回负值的范围有交集。不过，有一种情况，沿用这个惯例确实会出问题：系统调用 `fcntl()` 的 `F_GETOWN` 操作，会在 63.3 节加以描述。

图 3-1 以系统调用 `execve()` 为例，展示了上文述及事件的发生序列。在 Linux/x86-32 上，`execve()` 的系统调用号为 11（`_NR_execve`）。因此，在 `sys_call_table` 向量中，条目 11 包含了该系统调用的服务例程 `sys_execve()` 的地址。（在 Linux 中，系统调用服务例程的命名通常会采取 `sys_xyz()` 的形式，其中，`xyz()` 正是所论及的系统调用。）

若是单纯为了掌握本书的后续内容，这里的论述的确有些小题大做。但其要点在于即便对于一个简单的系统调用，仍要完成相当多的工作，因此系统调用的开销虽小，却也不容忽视。

可以以 `getppid()` 系统调用为例，研判一下发起系统调用的开销——该系统调用只是简单地返回调用进程的父进程 ID。在作者的一台运行 Linux 2.6.25 的 x86-32 系统上，调用 `getppid()` 一千万次大约需要 2.2 秒钟，每次调用大致需要 0.3 微秒。相形之下，在同一系统上，调用某个只返回整数的 C 语言函数一千万次，仅需 0.11 秒，约为调用 `getppid()` 耗费的时间的 1/20。当然，大多数系统调用的开销都明显高于 `getppid()`。

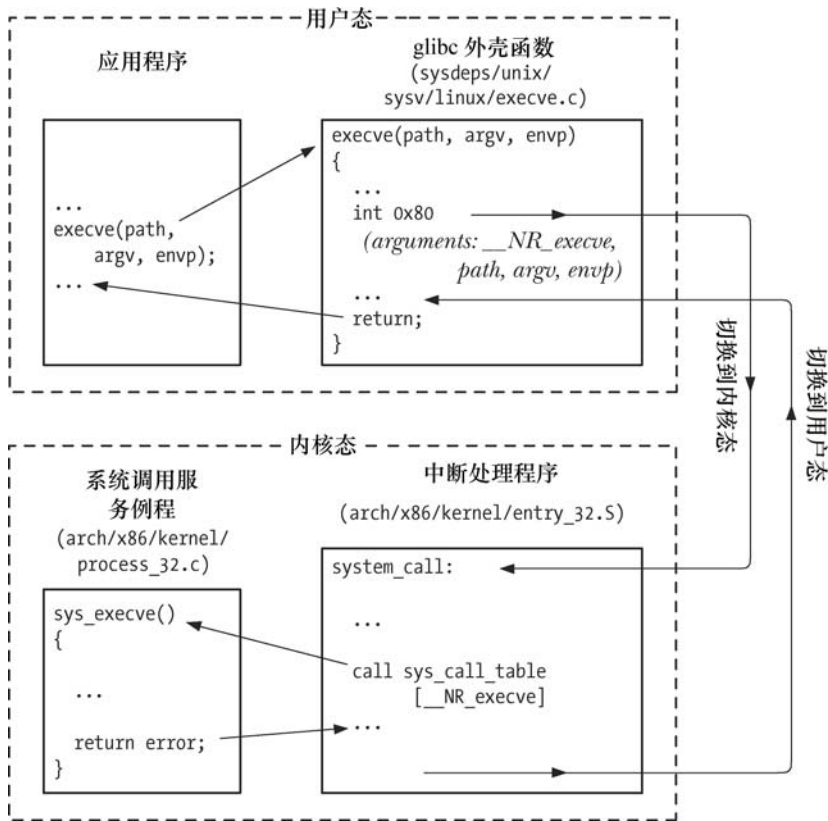


图 3-1: 系统调用的执行步骤

因此，从 C 语言编程的角度来看，调用 C 语言函数库的外壳（wrapper）函数等同于调用相应的系统调用服务例程，在本书后续内容中，“调用系统调用 xyz()”这类说法就意味着“调用外壳函数，由外壳函数去调用系统调用 xyz()”。

为调试程序，或是研究程序的运作机制，可使用附录 A 所介绍的 strace 命令，对程序发起的系统调用进行跟踪。

更多与 Linux 系统调用机制有关的信息请见[Love, 2010]、[Bovet & Cesati, 2005]以及 [Maxwell, 1999]。

3.2 库函数

一个库函数是构成标准 C 语言函数库的众多库函数之一。（出于简化，本书后文提到某具体函数时，通常将其称为“函数”而非“库函数”。）库函数的用途多种多样，可用来执行以下任务：打开文件、将时间转换为可读格式，以及进行字符串比较等。

许多库函数（比如，字符串操作函数）不会使用任何系统调用。另一方面，还有些库函数构建于系统调用层之上。例如，库函数 fopen()就利用系统调用 open()来执行打开文件的实际操作。往往，设计库函数是为了提供比底层系统调用更为方便的调用接口。例如，printf()函数可提供格式化输出和数据缓存功能，而 write()系统调用只能输出字节块。同理，与底层的 brk()系统调用相比，malloc()和 free()函数还执行了各种登记管理工作，内存的释放和分配也因此而容易许多。

3.3 标准 C 语言函数库；GNU C 语言函数库（glibc）

标准 C 语言函数库的实现随 UNIX 的实现而异。GNU C 语言函数库（glibc, <http://www.gnu.org/software/libc/>）是 Linux 上最常用的实现。

最初，Roland McGrath 是 GNU C 语言函数库的主要开发者和维护者。如今，Ulrich Drepper 挑起了这副重担。

Linux 同样支持各种其他 C 语言函数库，其中包括应用于嵌入式设备领域、受限内存条件下的 C 语言函数库。uClibc（<http://www.uclibc.org/>）和 dietlibc（<http://www.fefe.de/dietlibc/>）便是其中的两个例子。本书的讨论范围仅限于 glibc，因为为 Linux 开发的大多数应用程序都使用该函数库。

确定系统的 glibc 版本

有时，需要确定系统所安装的 glibc 版本。在 shell 中，可以直接运行 glibc 共享库文件——将其视为可执行文件——来获取 glibc 版本。这会输出各种文本信息，其中也包括了 glibc 的版本号。

```
$ /lib/libc.so.6
GNU C Library stable release version 2.10.1, by Roland McGrath et al.
Copyright (C) 2009 Free Software Foundation, Inc.
This is free software; see the source for copying conditions.
There is NO warranty; not even for MERCHANTABILITY or FITNESS FOR A
PARTICULAR PURPOSE.
Compiled by GNU CC version 4.4.0 20090506 (Red Hat 4.4.0-4).
Compiled on a Linux >>2.6.18-128.4.1.el5<< system on 2009-08-19.
Available extensions:
  The C stubs add-on version 2.1.2.
  crypt add-on version 2.1 by Michael Glad and others
  GNU Libidn by Simon Josefsson
  Native POSIX Threads Library by Ulrich Drepper et al
  BIND-8.2.3-T5B
  RT using linux kernel aio
For bug reporting instructions, please see:
<http://www.gnu.org/software/libc/bugs.html>.
```

在某些 Linux 发行版中，GNU C 语言函数库的路径名并非“/lib/libc.so.6”。确定该库所在位置的方法之一是：针对某个与 glibc 动态链接的可执行文件（大多数可执行文件都采用这种链接方式），运行 ldd（列出动态依赖性）程序。接下来，再检查输出的库依赖性列表，便能发现 glibc 共享库的位置：

```
$ ldd myprog | grep libc
libc.so.6 => /lib/tls/libc.so.6 (0x4004b000)
```

应用程序可通过测试常量和调用库函数这两种方法，来确定系统所安装的 glibc 版本。从版本 2.0 开始，glibc 定义了两个常量：__GLIBC__ 和 __GLIBC_MINOR__，供程序在编译时（在 #ifdef 语句中）测试使用。在安装有 glibc 2.12 版本的系统上，以上两个常量的值分别为 2 和 12。然而，如果程序在 A 系统上编译，而在 B 系统（安装了不同版本的 glibc）上运行，这两个常量作用就有限了。为应对这种可能，程序可以调用函数 gnu_get_libc_version()，来确定运行时的 glibc 版本。

```
#include <gnu/libc-version.h>

const char *gnu_get_libc_version(void);

Returns pointer to null-terminated, statically allocated string
containing GNU C library version number
```

函数 `gnu_get_libc_version()` 返回一个指针，指向诸如 “2.12” 的字符串。

获取 `glibc` 版本信息，还有一种方法：使用 `confstr()` 函数来获取（`glibc` 特有的）`_CS_GNU_LIBC_VERSION` 配置变量的值。这一调用会返回类似于 “`glibc 2.12`” 的字符串。

3.4 处理来自系统调用和库函数的错误

几乎每个系统调用和库函数都会返回某类状态值，用以表明调用成功与否。要了解调用是否成功，必须坚持对状态值进行检查。若调用失败，那么必须采取相应行动。至少，程序应该显示错误消息，警示有意想不到的事件发生。

不检查状态值，少敲几个字，听起来的确诱人（尤其是见识到了不检查状态值的 UNIX/Linux 程序以后），但实际却得不偿失。认定系统调用或库函数 “不可能失败”，不对状态返回值进行检查，这会浪费掉大把的程序调试时间。

少数几个系统函数在调用时从不失败。例如，`getpid()` 总能成功返回进程的 ID，而 `_exit()` 总能终止进程。无需对此类系统调用的返回值进行检查。

处理系统调用错误

每个系统调用的手册页记录有调用可能的返回值，并指出了哪些值表示错误。通常，返回值为 -1 表示出错。因此，可使用下列代码对系统调用进行检查：

```
fd = open(pathname, flags, mode);      /* system call to open a file */
if (fd == -1) {
    /* Code to handle the error */
}
...
if (close(fd) == -1) {
    /* Code to handle the error */
}
```

系统调用失败时，会将全局整形变量 `errno` 设置为一个正值，以标识具体的错误。程序应包含¹ `<errno.h>` 头文件，该文件提供了对 `errno` 的声明，以及一组针对各种错误编号而定义的常量。所有这些符号名都以字母 E 打头。在每个手册页内标题为 `ERRORS` 的章节内，都刊载有一份相应系统调用可能返回的 `errno` 值列表。以下便是利用 `errno` 诊断系统调用错误的一个简单示例：

```
cnt = read(fd, buf, numbytes);
if (cnt == -1) {
    if (errno == EINTR)
        fprintf(stderr, "read was interrupted by a signal\n");
    else {
        /* Some other error occurred */
    }
}
```

¹ 译者注：`#include`。

如果调用系统调用和库函数成功，`errno` 绝不会被重置为 0，故此，该变量值不为 0，可能是之前调用失败造成的。此外，SUSv3 允许在函数调用成功时，将 `errno` 设置为非零值（当然，几乎没有函数会这么做）。因此，在进行错误检查时，必须坚持首先检查函数的返回值是否表明调用出错，然后再检查 `errno` 确定错误原因。

少数系统调用（比如，`getpriority()`）在调用成功后，也会返回-1。要判断此类系统调用是否发生错误，应在调用前将 `errno` 置为 0，并在调用后对其进行检查（上述手法同样适用于某些库函数）。

系统调用失败后，常见的做法之一是根据 `errno` 值打印错误消息。提供库函数 `perror()` 和 `strerror()`，就是出于这一目的。

函数 `perror()` 会打印出其 `msg` 参数所指向的字符串，紧跟一条与当前 `errno` 值相对应的消息。

```
#include <stdio.h>

void perror(const char *msg);
```

以下是对系统调用错误进行处理的一种简单方式：

```
fd = open(pathname, flags, mode);
if (fd == -1) {
    perror("open");
    exit(EXIT_FAILURE);
}
```

函数 `strerror()` 会针对其 `errnum` 参数中所给定的错误号，返回相应的错误字符串。

```
#include <string.h>

char *strerror(int errnum);

Returns pointer to error string corresponding to errnum
```

由 `strerror()` 所返回的字符串可以是静态分配的，这意味着后续对 `strerror()` 的调用可能会覆盖该字符串。

若无法识别 `errnum` 所含的错误编号，则 `strerror()` 会返回 “Unknown error mnn.” 形式的字符串。在某些其他的实现中，在这种情况下，`strerror()` 会返回 NULL。

由于 `perror()` 和 `strerror()` 都属于对语言环境敏感（locale-sensitive）（参见 10.4 节）的函数，故而错误描述中使用的都是本地语言。

处理来自库函数的错误

不同的库函数在调用发生错误时，返回的数据类型和值也各不相同。（参见每个函数的手册页。）从错误处理的角度来说，可将库函数划分为以下几类。

- 某些库函数返回错误信息的方式与系统调用完全相同——返回值为-1，伴之以 `errno` 号来表示具体错误。`remove()` 便是其中一例，可使用该函数来删除文件（调用 `unlink()` 系统调用）或目录（调用 `rmdir()` 系统调用）。对此类函数所发生的错误进行诊断，其方式与系统调用完全相同。
- 某些库函数在出错时会返回-1 之外的其他值，但仍会设置 `errno` 来表明具体的出错情况。例如，`fopen()` 在出错时会返回一个 NULL 指针，还会根据出错的具体底层系统调

用来设置 `errno`。函数 `perror()` 和 `strerror()` 都用来诊断此类错误。

- 还有些函数根本不使用 `errno`。对此类函数来说，确定错误存在与否及其起因的方法各不相同，可见诸于相应函数的手册页中，不应使用 `errno`、`perror()` 或 `strerror()` 来诊断错误。

3.5 关于本书示例程序的注意事项

本节会就本书所载程序示例所普遍采用的各种惯例及特性加以介绍。

3.5.1 命令行选项及参数

本书所载的许多程序示例都会依照命令行选项及参数来决定其行为。

传统的 UNIX 命令行选项由一个连字符 (-)、表示选项的英文字母，以及一个可选参数组成。(GNU 实用工具则对选项语法有所扩展，以两个连字符开头 (--), 紧跟用来标识选项和可选参数的字符串。) 可使用标准库函数 `getopt()` (参见附录 B) 对命令行选项进行解析。

这些示例之中，但凡命令行语法颇为周正的，都为用户提供有一个简单的帮助工具：在以 `--help` 选项调用程序时，会显示用法信息，就命令行选项和参数的语法加以说明。

3.5.2 常用的函数及头文件

本书的大多数程序示例都包括有一个头文件，内含常用的各种定义。这些示例同样使用了一系列常用函数。本节会对这些头文件及函数进行讨论。

常用的头文件

程序清单 3-1 所列的头文件几乎为本书所有程序示例所使用。

程序清单 3-1：大多数程序示例所使用的头文件

```
----- lib/tlpi_hdr.h
#ifndef TLPI_HDR_H
#define TLPI_HDR_H      /* Prevent accidental double inclusion */

#include <sys/types.h> /* Type definitions used by many programs */
#include <stdio.h>     /* Standard I/O functions */
#include <stdlib.h>    /* Prototypes of commonly used library functions,
                       plus EXIT_SUCCESS and EXIT_FAILURE constants */
#include <unistd.h>    /* Prototypes for many system calls */
#include <errno.h>     /* Declares errno and defines error constants */
#include <string.h>    /* Commonly used string-handling functions */

#include "get_num.h"   /* Declares our functions for handling numeric
                       arguments (getInt(), getLong()) */

#include "error_functions.h" /* Declares our error-handling functions */

typedef enum { FALSE, TRUE } Boolean;

#define min(m,n) ((m) < (n) ? (m) : (n))
#define max(m,n) ((m) > (n) ? (m) : (n))

#endif
----- lib/tlpi_hdr.h
```

错误诊断函数

为简化本书程序示例中的错误处理，我们编写了错误诊断函数，对其的声明如程序清单 3-2 所示。

程序清单 3-2：常用错误处理函数的声明

```
----- lib/error_functions.h
#ifndef ERROR_FUNCTIONS_H
#define ERROR_FUNCTIONS_H

void errMsg(const char *format, ...);

#ifdef __GNUC__

/* This macro stops 'gcc -Wall' complaining that "control reaches
   end of non-void function" if we use the following functions to
   terminate main() or some other non-void function. */

#define NORETURN __attribute__((__noreturn__))
#else
#define NORETURN
#endif

void errExit(const char *format, ...) NORETURN ;

void err_exit(const char *format, ...) NORETURN ;

void errExitEN(int errnum, const char *format, ...) NORETURN ;

void fatal(const char *format, ...) NORETURN ;

void usageErr(const char *format, ...) NORETURN ;

void cmdLineErr(const char *format, ...) NORETURN ;

#endif
----- lib/error_functions.h
```

本书使用 `errMsg()`、`errExit()`、`err_exit()`以及 `errExitEN()`函数，以诊断调用系统调用和库函数时所发生的错误。

```
#include "tlpi_hdr.h"

void errMsg(const char *format, ...);
void errExit(const char *format, ...);
void err_exit(const char *format, ...);
void errExitEN(int errnum, const char *format, ...);
```

函数 `errMsg()`会在标准错误设备上打印消息。除了将一个终止换行符自动追加到输出字符串尾部以外，该函数的参数列表与 `printf()`所用相同。`errMsg()`函数会打印出与当前 `errno` 值相对应的错误文本，其中包括了错误名（比如，`EPERM`）以及由 `strerror()`返回的错误描述，外加由参数列表指定的格式化输出。

`errExit()`函数的操作方式与 `errMsg()`相似，只是还会以如下两种方式之一来终止程序。其一，调用 `exit()`退出。其二，若将环境变量 `EF_DUMPCORE` 定义为非空字符串，则调用 `abort()`退出，同时

生成核心转储 (core dump) 文件, 供调试器调试之用。(本书 22.1 节会对核心转储文件加以解释。)

函数 `err_exit()` 类似于 `errExit()`, 但存在两方面的差异。

- 打印错误消息之前, `err_exit()` 不会刷新标准输出。
- `err_exit()` 终止进程使用的是 `_exit()`, 而非 `exit()`。这一退出方式, 略去了对 `stdio` 缓冲区的刷新以及对退出处理程序 (exit handler) 的调用。

本书第 25 章描述了 `_exit()` 与 `exit()` 之间的区别, 还探讨了在 `fork()` 创建的子进程中对 `stdio` 缓冲区和退出处理程序的处理方式。阅读第 25 章, 会使上述 `err_exit()` 操作中的差异细节得以澄清。这里只是想提醒读者, 在编写的库函数创建了子进程, 且该子进程因发生错误而需要终止时, `err_exit()` 恰好能一显身手。它避免了对子进程继承自父进程 (即调用进程) 的 `stdio` 缓冲区副本进行刷新, 且不会调用由父进程所建立的退出处理程序。

在功能上, `errExitEN()` 函数与 `errExit()` 大体相同, 区别仅仅在于: 与 `errExit()` 打印与当前 `errno` 值相对应的错误文本不同, `errExitEN()` 只会打印与 `errno` 参数中给定的错误号 (error number) (这也是该函数后缀名 “EN” 的由来) 相对应的文本。

在本书中调用了 POSIX 线程 API 的程序示例中, 主要使用 `errExitEN()` 来处理错误。与传统的 UNIX 系统调用返回 -1 表示错误不同, POSIX 线程函数会在其结果中返回一个 (POSIX 线程函数返回 0 表示成功) 错误号 (正数, 类型为 `errno` 所专用)。

针对 POSIX 线程函数, 可使用如下代码来诊断错误:

```
errno = pthread_create(&thread, NULL, func, &arg);
if (errno != 0)
    errExit("pthread_create");
```

然而, 这一方法效率不高, 因为在线程程序中, `errno` 实际已被定义为宏, 展开后是返回可修改左值的一个函数调用。因此, 每次使用 `errno` 都会引发一次函数调用。使用 `errExitEN()` 改写上述代码, 功能相同, 但更为高效, 如下所示:

```
int s;

s = pthread_create(&thread, NULL, func, &arg);
if (s != 0)
    errExitEN(s, "pthread_create");
```

在 C 语言术语中, 左值是一个用来指代存储区域的表达式。左值最为常见的用法是作为一个变量的标识符。某些操作符也会产生左值。例如, 若 `p` 为指向某块存储区域的指针, 则 `*p` 便是一个左值。POSIX 线程 API 中, 将 `errno` 重新定义为一个函数¹, 该函数会返回一个指向线程专用存储区域的指针 (请参阅 31.3 节)。

诊断其他类型的错误时, 本书使用的是 `fatal()`、`usageErr()` 以及 `cmdLineErr()`。

```
#include "t1pi_hdr.h"

void fatal(const char *format, ...);
void usageErr(const char *format, ...);
void cmdLineErr(const char *format, ...);
```

函数 `fatal()` 用来诊断一般性错误, 其中包括未设置 `errno` 的库函数错误。除了将一个终止换行符自动追加到输出字符串尾部以外, `fatal()` 的参数列表与 `printf()` 基本相同。该函数会在标

¹ 译者注: 以宏的形式。

准错误上打印格式化输出，然后，像 `errExit()` 那样终止程序。

函数 `usageErr()` 用来诊断命令行参数使用方面的错误。其参数列表风格与 `printf()` 相同，并在标准错误上打印字符串 “Usage:”，随之以格式化输出，然后调用 `exit()` 终止程序。（本书的一些程序示例自行提供有对 `usageErr()` 的扩展版本，命名为 `usageError()`。）

函数 `cmdLineErr()` 酷似 `usageErr()`，但其错误诊断是针对于特定程序的命令行参数。

程序清单 3-3 列出的为本书错误诊断函数的实现。

程序清单 3-3：为本书所有程序所使用的错误处理函数

```
----- lib/error_functions.c
#include <stdarg.h>
#include "error_functions.h"
#include "tlpi_hdr.h"
#include "ename.c.inc"          /* Defines ename and MAX_ENAME */
#ifdef __GNUC__
__attribute__((__noreturn__))
#endif
static void
terminate(Boolean useExit3)
{
    char *s;

    /* Dump core if EF_DUMPCORE environment variable is defined and
       is a nonempty string; otherwise call exit(3) or _exit(2),
       depending on the value of 'useExit3'. */

    s = getenv("EF_DUMPCORE");

    if (s != NULL && *s != '\0')
        abort();
    else if (useExit3)
        exit(EXIT_FAILURE);
    else
        _exit(EXIT_FAILURE);
}
static void
outputError(Boolean useErr, int err, Boolean flushStdout,
            const char *format, va_list ap)
{
#define BUF_SIZE 500
    char buf[BUF_SIZE], userMsg[BUF_SIZE], errText[BUF_SIZE];

    vsnprintf(userMsg, BUF_SIZE, format, ap);

    if (useErr)
        snprintf(errText, BUF_SIZE, "[%s %s]",
                (err > 0 && err <= MAX_ENAME) ?
                ename[err] : "?UNKNOWN?", strerror(err));
    else
        snprintf(errText, BUF_SIZE, ":");

    snprintf(buf, BUF_SIZE, "ERROR%s %s\n", errText, userMsg);

    if (flushStdout)
        fflush(stdout);          /* Flush any pending stdout */
    fputs(buf, stderr);
}
```

```

    fflush(stderr);          /* In case stderr is not line-buffered */
}

void
errMsg(const char *format, ...)
{
    va_list argList;
    int savedErrno;

    savedErrno = errno;     /* In case we change it here */
    va_start(argList, format);
    outputError(TRUE, errno, TRUE, format, argList);
    va_end(argList);

    errno = savedErrno;
}

void
errExit(const char *format, ...)
{
    va_list argList;

    va_start(argList, format);
    outputError(TRUE, errno, TRUE, format, argList);
    va_end(argList);

    terminate(TRUE);
}

void
err_exit(const char *format, ...)
{
    va_list argList;

    va_start(argList, format);
    outputError(TRUE, errno, FALSE, format, argList);
    va_end(argList);

    terminate(FALSE);
}

void
errExitEN(int errnum, const char *format, ...)
{
    va_list argList;

    va_start(argList, format);
    outputError(TRUE, errnum, TRUE, format, argList);
    va_end(argList);

    terminate(TRUE);
}

void
fatal(const char *format, ...)
{
    va_list argList;

    va_start(argList, format);
    outputError(FALSE, 0, TRUE, format, argList);
}

```

```

    va_end(argList);

    terminate(TRUE);
}

void
usageErr(const char *format, ...)
{
    va_list argList;

    fflush(stdout);          /* Flush any pending stdout */

    fprintf(stderr, "Usage: ");
    va_start(argList, format);
    vfprintf(stderr, format, argList);
    va_end(argList);

    fflush(stderr);          /* In case stderr is not line-buffered */
    exit(EXIT_FAILURE);
}

void
cmdLineErr(const char *format, ...)
{
    va_list argList;

    fflush(stdout);          /* Flush any pending stdout */

    fprintf(stderr, "Command-line usage error: ");
    va_start(argList, format);
    vfprintf(stderr, format, argList);
    va_end(argList);

    fflush(stderr);          /* In case stderr is not line-buffered */
    exit(EXIT_FAILURE);
}

```

lib/error_functions.c

程序清单 3-4 列出了程序清单 3-3 所包含的文件 `enames.c.inc`。该文件定义了一个名为“`ename`”的字符串数组，其内容是与 `errno` 的各种可能值相对应的符号名称。本书所采用的错误处理函数会使用该数组，去打印与某个特定错误号相对应的符号名。之所以做如此变通，是为了应对以下两种实际情况：一方面，`strerror()`不会标识出与错误消息相对应的符号常量；而另一方面，手册页在描述错误时，使用的是符号名称。打印出符号名便于读者在手册页中查找错误原因。

由于 `errno` 值随 Linux 硬件架构的不同而有所变化，因此 `ename.c.inc` 文件的内容与特定的硬件架构相关。程序清单 3-4 所示的 `ename.c.inc` 文件版本专用于 Linux 2.6/x86-32 系统。构建该文件的脚本 `lib/Build_ename.sh`，包含于为本书发布的源码当中。可以使用该脚本，针对特定的硬件平台及内核版本，来构建 `ename.c.inc` 文件。

请注意，数组 `ename` 中的某些字符串为空。它们与未使用的错误值相对应。此外，其中的一些字符串包含了两个错误名称，之间以斜杠分隔，是对应两个符号错误名具有相同数值的情况。

从 `ename.c.inc` 文件中，可以看出错误 `EAGAIN` 和 `EWOULDBLOCK` 具有相同数值。`SUSv3` 明确允许这一做法，而且在大多数其他 UNIX 实现（并非全部）上，这些常量值均相同。系统调用返回此类错误的情况是：本应阻塞（亦即在完成调用前被强制等待），而调

用者要求系统调用返回错误。EAGAIN 源于 System V，是由实施 I/O 操作、信号操作、消息队列操作以及文件锁定操作（fcntl()）的系统调用所返回的错误。EWOULDBLOCK 则发源于 BSD，由文件锁定（flock()）以及与套接字相关的系统调用返回。

在 SUSv3 中，仅在与套接字相关的各种接口规范中提及 EWOULDBLOCK。对此类接口来说，SUSv3 允许非阻塞调用要么返回 EAGAIN，要么返回 EWOULDBLOCK。对于所有其他的非阻塞调用，SUSv3 只明确定义了 EAGAIN 错误。

程序清单 3-4: Linux 错误名 (x86-32 版)

```
lib/ename.c.inc
static char *ename[] = {
    /* 0 */ "",
    /* 1 */ "EPERM", "ENOENT", "ESRCH", "EINTR", "EIO", "ENXIO", "E2BIG",
    /* 8 */ "ENOEXEC", "EBADF", "ECHILD", "EAGAIN/EWOULDBLOCK", "ENOMEM",
    /* 13 */ "EACCES", "EFAULT", "ENOTBLK", "EBUSY", "EEXIST", "EXDEV",
    /* 19 */ "ENODEV", "ENODIR", "EISDIR", "EINVAL", "ENFILE", "EMFILE",
    /* 25 */ "ENOTTY", "ETXTBSY", "EFBIG", "ENOSPC", "ESPIPE", "EROFS",
    /* 31 */ "EMLINK", "EPIPE", "EDOM", "ERANGE", "EDEADLK/EDEADLOCK",
    /* 36 */ "ENAMETOOLONG", "ENOLCK", "ENOSYS", "ENOTEMPTY", "ELOOP", "",
    /* 42 */ "ENOMSG", "EIDRM", "ECHRNG", "EL2NSYNC", "EL3HLT", "EL3RST",
    /* 48 */ "ELNRNG", "EUNATCH", "ENOCSI", "EL2HLT", "EBADE", "EBADR",
    /* 54 */ "EXFULL", "ENOANO", "EBADRQC", "EBADSLT", "", "EBFONT", "ENOSTR",
    /* 61 */ "ENODATA", "ETIME", "ENOSR", "ENONET", "ENOPKG", "EREMOTE",
    /* 67 */ "ENOLINK", "EADV", "ESRMT", "ECOMM", "EPROTO", "EMULTIHOP",
    /* 73 */ "EDOTDOT", "EBADMSG", "E_OVERFLOW", "ENOTUNIQ", "EBADFD",
    /* 78 */ "EREMCHG", "ELIBACC", "ELIBBAD", "ELIBSCN", "ELIBMAX",
    /* 83 */ "ELIBEXEC", "EILSEQ", "ERESTART", "ESTRPIPE", "EUSERS",
    /* 88 */ "ENOTSOCK", "EDESTADDRREQ", "EMSGSIZE", "EPROTOTYPE",
    /* 92 */ "ENOPROTOOPT", "EPROTONOSUPPORT", "ESOCKTNOSUPPORT",
    /* 95 */ "EOPNOTSUPP/ENOTSUP", "EPFNOSUPPORT", "EAFNOSUPPORT",
    /* 98 */ "EADDRINUSE", "EADDRNOTAVAIL", "ENETDOWN", "ENETUNREACH",
    /* 102 */ "ENETRESET", "ECONNABORTED", "ECONNRESET", "ENOBUFS", "EISCONN",
    /* 107 */ "ENOTCONN", "ESHUTDOWN", "ETOOMANYREFS", "ETIMEDOUT",
    /* 111 */ "ECONNREFUSED", "EHOSTDOWN", "EHOSTUNREACH", "EALREADY",
    /* 115 */ "EINPROGRESS", "ESTALE", "EUCLEAN", "ENOTNAM", "EAVAIL",
    /* 120 */ "EISNAM", "EREMOTEIO", "EDQUOT", "ENOMEDIUM", "EMEDIUMTYPE",
    /* 125 */ "ECANCELED", "ENOKEY", "EKEYEXPIRED", "EKEYREVOKED",
    /* 129 */ "EKEYREJECTED", "EOWNERDEAD", "ENOTRECOVERABLE", "ERFKILL"
};

#define MAX_ENAME 132
lib/ename.c.inc
```

解析数值型命令行参数的函数

程序清单 3-5 中的头文件提供了两个函数声明，在本书中频繁用于解析整形命令行参数：getInt()和 getLong()。较之于 atoi()、atol()以及 strtol()，它们的主要优点在于针对数值型参数提供了一些基本的有效性检查。

```
#include "tspi_hdr.h"

int getInt(const char *arg, int flags, const char *name);
long getLong(const char *arg, int flags, const char *name);

Both return arg converted to numeric form
```

函数 `getInt()` 和 `getLong()` 分别将 `arg` 指向的字符串转换为 `int` 或 `long`。如果 `arg` 未包含一个有效的整数字符串（即仅包含数字以及字符“+”和“-”），那么这两个函数会打印一条错误消息，并终止程序。

若参数 `name` 非空，则所含内容应为一字符串，用于标识 `arg` 对应于命令行中相应参数的名称。在上述两函数中，无论打印任何错误消息，该字符串都是消息中的一部分。

可通过 `flags` 参数对 `getInt()` 和 `getLong()` 函数的操作施加一些控制。默认情况下，两个函数会处理包含有符号十进制整数的字符串。若将定义于程序清单 3-5 中的一个或多个 `GN_*` 系列常量与 `flags` 相或，则既可以选择其他的转换进制，也能将数值范围限制为非负或正整数。

虽然 `flags` 参数允许程序强制执行正文所述的范围检查，但在某些情况下，即便这么做看起来合理，程序示例也无需做类似检测。例如，程序清单 47-1 中就并未对参数 `init-value` 进行检查。这意味着，用户可将一个负数作为初始值赋给某个信号量，这会引发随后的 `semctl()` 系统调用返回错误（`ERANGE`），因为信号量不能为负值。在此类情况下省略对范围的检查，不但能够让读者体验对系统调用及库函数的正确使用，还能观察到输入无效参数时所发生的情形。通常，现实世界中的应用程序会对自身命令行参数施以更为严格的检查。

程序清单 3-6 给出了函数 `getInt()` 和 `getLong()` 的实现。

程序清单 3-5: `get_num.c` 的头文件

```
----- lib/get_num.h
#ifndef GET_NUM_H
#define GET_NUM_H

#define GN_NONNEG      01      /* Value must be >= 0 */
#define GN_GT_0       02      /* Value must be > 0 */

/* By default, integers are decimal */
#define GN_ANY_BASE   0100    /* Can use any base - like strtol(3) */
#define GN_BASE_8     0200    /* Value is expressed in octal */
#define GN_BASE_16    0400    /* Value is expressed in hexadecimal */

long getLong(const char *arg, int flags, const char *name);

int getInt(const char *arg, int flags, const char *name);

#endif
----- lib/get_num.h
```

程序清单 3-6: 解析数值型命令行参数的函数

```
----- lib/get_num.c

#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <limits.h>
#include <errno.h>
#include "get_num.h"
static void
gnFail(const char *fname, const char *msg, const char *arg, const char *name)
{
    fprintf(stderr, "%s error", fname);
}
```



```

    if (name != NULL)
        fprintf(stderr, " (in %s)", name);
    fprintf(stderr, ": %s\n", msg);
    if (arg != NULL && *arg != '\0')
        fprintf(stderr, "    offending text: %s\n", arg);

    exit(EXIT_FAILURE);
}

static long
getNum(const char *fname, const char *arg, int flags, const char *name)
{
    long res;
    char *endptr;
    int base;

    if (arg == NULL || *arg == '\0')
        gnFail(fname, "null or empty string", arg, name);

    base = (flags & GN_ANY_BASE) ? 0 : (flags & GN_BASE_8) ? 8 :
           (flags & GN_BASE_16) ? 16 : 10;

    errno = 0;
    res = strtol(arg, &endptr, base);
    if (errno != 0)
        gnFail(fname, "strtol() failed", arg, name);

    if (*endptr != '\0')
        gnFail(fname, "nonnumeric characters", arg, name);

    if ((flags & GN_NONNEG) && res < 0)
        gnFail(fname, "negative value not allowed", arg, name);

    if ((flags & GN_GT_0) && res <= 0)
        gnFail(fname, "value must be > 0", arg, name);

    return res;
}

long
getLong(const char *arg, int flags, const char *name)
{
    return getNum("getLong", arg, flags, name);
}

int
getInt(const char *arg, int flags, const char *name)
{
    long res;

    res = getNum("getInt", arg, flags, name);
    if (res > INT_MAX || res < INT_MIN)
        gnFail("getInt", "integer out of range", arg, name);

    return (int) res;
}

```

lib/get_num.c

3.6 可移植性问题

本节将探究可移植性系统编程方面的议题。除了会介绍特性测试宏，以及 SUSv3 所定义的标准系统数据类型之外，还会关注一些其他的可移植性问题。

3.6.1 特性测试宏

系统调用和库函数 API 的行为受各种标准（参见 1.3 节）的制约。这些标准中的一部分是由 Open Group（SUS）这样的标准机构来制定的，而另一部分则是由具有重要历史意义的两个 UNIX 实现 BSD 和 System V Release 4（以及相关的 System V 接口定义）来定义。

编写可移植性应用程序时，有时希望各个头文件只显露遵循特定标准的定义（常量、函数原型等）。要达到这一目的，在编译程序时需要定义下列一个或多个特性测试宏。方式之一是在程序源码包含¹任何头文件之前，定义如下宏：

```
#define _BSD_SOURCE 1
```

此外，还可以使用 C 编译器的 -D 选项：

```
$ cc -D_BSD_SOURCE prog.c
```

术语“特性测试宏”似乎易于让人产生误解，但只要从 UNIX 实现的角度看来，读者便会发现这一称谓其实颇有道理。通过测试（使用 #if）应用程序为宏所定义的值，实现可以决定应该让哪些（由头文件提供的）特性可见。

以下特性测试宏由相关标准定义而成，因而在支持这些标准的所有系统上，对这些宏的使用均是可移植的：

`_POSIX_SOURCE`

一经定义（任何值），头文件会显露符合 POSIX.1-1990 和 ISO C（1990）标准的定义。该宏已被 `_POSIX_C_SOURCE` 取代。

`_POSIX_C_SOURCE`

若定义为 1，效果与 `_POSIX_SOURCE` 相同。若将其值定义为大于等于 199309，头文件还会显露遵从 POSIX.1b（实时）标准的定义。若将其值定义为大于等于 199506，便会开启对 POSIX.1c（线程）定义的支持。若将其值定义为 200112，则开启对 POSIX.1-2001 基本规范（排除了 XSI 扩展）定义的支持。（2.3.3 版本之前，glibc 头文件对值为 200112 的 `_POSIX_C_SOURCE` 不做解释。）若将其值定义为 200809，便会开启对 POSIX.1-2008 基本规范定义的支持。（2.10 版本之前，glibc 头文件对值为 200809 的 `_POSIX_C_SOURCE` 不做解释。）

`_XOPEN_SOURCE`

一经定义（任何值），头文件会显露对 POSIX.1、POSIX.2 和 X/Open(XPG4)标准的定义。若将其值定义为大于等于 500，还会开启对 SUSv2 (UNIX 98 和 XPG5)扩展的支持。若将其值设置为大于等于 600，则又开启了对 SUSv3 XSI (UNIX 03)扩展和 C99 扩展的支持。（2.2 版本之前，glibc 头文件对值为 600 的 `_XOPEN_SOURCE` 不做解释。）若将其值设置为大于等于 700，便会开启对 SUSv4 XSI 扩展的支持（2.10 版本之前，glibc 头文件对值为 700 的 `_XOPEN_SOURCE` 不做解释）。之所以选择 500、600 和 700 作为取值，是因为 SUSv2、SUSv3 和 SUSv4 分别是

¹ 译者注：`#include`。

X/Open 规范的第 5 号、第 6 号和第 7 号。

以下列出为 glibc 专用的特性测试宏：

`_BSD_SOURCE`

一经定义（任何值），开启对 BSD 定义的支持。此外，只要定义了该宏，便以值 199506 定义了 `_POSIX_C_SOURCE`。极少数的情况下，当标准之间发生冲突时，显式设置该宏会导致系统向 BSD 定义倾斜。

`_SVID_SOURCE`

一经定义（任何值），头文件会显露符合 System V 接口规范（SVID）的定义。

`_GNU_SOURCE`

一经定义（任何值），头文件除了会显露符合前述所有标准的定义（通过设置前述所有宏来提供）外，还会开启对各种 GNU 扩展定义的支持。

在不带任何特殊选项调用 GNU C 编译器时，即默认定义了 `_POSIX_SOURCE`、`_POSIX_C_SOURCE=200809`（glibc 版本为 2.5-2.9 时，其值为 200112；glibc 版本低于 2.4 时，其值为 199506）、`_BSD_SOURCE` 以及 `_SVID_SOURCE`。

在对个别宏进行了定义，或以其标准模式之一去调用编译器时（比如，`cc -ansi` 或 `cc -std=c99`），只会按需提供定义。不过，有一个例外：若未对 `_POSIX_C_SOURCE` 另行定义，且未以标准模式之一去调用编译器，则 `_POSIX_C_SOURCE` 的值会被定义为 200809（glibc 版本为 2.4-2.9 时，其值为 200112；glibc 版本低于 2.4 时，其值为 199506）。

定义多个宏有叠加效应，故而缺省情况下所提供的宏设置，也可使用如下 `cc` 命令来明确选择：

```
$ cc -D_POSIX_SOURCE -D_POSIX_C_SOURCE=199506 \  
      -D_BSD_SOURCE -D_SVID_SOURCE prog.c
```

<features.h>头文件和 `feature_test_macros(7)`手册页，针对赋给每个特性测试宏的值，提供了更多精确信息。

`_POSIX_C_SOURCE`、`_XOPEN_SOURCE` 以及 POSIX.1/SUS

在 POSIX.1-2001/SUSv3 中，仅对 `_POSIX_C_SOURCE` 和 `_XOPEN_SOURCE` 特性测试宏进行了明确定义，应用程序要符合该标准，应分别将上述两宏的值定义为 200112 和 600。将 `_POSIX_C_SOURCE` 值定义为 200112，即表示应用程序符合 POSIX.1-2001 基本规范（即符合除 XSI 扩展规范以外的 POSIX 规范）。将 `_XOPEN_SOURCE` 值定义为 600，即表示应用程序符合 SUSv3 规范（即符合 XSI 规范基本规范加 XSI 扩展规范）。上述声明同样适用于 POSIX.1-2008/SUSv4，只是需要将上述两个特性测试宏的值分别定义为 200809 和 700。

SUSv3 明文规定将 `_XOPEN_SOURCE` 设置为 600 所提供的特性，就包含了将 `_POSIX_C_SOURCE` 设置为 200112 时所激活的所有特性。因此，为符合 SUSv3（即 XSI 规范），应用程序只需要定义 `_XOPEN_SOURCE`。SUSv4 做出了类似规定，将 `_XOPEN_SOURCE` 设置为 700 所提供的特性，包含了 `_POSIX_C_SOURCE` 值被设置为 200809 时所激活的所有特性。

函数原型及源码示例中的特性测试宏

手册页则描述了欲使某个特定常量定义或函数声明在头文件中可见，应该定义哪些特性测试宏。为本书编写的所有源码示例，编译时采用缺省的 GNU C 语言编译器选项或如下选项：

```
$ cc -std=c99 -D_XOPEN_SOURCE=600
```

对于在本书中出现的函数，为了能在以上述两种方式编译的程序中编译通过，在其原型处均注明了使用这些函数所必须定义的任何特性测试宏。手册页中，对于启用每一函数声明所需定义的特性测试宏，则有更为精确的描述。

3.6.2 系统数据类型

不同实现的数据类型，例如：进程 ID、用户 ID 以及文件偏移量，表示时均采用标准 C 语言类型。尽管也有可能使用 C 语言的基本类型，诸如 `int` 和 `long`，来声明存储此类信息的变量，但这一做法降低了不同 UNIX 系统间相互移植的难度，分析如下。

- 随着 UNIX 实现的不同（例如，`long` 型可能在系统 A 上长度为 4 字节，在系统 B 上为 8 字节），有时甚至是同一实现中编译环境的不同，这些基本类型的大小各不相同。更有甚者，不同实现可能会使用不同类型来表示相同信息。例如，进程 ID 在系统 A 上为 `int` 型，而在系统 B 上为 `long` 型。
- 即便是针对同一款 UNIX 实现，用以表征信息的类型在不同版本之间也会有所不同。Linux 上较为知名的例子是用户 ID 和组 ID。在 Linux 2.2 及其之前，这些值以 16 位表示。在 2.4 及其之后，则以 32 位表示。

为避免此类可移植性问题，SUSv3 规范了各种标准系统数据类型，并要求各个实现适当加以定义和使用。每种类型的定义均使用 C 语言的 `typedef` 特性。例如，`pid_t` 数据类型用以表示进程 ID，在 Linux/x86-32 上，其类型定义如下：

```
typedef int pid_t;
```

标准系统数据类型中的大多数，其命名均以 `_t` 结尾。其中的许多都声明于头文件 `<sys/types.h>` 中，余下的少量则定义于其他头文件中。

应用程序应采用这些类型定义来声明其使用的变量，才能保证可移植性。例如，如下声明将允许应用程序在任何符合 SUSv3 标准的系统上正确表示进程 ID。

```
pid_t mypid;
```

表 3-1 列出将在书中碰到的部分系统数据类型。对于表中的某些特定类型，SUSv3 要求以“运算类型（arithmetic type）”来加以实现。这意味着，实现所选择的底层类型，要么为整数类型，要么为浮点（实数或复数）型。

表 3-1: 系统数据类型选录

数据类型	SUSv3 类型需求	描述
<code>blkcnt_t</code>	有符号整型	文件块数量（15.1 节）
<code>blksize_t</code>	有符号整型	文件块大小（15.1 节）
<code>cc_t</code>	无符号整型	终端特殊字符（62.4 节）
<code>clock_t</code>	整型或浮点型实数	以时钟周期计量的系统时间（10.7 节）
<code>clockid_t</code>	运算类型之一	针对 POSIX.1b 时钟和定时器函数的时钟标识符
<code>comp_t</code>	SUSv3 未作规范	经由压缩处理的时钟周期（28.1 节）
<code>dev_t</code>	运算类型之一	设备号，包含主、次设备号（15.1 节）
<code>DIR</code>	无类型要求	目录流（18.8 节）
<code>fd_set</code>	结构类型	<code>select()</code> （63.2.1 节）中的文件描述符集合
<code>fsblkcnt_t</code>	无符号整型	文件系统块数量（14.11 节）

数据类型	SUSv3 类型需求	描 述
fsfilcnt_t	无符号整型	文件数量 (14.11 节)
gid_t	整型	数值型组标识符 (8.3 节)
id_t	整型	用以存放标识符的通用类型, 其大小至少可放置 pid_t、uid_t 和 gid_t 类型
in_addr_t	32 位无符号整型	IPv4 地址 (59.4 节)
in_port_t	16 位无符号整型	IP 端口号 (59.4 节)
ino_t	无符号整型	文件 i-node 号 (15.1 节)
key_t	运算类型之一	System V IPC 键 (45.2 节)
mode_t	整型	文件权限及类型 (15.1 节)
mqd_t	无类型要求, 但不能为数组类型	POSIX 消息队列描述符
msglen_t	无符号整型	System V 消息队列所允许的字节数 (46.4 节)
msgqnum_t	无符号整型	System V 消息队列中的消息数量 (46.4 节)
nfds_t	无符号整型	poll() (63.2.2 节) 中的文件描述符数量
nlink_t	整型	文件的 (硬) 连接数量 (15.1 节)
off_t	有符号整型	文件偏移量或大小 (4.7 节及 15.1 节)
pid_t	有符号整型	进程 ID、进程组 ID 或会话 ID (6.3 节、34.2 节、34.3 节)
ptrdiff_t	有符号整型	两指针差值, 为有符号整型
rlim_t	无符号整型	资源限制 (36.2 节)
sa_family_t	无符号整型	套接字地址族 (56.4 节)
shmatt_t	无符号整型	与 System V 共享内存段相连的进程数量
sig_atomic_t	整型	可进行原子访问的数据类型 (21.1.3 节)
siginfo_t	结构类型	信号起源的相关信息 (21.4 节)
sigset_t	整形或结构类型	信号集合 (20.9 节)
size_t	无符号整型	对象大小 (以字节数计)
socklen_t	至少 32 位的整型	套接字地址结构大小 (以字节数计) (56.3 节)
speed_t	无符号整型	终端线速度 (62.7 节)
ssize_t	有符号整型	字节数或 (为负时) 标识错误
stack_t	结构类型	对备选信号栈的描述 (21.3 节)
suseconds_t	有符号整型, 范围为 [-1,1000000]	微秒级的时间间隔 (10.1 节)
tflag_t	无符号整型	终端模式标志位的位掩码 (62.2 节)
time_t	整型或浮点型实数	自所谓纪元 (Epoch) (10.1 节) 始, 以秒计的日历时间
timer_t	运算类型之一	POSIX.1b 间隔定时器函数 (23.6 节) 的定时器标识符
uid_t	整型	数值型用户标识符 (8.1 节)

在后续章节中论及表 3-1 的数据类型时，常会作如下表述：某类型“为一整数类型（由 SUSv3 所规定）”。这是指 SUSv3 要求以整型来定义该类型，但不要求使用某一特定的原生（native）数据类型（例如：short、int 或 long）。（通常，针对 Linux 中的每种系统数据类型，书中不会言明实际会使用哪种原生数据类型加以表示，因为编写可移植应用程序时无需关注这点。）

打印系统数据类型值

当需要打印表 3-1 所列数值型系统数据类型（例如：pid_t 和 uid_t）的值时，调用 printf() 应留意不要引入对表现形式的依赖。这一依赖是由于 C 语言的（升级型）自动类型转换造成的。该转换会将 short 型转换为 int 型，而对于 int 型和 long 型则置之不问。这就意味着传入 printf() 的要么为 int 型，要么为 long 型。然而，因为 printf() 在运行时无从判定其参数类型，调用者必须明确其格式限定符为 %d 还是 %ld。问题在于在 printf() 中仅就一种限定符进行编码会导致对实现的依赖。常见的应对策略是强制转换相应值为 long 型后，再使用 %ld 的限定符，如下所示：

```
pid_t mypid;

mypid = getpid();          /* Returns process ID of calling process */
printf("My PID is %ld\n", (long) mypid);
```

上述技术有个例外。因为在一些编译环境中数据类型 off_t 大小与 long long 相当，所以会将 off_t 强制转换为该类型并使用 %lld 的限定符，如 5.10 节所述。

C99 标准为 printf 定义了名为 z 的长度修饰符，以表明紧随其后的整型转换是与 size_t 或 ssize_t 类型相对应的。因而，要对付这些类型，就可以用 %zd 来取代 %ld 外加类型转换的方法了。尽管 glibc 支持该限定符，但其并未获得所有 UNIX 实现的支持，故而本书也避免采用这一做法。

C99 标准还定义有名为 j 的长度修饰符，并规定其相应参数的类型为 intmax_t（或 uintmax_t），该类型之大，足以用其表示任何类型的整数。最终，使用 intmax_t 强制转换外加 %jd 限定符的方案应取代 long 强制转换外加 %ld 限定符的方案，成为打印数值型系统数据类型的首选，因为前者可以处理 long long 型值和诸如 int128_t 之类的任一可扩展整数类型。然而，出于相同的原因（未获得所有 UNIX 实现的支持），本书没有采用这一技术。

3.6.3 其他的可移植性问题

本节所讨论的，是进行系统编程时可能会遇到的一些其他可移植性问题。

初始化操作和使用结构

每种 UNIX 实现，都明确定义了一系列标准结构，用于各种系统调用及库函数。现试举一例，请考虑用来表示信号量操作（通过 semop() 系统调用去执行）的结构 sembuf：

```
struct sembuf {
    unsigned short sem_num;    /* Semaphore number */
    short          sem_op;     /* Operation to be performed */
    short          sem_flg;    /* Operation flags */
};
```

尽管 SUSv3 定义了诸如 sembuf 之类的结构，但意识到如下两点尤为重要。

- 总体而言，未对此类结构内部的字段顺序作出规范。
- 一些情况下，此类结构内会包含额外的、与实现相关的字段。

因此，以如下方式对数据结构进行初始化，其代码是无法移植的：

```
struct sembuf s = { 3, -1, SEM_UNDO };
```

这段初始化程序尽管在 Linux 上运行没有问题，但在其他一些实现上，由于对 sembuf 结构的定义中顺序会有所不同，故而将无法工作。要在初始化时消除此类移植问题，必须明确采用如下赋值语句：

```
struct sembuf s;  
  
s.sem_num = 3;  
s.sem_op = -1;  
s.sem_flg = SEM_UNDO;
```

如果采用的是 C99 语言标准，可以利用该语言针对结构初始化的新语法，写出等价代码：

```
struct sembuf s = { .sem_num = 3, .sem_op = -1, .sem_flg = SEM_UNDO };
```

若需要将标准结构的内容转储到文件时，标准结构的成员顺序也要加以考虑。此时要消除可移植性问题，简单地将结构以二进制形式写入是无济于事的。相反，必须将结构内字段以特定顺序逐一加以记录（可能是以文本形式）。

使用未见诸于所有实现的宏

有时，未必所有的 UNIX 实现都对一个宏做了定义。例如，WCOREDUMP()宏（用于检测子进程是否生成了核心转储文件）的使用非常广泛，但 SUSv3 却并未对其进行规范。因此，在某些 UNIX 实现上，该宏并不存在。要妥善处理此类潜在的可移植性问题，可以使用 C 语言的预编译指令#ifdef，如下所示：

```
#ifdef WCOREDUMP  
    /* Use WCOREDUMP() macro */  
#endif
```

不同实现间所需头文件的变化

有些情况下，包含各种系统调用和库函数原型的头文件，在不同 UNIX 实现之间会有所不同。本书仅展示 Linux 的需求，并注明其与 SUSv3 间的各种变化。

本书论及部分函数时，会引入特定头文件，伴之以“/*出于可移植型考虑*/”形式的注解。这表明 Linux 和 SUSv3 都不需要此头文件，但由于某些其他（尤其是老一些的）实现可能需要，在可移植程序中应将其纳入。

POSIX.1-1990 曾规定，编程时如需使用其所规范的许多函数，在包含与该函数相关的任何其他头文件之前，必须包含头文件<sys/types.h>。可是，这一要求不久就成了多此一举。绝大多数现代 UNIX 实现中的应用程序，在使用这些函数时都无需包含此头文件。SUSv1 因而删除了这一要求。然而，在编写可移植程序时，将以其为首的头文件包括在内，仍不失为明智之举。（不过，本书的程序示例中省去这一头文件，因为在 Linux 平台下此举实属多余，而程序代码也因此少了一行。）

3.7 总结

系统调用允许进程向内核请求服务。与用户空间的函数调用相比，哪怕是最简单的系统调用都会产生显著的开销，其原因是为了执行系统调用，系统需要临时性地切换到

核心态，此外，内核还需验证系统调用的参数、用户内存和内核内存之间也有数据需要传递。

标准的 C 语言函数库提供了大量库函数，功能五花八门。有些库函数会利用系统调用来完成工作，而另一些库函数则完全在用户空间内执行任务。在 Linux 上，一般情况下，使用 glibc 作为 C 语言标准库的实现。

大多数系统调用和库函数都会返回一个状态值，以表明调用成功与否。对这一返回状态进行检查是一条编程铁律。

为本书的程序示例还实现有一批函数。其所执行的任务包括诊断错误和解析命令行参数。

本章也提供了一系列指南及技术，以帮助读者编写可移植的系统程序，此类程序可在任何符合标准的系统上运行。

编译应用程序时，可定义不同的特性测试宏，以控制头文件显露对特定标准的定义。当希望确保程序符合某些正式或由实现定义的标准时，上述做法可谓是非常实用。

利用定义于各个标准中（而非原生 C 语言类型）的系统数据类型，能够改善系统编程的可移植性。SUSv3 定义有大量系统数据类型，UNIX 实现应加以支持，应用程序应予以采用。

3.8 练习

- 3-1** 使用 Linux 专有的 `reboot()` 系统调用重启系统时，必须将第二个参数 `MAGIC2` 定义为一组 `MAGIC` 号之一（例如，`LINUX_REBOOT_MAGIC2`）。这些 `MAGIC` 号有何意义？（将 `MAGIC` 号转换为十六进制数，对解题会有所帮助。）

第 4 章

文件 I/O：通用的 I/O 模型

现在，我们开始深入研究系统调用 API。作为 UNIX 系统设计思想的核心理念，文件（file）是一个不错的起点。本章重点介绍用于文件输入/输出的系统调用。

本章开篇会讨论文件描述符的概念，随后会逐一讲解构成通用 I/O 模型的系统调用，其中包括：打开文件、关闭文件、从文件中读数据和向文件中写数据。

本章所关注的是磁盘文件的 I/O 操作。然而，鉴于可以采用相同的系统调用对诸如管道、终端等所有类型的文件施以输入/输出操作，故而本章的大部分内容会与后续章节有关。

第 5 章会在本章基础上对文件 I/O 做深入探讨。缓冲（buffering）是文件 I/O 的另一要点，其复杂程度足以专辟一章讲述。第 13 章就涵盖了内核和 stdio 库中的 I/O 缓冲。

4.1 概述

所有执行 I/O 操作的系统调用都以文件描述符，一个非负整数（通常是小整数），来指代打开的文件。文件描述符用以表示所有类型的已打开文件，包括管道（pipe）、FIFO、socket、终端、设备和普通文件。针对每个进程，文件描述符都自成一套。

按照惯例，大多数程序都期望能够使用 3 种标准的文件描述符，见表 4-1。在程序开始运行之前，shell 代表程序打开这 3 个文件描述符。更确切地说，程序继承了 shell 文件描述符的副本——在 shell 的日常操作中，这 3 个文件描述符始终是打开的。（在交互式 shell 中，这 3 个文件描述符通常指向 shell 运行所在的终端。）如果命令行指定对输入/输出进行重定向操作，那么 shell 会对文件描述符做适当修改，然后再启动程序。

表 4-1：标准文件描述符

文件描述符	用途	POSIX 名称	stdio 流
0	标准输入	STDIN_FILENO	stdin
1	标准输出	STDOUT_FILENO	stdout
2	标准错误	STDERR_FILENO	stderr

在程序中指代这些文件描述符时，可以使用数字（0、1、2）表示，或者采用<unistd.h>所定义的 POSIX 标准名称——此方法更为可取。

虽然 `stdin`、`stdout` 和 `stderr` 变量在程序初始化时用于指代进程的标准输入、标准输出和标准错误，但是调用 `freopen()` 库函数可以使这些变量指代其他任何文件对象。作为其操作的一部分，`freopen()` 可以在将流（stream）重新打开之际一并更换隐匿其中的文件描述符。换言之，针对 `stdout` 调用 `freopen()` 函数后，无法保证 `stdout` 变量值仍然为 1。

下面介绍执行文件 I/O 操作的 4 个主要系统调用（编程语言和软件包通常会利用 I/O 函数库对它们进行间接调用）。

- `fd = open(pathname, flags, mode)` 函数打开 `pathname` 所标识的文件，并返回文件描述符，用以在后续函数调用中指代打开的文件。如果文件不存在，`open()` 函数可以创建之，这取决于对位掩码参数 `flags` 的设置。`flags` 参数还可指定文件的打开方式：只读、只写亦或是读写方式。`mode` 参数则指定了由 `open()` 调用创建文件的访问权限，如果 `open()` 函数并未创建文件，那么可以忽略或省略 `mode` 参数。
- `numread = read(fd, buffer, count)` 调用从 `fd` 所指代的打开文件中读取至多 `count` 字节的数据，并存储到 `buffer` 中。`read()` 调用的返回值为实际读取到的字节数。如果再无字节可读（例如：读到文件结尾符 EOF 时），则返回值为 0。
- `numwritten = write(fd, buffer, count)` 调用从 `buffer` 中读取多达 `count` 字节的数据写入由 `fd` 所指代的已打开文件中。`write()` 调用的返回值为实际写入文件中的字节数，且有可能小于 `count`。
- `status = close(fd)` 在所有输入/输出操作完成后，调用 `close()`，释放文件描述符 `fd` 以及与之相关的内核资源。

在详细说明这些系统调用之前，程序清单 4-1 简要展示了它们的使用方法。该程序实现了一个简版的 `cp(1)` 命令，将源文件内容复制到新文件中。在命令行中，程序的第一个参数代表已存在的源文件，第二个参数则代表新文件。程序清单 4-1 如下所示：

```
$ ./copy oldfile newfile
```

程序清单 4-1：使用 I/O 系统调用

```
----- fileio/copy.c
#include <sys/stat.h>
#include <fcntl.h>
#include "tlpi_hdr.h"

#ifdef BUF_SIZE /* Allow "cc -D" to override definition */
#define BUF_SIZE 1024
#endif

int
main(int argc, char *argv[])
{
    int inputFd, outputFd, openFlags;
    mode_t filePerms;
    ssize_t numRead;
    char buf[BUF_SIZE];

    if (argc != 3 || strcmp(argv[1], "--help") == 0)
        usageErr("%s old-file new-file\n", argv[0]);
```

```

/* Open input and output files */

inputFd = open(argv[1], O_RDONLY);
if (inputFd == -1)
    errExit("opening file %s", argv[1]);

openFlags = O_CREAT | O_WRONLY | O_TRUNC;
filePerms = S_IRUSR | S_IWUSR | S_IRGRP | S_IWGRP |
            S_IROTH | S_IWOTH; /* rw-rw-rw- */
outputFd = open(argv[2], openFlags, filePerms);
if (outputFd == -1)
    errExit("opening file %s", argv[2]);

/* Transfer data until we encounter end of input or an error */

while ((numRead = read(inputFd, buf, BUF_SIZE)) > 0)
    if (write(outputFd, buf, numRead) != numRead)
        fatal("couldn't write whole buffer");
if (numRead == -1)
    errExit("read");

if (close(inputFd) == -1)
    errExit("close input");
if (close(outputFd) == -1)
    errExit("close output");

exit(EXIT_SUCCESS);
}

```

fileio/copy.c

4.2 通用 I/O

UNIX I/O 模型的显著特点之一是其输入/输出的通用性概念。这意味着使用 4 个同样的系统调用 `open()`、`read()`、`write()` 和 `close()` 可以对所有类型的文件执行 I/O 操作，包括终端之类的设备。因此，仅使用这些系统调用编写的程序，将对任何类型的文件有效。例如，针对程序清单 4-1 中的程序，如下操作都是有效的：

```

$ ./copy test test.old      Copy a regular file
$ ./copy a.txt /dev/tty     Copy a regular file to this terminal
$ ./copy /dev/tty b.txt     Copy input from this terminal to a regular file
$ ./copy /dev/pts/16 /dev/tty Copy input from another terminal

```

要实现通用 I/O，就必须确保每一文件系统和设备驱动程序都实现了相同的 I/O 系统调用集。由于文件系统或设备所特有的操作细节在内核中处理，在编程时通常可以忽略设备专有的因素。一旦应用程序需要访问文件系统或设备的专有功能时，可以选择瑞士军刀般的 `ioctl()` 系统调用（4.8 节），该调用为通用 I/O 模型之外的专有特性提供了访问接口。

4.3 打开一个文件：open()

`open()` 调用既能打开一个业已存在的文件，也能创建并打开一个新文件。

```

#include <sys/stat.h>
#include <fcntl.h>

```

```
int open(const char *pathname, int flags, ... /* mode_t mode */);
```

Returns file descriptor on success, or -1 on error

要打开的文件由参数 `pathname` 来标识。如果 `pathname` 是一符号链接，会对其进行解引用。如果调用成功，`open()` 将返回一文件描述符，用于在后续函数调用中指代该文件。若发生错误，则返回 -1，并将 `errno` 置为相应的错误标志。

参数 `flags` 为位掩码，用于指定文件的访问模式，可选择表 4-2 所示的常量之一。

早期的 UNIX 实现中使用数字 0、1、2，而非表 4-2 中所列的常量名称。大多数现代 UNIX 实现将这些常量定义为上述相应数字（以期与早期系统保持兼容）。由此可见，`O_RDWR` 并不等同于 `O_RDONLY | O_WRONLY`，后者（或组合）属于逻辑错误。

当调用 `open()` 创建新文件时，位掩码参数 `mode` 指定了文件的访问权限。（SUSv3 规定，`mode` 的数据类型 `mode_t` 属于整数类型。）如果 `open()` 并未指定 `O_CREAT` 标志，则可以省略 `mode` 参数。

表 4-2: 文件访问模式

访问模式	描述
<code>O_RDONLY</code>	以只读方式打开文件
<code>O_WRONLY</code>	以只写方式打开文件
<code>O_RDWR</code>	以读写方式打开文件

15.4 节将详细描述文件权限。之后，读者会了解到新建文件的访问权限不仅仅依赖于参数 `mode`，而且受到进程的 `umask` 值（15.4.6 节）和（可能存在的）父目录的默认访问控制列表（17.6 节）影响。与此同时，需要注意 `mode` 参数可以指定为数字（通常为八进制数），更为可取的做法是对 0 个或多个表 15-4（15.4.1 节）中所列位掩码常量进行逻辑或（`|`）操作。

程序清单 4-2 展示了 `open()` 调用的几个使用实例，其中有些调用用到了其他标志位，后续将会加以介绍。

程序清单 4-2: `open` 函数使用的例子

```
/* Open existing file for reading */

fd = open("startup", O_RDONLY);
if (fd == -1)
    errExit("open");

/* Open new or existing file for reading and writing, truncating to zero
bytes; file permissions read+write for owner, nothing for all others */

fd = open("myfile", O_RDWR | O_CREAT | O_TRUNC, S_IRUSR | S_IWUSR);
if (fd == -1)
    errExit("open");

/* Open new or existing file for writing; writes should always
append to end of file */

fd = open("w.log", O_WRONLY | O_CREAT | O_TRUNC | O_APPEND,
          S_IRUSR | S_IWUSR);
```

```

if (fd == -1)
    errExit("open");

```

open()调用所返回的文件描述符数值

SUSv3 规定，如果调用 `open()` 成功，必须保证其返回值为进程未用文件描述符中数值最小者。可以利用该特性以特定文件描述符打开某一文件。例如，如下代码序列就会确保使用标准输入（文件描述符 0）打开一文件。

```

if (close(STDIN_FILENO) == -1)    /* Close file descriptor 0 */
    errExit("close");

fd = open(pathname, O_RDONLY);
if (fd == -1)
    errExit("open");

```

由于文件描述符 0 未用，所以 `open()` 调用势必使用此描述符打开文件。5.5 节中所论及的 `dup2()` 和 `fcntl()` 也可实现类似功能，但对于文件描述符的控制更加灵活。该节还将举例说明对于业已打开的文件，控制其描述符为何大有益处。

4.3.1 open()调用中的 flags 参数

在程序清单 4-2 展示的一些 `open()` 调用例子中，`flags` 参数除了使用文件访问标志外，还使用了其他操作标志（`O_CREAT`、`O_TRUNC` 和 `O_APPEND`）。现在将详细介绍 `flags` 参数。表 4-3 总结了可参与 `flags` 参数逐位或运算（|）的一整套常量。最后一列显示常量标准化于 SUSv3 还是 SUSv4。

表 4-3: `open()` 系统调用的 `flags` 参数值介绍

标 志	用 途	统一 UNIX 规范版本
<code>O_RDONLY</code>	以只读方式打开	v3
<code>O_WRONLY</code>	以只写方式打开	v3
<code>O_RDWR</code>	以读写方式打开	v3
<code>O_CLOEXEC</code>	设置 <code>close-on-exec</code> 标志（自 Linux 2.6.23 版本开始）	v4
<code>O_CREAT</code>	若文件不存在则创建之	v3
<code>O_DIRECT</code>	无缓冲的输入/输出	
<code>O_DIRECTORY</code>	如果 <code>pathname</code> 不是目录，则失败	v4
<code>O_EXCL</code>	结合 <code>O_CREAT</code> 参数使用，专门用于创建文件	v3
<code>O_LARGEFILE</code>	在 32 位系统中使用此标志打开大文件	
<code>O_NOATIME</code>	调用 <code>read()</code> 时，不修改文件最近访问时间（自 Linux 2.6.8 版本开始）	
<code>O_NOCTTY</code>	不要让 <code>pathname</code> （所指向的终端设备）成为控制终端	v3
<code>O_NOFOLLOW</code>	对符号链接不予解引用	v4
<code>O_TRUNC</code>	截断已有文件，使其长度为零	v3

标志	用途	统一 UNIX 规范版本
O_APPEND	总在文件尾部追加数据	v3
O_ASYNC	当 I/O 操作可行时，产生信号 (signal) 通知进程	
O_DSYNC	提供同步的 I/O 数据完整性 (自 Linux 2.6.33 版本开始)	v3
O_NONBLOCK	以非阻塞方式打开	v3
O_SYNC	以同步方式写入文件	v3

表 4-3 中常量分为如下几组。

- 文件访问模式标志：先前描述的 O_RDONLY、O_WRONLY 和 O_RDWR 标志均在此列，调用 open() 时，上述三者 in flags 参数中不能同时使用，只能指定其中一种。调用 fcntl() 的 F_GETFL 操作能够检索文件的访问模式 (见 5.3 节)。
- 文件创建标志：这些标志在表 4-3 中位于第二部分，其控制范围不拘于 open() 调用行为的方方面面，还涉及后续 I/O 操作的各个选项。这些标志不能检索，也无法修改。
- 已打开文件的状态标志：这些标志是表 4-3 中的剩余部分，使用 fcntl() 的 F_GETFL 和 F_SETFL 操作可以分别检索和修改此类标志。有时干脆将其称之为文件状态标志。

始于内核版本 2.6.22，读取位于 /proc/PID/fdinfo 目录下的 linux 系统专有文件，可以获取系统内任一进程中文件描述符的相关信息。针对进程中每一个已打开的文件描述符，该目录下都有相应文件，以对应文件描述符的数值命名。文件中的 pos 字段表示当前的文件偏移量 (4.7 节)。而 flags 字段则为一个八进制数，表征文件访问标志和已打开文件的状态标志。(该数字的解码需要参考这些标志在 C 语言函数库头文件中所定义的数值。)

如下是 flags 常量的详细描述。

O_APPEND

总是在文件尾部追加数据，5.1 节将讨论此标志的意义。

O_ASYNC

当对于 open() 调用所返回的文件描述符可以实施 I/O 操作时，系统会产生一个信号通知进程。这一特性，也被称为信号驱动 I/O，仅对特定类型的文件有效，诸如终端、FIFOS 及 socket。(在 SUSv3 中并未规定 O_ASYNC 标志，但大多数 UNIX 实现都支持此标志或者老版本中与其等效的 FASYNC 标志。) 在 Linux 中，调用 open() 时指定 O_ASYNC 标志没有任何实质效果。要启用信号驱动 I/O 特性，必须调用 fcntl() 的 F_SETFL 操作来设置 O_ASYNC 标志 (见 5.3 节)。(其他一些 UNIX 系统的实现有类似行为。) 关于 O_ASYNC 标志的更多内容请参考 63.3 节。

O_CLOEXEC (自 Linux 2.6.23 版本开始支持)

为新 (创建) 的文件描述符启用 close-on-flag 标志 (FD_CLOEXEC)。27.4 节将描述 FD_CLOEXEC 标志。使用 O_CLOEXEC 标志 (打开文件)，可以免去程序执行 fcntl() 的 F_GETFD 和 F_SETFD 操作来设置 close-on-exec 标志的额外工作。在多线程程序中执行 fcntl() 的 F_GETFD 和 F_SETFD 操作有可能导致竞争状态，而使用 O_CLOEXEC 标志则能够避免这一点。可能引发竞争的场景是：线程某甲打开一文件描述符，尝试为该描述符标记 close-on-exec 标志，于此同时，线程某乙执行 fork() 调用，然后调用 exec() 执行任意一个程序。(假设在某甲打开文件描述符和调用

`fcntl()`设置 `close-on-exec` 标志之间, 某乙成功地执行了 `fork()`和 `exec()`操作。) 此类竞争可能会在无意间将打开的文件描述符泄露给不安全的程序。(更多关于竞争状态的内容请参考 5.1 节。)

O_CREAT

如果文件不存在, 将创建一个新的空文件。即使文件以只读方式打开, 此标志依然有效。如果在 `open()`调用中指定 `O_CREAT` 标志, 那么还需要提供 `mode` 参数, 否则, 会将新文件的权限设置为栈中的某个随机值。

O_DIRECT

无系统缓冲的文件 I/O 操作。该特性将在 13.6 节中详述。为使 `O_DIRECT` 标志的常量定义在 `<fcntl.h>`中有效, 必须定义 `_GNU_SOURCE` 功能测试宏。

O_DIRECTORY

如果 `pathname` 参数并非目录, 将返回错误 (错误号 `errno` 为 `ENOTDIR`)。这一标志是专为实现 `opendir()`函数 (18.8 节) 而设计的扩展标志。为使 `O_DIRECTORY` 标志的常量定义在 `<fcntl.h>`中有效, 必须定义 `_GNU_SOURCE` 功能测试宏。

O_DSYNC (自 Linux 2.6.33 版本开始支持)

根据同步 I/O 数据完整性的完成要求¹来执行文件写操作。参见 13.3 节中关于内核 I/O 缓冲的讨论。

O_EXCL

此标志与 `O_CREAT` 标志结合使用表明如果文件已经存在, 则不会打开文件, 且 `open()`调用失败, 并返回错误, 错误号 `errno` 为 `EEXIST`。换言之, 此标志确保了调用者 (`open()`的调用进程) 就是创建文件的进程。检查文件存在与否和创建文件这两步属于同一原子操作。5.1 节将讨论原子操作的概念。如果在 `flags` 参数中同时指定了 `O_CREAT` 和 `O_EXCL` 标志, 且 `pathname` 参数是符号链接, 则 `open()`函数调用失败 (错误号 `errno` 为 `EEXIST`)。SUSv3 之所以如此规定, 是要求有特权的应用程序在已知目录下创建文件, 从而消除了如下安全隐患, 使用符号链接打开文件会导致在另一位置创建文件 (例如, 系统目录)。

O_LARGEFILE

支持以大文件方式打开文件。在 32 位操作系统中使用此标志, 以支持大文件操作。尽管在 SUSv3 中没有规定这一标志, 但其他一些 UNIX 实现都支持这一特性。此标志在诸如 Alpha、IA-64 之类的 64 位 Linux 实现中是无效的。更多的内容将在 5.10 节中讨论。

O_NOATIME (自 Linux 2.6.8 版本开始)

在读文件时, 不更新文件的最近访问时间 (15.1 节中所描述的 `st_atime` 属性)。要使用该标志, 要么调用进程的有效用户 ID 必须与文件的拥有者相匹配, 要么进程需要拥有特权 (`CAP_FOWNER`)。否则, `open()`调用失败, 并返回错误, 错误号 `errno` 为 `EPERM`。(事实上, 如 9.5 节所述, 对于非特权进程, 当以 `O_NOATIME` 标志打开文件时, 与文件用户 ID 必须匹配的是进程的文件系统用户 ID, 而非进程的有效用户 ID。) 此标志是 Linux 特有的非标准扩展。要从 `<fcntl.h>`中启用此标志, 必须定义 `_GNU_SOURCE` 功能测试宏。 `O_NOATIME` 标志的设计旨在为索引和备份程序服务。该标志的使用能够显著减少磁盘的活动量, 省却了既

¹ 译者注: 所谓 `synchronized I/O data integration completion` 在 SUS 的 `base definition 3.374` 中有详细定义, 但学究气十足, 语焉不详。建议参考《UNIX 环境高级编程》v2 一书 (后续译注中简称为 APUEv2) 3.3 节关于 `O_DSYNC` 的描述。

要读取文件内容，又要更新文件 i-node 结构中最近访问时间的繁琐，进而节省了磁头在磁盘上的反复寻道时间（14.4 节）。`mount()`函数中 `MS_NOATIME` 标志（14.8.1 节）和 `FS_NOATIME_FL` 标志（15.5 节）与 `O_NOATIME` 标志功能相似。

`O_NOCTTY`

如果正在打开的文件属于终端设备，`O_NOCTTY` 标志防止其成为控制终端。34.4 节将讨论控制终端。如果正在打开的文件不是终端设备，则此标志无效。

`O_NOFOLLOW`

通常，如果 `pathname` 参数是符号链接，`open()`函数将对 `pathname` 参数进行解引用。一旦在 `open()`函数中指定了 `O_NOFOLLOW` 标志，且 `pathname` 参数属于符号链接，则 `open()`函数将返回失败（错误号 `errno` 为 `ELOOP`）。此标志在特权程序中极为有用，能够确保 `open()`函数不对符号链接进行解引用。为使 `O_NOFOLLOW` 标志在 `<fcntl.h>`中有效，必须定义 `_GNU_SOURCE` 功能测试宏。

`O_NONBLOCK`

以非阻塞方式打开文件，参照 5.9 节。

`O_SYNC`

以同步 I/O 方式打开文件，参见 13.3 节针对内核 I/O 缓冲的讨论。

`O_TRUNC`

如果文件已经存在且为普通文件，那么将清空文件内容，将其长度置 0。在 Linux 下使用此标志，无论以读、写方式打开文件，都可清空文件内容（在这两种情况下，都必须拥有对文件的写权限）。SUSv3 对 `O_RDONLY` 与 `O_TRUNC` 标志的组合未作规定，但多数其他 UNIX 实现与 Linux 的处理方式相同。

4.3.2 `open()`函数的错误

若打开文件时发生错误，`open()`将返回 -1，错误号 `errno` 标识错误原因。以下是一些可能发生的错误（除了在上节参数描述中已经提及的错误之外）。

`EACCES`

文件权限不允许调用进程以 `flags` 参数指定的方式打开文件。无法访问文件，其可能的原因有目录权限的限制、文件不存在并且也无法创建该文件。

`EISDIR`

所指定的文件属于目录，而调用者企图打开该文件进行写操作。不允许这种用法。（另一方面，在某些场合中，打开目录进行读操作是必要的。18.11 节将举例说明。）

`EMFILE`

进程已打开的文件描述符数量达到了进程资源限制所设定的上限（在 36.3 节将描述 `RLIMIT_NOFILE` 参数）。

`ENFILE`

文件打开数量已经达到系统允许的上限。

`ENOENT`

要么文件不存在且未指定 `O_CREAT` 标志，要么指定了 `O_CREAT` 标志，但 `pathname` 参

数所指定路径的目录之一不存在，或者 `pathname` 参数为符号链接，而该链接指向的文件不存在（空链接）。

EROFS

所指定的文件隶属于只读文件系统，而调用者企图以写方式打开文件。

ETXTBSY

所指定的文件为可执行文件（程序），且正在运行。系统不允许修改正在运行的程序（比如以写方式打开文件）。（必须首先终止程序运行，然后方可修改可执行文件。）

后续在描述其他系统调用或库函数时，一般不会再以上述方式展现可能发生的一系列错误。（每个系统调用或库函数的错误列表可从相关操作手册中查询获得。）采用上述方式原因有二，一是因为 `open()` 是本书详细描述的首个系统调用，而上述列表表明任一原因都有可能导导致系统调用或库函数的调用失败。二是 `open()` 调用失败的具体原因列表本身就颇为值得玩味，它展示了影响文件访问的若干因素，以及访问文件时系统所执行的一系列检查。（上述错误列表并不完整，更多 `open()` 调用失败的错误原因请查看 `open(2)` 的操作手册。）

4.3.3 `creat()` 系统调用

在早期的 UNIX 实现中，`open()` 只有两个参数，无法创建新文件，而是使用 `creat()` 系统调用来创建并打开一个新文件。

```
#include <fcntl.h>

int creat(const char *pathname, mode_t mode);

Returns file descriptor, or -1 on error
```

`creat()` 系统调用根据 `pathname` 参数创建并打开一个文件，若文件已存在，则打开文件，并清空文件内容，将其长度清 0。`creat()` 返回一文件描述符，供后续系统调用使用。`creat()` 系统调用等价于如下 `open()` 调用：

```
fd = open(pathname, O_WRONLY | O_CREAT | O_TRUNC, mode);
```

尽管 `creat()` 在一些老旧程序的代码中还时有所见，但由于 `open()` 的 `flags` 参数能对文件打开方式提供更多控制（例如：可以指定 `O_RDWR` 标志，代替 `O_WRONLY` 标志），对 `creat()` 的使用现在已不多见。

4.4 读取文件内容：`read()`

`read()` 系统调用从文件描述符 `fd` 所指代的打开文件中读取数据。

```
#include <unistd.h>

ssize_t read(int fd, void *buffer, size_t count);

Returns number of bytes read, 0 on EOF, or -1 on error
```

`count` 参数指定最多能读取的字节数。（`size_t` 数据类型属于无符号整数类型。）`buffer` 参数提供用来存放输入数据的内存缓冲区地址。缓冲区至少应有 `count` 个字节。

系统调用不会分配内存缓冲区用以返回信息给调用者。所以，必须预先分配大小合适的缓冲区并将缓冲区指针传递给系统调用。与此相反，有些库函数却会分配内存缓冲区用以返回信息给调用者。

如果 `read()`调用成功，将返回实际读取的字节数，如果遇到文件结束（EOF）则返回 0，如果出现错误则返回-1。`ssize_t` 数据类型属于有符号的整数类型，用来存放（读取的）字节数或-1（表示错误）。

一次 `read()`调用所读取的字节数可以小于请求的字节数。对于普通文件而言，这有可能是因为当前读取位置靠近文件尾部。

当 `read()`应用于其他文件类型时，比如管道、FIFO、`socket` 或者终端，在不同环境下也会出现 `read()`调用读取的字节数小于请求字节数的情况。例如，默认情况下从终端读取字符，一遇到换行符（`\n`），`read()`调用就会结束。在后续章节论及其他类型文件时，会再次针对这些情况进行探讨。

使用 `read()`从终端读取一连串字符，我们也许期望下面的代码会起作用：

```
#define MAX_READ 20
char buffer[MAX_READ];

if (read(STDIN_FILENO, buffer, MAX_READ) == -1)
    errExit("read");
printf("The input data was: %s\n", buffer);
```

这段代码的输出可能会很奇怪，因为输出结果除了实际输入的字符串外还会包括其他字符。这是因为 `read()`调用没有在 `printf()`函数打印的字符串尾部添加一个表示终止的空字符。思索片刻就会意识到这肯定是症结所在，因为 `read()`能够从文件中读取任意序列的字节。有些情况下，输入信息可能是文本数据，但在其他情况下，又可能是二进制整数或者二进制形式的 C 语言数据结构。`read()`无从区分这些数据，故而也无法遵从 C 语言对字符串处理的约定，在字符串尾部追加标识字符串结束的空字符。如果输入缓冲区的结尾处需要一个表示终止的空字符，必须显式追加。

```
char buffer[MAX_READ + 1];
ssize_t numRead;

numRead = read(STDIN_FILENO, buffer, MAX_READ);
if (numRead == -1)
    errExit("read");

buffer[numRead] = '\0';
printf("The input data was: %s\n", buffer);
```

由于表示字符串终止的空字符需要一个字节的内存，所以缓冲区的大小至少要比预计读取的最大字符串长度多出 1 个字节。

4.5 数据写入文件：write()

`write()`系统调用将数据写入一个已打开的文件中。

```
#include <unistd.h>

ssize_t write(int fd, void *buffer, size_t count);

Returns number of bytes written, or -1 on error
```

`write()`调用的参数含义与 `read()`调用相类似。`buffer` 参数为要写入文件中数据的内存地址，`count` 参数为欲从 `buffer` 写入文件的数据字节数，`fd` 参数为一文件描述符，指代数据要写入的文件。

如果 `write()`调用成功，将返回实际写入文件的字节数，该返回值可能小于 `count` 参数值。这被称为“部分写”。对磁盘文件来说，造成“部分写”的原因可能是由于磁盘已满，或是因为进程资源对文件大小的限制。（相关的限制为 `RLIMIT_FSIZE`，将在 36.3 节描述。）

对磁盘文件执行 I/O 操作时，`write()`调用成功并不能保证数据已经写入磁盘。因为为了减少磁盘活动量和加快 `write()`系统调用，内核会缓存磁盘的 I/O 操作，第 13 章将会详加介绍。

4.6 关闭文件：close()

`close()`系统调用关闭一个打开的文件描述符，并将其释放回调用进程，供该进程继续使用。当一进程终止时，将自动关闭其已打开的所有文件描述符。

```
#include <unistd.h>

int close(int fd);

Returns 0 on success, or -1 on error
```

显式关闭不再需要的文件描述符往往是良好的编程习惯，会使代码在后续修改时更具可读性，也更可靠。进而言之，文件描述符属于有限资源，因此文件描述符关闭失败可能会导致一个进程将文件描述符资源消耗殆尽。在编写需要长期运行并处理大量文件的程序时，比如 `shell` 或者网络服务器软件，需要特别加以关注。

像其他所有系统调用一样，应对 `close()`的调用进行错误检查，如下所示：

```
if (close(fd) == -1)
    errExit("close");
```

上述代码能够捕获的错误有：企图关闭一个未打开的文件描述符，或者两次关闭同一文件描述符，也能捕获特定文件系统在关闭操作中诊断出的错误条件。

针对特定文件系统的错误，NFS（网络文件系统）就是一例。如果 NFS 出现提交失败，这意味着数据没有抵达远程磁盘，随之将这一错误作为 `close()`调用失败的原因传递给应用系统。

4.7 改变文件偏移量：lseek()

对于每个打开的文件，系统内核会记录其文件偏移量，有时也将文件偏移量称为读写偏移量或指针。文件偏移量是指执行下一个 `read()`或 `write()`操作的文件起始位置，会以相对于文件头部起始点的文件当前位置来表示。文件第一个字节的偏移量为 0。

文件打开时，会将文件偏移量设置为指向文件开始，以后每次 `read()`或 `write()`调用将自动对其进行调整，以指向已读或已写数据后的下一字节。因此，连续的 `read()`和 `write()`调用将按顺序递进，对文件进行操作。

针对文件描述符 `fd` 参数所指代的已打开文件，`lseek()`系统调用依照 `offset` 和 `whence` 参数值调整该文件的偏移量。

```
#include <unistd.h>

off_t lseek(int fd, off_t offset, int whence);

Returns new file offset if successful, or -1 on error
```

`offset` 参数指定了一个以字节为单位的数值。（SUSv3 规定 `off_t` 数据类型为有符号整型数。）`whence` 参数则表明应参照哪个基点来解释 `offset` 参数，应为下列其中之一：

SEEK_SET

将文件偏移量设置为从文件头部起始点开始的 `offset` 个字节。

SEEK_CUR

相对于当前文件偏移量，将文件偏移量调整 `offset` 个字节¹。

SEEK_END

将文件偏移量设置为起始于文件尾部的 `offset` 个字节。也就是说，`offset` 参数应该从文件最后一个字节之后的下一个字节算起。

图 4-1 展示了 `whence` 参数的含义。

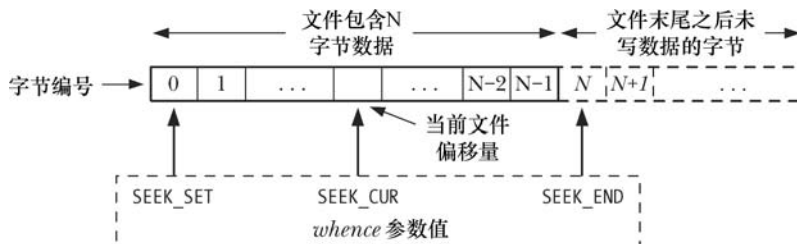


图 4-1: 解释 `lseek()`函数中 `whence` 参数

在早期的 UNIX 实现中，`whence` 参数用整数 0、1、2 来表示，而非正文中显示的 `SEEK_*`常量。BSD 的早期版本使用另一套命名：`L_SET`、`L_INCR` 和 `L_XTND` 来表示 `whence` 参数。

如果 `whence` 参数值为 `SEEK_CUR` 或 `SEEK_END`，`offset` 参数可以为正数也可以为负数；如果 `whence` 参数值为 `SEEK_SET`，`offset` 参数值必须为非负数。

`lseek()`调用成功会返回新的文件偏移量。下面的调用只是获取文件偏移量的当前位置，并没有修改它。

```
curr = lseek(fd, 0, SEEK_CUR);
```

有些 UNIX 系统（Linux 不在此列）实现了非标准的 `tell(fd)`函数，其调用目的与上述 `lseek()`相同。

¹ 译者注：简而言之，相对于文件头部的绝对偏移量=当前文件偏移量+`offset`。

这里给出了 `lseek()`调用的其他一些例子，在注释中说明了将文件偏移量移到的具体位置。

```
lseek(fd, 0, SEEK_SET);      /* Start of file */
lseek(fd, 0, SEEK_END);     /* Next byte after the end of the file */
lseek(fd, -1, SEEK_END);    /* Last byte of file */
lseek(fd, -10, SEEK_CUR);   /* Ten bytes prior to current location */
lseek(fd, 10000, SEEK_END); /* 10001 bytes past last byte of file */
```

`lseek()`调用只是调整内核中与文件描述符相关的文件偏移量记录，并没有引起对任何物理设备的访问。

5.4 节将进一步描述文件偏移量、文件描述符、已打开文件三者之间的关系。

`lseek()`并不适用于所有类型的文件。不允许将 `lseek()`应用于管道、FIFO、socket 或者终端。一旦如此，调用将会失败，并将 `errno` 置为 `ESPIPE`。另一方面，只要合情合理，也可以将 `lseek()`应用于设备。例如，在磁盘或者磁带上查找一处具体位置。

`lseek()`调用名中的 `l` 源于这样一个事实：`offset` 参数和调用返回值的类型起初都是 `long` 型。早期的 UNIX 系统还提供了 `seek()`系统调用，当时这两个值的类型为 `int` 型。

文件空洞

如果程序的文件偏移量已然跨越了文件结尾，然后再执行 I/O 操作，将会发生什么情况？`read()`调用将返回 0，表示文件结尾。有点令人惊讶的是，`write()`函数可以在文件结尾后的任意位置写入数据。

从文件结尾后到新写入数据间的这段空间被称为文件空洞。从编程角度看，文件空洞中是存在字节的，读取空洞将返回以 0（空字节）填充的缓冲区。

然而，文件空洞不占用任何磁盘空间。直到后续某个时点，在文件空洞中写入了数据，文件系统才会为之分配磁盘块。文件空洞的主要优势在于，与为实际需要的空字节分配磁盘块相比，稀疏填充的文件会占用较少的磁盘空间。核心转储文件（`core dump`）（见 22.1 节）是包含空洞文件的常见例子。

对于文件空洞不占用磁盘空间的说法需要稍微限定一下。在大多数文件系统中，文件空间的分配是以块为单位的（14.3 节）。块的大小取决于文件系统，通常是 1024 字节、2048 字节、4096 字节。如果空洞的边界落在块内，而非恰好落在块边界上，则会分配一个完整的块来存储数据，块中与空洞相关的部分则以空字节填充。

大多数“原生”UNIX 文件系统都支持文件空洞的概念，但很多“非原生”文件系统（比如，微软的 VFAT）并不支持这一概念。不支持文件空洞的文件系统会显式地将空字节写入文件。

空洞的存在意味着一个文件名义上的大小可能要比其占用的磁盘存储总量要大（有时会大出许多）。向文件空洞中写入字节，内核需要为其分配存储单元，即使文件大小不变，系统的可用磁盘空间也将减少。这种情况并不常见，但也需要了解。

SUSv3 的函数 `posix_fallocate(fd, offset, len)`规定，针对文件描述符 `fd` 所指代的文件，能确保按照由 `offset` 参数和 `len` 参数所确定的字节范围为其在磁盘上分配存储空间。这样，应用程序对文件的后续 `write()`调用不会因磁盘空间耗尽而失败（否则，当文件中一个空洞被填满后，或者因其他应用程序消耗了磁盘空间时，都可能因磁盘空间耗尽而引发此类错误）。在过去，`glibc` 库在实现 `posix_fallocate()`函数时，通过向指定范围内的每个块写入一个值为 0 的字节以达到预期结果。自内核版本 2.6.23 开始，Linux 系统提供了 `fallocate()`系

统调用，能更为高效地确保所需存储空间的分配。当 `fallocate()`调用可用时，`glibc` 库会利用其来实现 `posix_fallocate()`函数的功能。

14.4 节将描述空洞在文件中的表示方式。15.1 节将描述 `stat()`系统调用，该调用能够提供文件当前大小和实际分配给文件的块数量等信息。

示例程序

程序清单 4-3 演示了 `lseek()`与 `read()`、`write()`的协作使用。该程序的第一个命令行参数为将要打开的文件名称，余下的参数则指定了在文件上执行的输入/输出操作。每个表示操作的参数都以一个字母开头，紧跟以相关值（中间无空格分隔）。

- `soffset`: 从文件开始检索到 `offset` 字节位置。
- `rlength`: 在当前文件偏移量处，从文件中读取 `length` 字节数据，并以文本形式显示。
- `Rlength`: 在当前文件偏移量处，从文件中读取 `length` 字节数据，并以十六进制形式显示。
- `wstr`: 在当前文件偏移量处，向文件写入由 `str` 指定的字符串。

程序清单 4-3: `read()`、`write()`和 `lseek()`的使用示范

```
----- fileio/seek_io.
#include <sys/stat.h>
#include <fcntl.h>
#include <ctype.h>
#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    size_t len;
    off_t offset;
    int fd, ap, j;
    char *buf;
    ssize_t numRead, numWritten;

    if (argc < 3 || strcmp(argv[1], "--help") == 0)
        usageErr("%s file {r<length>|R<length>|w<string>|s<offset>}...\n",
                 argv[0]);

    fd = open(argv[1], O_RDWR | O_CREAT,
              S_IRUSR | S_IWUSR | S_IRGRP | S_IWGRP |
              S_IROTH | S_IWOTH);          /* rw-rw-rw- */
    if (fd == -1)
        errExit("open");

    for (ap = 2; ap < argc; ap++) {
        switch (argv[ap][0]) {
            case 'r': /* Display bytes at current offset, as text */
            case 'R': /* Display bytes at current offset, in hex */
                len = getLong(&argv[ap][1], GN_ANY_BASE, argv[ap]);
                buf = malloc(len);
                if (buf == NULL)
                    errExit("malloc");

                numRead = read(fd, buf, len);
                if (numRead == -1)
```

```

        errExit("read");

    if (numRead == 0) {
        printf("%s: end-of-file\n", argv[ap]);
    } else {
        printf("%s: ", argv[ap]);
        for (j = 0; j < numRead; j++) {
            if (argv[ap][0] == 'r')
                printf("%c", isprint((unsigned char) buf[j]) ?
                    buf[j] : '?');
            else
                printf("%02x ", (unsigned int) buf[j]);
        }
        printf("\n");
    }

    free(buf);
    break;

case 'w': /* Write string at current offset */
    numWritten = write(fd, &argv[ap][1], strlen(&argv[ap][1]));
    if (numWritten == -1)
        errExit("write");
    printf("%s: wrote %ld bytes\n", argv[ap], (long) numWritten);
    break;

case 's': /* Change file offset */
    offset = getLong(&argv[ap][1], GN_ANY_BASE, argv[ap]);
    if (lseek(fd, offset, SEEK_SET) == -1)
        errExit("lseek");
    printf("%s: seek succeeded\n", argv[ap]);
    break;

default:
    cmdlineErr("Argument must start with [rRws]: %s\n", argv[ap]);
}

}

exit(EXIT_SUCCESS);
}

```

fileio/seek_io.c

下面的 shell 会话演示了程序清单 4-3 程序的使用，还显示了从文件空洞中读取字节时的情况：

```

$ touch tfile Create new, empty file
$ ./seek_io tfile s100000 wabc Seek to offset 100,000, write "abc"
s100000: seek succeeded
wabc: wrote 3 bytes
$ ls -l tfile Check size of file
-rw-r--r-- 1 mtk users 100003 Feb 10 10:35 tfile
$ ./seek_io tfile s10000 R5 Seek to offset 10,000, read 5 bytes from hole
s10000: seek succeeded
R5: 00 00 00 00 00 Bytes in the hole contain 0

```

4.8 通用 I/O 模型以外的操作：ioctl()

在本章上述通用 I/O 模型之外，ioctl() 系统调用又为执行文件和设备操作提供了一种多用途机制。

```
#include <sys/ioctl.h>
```

```
int ioctl(int fd, int request, ... /* argp */);
```

Value returned on success depends on *request*, or -1 on error

`fd` 参数为某个设备或文件已打开的文件描述符，`request` 参数指定了将在 `fd` 上执行的控制操作。具体设备的头文件定义了可传递给 `request` 参数的常量。

`ioctl()`调用的第三个参数采用了标准 C 语言的省略符号 (...) 来表示（称之为 `argp`），可以是任意数据类型。`ioctl()`根据 `request` 的参数值来确定 `argp` 所期望的类型。通常情况下，`argp` 是指向整数或结构的指针，有些情况下，不需要使用 `argp`。

后面各章中将会有许多 `ioctl()`的用法展示（例如 15.5 节）。

SUSv3 为 `ioctl()`制定的唯一规定是针对流（STREAM）设备的控制操作。（流是 System V 操作系统中的特性。尽管为其开发有一些插件，主流的 Linux 内核并不支持该特性。）本书述及的 `ioctl()`的其他操作都不在 SUSv3 的规范之列。然而，从早期版本开始，`ioctl()`调用就是 UNIX 系统的一部分，因此本书所描述的几个 `ioctl()`操作在许多其他 UNIX 系统中都已实现。在讨论 `ioctl()`调用的各个操作时，会点出存在的可移植性问题。

4.9 总结

为了对普通文件执行 I/O 操作，首先必须调用 `open()`以获得一个文件描述符。随之使用 `read()`和 `write()`执行文件的 I/O 操作，然后应使用 `close()`释放文件描述符及相关资源。这些系统调用可对所有类型的文件执行 I/O 操作。

所有类型的文件和设备驱动都实现了相同的 I/O 接口，这保证了 I/O 操作的通用性，同时也意味着在无需针对特定文件类型编写代码的情况下，程序通常就能操作所有类型的文件。

对于已打开的每个文件，内核都维护有一个文件偏移量，这决定了下一次读或写操作的起始位置。读和写操作会隐式修改文件偏移量。使用 `lseek()`函数可以显式地将文件偏移量置为文件中或文件结尾后的任一位置。在文件原结尾处之后的某一位置写入数据将导致文件空洞。从文件空洞处读取文件将返回全 0 字节。

对于未纳入标准 I/O 模型的所有设备和文件操作而言，`ioctl()`系统调用是个“百宝箱”。

4.10 练习

- 4-1. `tee` 命令是从标准输入中读取数据，直至文件结尾，随后将数据写入标准输出和命令行参数所指定的文件。（44.7 节讨论 FIFO 时，会展示使用 `tee` 命令的一个例子。）请使用 I/O 系统调用实现 `tee` 命令。默认情况下，若已存在与命令行参数指定文件同名的文件，`tee` 命令会将其覆盖。如文件已存在，请实现 `-a` 命令行选项（`tee -a file`）在文件结尾处追加数据。（请参考附录 B 中对 `getopt()`函数的描述来解析命令行选项。）
- 4-2. 编写一个类似于 `cp` 命令的程序，当使用该程序复制一个包含空洞（连续的空字节）的普通文件时，要求目标文件的空洞与源文件保持一致。

第 5 章

深入探究文件 I/O

本章将延续上一章的讨论，进一步探究文件 I/O。

在后续的关于 `open()` 系统调用的探讨中，将引入原子（atomicity）操作的概念——将某一系统调用所要完成的各个动作作为不可中断的操作，一次性加以执行。原子操作是许多系统调用得以正确执行的必要条件。

本章还将介绍另一个与文件操作相关的系统调用：多用途的 `fcntl()`，并展示其应用之一读取和设置打开文件的状态标志。

随后，本章将审视用于表示文件描述符和已打开文件的内核数据结构。后续各章将探讨文件 I/O 的某些微妙之处，理解这些数据结构之间的关系对此将有所助益。基于这一模型，本章还将解释如何复制文件描述符。

之后，本章将讨论一些支持扩展读写功能的系统调用。此类调用可以在不改变文件当前偏移量的情况下，在文件的特定位置处进行读写操作，以及对程序中多个缓冲区进行数据（双向）传输。

最后，将简要介绍非阻塞 I/O 的概念，并述及一些用于读写大文件的扩展接口。

此外，因为临时文件在许多系统程序中有广泛的应用，所以本章也会介绍一些相关库函数：在保证随机生成唯一文件名称的同时，用于创建和操作临时文件。

5.1 原子操作和竞争条件

在探究系统调用时会反复涉及原子操作的概念。所有系统调用都是以原子操作方式执行的。之所以这么说，是指内核保证了某系统调用中的所有步骤会作为独立操作而一次性加以执行，其间不会为其他进程或线程所中断。

原子性是某些操作得以圆满成功的关键所在。特别是它规避了竞争状态（race conditions）（有时也称为竞争冒险）。竞争状态是这样一种情形：操作共享资源的两个进程（或线程），其结果取决于一个无法预期的顺序，即这些进程¹获得 CPU 使用权的先后相对顺序。

接下来，将讨论涉及文件 I/O 的两种竞争状态，并展示了如何使用 `open()` 的标志位，来保

¹ 译者注：或线程。

证相关文件操作的原子性，从而消除这些竞争状态。

22.9 节将介绍 `sigsuspend()` 系统调用。24.4 节将介绍 `fork()` 调用，届时将再次探讨竞争状态。

以独占方式创建一个文件

4.3.1 节曾提及：当同时指定 `O_EXCL` 与 `O_CREAT` 作为 `open()` 的标志位时，如果要打开的文件已然存在，则 `open()` 将返回一个错误。这提供了一种机制，保证进程是打开文件的创建者。对文件是否存在的检查和创建文件属于同一原子操作。要理解这一点的重要性，请思考程序清单 5-1 所示代码，该段代码中并未使用 `O_EXCL` 标志。（在此，为了对执行该程序的不同进程加以区分，在输出信息中打印有通过调用 `getpid()` 所返回的进程号。）

程序清单 5-1：试图以独占方式打开文件的错误代码

```
----- from fileio/bad_exclusive_open.c
fd = open(argv[1], O_WRONLY);      /* Open 1: check if file exists */
if (fd != -1) {                    /* Open succeeded */
    printf("[PID %ld] File \"%s\" already exists\n",
           (long) getpid(), argv[1]);
    close(fd);
} else {
    if (errno != ENOENT) {          /* Failed for unexpected reason */
        errExit("open");
    } else {
        /* WINDOW FOR FAILURE */
        fd = open(argv[1], O_WRONLY | O_CREAT, S_IRUSR | S_IWUSR);
        if (fd == -1)
            errExit("open");

        printf("[PID %ld] Created file \"%s\" exclusively\n",
               (long) getpid(), argv[1]);      /* MAY NOT BE TRUE! */
    }
}
----- from fileio/bad_exclusive_open.c
```

程序清单 5-1 中所示的代码，除了要啰啰嗦嗦地调用 `open()` 两次外，还潜伏着一个 bug。假设如下情况：当第一次调用 `open()` 时，希望打开的文件还不存在，而当第二次调用 `open()` 时，其他进程已经创建了该文件。如图 5-1 所示，若内核调度器判断出分配给 A 进程的时间片已经耗尽，并将 CPU 使用权交给 B 进程，就可能会发生这种问题。再比如两个进程在一个多 CPU 系统上同时运行时，也会出现这种情况。图 5-1 展示了两个进程同时执行程序清单 5-1 中代码的情形。在这一场景下，进程 A 将得出错误的结论：目标文件是由自己创建的。因为无论目标文件存在与否，进程 A 对 `open()` 的第二次调用都会成功。

虽然进程将自己误认为文件创建者的可能性相对较小，但毕竟是存在的，这已然将此段代码置于不可靠的境地。操作的结果将依赖于对两个进程的调度顺序，这一事实也就意味着出现了竞争状态。

为了说明这段代码的确存在问题，可以用一段代码替换程序清单 5-1 中的注释行“处理文件不存在的情况”，在检查文件是否存在与创建文件这两个动作之间人为制造一个长时间的等待。

```
printf("[PID %ld] File \"%s\" doesn't exist yet\n", (long) getpid(), argv[1]);
if (argc > 2) {                    /* Delay between check and create */
    sleep(5);                      /* Suspend execution for 5 seconds */
    printf("[PID %ld] Done sleeping\n", (long) getpid());
}
```

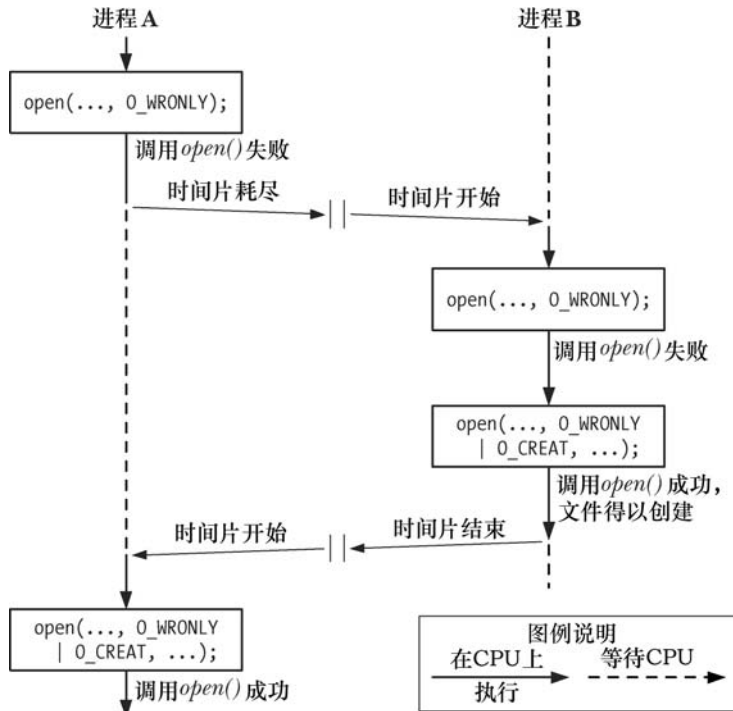


图 5-1: 未能以独占方式创建文件

sleep()库函数可将当前执行的进程挂起指定的秒数。23.4 节将讨论该函数。

如果同时运行程序清单 5-1 中程序的两个实例，两个进程都会声称自己以独占方式创建了文件。

```

$ ./bad_exclusive_open tfile sleep &
[PID 3317] File "tfile" doesn't exist yet
[1] 3317
$ ./bad_exclusive_open tfile
[PID 3318] File "tfile" doesn't exist yet
[PID 3318] Created file "tfile" exclusively
$ [PID 3317] Done sleeping
[PID 3317] Created file "tfile" exclusively           Not true

```

从上面输出的倒数第二行可以发现，shell 提示符里夹杂了第一个实例的输出信息。

由于第一个进程在检查文件是否存在和创建文件之间发生了中断，造成两个进程都声称自己是文件的创建者。结合 O_CREAT 和 O_EXCL 标志来一次性地调用 open()可以防止这种情况，因为这确保了检查文件和创建文件的步骤属于一个单一的原子（即不可中断的）操作。

向文件尾部追加数据

用以说明原子操作必要性的第二个例子是：多个进程同时向同一个文件（例如，全局日志文件）尾部添加数据。为了达到这一目的，也许可以考虑在每个写进程中使用如下代码。

```

if (lseek(fd, 0, SEEK_END) == -1)
    errExit("lseek");
if (write(fd, buf, len) != len)
    fatal("Partial/failed write");

```

但是，这段代码存在的缺陷与前一个例子如出一辙。如果第一个进程执行到 `lseek()` 和 `write()` 之间，被执行相同代码的第二个进程所中断，那么这两个进程会在写入数据前，将文件偏移量设为相同位置，而当第一个进程再次获得调度时，会覆盖第二个进程已写入的数据。此时再次出现了竞争状态，因为执行的结果依赖于内核对两个进程的调度顺序。

要规避这一问题，需要将文件偏移量的移动与数据写操作纳入同一原子操作。在打开文件时加入 `O_APPEND` 标志就可以保证这一点。

有些文件系统（例如 NFS）不支持 `O_APPEND` 标志。在这种情况下，内核会选择按如上代码所示的方式，施之以非原子操作的调用序列，从而可能导致上述的文件脏写入问题。

5.2 文件控制操作：fcntl()

`fcntl()` 系统调用对一个打开的文件描述符执行一系列控制操作。

```
#include <fcntl.h>

int fcntl(int fd, int cmd, ...);

Return on success depends on cmd, or -1 on error
```

`cmd` 参数所支持的操作范围很广。本章随后各节会对其中的部分操作加以研讨，剩下的操作将在后续各章中进行论述。

`fcntl()` 的第三个参数以省略号来表示，这意味着可以将其设置为不同的类型，或者加以省略。内核会依据 `cmd` 参数（如果有的话）的值来确定该参数的数据类型。

5.3 打开文件的状态标志

`fcntl()` 的用途之一是针对一个打开的文件，获取或修改其访问模式和状态标志（这些值是通过指定 `open()` 调用的 `flag` 参数来设置的）。要获取这些设置，应将 `fcntl()` 的 `cmd` 参数设置为 `F_GETFL`。

```
int flags, accessMode;

flags = fcntl(fd, F_GETFL);          /* Third argument is not required */
if (flags == -1)
    errExit("fcntl");
```

在上述代码之后，可以以如下代码测试文件是否以同步写方式打开：

```
if (flags & O_SYNC)
    printf("writes are synchronized\n");
```

SUSv3 规定：针对一个打开的文件，只有通过 `open()` 或后续 `fcntl()` 的 `F_SETFL` 操作，才能对该文件的状态标志进行设置。然而在如下方面，Linux 实现与标准有所偏离：如果一个程序编译时采用了 5.10 节所提及的打开大文件技术，那么当使用 `F_GETFL` 命令获取文件状态标志时，标志中将总是包含 `O_LARGEFILE` 标志。

判定文件的访问模式有一点复杂，这是因为 `O_RDONLY(0)`、`O_WRONLY(1)`和 `O_RDWR(2)` 这 3 个常量并不与打开文件状态标志中的单个比特位对应。因此，要判定访问模式，需使用掩码 `O_ACCMODE` 与 `flag` 相与，将结果与 3 个常量进行比对，示例代码如下：

```
accessMode = flags & O_ACCMODE;
if (accessMode == O_WRONLY || accessMode == O_RDWR)
    printf("file is writable\n");
```

可以使用 `fcntl()` 的 `F_SETFL` 命令来修改打开文件的某些状态标志。允许更改的标志有 `O_APPEND`、`O_NONBLOCK`、`O_NOATIME`、`O_ASYNC` 和 `O_DIRECT`。系统将忽略对其他标志的修改操作。（有些其他的 UNIX 实现允许 `fcntl()` 修改其他标志，如 `O_SYNC`。）

使用 `fcntl()` 修改文件状态标志，尤其适用于如下场景。

- 文件不是由调用程序打开的，所以程序也无法使用 `open()` 调用来控制文件的状态标志（例如，文件是 3 个标准输入输出描述符中的一员，这些描述符在程序启动之前就被打开）。
- 文件描述符的获取是通过 `open()` 之外的系统调用。比如 `pipe()` 调用，该调用创建一个管道，并返回两个文件描述符分别对应管道的两端。再比如 `socket()` 调用，该调用创建一个套接字并返回指向该套接字的文件描述符。

为了修改打开文件的状态标志，可以使用 `fcntl()` 的 `F_GETFL` 命令来获取当前标志的副本，然后修改需要变更的比特位，最后再次调用 `fcntl()` 函数的 `F_SETFL` 命令来更新此状态标志。因此，为了添加 `O_APPEND` 标志，可以编写如下代码：

```
int flags;

flags = fcntl(fd, F_GETFL);
if (flags == -1)
    errExit("fcntl");
flags |= O_APPEND;
if (fcntl(fd, F_SETFL, flags) == -1)
    errExit("fcntl");
```

5.4 文件描述符和打开文件之间的关系

到目前为止，文件描述符和打开的文件之间似乎呈现出一一对应的关系。然而，实际并非如此。多个文件描述符指向同一打开文件，这既有可能，也属必要。这些文件描述符可在相同或不同的进程中打开。

要理解具体情况如何，需要查看由内核维护的 3 个数据结构。

- 进程级的文件描述符表。
- 系统级的打开文件表。
- 文件系统的 `i-node` 表。

针对每个进程，内核为其维护打开文件的描述符（`open file descriptor`）表。该表的每一条目都记录了单个文件描述符的相关信息，如下所示。

- 控制文件描述符操作的一组标志。（目前，此类标志仅定义了一个，即 `close-on-exec` 标志，将在 27.4 节予以介绍。）
- 对打开文件句柄的引用。

内核对所有打开的文件维护有一个系统级的描述表格（`open file description table`）。有时，也

称之为打开文件表（open file table），并将表中各条目称为打开文件句柄（open file handle）¹。一个打开文件句柄存储了与一个打开文件相关的全部信息，如下所示。

- 当前文件偏移量（调用 read()和 write()时更新，或使用 lseek()直接修改）。
- 打开文件时所使用的状态标志（即，open()的 flags 参数）。
- 文件访问模式（如调用 open()时所设置的只读模式、只写模式或读写模式）。
- 与信号驱动 I/O 相关的设置（见 63.3 节）。
- 对该文件 i-node 对象的引用。

每个文件系统都会为驻留其上的所有文件建立一个 i-node 表。第 14 章将详细讨论 i-node 结构和文件系统的总体结构。这里只是列出每个文件的 i-node 信息，具体如下。

- 文件类型（例如，常规文件、套接字或 FIFO）和访问权限。
- 一个指针，指向该文件所持有的锁的列表。
- 文件的各种属性，包括文件大小以及与不同类型操作相关的时间戳。

此处将忽略 i-node 在磁盘和内存中的表示差异。磁盘上的 i-node 记录了文件的固有属性，诸如：文件类型、访问权限和时间戳。访问一个文件时，会在内存中为 i-node 创建一个副本，其中记录了引用该 i-node 的打开文件句柄数量以及该 i-node 所在设备的主、从设备号，还包括一些打开文件时与文件相关的临时属性，例如：文件锁。

图 5-2 展示了文件描述符、打开的文件句柄以及 i-node 之间的关系。在下图中，两个进程拥有诸多打开的文件描述符。

在进程 A 中，文件描述符 1 和 20 都指向同一个打开的文件句柄（标号为 23）。这可能是通过调用 dup()、dup2()或 fcntl()而形成的（参见 5.5 节）。

进程 A 的文件描述符 2 和进程 B 的文件描述符 2 都指向同一个打开的文件句柄（标号为 73）。这种情形可能在调用 fork()后出现（即，进程 A 与进程 B 之间是父子关系），或者当某进程通过 UNIX 域套接字将一个打开的文件描述符传递给另一进程时，也会发生（参见 61.13.3 节）。

此外，进程 A 的描述符 0 和进程 B 的描述符 3 分别指向不同的打开文件句柄，但这些句柄均指向 i-node 表中的相同条目（1976），换言之，指向同一文件。发生这种情况是因为每个进程各自对同一文件发起了 open()调用。同一个进程两次打开同一文件，也会发生类似情况。

上述讨论揭示出如下要点。

- 两个不同的文件描述符，若指向同一打开文件句柄，将共享同一文件偏移量。因此，如果通过其中一个文件描述符来修改文件偏移量（由调用 read()、write()或 lseek()所致），那么从另一文件描述符中也会观察到这一变化。无论这两个文件描述符分属于不同进程，还是同属于一个进程，情况都是如此。
- 要获取和修改打开的文件标志（例如，O_APPEND、O_NONBLOCK 和 O_ASYNC），可执行 fcntl()的 F_GETFL 和 F_SETFL 操作，其对作用域的约束与上一条颇为类似。
- 相形之下，文件描述符标志（亦即，close-on-exec 标志）为进程和文件描述符所私有。对这一标志的修改将不会影响同一进程或不同进程中的其他文件描述符。

¹ 译者注：为避免混淆，译文将原文中的 open file description table 和 open file description 分别以“打开文件表”和“打开文件句柄”替换。但在给译者的回信中，作者尽管承认这一表述方式容易使读者产生混淆，但仍坚持 open file description table 和 open file description 的称谓，原因有二：一，open file description 是相关标准所采用的术语，而与标准保持一致实属必要；二，handle 通常用于引用用户空间中的应用对象，而此处的 open file description 则无法由用户空间中的应用直接访问。

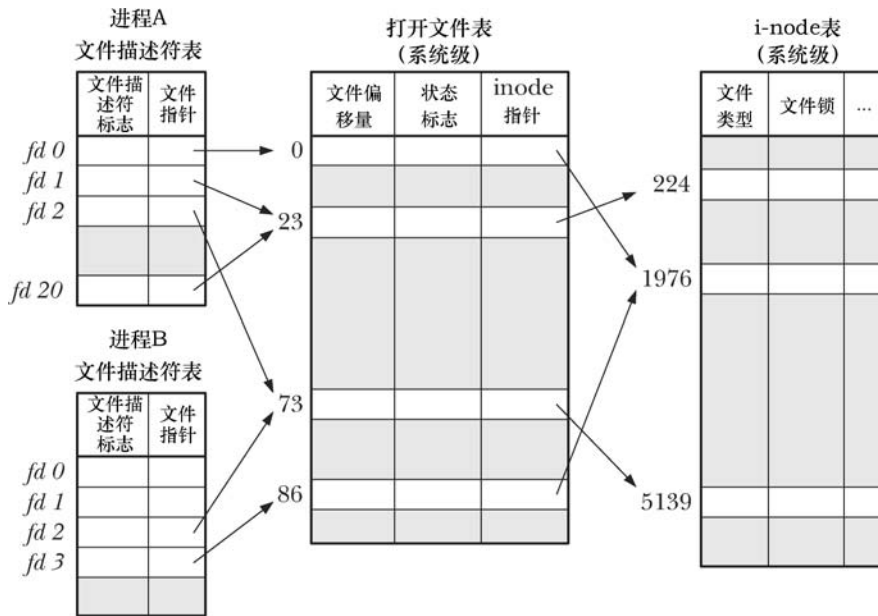


图 5-2: 文件描述符、打开的文件句柄和 i-node 之间的关系

5.5 复制文件描述符

Bourne shell 的 I/O 重定向语法 `2>&1`, 意在通知 shell 把标准错误 (文件描述符 2) 重定向到标准输出 (文件描述符 1)。因此, 下列命令将把 (因为 shell 按从左至右的顺序处理 I/O 重定向语句) 标准输出和标准错误写入 `result.log` 文件:

```
$ ./myscript > results.log 2>&1
```

shell 通过复制文件描述符 ¹ 实现了标准错误的重定向操作, 因此文件描述符 2 与文件描述符 1 指向同一个打开文件句柄 (类似于图 5-2 中进程 A 的描述符 1 和 20 指向同一打开文件句柄的情况)。可以通过调用 `dup()` 和 `dup2()` 来实现此功能。

请注意, 要满足 shell 的这一要求, 仅仅简单地打开 `results.log` 文件两次是远远不够的 (第一次在描述符 1 上打开, 第二次在描述符 2 上打开)。首先两个文件描述符不能共享相同的文件偏移量指针, 因此有可能导致相互覆盖彼此的输出。再者打开的文件不一定是磁盘文件。在如下命令中, 标准错误就将和标准输出一起送达同一管道:

```
$ ./myscript 2>&1 | less
```

`dup()` 调用复制一个打开的文件描述符 `oldfd`, 并返回一个新描述符, 二者都指向同一打开的文件句柄。系统会保证新描述符一定是编号值最低的未用文件描述符。

```
#include <unistd.h>

int dup(int oldfd);

Returns (new) file descriptor on success, or -1 on error
```

假设发起如下调用:

¹ 译者注: 将文件描述符 1 复制到文件描述符 2。

```
newfd = dup(1);
```

再假定在正常情况下，shell 已经代表程序打开了文件描述符 0、1 和 2，且没有其他描述符在用，dup()调用会创建文件描述符 1 的副本，返回的文件描述符编号值为 3。

如果希望返回文件描述符 2，可以使用如下技术：

```
close(2);          /* Frees file descriptor 2 */
newfd = dup(1);    /* Should reuse file descriptor 2 */
```

只有当描述符 0 已经打开时，这段代码方可工作。如果想进一步简化上述代码，同时总是能获得所期望的文件描述符，可以调用 dup2()。

```
#include <unistd.h>

int dup2(int oldfd, int newfd);

Returns (new) file descriptor on success, or -1 on error
```

dup2()系统调用会为 oldfd 参数所指定的文件描述符创建副本，其编号由 newfd 参数指定。如果由 newfd 参数所指定编号的文件描述符之前已经打开，那么 dup2()会首先将其关闭。(dup2()调用会默然忽略 newfd 关闭期间出现的任何错误。故此，编码时更为安全的做法是：在调用 dup2()之前，若 newfd 已经打开，则应显式调用 close()将其关闭。)

前述调用 close()和 dup()的代码可以简化为：

```
dup2(1, 2);
```

若调用 dup2()成功，则将返回副本的文件描述符编号（即 newfd 参数指定的值）。

如果 oldfd 并非有效的文件描述符，那么 dup2()调用将失败并返回错误 EBADF，且不关闭 newfd。如果 oldfd 有效，且与 newfd 值相等，那么 dup2()将什么也不做，不关闭 newfd，并将其作为调用结果返回。

fcntl()的 F_DUPFD 操作是复制文件描述符的另一接口，更具灵活性。

```
newfd = fcntl(oldfd, F_DUPFD, startfd);
```

该调用为 oldfd 创建一个副本，且将使用大于等于 startfd 的最小未用值作为描述符编号。该调用还能保证新描述符（newfd）编号落在特定的区间范围内。总是能将 dup()和 dup2()调用改写为对 close()和 fcntl()的调用，虽然前者更为简洁。（还需注意，正如手册页中所描述的，dup2()和 fcntl()二者返回的 errno 错误码存在一些差别。）

由图 5-2 可知，文件描述符的正、副本之间共享同一打开文件句柄所含的文件偏移量和状态标志。然而，新文件描述符有其自己的一套文件描述符标志，且其 close-on-exec 标志 (FD_CLOEXEC) 总是处于关闭状态。下面将要介绍的接口，可以直接控制新文件描述符的 close-on-exec 标志。

dup3()系统调用完成的工作与 dup2()相同，只是新增了一个附加参数 flag，这是一个可以修改系统调用行为的位掩码。

```
#define _GNU_SOURCE
#include <unistd.h>

int dup3(int oldfd, int newfd, int flags);

Returns (new) file descriptor on success, or -1 on error
```

目前，dup3()只支持一个标志 O_CLOEXEC，这将促使内核为新文件描述符设置 close-on-exec

标志 (FD_CLOEXEC)。设计该标志的缘由，类似于 4.3.1 节对 `open()` 调用中 `O_CLOEXEC` 标志的描述。

`dup3()` 系统调用始见于 Linux 2.6.27，为 Linux 所特有。

Linux 从 2.6.24 开始支持 `fcntl()` 用于复制文件描述符的附加命令：`F_DUPFD_CLOEXEC`。该标志不仅实现了与 `F_DUPFD` 相同的功能，还为新文件描述符设置 `close-on-exec` 标志。同样，此命令之所以得以一显身手，其原因也类似于 `open()` 调用中的 `O_CLOEXEC` 标志。SUSv3 并未提及 `F_DUPFD_CLOEXEC` 标志，但 SUSv4 对其作了规范。

5.6 在文件特定偏移量处的 I/O: `pread()` 和 `pwrite()`

系统调用 `pread()` 和 `pwrite()` 完成与 `read()` 和 `write()` 相类似的工作，只是前两者会在 `offset` 参数所指定的位置进行文件 I/O 操作，而非始于文件的当前偏移量处，且它们不会改变文件的当前偏移量。

```
#include <unistd.h>

ssize_t pread(int fd, void *buf, size_t count, off_t offset);
           Returns number of bytes read, 0 on EOF, or -1 on error
ssize_t pwrite(int fd, const void *buf, size_t count, off_t offset);
           Returns number of bytes written, or -1 on error
```

`pread()` 调用等同于将如下调用纳入同一原子操作：

```
off_t orig;
```

```
orig = lseek(fd, 0, SEEK_CUR);   /* Save current offset */
lseek(fd, offset, SEEK_SET);
s = read(fd, buf, len);
lseek(fd, orig, SEEK_SET);      /* Restore original file offset */
```

对 `pread()` 和 `pwrite()` 而言，`fd` 所指代的文件必须是可定位的（即允许对文件描述符执行 `lseek()` 调用）。

多线程应用为这些系统调用提供了用武之地。正如第 29 章所述，进程下辖的所有线程将共享同一文件描述符表。这也意味着每个已打开文件的文件偏移量为所有线程所共享。当调用 `pread()` 或 `pwrite()` 时，多个线程可同时对同一文件描述符执行 I/O 操作，且不会因其他线程修改文件偏移量而受到影响。如果还试图使用 `lseek()` 和 `read()`（或 `write()`）来代替 `pread()`（或 `pwrite()`），那么将引发竞争状态，这类似于 5.1 节讨论 `O_APPEND` 标志时的描述（当多个进程的文件描述符指向相同的打开文件句柄时，使用 `pread()` 和 `pwrite()` 系统调用同样能够避免进程间出现竞争状态）。

如果需要反复执行 `lseek()`，并伴之以文件 I/O，那么 `pread()` 和 `pwrite()` 系统调用在某些情况下是具有性能优势的。这是因为执行单个 `pread()`（或 `pwrite()`）系统调用的成本要低于执行 `lseek()` 和 `read()`（或 `write()`）两个系统调用。然而，较之于执行 I/O 实际所需的时间，系统调用的开销就有些相形见绌了¹。

¹ 译者注：执行实际 I/O 的开销要远大于执行系统调用，系统调用的性能优势作用有限。

5.7 分散输入和集中输出 (Scatter-Gather I/O): readv() 和 writev()

readv()和 writev()系统调用分别实现了分散输入和集中输出的功能。

```
#include <sys/uio.h>

ssize_t readv(int fd, const struct iovec *iov, int iovcnt);
        Returns number of bytes read, 0 on EOF, or -1 on error

ssize_t writev(int fd, const struct iovec *iov, int iovcnt);
        Returns number of bytes written, or -1 on error
```

这些系统调用并非只对单个缓冲区进行读写操作，而是一次即可传输多个缓冲区的数据。数组 `iov` 定义了一组用来传输数据的缓冲区。整型数 `iovcnt` 则指定了 `iov` 的成员个数。`iov` 中的每个成员都是如下形式的数据结构。

```
struct iovec {
    void *iov_base;      /* Start address of buffer */
    size_t iov_len;     /* Number of bytes to transfer to/from buffer */
};
```

SUSv3 标准允许系统实现对 `iov` 中的成员个数加以限制。系统实现可以通过定义 `<limits.h>` 文件中 `IOV_MAX` 来通告这一限额，程序也可以在系统运行时调用 `sysconf(_SC_IOV_MAX)` 来获取这一限额。(11.2 节将介绍 `sysconf()`。) SUSv3 要求该限额不得少于 16。Linux 将 `IOV_MAX` 的值定义为 1024，这是与内核对该向量大小的限制（由内核常量 `UIO_MAXIOV` 定义）相对应的。

然而，`glibc` 对 `readv()`和 `writev()`的封装函数¹还悄悄做了些额外工作。若系统调用因 `iovcnt` 参数值过大而失败，外壳函数将临时分配一块缓冲区，其大小足以容纳 `iov` 参数所有成员所描述的数据缓冲区，随后再执行 `read()`或 `write()`调用（参见后文对使用 `write()`实现 `writev()`功能的讨论）。

图 5-3 展示的是一个关于 `iov`、`iovcnt` 以及 `iov` 指向缓冲区之间关系的示例。

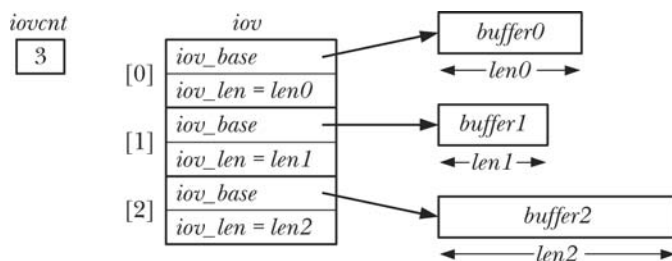


图 5-3: iovec 数组及其相关缓冲区的示例

¹ 译者注：又称外壳函数。

分散输入

`readv()`系统调用实现了分散输入的功能：从文件描述符 `fd` 所指代的文件中读取一片连续的字节，然后将其散置（“分散放置”）于 `iov` 指定的缓冲区中。这一散置动作从 `iov[0]`开始，依次填满每个缓冲区。

原子性是 `readv()`的重要属性。换言之，从调用进程的角度来看，当调用 `readv()`时，内核在 `fd` 所指代的文件与用户内存之间一次性地完成了数据转移。这意味着，假设即使有另一进程（或线程）与其共享同一文件偏移量，且在调用 `readv()`的同时企图修改文件偏移量，`readv()`所读取的数据仍将是连续的。

调用 `readv()`成功将返回读取的字节数，若文件结束¹将返回 0。调用者必须对返回值进行检查，以验证读取的字节数是否满足要求。若数据不足以填充所有缓冲区，则只会占用²部分缓冲区，其中最后一个缓冲区可能只存有部分数据。

程序清单 5-2 展示了 `readv()`的用法。

在本书中，当以函数名称冠以“t_”来命名示例程序时（例如：程序清单 5-2 中的程序 `t_readv.c`），意在表明该程序主要用于展示单个系统调用或库函数的用法。

程序清单 5-2：使用 `readv()`执行分散输入

```
----- fileio/t_readv.c

#include <sys/stat.h>
#include <sys/uio.h>
#include <fcntl.h>
#include "tlpi_hdr.h"

int
main(int argc, char *argv[])
{
    int fd;
    struct iovec iov[3];
    struct stat myStruct;      /* First buffer */
    int x;                    /* Second buffer */
#define STR_SIZE 100
    char str[STR_SIZE];       /* Third buffer */
    ssize_t numRead, totRequired;

    if (argc != 2 || strcmp(argv[1], "--help") == 0)
        usageErr("%s file\n", argv[0]);

    fd = open(argv[1], O_RDONLY);
    if (fd == -1)
        errExit("open");

    totRequired = 0;

    iov[0].iov_base = &myStruct;
    iov[0].iov_len = sizeof(struct stat);
```

1 译者注：EOF。

2 译者注：按 `iov` 数组顺序。

```

totRequired += iov[0].iov_len;

iov[1].iov_base = &x;
iov[1].iov_len = sizeof(x);
totRequired += iov[1].iov_len;

iov[2].iov_base = str;
iov[2].iov_len = STR_SIZE;
totRequired += iov[2].iov_len;

numRead = readv(fd, iov, 3);
if (numRead == -1)
    errExit("readv");

if (numRead < totRequired)
    printf("Read fewer bytes than requested\n");

printf("total bytes requested: %ld; bytes read: %ld\n",
       (long) totRequired, (long) numRead);
exit(EXIT_SUCCESS);
}

```

fileio/t_readv.c

集中输出

`writev()`系统调用实现了集中输出：将 `iov` 所指定的所有缓冲区中的数据拼接（“集中”）起来，然后以连续的字节序列写入文件描述符 `fd` 指代的文件中。对缓冲区中数据的“集中”始于 `iov[0]`所指定的缓冲区，并按数组顺序展开。

像 `readv()`调用一样，`writev()`调用也属于原子操作，即所有数据将一次性地从用户内存传输到 `fd` 指代的文件中。因此，在向普通文件写入数据时，`writev()`调用会把所有的请求数据连续写入文件，而不会在其他进程（或线程）写操作的影响下¹分散地写入文件²。

如同 `write()`调用，`writev()`调用也可能存在部分写的问题。因此，必须检查 `writev()`调用的返回值，以确定写入的字节数是否与要求相符。

`readv()`调用和 `writev()`调用的主要优势在于便捷。如下两种方案，任选其一都可替代对 `writev()`的调用。

- 编码时，首先分配一个大缓冲区，随即再从进程地址空间的其他位置将数据复制过来，最后调用 `write()`输出其中的所有数据。
- 发起一系列 `write()`调用，逐一输出每个缓冲区中的数据。

尽管方案一在语义上等同于 `writev()`调用，但需要在用户空间内分配缓冲区，进行数据复制，很不方便（效率也低）。

方案二在语义上就不同于单次的 `writev()`调用，因为发起多次 `write()`调用将无法保证原子性。更何况，执行一次 `writev()`调用比执行多次 `write()`调用开销要小（参见 3.1 节关于系统调用的讨论）。

在指定的文件偏移量处执行分散输入/集中输出

Linux 2.6.30 版本新增了两个系统调用：`preadv()`、`pwritev()`，将分散输入/集中输出和于指定文件偏移量处的 I/O 二者集于一身。它们并非标准的系统调用，但获得了现代 BSD 的支持。

¹ 译者注：即不受其他进（线）程改变文件偏移量的影响。

² 译者注：应当指出，`readv()`和 `writev()`会改变打开文件句柄的当前文件偏移量。

```
#define _BSD_SOURCE
#include <sys/uio.h>

ssize_t preadv(int fd, const struct iovec *iov, int iovcnt, off_t offset);
    Returns number of bytes read, 0 on EOF, or -1 on error
ssize_t pwritev(int fd, const struct iovec *iov, int iovcnt, off_t offset);
    Returns number of bytes written, or -1 on error
```

`preadv()`和 `pwritev()`系统调用所执行的任务与 `readv()`和 `writev()`相同，但执行 I/O 的位置将由 `offset` 参数指定（类似于 `pread()`和 `pwrite()`系统调用）¹。

对于那些既想从分散-集中 I/O 中受益，又不愿受制于当前文件偏移量的应用程序（比如，多线程的应用程序）而言，这些系统调用恰好可以派上用场。

5.8 截断文件：truncate()和 ftruncate()系统调用

`truncate()`和 `ftruncate()`系统调用将文件大小设置为 `length` 参数指定的值。

```
#include <unistd.h>

int truncate(const char *pathname, off_t length);
int ftruncate(int fd, off_t length);

    Both return 0 on success, or -1 on error
```

若文件当前长度大于参数 `length`，调用将丢弃超出部分，若小于参数 `length`，调用将在文件尾部添加一系列空字节或是一个文件空洞。

两个系统调用之间的差别在于如何指定操作文件。`truncate()`以路径字符串来指定文件，并要求可访问该文件²，且对文件拥有写权限。若文件名为符号链接，那么调用将对其进行解引用。而调用 `ftruncate()`之前，需以可写方式打开操作文件，获取其文件描述符以指代该文件，该系统调用不会修改文件偏移量。

若 `ftruncate()`的 `length` 参数值超出文件的当前大小，SUSv3 允许两种行为：要么扩展该文件（如 Linux），要么返回错误。而符合 XSI 标准的系统则必须采取前一种行为。相同的情况，对于 `truncate()`系统调用，SUSv3 则要求总是能扩展文件。

`truncate()`无需先以 `open()`（或是一些其他方法）来获取文件描述符，却可修改文件内容，在系统调用中可谓独树一帜。

5.9 非阻塞 I/O

在打开文件时指定 `O_NONBLOCK` 标志，目的有二。

- 若 `open()`调用未能立即打开文件，则返回错误，而非陷入阻塞。有一种情况属于例外，调用 `open()`操作 FIFO 可能会陷入阻塞（参见 44.7 节）。

¹ 译者注：作者于此处暗示，这两个系统调用所执行的 I/O 将不影响文件的当前偏移量。

² 译者注：即对组成路径名的各目录拥有可执行（x）权限。

- 调用 `open()` 成功后，后续的 I/O 操作也是非阻塞的。若 I/O 系统调用未能立即完成，则可能会只传输部分数据，或者系统调用失败，并返回 `EAGAIN` 或 `EWOULDBLOCK` 错误。具体返回何种错误将依赖于系统调用。Linux 系统与许多 UNIX 实现一样，将两个错误常量视为同义。

管道、FIFO、套接字、设备（比如终端、伪终端）都支持非阻塞模式。（因为无法通过 `open()` 来获取管道和套接字的文件描述符，所以要启用非阻塞标志，就必须使用 5.3 节所述 `fcntl()` 的 `F_SETFL` 命令。）

正如 13.1 节所述，由于内核缓冲区保证了普通文件 I/O 不会陷入阻塞，故而打开普通文件时一般会忽略 `O_NONBLOCK` 标志。然而，当使用强制文件锁时（55.4 节），`O_NONBLOCK` 标志对普通文件也是起作用的。

更多关于非阻塞 I/O 的信息请参见 44.9 节和第 63 章。

历史上，派生自 System V 的系统提供有 `O_NDELAY` 标志，语义上类似于 `O_NONBLOCK` 标志。二者主要的区别在于：在 System V 系统中，若非阻塞的 `write()` 调用未能完成写操作，或者非阻塞的 `read()` 调用无输入数据可读时，则两个调用将返回 0。这对于 `read()` 调用来说会有问题，因为程序将无法区分返回 0 的 `read()` 到底是没有可用的输入数据，还是遇到了文件结尾¹。故而 POSIX.1 标准在初版中引入了 `O_NONBLOCK` 标志。有些 UNIX 实现一直还在支持旧语义的 `O_NDELAY` 标志。Linux 系统虽然也定义了 `O_NDELAY` 常量，但其与 `O_NONBLOCK` 标志同义。

5.10 大文件 I/O

通常将存放文件偏移量的数据类型 `off_t` 实现为一个有符号的长整型。（之所以采用有符号数据类型，是要以 -1 来表示错误情况。）在 32 位体系架构中（比如 x86-32），这将文件大小置于 $2^{31}-1$ 个字节（即 2GB）的限制之下。

然而，磁盘驱动器的容量早已超出这一限制，因此 32 位 UNIX 实现有处理超过 2GB 大小文件的需求，这也在情理之中。由于问题较为普遍，UNIX 厂商联盟在大型文件峰会（Large File Summit）上就此进行了协商，并针对必需的大文件访问功能，形成了对 SUSv2 规范的扩展。本节将概述 LFS 的增强特性。（完整的 LFS 规范定稿于 1996 年，可通过 <http://opengroup.org/platform/lfs.html> 访问。）

始于内核版本 2.4，32 位 Linux 系统开始提供对 LFS 的支持（glibc 版本必须为 2.2 或更高）。另一个前提是，相应的文件系统也必须支持大文件操作。大多数“原生”Linux 文件系统提供了 LFS 支持，但一些“非原生”文件系统则未提供该功能（微软的 VFAT 和 NFSv2 系统是其中较为知名的范例，无论系统是否启用了 LFS 扩展功能，2GB 的文件大小限制都是硬杠杠）。

由于 64 位系统架构（例如，Alpha、IA-64）的长整型类型长度为 64 位，故而 LFS 增强特性所要突破的限制对其而言并不是问题。然而，即便在 64 位系统中，一些“原生”Linux 文件系统的实现细节还是将文件大小的理论值默认为不会超过 $2^{63}-1$ 个字节。在大多数情况下，此限额远远超出了目前的磁盘容量，故而这一对文件大小的限制并无实际意义。

¹ 译者注：此处所谓非阻塞意指 `O_NDELAY`。另外，原文表述似有错误，酌改，请参见 APUEv2 第 14.2 节。

应用程序可使用如下两种方式之一以获得 LFS 功能。

- 使用支持大文件操作的备选 API。该 API 由 LFS 设计，意在作为 SUS 规范的“过渡型扩展”。因此，尽管大部分系统都支持这一 API，但这对于符合 SUSv2 或 SUSv3 规范的系统其实并非必须。这一方法现已过时。
- 在编译应用程序时，将宏 `_FILE_OFFSET_BITS` 的值定义为 64。这一方法更为可取，因为符合 SUS 规范的应用程序无需修改任何源码即可获得 LFS 功能。

过渡型 LFS API

要使用过渡型的 LFS API，必须在编译程序时定义 `_LARGEFILE64_SOURCE` 功能测试宏，该定义可以通过命令行指定，也可以定义于源文件中包含所有头文件之前的位置。该 API 所属函数具有处理 64 位文件大小和文件偏移量的能力。这些函数与其 32 位版本命名相同，只是尾部缀以 64 以示区别。其中包括：`fopen64()`、`open64()`、`lseek64()`、`truncate64()`、`stat64()`、`mmap64()`和 `setrlimit64()`。（针对这些函数的 32 位版本，本书前面已然讨论了一部分，还有一些将在后续章节中描述。）

要访问大文件，可以使用这些函数的 64 位版本。例如，打开大文件的编码示例如下：

```
fd = open64(name, O_CREAT | O_RDWR, mode);
if (fd == -1)
    errExit("open");
```

调用 `open64()`，相当于在调用 `open()` 时指定 `O_LARGEFILE` 标志。若调用 `open()` 时未指定此标志，且欲打开的文件大小大于 2GB，那么调用将返回错误。

另外，除去上述提及的函数之外，过渡型 LFS API 还增加了一些新的数据类型，如下所示。

- `struct stat64`：类似于 `stat` 结构（参见 15.1 节），支持大文件尺寸。
- `off64_t`：64 位类型，用于表示文件偏移量。

如程序清单 5-3 所示，除去使用了该 API 中的其他 64 位函数之外，`lseek64()` 就用到了数据类型 `off64_t`。该程序接受两个命令行参数：欲打开的文件名称和给定的文件偏移量（整型）值。程序首先打开指定的文件，然后检索至给定的文件偏移量处，随即写入一串字符。如下所示的 shell 会话中，程序检索到一个超大的文件偏移量处（超过 10GB），再写入一些字节：

```
$ ./large_file x 10111222333
$ ls -l x                                     Check size of resulting file
-rw-----  1 mtk      users   10111222337 Mar  4 13:34 x
```

程序清单 5-3：访问大文件

```
fileio/large_file.c

#define _LARGEFILE64_SOURCE
#include <sys/stat.h>
#include <fcntl.h>
#include "tlpi_hdr.h"

int
main(int argc, char *argv[])
{
    int fd;
    off64_t off;
    if (argc != 3 || strcmp(argv[1], "--help") == 0)
        usageErr("%s pathname offset\n", argv[0]);

    fd = open64(argv[1], O_RDWR | O_CREAT, S_IRUSR | S_IWUSR);
```

```

    if (fd == -1)
        errExit("open64");

    off = atoll(argv[2]);
    if (lseek64(fd, off, SEEK_SET) == -1)
        errExit("lseek64");

    if (write(fd, "test", 4) == -1)
        errExit("write");
    exit(EXIT_SUCCESS);
}

```

fileio/large_file.c

__FILE_OFFSET_BITS 宏

要获取 LFS 功能，推荐的作法是：在编译程序时，将宏 `__FILE_OFFSET_BITS` 的值定义为 64。做法之一是利用 C 语言编译器的命令行选项：

```
$ cc -D__FILE_OFFSET_BITS=64 prog.c
```

另外一种方法，是在 C 语言的源文件中，在包含所有头文件之前添加如下宏定义：

```
#define __FILE_OFFSET_BITS 64
```

所有相关的 32 位函数和数据类型将自动转换为 64 位版本。因而，例如，实际会将 `open()` 转换为 `open64()`，数据类型 `off_t` 的长度也将转而定义为 64 位。换言之，无需对源码进行任何修改，只要对已有程序进行重新编译，就能够实现大文件操作。

显然，使用宏 `__FILE_OFFSET_BITS` 要比采用过渡型的 LFS API 更为简单，但这也要求应用程序的代码编写必须规范（例如，声明用于放置文件偏移量的变量，应正确地使用 `off_t`，而不能使用“原生”的 C 语言整型）。

LFS 规范对于支持 `__FILE_OFFSET_BITS` 宏未作硬性规定，仅仅提及将该宏作为指定数据类型 `off_t` 大小的可选方案。一些 UNIX 实现使用不同的特性测试宏来获取此功能。

若试图使用 32 位函数访问大文件（即在编译程序时，未将宏 `__FILE_OFFSET_BITS` 的值设置为 64），调用可能会返回 `E_OVERFLOW` 错误。例如，为获取大小超过 2G 文件的信息，若使用 `stat` 的 32 位版本时就会遇到这一错误。

向 printf() 调用传递 off_t 值

LFS 扩展功能没有解决的问题之一是，如何向 `printf()` 调用传递 `off_t` 值。3.6.2 节曾特别指出，对于预定义的系统数据类型（诸如 `pid_t`、`uid_t`），展示其值的可移植方法是将该值强制转换为 `long` 型，并在 `printf()` 中使用限定符 `%ld`。然而，一旦使用了 LFS 扩展功能，`%ld` 将不足以处理 `off_t` 数据类型，因为对该数据类型的定义可能会超出 `long` 类型的范围，一般为 `long long` 类型。据此，若要显示 `off_t` 类型的值，则先要将其强制转换为 `long long` 类型，然后使用 `printf()` 函数的 `%lld` 限定符显示，如下所示：

```

#define __FILE_OFFSET_BITS 64

off_t offset;          /* Will be 64 bits, the size of 'long long' */

/* Other code assigning a value to 'offset' */

printf("offset=%lld\n", (long long) offset);

```


在处理 `stat` 结构所使用的 `blkcnt_t` 数据类型时，也应予以类似关注（参见 15.1 节的描述）。

如需在独立的编译模块之间传递 `off_t` 或 `stat` 类型的参数值，则需确保在所有模块中，这些数据类型的大小相同（即编译这些模块时，要么将宏 `_FILE_OFFSET_BITS` 的值都定义为 64，要么都不做定义）。

5.11 /dev/fd 目录

对于每个进程，内核都提供一个特殊的虚拟目录 `/dev/fd`。该目录中包含 “`/dev/fd/n`” 形式的文件名，其中 `n` 是与进程中的打开文件描述符相对应的编号。因此，例如，`/dev/fd/0` 就对应于进程的标准输入。（SUSv3 对 `/dev/fd` 特性未做规定，但有些其他的 UNIX 实现也提供了这一特性。）

打开 `/dev/fd` 目录中的一个文件等同于复制相应的文件描述符，所以下列两行代码是等价的：

```
fd = open("/dev/fd/1", O_WRONLY);
fd = dup(1);          /* Duplicate standard output */
```

在为 `open()` 调用设置 `flag` 参数时，需要注意将其设置为与原描述符相同的访问模式。这一场景下，在 `flag` 标志的设置中引入其他标志，诸如 `O_CREAT`，是毫无意义的（系统会将其忽略）。

`/dev/fd` 实际上是一个符号链接，链接到 Linux 所专有的 `/proc/self/fd` 目录。后者又是 Linux 特有的 `/proc/PID/fd` 目录族的特例之一，此目录族中的每一目录都包含有符号链接，与一进程所打开的所有文件相对应。

程序中很少会使用 `/dev/fd` 目录中的文件。其主要用途在 shell 中。许多用户级 shell 命令将文件名作为参数，有时需要将命令输出至管道，并将某个参数替换为标准输入或标准输出。出于这一目的，有些命令（例如，`diff`、`ed`、`tar` 和 `comm`）提供了一个解决方法，使用 “-” 符号作为命令的参数之一，用以表示标准输入或输出（视情况而定）。所以，要比较 `ls` 命令输出的文件名列表与之前生成的文件名列表，命令就可以写成：

```
$ ls | diff - oldfilelist
```

这种方法有不少问题。首先，该方法要求每个程序都对 “-” 符号做专门处理，但是许多程序并未实现这样的功能，这些命令只能处理文件，不支持将标准输入或输出作为参数。其次，有些程序还将单个 “-” 符解释为表征命令行选项结束的分隔符。

使用 `/dev/fd` 目录，上述问题将迎刃而解，可以把标准输入、标准输出和标准错误作为文件名参数传递给任何需要它们的程序。所以，可以将前一个 shell 命令改写成如下形式：

```
$ ls | diff /dev/fd/0 oldfilelist
```

方便起见，系统还提供了 3 个符号链接：`/dev/stdin`、`/dev/stdout` 和 `/dev/stderr`，分别链接到 `/dev/fd/0`、`/dev/fd/1` 和 `/dev/fd/2`。

5.12 创建临时文件

有些程序需要创建一些临时文件，仅供其在运行期间使用，程序终止后即行删除。例如，

很多编译器程序会在编译过程中创建临时文件。GNU C 语言函数库为此而提供了一系列库函数。（之所以有“一系列”的库函数，部分原因是由于这些函数分别继承自各种 UNIX 实现。）本节将介绍其中的两个函数：`mkstemp()`和 `tmpfile()`。

基于调用者提供的模板，`mkstemp()`函数生成一个唯一文件名并打开该文件，返回一个可用于 I/O 调用的文件描述符。

```
#include <stdlib.h>

int mkstemp(char *template);

Returns file descriptor on success, or -1 on error
```

模板参数采用路径名形式，其中最后 6 个字符必须为 `XXXXXX`。这 6 个字符将被替换，以保证文件名的唯一性，且修改后的字符串将通过 `template` 参数传回。因为会对传入的 `template` 参数进行修改，所以必须将其指定为字符数组，而非字符串常量。

文件所有者对 `mkstemp()`函数建立的文件拥有读写权限（其他用户则没有任何操作权限），且打开文件时使用了 `O_EXCL` 标志，以保证调用者以独占方式访问文件。

通常，打开临时文件不久，程序就会使用 `unlink` 系统调用（参见 18.3 节）将其删除。故而，`mkstemp()`函数的示例代码如下所示：

```
int fd;
char template[] = "/tmp/somestringXXXXXX";

fd = mkstemp(template);
if (fd == -1)
    errExit("mkstemp");
printf("Generated filename was: %s\n", template);
unlink(template); /* Name disappears immediately, but the file
                  is removed only after close() */

/* Use file I/O system calls - read(), write(), and so on */

if (close(fd) == -1)
    errExit("close");
```

使用 `tmpnam()`、`tempnam()`和 `mktemp()`函数也能生成唯一的文件名。然而，由于这会导致应用程序出现安全漏洞，应当避免使用这些函数。关于这些函数的进一步细节请参考手册页。

`tmpfile()`函数会创建一个名称唯一的临时文件，并以读写方式将其打开。（打开该文件时使用了 `O_EXCL` 标志，以防一个可能性极小的冲突，即另一个进程已经创建了一个同名文件。）

```
#include <stdio.h>

FILE *tmpfile(void);

Returns file pointer on success, or NULL on error
```

`tmpfile()`函数执行成功，将返回一个文件流供 `stdio` 库函数使用。文件流关闭后将自动删除临时文件。为达到这一目的，`tmpfile()`函数会在打开文件后，从内部立即调用 `unlink()`来删除该文件名¹。

¹ 译者注：进程终止时会关闭所有打开的文件描述符，关闭文件就会删除这些临时文件（参考 `mkstemp` 代码示例中的注释），由此可以推导出，进程退出时将自动删除临时文件。

5.13 总结

本章介绍了原子操作的概念，这对于一些系统调用的正确操作至关重要。特别是，指定 `O_EXCL` 标志调用 `open()`，这确保了调用者就是文件的创建者。而指定 `O_APPEND` 标志来调用 `open()`，还确保了多个进程在对同一文件追加数据时不会覆盖彼此的输出。

系统调用 `fcntl()` 可以执行许多文件控制操作，其中包括：修改打开文件的状态标志、复制文件描述符。使用 `dup()` 和 `dup2()` 系统调用也能实现文件描述符的复制功能。

本章接着研究了文件描述符、打开文件句柄和文件 `i-node` 之间的关系，并特别指出这 3 个对象各自包含的不同信息。文件描述符及其副本指向同一个打开文件句柄，所以也将共享打开文件的状态标志和文件偏移量。

之后描述的诸多系统调用，是对常规 `read()` 和 `write()` 系统调用的功能扩展。`pread()` 和 `pwrite()` 系统调用可在文件的指定位置处执行 I/O 功能，且不会修改文件偏移量。`readv()` 和 `writelv()` 系统调用实现了分散输入和集中输出的功能。`preadv()` 和 `pwritev()` 系统调用则集上述两对系统调用的功能于一身。

使用 `truncate()` 和 `ftruncate()` 系统调用，既可以丢弃多余的字节以缩小文件大小，又能使用填充为 0 的文件空洞来增加文件大小。

本章还简单介绍了非阻塞 I/O 的概念，后续章节中还将继续讨论。

LFS 规范定义了一套扩展功能，允许在 32 位系统中运行的进程来操作无法以 32 位表示的大文件。

运用虚拟目录 `/dev/fd` 中的编号文件，进程就可以通过文件描述符编号来访问自己打开的文件，这在 `shell` 命令中尤其有用。

`mkstemp()` 和 `tmpfile()` 函数允许应用程序去创建临时文件。

5.14 练习

- 5-1. 请使用标准文件 I/O 系统调用 (`open()` 和 `lseek()`) 和 `off_t` 数据类型修改程序清单 5-3 中的程序。将宏 `_FILE_OFFSET_BITS` 的值设置为 64 进行编译，并测试该程序是否能够成功创建一个文件。
- 5-2. 编写一个程序，使用 `O_APPEND` 标志并以写方式打开一个已存在的文件，且将文件偏移量置于文件起始处，再写入数据。数据会显示在文件中的哪个位置？为什么？
- 5-3. 本习题的设计目的在于展示为何以 `O_APPEND` 标志打开文件来保障操作的原子性是必要的。请编写一程序，可接收多达 3 个命令行参数：

```
$ atomic_append filename num-bytes [x]
```

该程序应打开所指定的文件（如有必要，则创建之），然后以每次调用 `write()` 写入一个字节的方式，向文件尾部追加 `num-bytes` 个字节。缺省情况下，程序使用 `O_APPEND` 标志打开文件，但若存在第三个命令行参数 (`x`)，那么打开文件时将不再使用 `O_APPEND` 标志，代之以在每次调用 `write()` 前调用 `lseek(fd,0,SEEK_END)`。同时运行该程序的两个实例，不带 `x` 参数，将 100 万个字节写入同一文件：

```
$ atomic_append f1 1000000 & atomic_append f1 1000000
```

重复上述操作，将数据写入另一文件，但运行时加入 x 参数：

```
$ atomic_append f2 1000000 x & atomic_append f2 1000000 x
```

使用 ls-l 命令检查文件 f1 和 f2 的大小，并解释两文件大小不同的原因。

- 5-4.** 使用 `fcntl()` 和 `close()`（若有必要）来实现 `dup()` 和 `dup2()`。（对于某些错误，`dup2()` 和 `fcntl()` 返回的 `errno` 值并不相同，此处可不予考虑。）务必牢记 `dup2()` 需要处理的一种特殊情况，即 `oldfd` 与 `newfd` 相等。这时，应检查 `oldfd` 是否有效，测试 `fcntl(oldfd, F_GETFL)` 是否成功就能达到这一目的。若 `oldfd` 无效，则 `dup2()` 将返回 -1，并将 `errno` 置为 `EBADF`。

- 5-5.** 编写一程序，验证文件描述符及其副本是否共享了文件偏移量和打开文件的状态标志。

- 5-6.** 说明下列代码中每次执行 `write()` 后，输出文件的内容是什么，为什么。

```
fd1 = open(file, O_RDWR | O_CREAT | O_TRUNC, S_IRUSR | S_IWUSR);
fd2 = dup(fd1);
fd3 = open(file, O_RDWR);
write(fd1, "Hello,", 6);
write(fd2, "world", 6);
lseek(fd2, 0, SEEK_SET);
write(fd1, "HELLO,", 6);
write(fd3, "Giddy", 6);
```

- 5-7.** 使用 `read()`、`write()` 以及 `malloc` 函数包（见 7.1.2 节）中的必要函数以实现 `readv()` 和 `writenv()` 功能。

第 6 章

进 程

本章将研究进程结构，并将重点关注进程虚拟内存的布局及内容。同时，还会对进程的某些属性进行考察。后续章节会对进程属性做进一步的探究（例如第 9 章的进程凭证，第 35 章的进程优先权及进程调度）。第 24 章至第 27 章将讨论如何创建、终止进程，以及进程如何执行新的程序。

6.1 进程和程序

进程（process）是一个可执行程序（program）的实例。本节将阐述进程定义，并澄清其与程序之间的区别。

程序是包含了一系列信息的文件，这些信息描述了如何在运行时创建一个进程，所包括的内容如下所示。

- 二进制格式标识：每个程序文件都包含用于描述可执行文件格式的元信息（metainformation）。内核（kernel）利用此信息来解释文件中的其他信息。历史上，UNIX 可执行文件曾有两种广泛使用的格式，分别为最初的 a.out（汇编程序输出）和更加复杂的 COFF（通用对象文件格式）。现在，大多数 UNIX 实现（包括 Linux）采用可执行连接格式（ELF），这一文件格式比老版本格式具有更多优点。
- 机器语言指令：对程序算法进行编码。
- 程序入口地址：标识程序开始执行时的起始指令位置。
- 数据：程序文件包含的变量初始值和程序使用的字面常量（literal constant）值（比如字符串）。
- 符号表及重定位表：描述程序中函数和变量的位置及名称。这些表格有多种用途，其中包括调试和运行时的符号解析（动态链接）。
- 共享库和动态链接信息：程序文件所包含的一些字段，列出了程序运行时需要使用的共享库，以及加载共享库的动态链接器的路径名。
- 其他信息：程序文件还包含许多其他信息，用以描述如何创建进程。

可以用一个程序来创建许多进程，或者反过来说，许多进程运行的可以是同一程序。在此将本节开始时给出的进程定义重新改写为，进程是由内核定义的抽象的实体，并为

该实体分配用以执行程序的各项系统资源。

从内核角度看，进程由用户内存空间（user-space memory）和一系列内核数据结构组成，其中用户内存空间包含了程序代码及代码所使用的变量，而内核数据结构则用于维护进程状态信息。记录在内核数据结构中的信息包括许多与进程相关的标识号（IDs）、虚拟内存表、打开文件的描述符表、信号传递及处理的有关信息、进程资源使用及限制、当前工作目录和大量的其他信息。

6.2 进程号和父进程号

每个进程都有一个进程号（PID），进程号是一个正数，用以唯一标识系统中的某个进程。对各种系统调用而言，进程号有时可以作为传入参数，有时可以作为返回值。比如，系统调用 `kill()`（20.5 节）允许调用者向拥有特定进程号的进程发送一个信号。当需要创建一个对某进程而言唯一的标识符时，进程号就会派上用场。常见的例子是将进程号作为与进程相关文件名的一部分。

系统调用 `getpid()` 返回调用进程的进程号。

```
#include <unistd.h>

pid_t getpid(void);
```

Always successfully returns process ID of caller

`getpid()` 返回值的数据类型为 `pid_t`，该类型是由 SUSv3 所规定的整数类型，专用于存储进程号。

除了少数系统进程外，比如 `init` 进程（进程号为 1），程序与运行该程序进程的进程号之间没有固定关系。

Linux 内核限制进程号需小于等于 32767。新进程创建时，内核会按顺序将下一个可用的进程号分配给其使用。每当进程号达到 32767 的限制时，内核将重置进程号计数器，以便从小整数开始分配。

一旦进程号达到 32767，会将进程号计数器重置为 300，而不是 1。之所以如此，是因为低数值的进程号为系统进程和守护进程所长期占用，在此范围内搜索尚未使用的进程号只会是浪费时间。

在 Linux 2.4 版本及更早版本中，进程号的上限 32767，由内核常量 `PID_MAX` 所定义。在 Linux 2.6 版本中，情况有所改变。尽管进程号的默认上限仍是 32767，但可以通过 Linux 系统特有的 `/proc/sys/kernel/pid_max` 文件来进行调整（其值=最大进程号+1）。在 32 位平台中，`pid_max` 文件的最大值为 32768，但在 64 位平台中，该文件的最大值可以高达 2^{22} （约 400 万），系统可能容纳的进程数量会非常庞大。

每个进程都有一个创建自己的父进程。使用系统调用 `getppid()` 可以检索到父进程的进程号。

```
#include <unistd.h>

pid_t getppid(void);
```

Always successfully returns process ID of parent of caller

实际上，每个进程的父进程号属性反映了系统上所有进程间的树状关系。每个进程的父进程又有自己的父进程，以此类推，回溯到 1 号进程——`init` 进程，即所有进程的始祖。使用

ps-tree(1)命令可以查看到这一“家族树”(family tree)。

如果子进程的父进程终止，则子进程就会变成“孤儿”，init 进程随即将收养该进程，子进程后续对 getppid()的调用将返回进程号 1 (参照 26.2 节)。

通过查看由 Linux 系统所特有的 /proc/PID/status 文件所提供的 PPid 字段，可以获知每个进程的父进程。

6.3 进程内存布局

每个进程所分配的内存由很多部分组成，通常称之为“段(segment)”。如下所示。

- 文本段包含了进程运行的程序机器语言指令。文本段具有只读属性，以防止进程通过错误指针意外修改自身指令。因为多个进程可同时运行同一程序，所以又将文本段设为可共享，这样，一份程序代码的拷贝可以映射到所有这些进程的虚拟地址空间中。
- 初始化数据段包含显式初始化的全局变量和静态变量。当程序加载到内存时，从可执行文件中读取这些变量的值。
- 未初始化数据段包含了未进行显式初始化的全局变量和静态变量。程序启动之前，系统将本段内所有内存初始化为 0。出于历史原因，此段常被称为 BSS 段，这源于老版本的汇编语言助记符“block started by symbol”。将经过初始化的全局变量和静态变量与未经初始化的全局变量和静态变量分开存放，其主要原因在于程序在磁盘上存储时，没有必要为未经初始化的变量分配存储空间。相反，可执行文件只需记录未初始化数据段的位置及所需大小，直到运行时再由程序加载器来分配这一空间。
- 栈(stack)是一个动态增长和收缩的段，由栈帧(stack frames)组成。系统会为每个当前调用的函数分配一个栈帧。栈帧中存储了函数的局部变量(所谓自动变量)、实参和返回值。6.5 节将深入讨论栈帧。
- 堆(heap)是可在运行时(为变量)动态进行内存分配的一块区域。堆顶端称作 program break。

对于初始化和未初始化的数据段而言，不太常用、但表述更清晰的称谓分别是用户初始化数据段(user-initialized data segment)和零初始化数据段(zero-initialized data segment)。

size(1)命令可显示二进制可执行文件的文本段、初始化数据段、非初始化数据段(bss)的段大小。

正文中使用的术语“段(segment)”不应与一些硬件体系架构，比如 x86-32 中使用的硬件分段(segmentation)相混淆。相反，本文中的段是对 UNIX 系统中进程虚拟内存的逻辑划分。有时，会使用术语“区(section)”来替代段，因为在当下风行的可执行文件格式(ELF)规范中，采用的术语与“区”更趋一致。

本书会在多处涉及这种情况：库函数返回的指针指向静态分配的内存。这意味着，该内存既可在初始化数据段中分配，也可在非初始化数据段中分配。(某些情况下，库函数转而在堆上对内存做一次性动态分配，然而，这一实现细节与这里所要表达的意思无关。)库函数有时会通过静态分配的内存来返回信息，了解这一情况至关重要，因为这片内存的存在独立于函数调用，后续对同一函数的调用可能会将其覆盖(有时，后续对相关函数的调用也有相同的效应)。使用静态分配的内存会使函数不可重入(nonreentrant)。21.1.2 节和 31.1 节将深入讨论重入(reentrancy)问题。

程序清单 6-1 展示了不同类型的 C 语言变量，并以注释说明每种变量分属于哪个段。这些说明**正确的前提是假定使用了非优化的编译器**，且在应用程序二进制接口（ABI）中，是通过栈来传递所有参数的。实际上，优化编译器会将频繁使用的变量分配于寄存器中，或者索性将变量彻底剔除¹。此外，一些 ABI 需要通过寄存器，而不是栈，来传递函数实参和结果。尽管如此，本例只是意在展示 C 语言变量和进程各段间的映射关系。

程序清单 6-1：程序变量在进程内存各段中的位置

```
----- proc/mem_segments.c
#include <stdio.h>
#include <stdlib.h>

char globBuf[65536];          /* Uninitialized data segment */
int primes[] = { 2, 3, 5, 7 }; /* Initialized data segment */

static int
square(int x)                /* Allocated in frame for square() */
{
    int result;              /* Allocated in frame for square() */

    result = x * x;
    return result;          /* Return value passed via register */
}

static void
doCalc(int val)              /* Allocated in frame for doCalc() */
{
    printf("The square of %d is %d\n", val, square(val));

    if (val < 1000) {
        int t;               /* Allocated in frame for doCalc() */

        t = val * val * val;
        printf("The cube of %d is %d\n", val, t);
    }
}

int
main(int argc, char *argv[]) /* Allocated in frame for main() */
{
    static int key = 9973;    /* Initialized data segment */
    static char mbuf[10240000]; /* Uninitialized data segment */
    char *p;                 /* Allocated in frame for main() */

    p = malloc(1024);        /* Points to memory in heap segment */

    doCalc(key);

    exit(EXIT_SUCCESS);
}
----- proc/mem_segments.c
```

应用程序二进制接口（ABI）是一套规则，规定了二进制可执行文件在运行时应如何与某些服务（诸如内核或函数库所提供的服务）交换信息。ABI 特别规定了使用哪些寄存器

¹ 译者注：例如，以寄存器取代变量。

和栈地址来交换信息以及所交换值的含义，一旦针对某个特定 ABI 进行了编译，其二进制可执行文件应能在 ABI 相同的任何系统上运行。与之相反，标准化的 API（如 SUSv3）仅能通过编译源代码来保证应用程序的可移植性。

虽然 SUSv3 未作规定，但在大多数 UNIX 实现（包括 Linux）中 C 语言编程环境提供了 3 个全局符号（symbol）：`etext`、`edata` 和 `end`，可在程序内使用这些符号以获取相应程序文本段、初始化数据段和非初始化数据段结尾处下一字节的地址。使用这些符号，必须显式声明如下：

```
extern char etext, edata, end;
/* For example, &etext gives the address of the end
of the program text / start of initialized data */
```

图 6-1 展示了各种内存段在 x86-32 体系结构中的布局，该图的顶部标记为 `argv`、`environ` 的空间用来存储程序命令行实参（通过 C 语言中 `main()` 函数的 `argv` 参数获得）和进程环境列表（稍后讨论），图中十六进制的地址会因内核配置和程序链接选项差异而有所不同。图中标灰的区域表示这些范围在进程虚拟地址空间中不可用，也就是说，没有为这些区域创建页表（`page table`）（参考以下关于虚拟内存管理的讨论）。

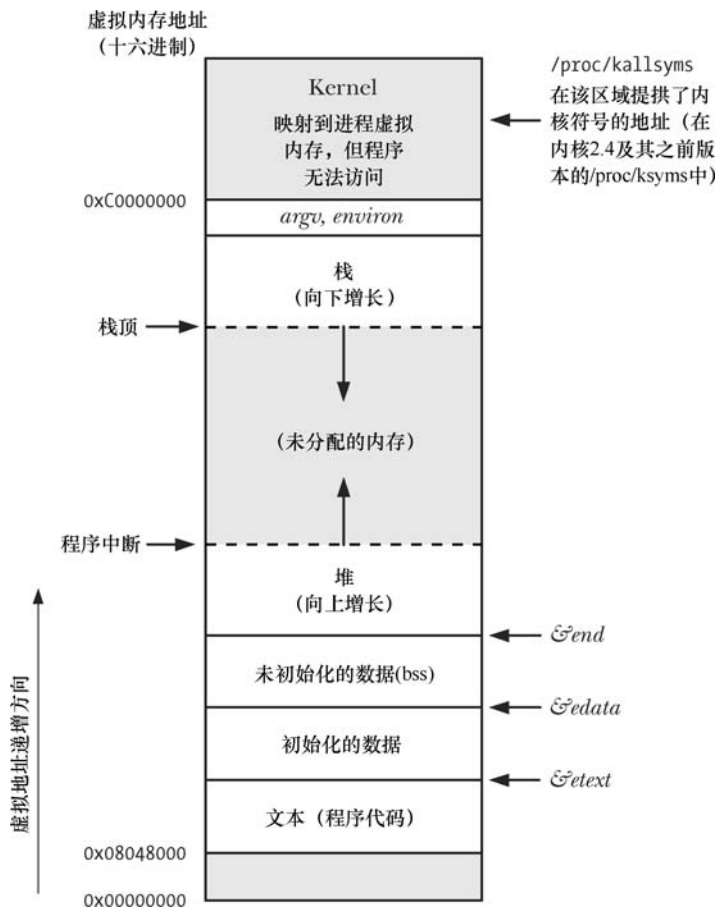


图 6-1：在 Linux/x86-32 中典型的进程内存结构

48.5 节将更为详细地重新讨论进程内存布局的课题，还将论及共享内存和共享库在进程虚拟内存中的放置位置。

6.4 虚拟内存管理

上述关于进程内存布局的讨论忽略了一个事实：这一布局存在于虚拟内存中。因为对虚拟内存的理解将有助于后续对诸如 `fork()` 系统调用、共享内存和映射文件之类主题的阐述，所以这里将探讨一些有关虚拟内存的详细内容。

Linux，像多数现代内核一样，采用了虚拟内存管理技术。该技术利用了大多数程序的一个典型特征，即访问局部性（locality of reference），以求高效使用 CPU 和 RAM（物理内存）资源。大多数程序都展现了两种类型的局部性。

- 空间局部性（Spatial locality）：是指程序倾向于访问在最近访问过的内存地址附近的内存（由于指令是顺序执行的，且有时会按顺序处理数据结构）。
- 时间局部性（Temporal locality）：是指程序倾向于在不久的将来再次访问最近刚访问过的内存地址（由于循环）。

正是由于访问局部性特征，使得程序即便仅有部分地址空间存在于 RAM 中，依然可能得以执行。

虚拟内存的规划之一是将每个程序使用的内存切割成小型的、固定大小的“页”（page）单元。相应地，将 RAM 划分成一系列与虚存页尺寸相同的页帧。任一时刻，每个程序仅有部分页需要驻留在物理内存页帧中。这些页构成了所谓驻留集（resident set）。程序未使用的页拷贝保存在交换区（swap area）内——这是磁盘空间中的保留区域，作为计算机 RAM 的补充——仅在需要时才会载入物理内存。若进程欲访问的页面目前并未驻留在物理内存中，将会发生页面错误（page fault），内核即刻挂起进程的执行，同时从磁盘中将该页面载入内存。

在 x86-32 中，页面大小为 4096 个字节。其他一些 Linux 实现使用的页面比 4096 个字节更大。例如，Alpha 使用的页面大小为 8192 个字节，IA-64 使用的页面大小是可变的，默认为 16384 个字节。程序可调用 `sysconf(_SC_PAGESIZE)` 来获取系统虚拟内存的页面大小，具体参见 11.2 节的描述。

为支持这一组织方式，内核需要为每个进程维护一张页表（page table）（见图 6-2）。该页表描述了每页在进程虚拟地址空间（virtual address space）中的位置（可为进程所用的所有虚拟内存页面的集合）。页表中的每个条目要么指出一个虚拟页面在 RAM 中的所在位置，要么表明其当前驻留在磁盘上。

在进程虚拟地址空间中，并非所有的地址范围都需要页表条目。通常情况下，由于可能存在大段的虚拟地址空间并未投入使用，故而也无必要为其维护相应的页表条目。若进程试图访问的地址并无页表条目与之对应，那么进程将收到一个 SIGSEGV 信号。

由于内核能够为进程分配和释放页（和页表条目），所以进程的有效虚拟地址范围在其生命周期中可以发生变化。这可能会发生于如下场景。

- 由于栈向下增长超出之前曾达到的位置。
- 当在堆中分配或释放内存时，通过调用 `brk()`、`sbrk()` 或 `malloc` 函数族（第 7 章）来提升 program break 的位置。
- 当调用 `shmat()` 连接 System V 共享内存区时，或者当调用 `shmdt()` 脱离共享内存区时（第 48 章）。

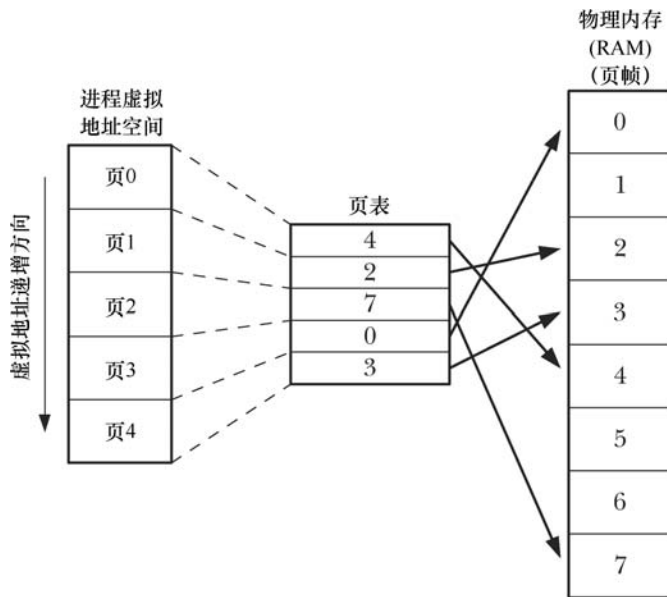


图 6-2: 虚拟内存概览

- 当调用 `mmap()` 创建内存映射时，或者当调用 `munmap()` 解除内存映射时（第 49 章）。

虚拟内存的实现需要硬件中分页内存管理单元（PMMU）的支持。PMMU 把要访问的每个虚拟内存地址转换成相应的物理内存地址，当特定虚拟内存地址所对应的页没有驻留于 RAM 中时，将以页面错误通知内核。

虚拟内存管理使进程的虚拟地址空间与 RAM 物理地址空间隔离开来，这带来许多优点。

- 进程与进程、进程与内核相互隔离，所以一个进程不能读取或修改另一进程或内核的内存。这是因为每个进程的页表条目指向 RAM（或交换区）中截然不同的物理页面集合。
- 适当情况下，两个或者更多进程能够共享内存。这是由于内核可以使不同进程的页表条目指向相同的 RAM 页。内存共享常发生于如下两种场景。
 - 执行同一程序的多个进程，可共享一份（只读的）程序代码副本。当多个程序执行相同的程序文件（或加载相同的共享库）时，会隐式地实现这一类型的共享。
 - 进程可以使用 `shmget()` 和 `mmap()` 系统调用显式地请求与其他进程共享内存区。这么做是出于进程间通信的目的。
- 便于实现内存保护机制；也就是说，可以对页表条目进行标记，以表示相关页面内容是可读、可写、可执行亦或是这些保护措施的组合。多个进程共享 RAM 页面时，允许每个进程对内存采取不同的保护措施。例如，一个进程可能以只读方式访问某页面，而另一进程则以读写方式访问同一页面。
- 程序员和编译器、链接器之类的工具无需关注程序在 RAM 中的物理布局。
- 因为需要驻留在内存中的仅是程序的一部分，所以程序的加载和运行都很快。而且，一个进程所占用的内存（即虚拟内存大小）能够超出 RAM 容量。

虚拟内存管理的最后一个优点是：由于每个进程使用的 RAM 减少了，RAM 中同时可以容纳的进程数量就增多了。这增大了如下事件的概率：在任一时刻，CPU 都可执行至少一个进程，因而往往也会提高 CPU 的利用率。

6.5 栈和栈帧

函数的调用和返回使栈的增长和收缩呈线性。X86-32 体系架构之上的 Linux（和多数其他 Linux 和 UNIX 实现），栈驻留在内存的高端并向下增长（朝堆的方向）。专用寄存器——栈指针（stack pointer），用于跟踪当前栈顶。每次调用函数时，会在栈上新分配一帧，每当函数返回时，再从栈上将此帧移去。

虽然栈向下增长，但仍将栈的增长端称为栈顶，因为抽象地说来，情况本就如此。栈的实际增长方向是个（属于硬件范畴的）实现细节。在 HP PA-RISC 的 Linux 实现中，栈的增长方向就是向上的。

就虚拟内存而言，分配栈帧后，栈段的大小将会增长，但在大多数（Linux）实现中，释放这些栈帧后，栈的大小并未减少（在分配新的栈帧时，会对这些内存重新加以利用）。当谈论栈段的增长和收缩时，只是从逻辑视角来看待栈帧在栈中的增减情况。

有时，会用用户栈（user stack）来表示此处所讨论的栈，以便与内核栈区分开来。内核栈是每个进程保留在内核内存中的内存区域，在执行系统调用的过程中供（内核）内部函数调用使用。（由于用户栈驻留在不受保护的用户内存中，所以内核无法利用用户栈来达成这一目的。）

每个（用户）栈帧包括如下信息。

- 函数实参和局部变量：由于这些变量都是在调用函数时自动创建的，因此在 C 语言中称其为自动变量。函数返回时将自动销毁这些变量（因为栈帧会被释放），这也是自动变量与静态（以及全局）变量主要的语义区别：后者与函数执行无关，且长期存在。
- （函数）调用的链接信息：每个函数都会用到一些 CPU 寄存器，比如程序计数器，其指向下一条将要执行的机器语言指令。每当一函数调用另一函数时，会在被调用函数的栈帧中保存这些寄存器的副本，以便函数返回时能为函数调用者将寄存器恢复原状。

因为函数能够嵌套调用，所以栈中可能有多个栈帧。（若一函数递归调用自身，则该函数在栈中将有多个栈帧。）参考程序清单 6-1，在 `square()` 函数执行期间，栈中包含的帧如图 6-3 所示。

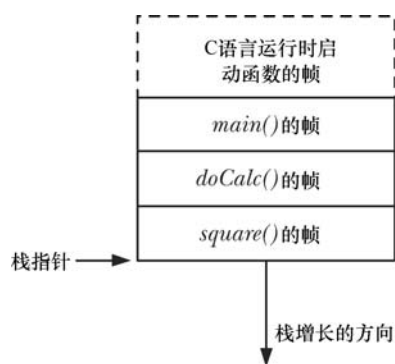


图 6-3：一个进程栈的示例

6.6 命令行参数（argc, argv）

每个 C 语言程序都必须有一个称为 `main()` 的函数，作为程序启动的起点。当执行程序时，命令行参数（command-line argument）（由 shell 逐一解析）通过两个入参提供给 `main()` 函数。第一个参数 `int argc`，表示命令行参数的个数。第二个参数 `char *argv[]`，是一个指向命令行参数的指针数组，每一参数又都是以空字符（`null`）¹ 结尾的字符串。第一个字符串，亦即 `argv[0]`

¹ 译者注：\0。

指向的，(通常)是该程序的名称。`argv` 中的指针列表以 `NULL` 指针结尾 (即 `argv[argc]` 为 `NULL`)。

`argv[0]` 包含了调用程序的名称，可以利用这一特性玩个实用的小技巧。首先为同一程序创建多个链接 (即名称不同)，然后让该程序查看 `argv[0]`，并根据调用程序的名称来执行不同任务。`gzip(1)`、`gunzip(1)` 和 `zcat(1)` 命令是该技术应用的一个例子，这些命令链接的都是同一可执行文件。(使用该技术，必须小心处理如下情况：用户通过链接调用程序，但链接名又在该程序的意料之外。)

图 6-4 展示了执行程序清单 6-2 中程序所传入参 `argc` 和 `argv` 的数据结构。该图使用 C 语言符号 “\0” 来表示每个字符串末尾的终止空字节。

程序清单 6-2 中的程序回显了其命令行参数，逐一按行输出，前面还冠以要显示的 `argv` 成员名称。

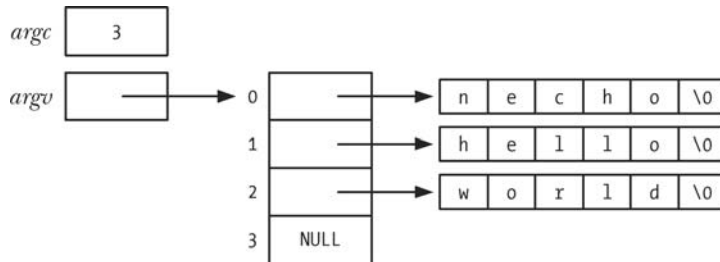


图 6-4: 命令 “necho hello world” 的 `argc` 和 `argv` 值

程序清单 6-2: 回显命令行参数

```
----- proc/necho.c
#include "tldpi_hdr.h"

int
main(int argc, char *argv[])
{
    int j;

    for (j = 0; j < argc; j++)
        printf("argv[%d] = %s\n", j, argv[j]);

    exit(EXIT_SUCCESS);
}
----- proc/necho.c
```

因为 `argv` 列表以 `NULL` 值终止，所以可以将程序清单 6-2 中的程序主体改写如下，且每行只输出一个命令行实参：

```
char **p;

for (p = argv; *p != NULL; p++)
    puts(*p);
```

`argc/argv` 参数机制的局限之一在于这些变量仅对 `main()` 函数可用。在保证可移植性的同时，为使这些命令行参数能为其他函数所用，必须把 `argv` 以参数形式传递给这些函数，或是设置一个指向 `argv` 的全局变量。

要想从程序内任一位置访问这些信息的部分或者全部内容，还有两个方法，但是会破坏程序的可移植性。

- 通过 linux 系统专有的 `/proc/PID/cmdline` 文件可以读取任一进程的命令行参数，每个参数

都以空（null）字节终止。（程序可以通过`/proc/self/cmdline`文件访问自己的命令行参数。）

- GNU C 语言库提供有两个全局变量，可在程序内任一位置使用以获取调用该程序时的程序名称（即命令行的第一个参数）。第一个全局变量 `program_invocation_name`，提供了用于调用该程序的完整路径名。第二个全局变量 `program_invocation_short_name`，提供了不含目录的程序名称，即路径名的基本名称（`basename`）部分，定义 `_GNU_SOURCE` 宏后即可从 `<errno.h>` 中获得对这两个全局变量的声明。

正如图 6-1 所示，`argv` 和 `environ` 数组，以及这些参数最初指向的字符串，都驻留在进程栈之上的一个单一、连续的内存区域。（下一节将描述 `environ` 参数，该参数用于存储程序的环境列表。）此区域可存储的字节数有上限要求，SUSv3 规定使用 `ARG_MAX` 常量（定义于 `<limits.h>`）或者调用 `sysconf(_SC_ARG_MAX)` 函数以确定该上限值（将在 11.2 节描述 `sysconf()` 函数），并且 SUSv3 还要求 `ARG_MAX` 常量的下限为 `_POSIX_ARG_MAX`（4096）个字节，而大多数 UNIX 实现的限制都远高于此。但 SUSv3 并未规定对 `ARG_MAX` 限制的实现中是否要将一些开销字节计算在内（比如终止空字符、字节对齐、`argv` 和 `environ` 指针数组）。

Linux 中的 `ARG_MAX` 参数值曾一度固定为 32 个页面（在 Linux/x86-32 中即为 131072 个字节），且包含了开销字节。自内核 2.6.23 版本开始，可以通过资源限制 `RLIMIT_STACK` 来控制 `argv` 和 `environ` 参数所使用的空间总量上限，在这种情况下，允许 `argv` 和 `environ` 参数使用的空间上限要比以前大出许多，具体限额为资源软限制 `RLIMIT_STACK` 的四分之一，`RLIMIT_STACK` 在调用 `execve()` 时已经生效。更多详细信息请参照 `execve(2)` 手册页。

许多程序（包括本书中的几个例子）使用 `getopt()` 库函数解析命令行选项（即以“-”符号开头的参数）。附录（Appendix）B 将描述 `getopt()` 函数。

6.7 环境列表

每一个进程都有与其相关的称之为环境列表（`environment list`）的字符串数组，或简称为环境（`environment`）。其中每个字符串都以名称=值（`name=value`）形式定义。因此，环境是“名称-值”的成对集合，可存储任何信息。常将列表中的名称称为环境变量（`environment variables`）。

新进程在创建之时，会继承其父进程的环境副本。这是一种原始的进程间通信方式，却颇为常用。环境（`environment`）提供了将信息从父进程传递给子进程的方法。由于子进程只有在创建时才能获得其父进程的环境副本，所以这一信息传递是单向的、一次性的。子进程创建后，父、子进程均可更改各自的环境变量，且这些变更对对方而言不再可见。

环境变量的常见用途之一是在 `shell` 中。通过在自身环境中放置变量值，`shell` 就可确保把这些值传递给其所创建的进程，并以此来执行用户命令。例如，环境变量 `SHELL` 被设置为 `shell` 程序本身的路径名，如果程序需要执行 `shell` 时，大多会将此变量视为需要执行的 `shell` 名称。

可以通过设置环境变量来改变一些库函数的行为。正因如此，用户无需修改程序代码或者重新链接相关库，就能控制调用该函数的应用程序行为。`getopt()` 函数就是其中一例（附录 B），可通过设置 `POSIXLY_CORRECT` 环境变量来改变此函数的行为。

大多数 `shell` 使用 `export` 命令向环境中添加变量值。

```
$ SHELL=/bin/bash          Create a shell variable
$ export SHELL              Put variable into shell process's environment
```

在 bash shell 和 Korn shell 中，可以简写为：

```
$ export SHELL=/bin/bash
```

在 C shell 中，使用的则是 setenv 命令：

```
% setenv SHELL /bin/bash
```

上述命令把一个值永久地添加到 shell 环境中，此后这个 shell 创建的所有子进程都将继承此环境。在任一时刻，可以使用 unset 命令撤销一个环境变量（在 C shell 中则使用 unsetenv 命令）。

在 Bourne shell 和其衍生 shell（诸如 bash shell 和 Korn shell）中，可使用下列语法向执行某应用程序的环境中添加一个变量值，而不影响其父 shell（和后续命令）：

```
$ NAME=value program
```

此命令仅向执行特定程序的子进程环境添加了一个（环境变量）定义。如果希望（多个变量对该程序有效），可以在 program 前放置多对赋值（以空格分隔）。

env 命令在运行程序时使用了一份经过修改的 shell 环境列表副本。可同时为 shell 环境列表副本增加和移除环境变量定义，以修改此环境列表。详细内容请参阅 env(1)手册。

printenv 命令显示当前的环境列表，此处是其输出的一例：

```
$ printenv
LOGNAME=mtk
SHELL=/bin/bash
HOME=/home/mtk
PATH=/usr/local/bin:/usr/bin:/bin:.
TERM=xterm
```

后续章节中将适时描述大多数上述环境变量的用途（也可参阅 environ(7)手册）。

由以上输出可知，环境列表的排列是无序的，列表中的字符串顺序不过是最易于实现的排列形式。一般而言，无序的环境列表不是问题，因为通常都是访问单个的环境变量，而非环境列表中按序排列的一串。

通过 Linux 专有的 /proc/PID/environ 文件检查任一进程的环境列表，每一个“NAME=value”对都以空字节终止。

从程序中访问环境

在 C 语言程序中，可以使用全局变量 `char **environ` 访问环境列表。（C 运行时启动代码定义了该变量并以环境列表位置为其赋值。）`environ` 与 `argv` 参数类似，指向一个以 NULL 结尾的指针列表，每个指针又指向一个以空字节终止的字符串。图 6-5 所示为与上述 `printenv` 命令输出环境相对应的环境列表数据结构。

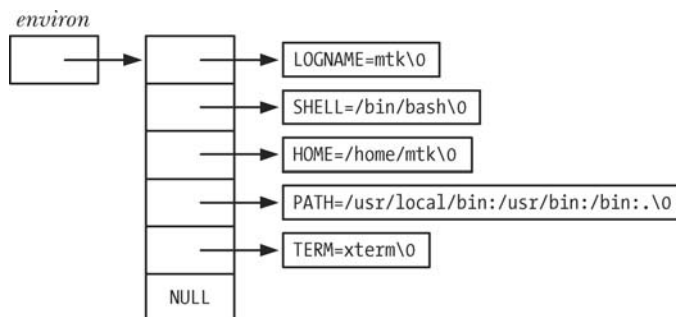


图 6-5：进程环境列表数据结构的示例

程序清单 6-3 中的程序通过访问 `environ` 变量来展示该进程环境中的所有值。该程序的输出结果与 `printenv` 命令的输出结果相同。程序中的循环利用指针来遍历 `environ` 变量。虽然可以把 `environ` 当成数组来使用（正如程序清单 6-2 中 `argv` 的用法），但这多少有些生硬，因为环境列表中各项的排列不分先后，而且也没有变量（相当于 `argc`）用来指定环境列表的长度。（出于同样原因，也没有对图 6-5 中的 `environ` 数组诸元素进行编号。）

程序清单 6-3：显示进程环境

```
----- proc/display_env.c
#include "tspi_hdr.h"

extern char **environ;

int
main(int argc, char *argv[])
{
    char **ep;

    for (ep = environ; *ep != NULL; ep++)
        puts(*ep);

    exit(EXIT_SUCCESS);
}
----- proc/display_env.c
```

另外，还可以通过声明 `main()` 函数中的第三个参数来访问环境列表：

```
int main(int argc, char *argv[], char *envp[])
```

该参数随即可被视为 `environ` 变量来使用，所不同的是，该参数的作用域在 `main()` 函数内。虽然 UNIX 系统普遍实现了这一特性，但还是要避免使用，因为除了局限于作用域限制外，该特性也不在 SUSv3 的规范之列。

`getenv()` 函数能够从进程环境中检索单个值。

```
#include <stdlib.h>

char *getenv(const char *name);

Returns pointer to (value) string, or NULL if no such variable
```

向 `getenv()` 函数提供环境变量名称，该函数将返回相应字符串指针。因此，就前面所示的环境（列表）示例来看，如果指定 `SHELL` 为参数 `name`，那么将返回 `/bin/bash`。如果不存在指定名称的环境变量，那么 `getenv()` 函数将返回 `NULL`。

以下是使用 `getenv()` 函数时可移植性方面的注意事项。

- SUSv3 规定应用程序不应修改 `getenv()` 函数返回的字符串，这是由于（在大多数 UNIX 实现中）该字符串实际上属于环境的一部分（即 `name=value` 字符串的 `value` 部分）。若需要改变一个环境变量的值，可以使用 `setenv()` 函数或 `putenv()` 函数（见下文）。
- SUSv3 允许 `getenv()` 函数的实现使用静态分配的缓冲区返回执行结果，后续对 `getenv()`、`setenv()`、`putenv()` 或者 `unsetenv()` 的函数调用可以重写该缓冲区。虽然 `glibc` 库的 `getenv()` 函数实现并未这样使用静态缓冲区，但具备可移植性的程序如需保留 `getenv()` 调用返回的字符串，就应先将返回字符串复制到其他位置，之后方可对上述函数发起调用。

修改环境

有时，对进程来说，修改其环境很有用处。原因之一是这一修改对该进程后续创建的所有子进程均可见。另一个可能的原因在于设定某一变量，以求对于将要载入进程内存的新程序（“execed”）可见。从这个意义上讲，环境不仅是一种进程间通信的形式，还是程序间通信的方法。（第 27 章将深入描述这一点，还将解释在同一进程中 `exec()` 函数如何使当前程序被一新程序所替代。）

`putenv()` 函数向调用进程的环境中添加一个新变量，或者修改一个已经存在的变量值。

```
#include <stdlib.h>

int putenv(char *string);

Returns 0 on success, or nonzero on error
```

参数 `string` 是一指针，指向 `name=value` 形式的字符串。调用 `putenv()` 函数后，该字符串就成为环境的一部分，换言之，`putenv` 函数将设定 `environ` 变量中某一元素的指向与 `string` 参数的指向位置相同，而非 `string` 参数所指向字符串的复制副本。因此，如果随后修改 `string` 参数所指的内容，这将影响该进程的环境。出于这一原因，`string` 参数不应为自动变量（即在栈中分配的字符数组¹），因为定义此变量的函数一旦返回，就有可能重写这块内存区域。

注意，`putenv()` 函数调用失败将返回非 0 值，而非 -1。

`putenv()` 函数的 `glibc` 库实现还提供了一个非标准扩展。如果 `string` 参数内容不包含一个等号 (=)，那么将从环境列表中移除以 `string` 参数命名的环境变量。

`setenv()` 函数可以代替 `putenv()` 函数，向环境中添加一个变量。

```
#include <stdlib.h>

int setenv(const char *name, const char *value, int overwrite);

Returns 0 on success, or -1 on error
```

`setenv()` 函数为形如 `name=value` 的字符串分配一块内存缓冲区，并将 `name` 和 `value` 所指向的字符串复制到此缓冲区，以此来创建一个新的环境变量。注意，不需要（实际上，是绝对不要）在 `name` 的结尾处或者 `value` 的开始处提供一个等号字符，因为 `setenv()` 函数会在向环境添加新变量时添加等号字符。

若以 `name` 标识的变量在环境中已经存在，且参数 `overwrite` 的值为 0，则 `setenv()` 函数将不改变环境，如果参数 `overwrite` 的值为非 0，则 `setenv()` 函数总是改变环境。

这一事实——`setenv()` 函数复制其参数（到环境中）——意味着与 `putenv()` 函数不同，之后对 `name` 和 `value` 所指字符串内容的修改将不会影响环境。此外，使用自动变量作为 `setenv()` 函数的参数也不会有任何问题。

`unsetenv()` 函数从环境中移除由 `name` 参数标识的变量。

```
#include <stdlib.h>

int unsetenv(const char *name);

Returns 0 on success, or -1 on error
```

¹ 译者注：请读者仔细思考 C 语言中数组与指针的对等关系。

同 `setenv()` 函数一样，参数 `name` 不应包含等号字符。

`setenv()` 函数和 `unsetenv()` 函数均来自 BSD，不如 `putenv()` 函数使用普遍。尽管起初的 POSIX.1 标准和 SUSv2 并未定义这两个函数，但 SUSv3 已将其纳入规范。

在 glibc 2.2.2 之前版本中，`unsetenv()` 函数原型的返回值为 `void` 类型，这与最初的 BSD 实现中 `unsetenv` 的函数原型相同，一些 UNIX 实现目前仍然沿用 BSD 原型。

有时，需要清除整个环境，然后以所选值进行重建。例如，为了以安全方式执行 `set-user-ID` 程序（38.8 节），就需要这样做。可以通过将 `environ` 变量赋值为 `NULL` 来清除环境。

```
environ = NULL;
```

这也正是 `clearenv()` 库函数的工作内容。

```
#define _BSD_SOURCE          /* Or: #define _SVID_SOURCE */
#include <stdlib.h>

int clearenv(void)

Returns 0 on success, or a nonzero on error
```

在某些情况下，使用 `setenv()` 函数和 `clearenv()` 函数可能会导致程序内存泄露。前面已然提及：`setenv()` 函数所分配的一块内存缓冲区，随之会成为进程环境的一部分。而调用 `clearenv()` 时则没有释放该缓冲区（`clearenv()` 调用并不知晓该缓冲区的存在，故而也无法将其释放）。反复调用这两个函数的程序，会不断产生内存泄露。实际上，这不大可能成为一个问题，因为程序通常仅在启动时调用 `clearenv()` 函数一次，用于移除继承自其父进程（即调用 `exec()` 函数来启动当前程序的程序）环境中的所有条目。

许多 UNIX 实现都支持 `clearenv()` 函数，但是 SUSv3 没有对此函数进行规范。SUSv3 规定如果应用程序直接修改 `environ` 变量，正如 `clearenv()` 函数所做的那样，则不对 `setenv()`、`unsetenv()` 和 `getenv()` 的行为进行定义。（这一作法的根本原因在于禁止符合 SUSv3 标准的应用程序直接修改环境，意在使 UNIX 实现能完全控制其实现环境变量时所采用的数据结构。）SUSv3 允许应用程序清空自身环境的唯一方法是首先获取所有环境变量的列表（通过 `environ` 变量获得所有环境变量的名称），然后逐一调用 `unsetenv()` 移除每个环境变量。

程序示例

程序清单 6-4 展示了本节讨论的所有函数的用法。该应用程序首先清空环境，然后向环境中逐一添加命令行参数所提供的环境变量定义；之后，如果环境中尚无名为 `GREET` 的变量，就向环境中添加该变量；接着，从环境中移除名为 `BYE` 的变量；最后打印当前环境列表。此处为该程序运行时输出结果的一例：

```
$./modify_env "GREET=Guten Tag" SHELL=/bin/bash BYE=Ciao
GREET=Guten Tag
SHELL=/bin/bash
$./modify_env SHELL=/bin/sh BYE=byebye
SHELL=/bin/sh
GREET>Hello world
```

如果将 `environ` 参数赋值为 `NULL`（正如程序清单 6-4 中 `clearenv()` 函数调用的所作所为），那么可以预见如下形式的循环（如程序清单 6-4 中使用的循环）将失败，因为 `*environ` 是无效的。

```
for (ep = environ; *ep != NULL; ep++)
    puts(*ep);
```

然而，如果 `setenv()` 函数和 `putenv()` 函数发现 `environ` 参数为 `NULL`，则会创建一个新的环境列表，并使 `environ` 参数指向此列表，结果上面的循环操作又将正确运行。

程序清单 6-4：修改进程环境

```
----- proc/modify_env.c
#define _GNU_SOURCE /* To get various declarations from <stdlib.h> */
#include <stdlib.h>
#include "tspi_hdr.h"

extern char **environ;

int
main(int argc, char *argv[])
{
    int j;
    char **ep;

    clearenv(); /* Erase entire environment */

    for (j = 1; j < argc; j++)
        if (putenv(argv[j]) != 0)
            errExit("putenv: %s", argv[j]);

    if (setenv("GREET", "Hello world", 0) == -1)
        errExit("setenv");

    unsetenv("BYE");

    for (ep = environ; *ep != NULL; ep++)
        puts(*ep);

    exit(EXIT_SUCCESS);
}
----- proc/modify_env.c
```

6.8 执行非局部跳转：setjmp()和 longjmp()

使用库函数 `setjmp()` 和 `longjmp()` 可执行非局部跳转（nonlocal goto）。术语“非局部（nonlocal）”是指跳转的目标为当前执行函数之外的某个位置。

C 语言，像许多其他编程语言一样，包含 `goto` 语句。这就好比打开了潘多拉的魔盒。若无止境的滥用，将使程序难以阅读和维护。不过偶尔也能一显身手，令程序更简单、更快速，或是兼而有之。

C 语言的 `goto` 语句存在一个限制，即不能从当前函数跳转到另一函数。然而，偶尔还是需要这一功能的。考虑错误处理中经常出现的如下场景：**在一个深度嵌套的函数调用中发生了错误，需要放弃当前任务，从多层函数调用中返回**，并在较高层级的函数中继续执行（也许甚至是在 `main()` 中）。要做到这一点，可以让每个函数都返回一个状态值，由函数的调用者检查并做相应处理。这一方法完全有效，而且，在许多情况下，是处理这类场景的理想方法。然而，有时候如果能从嵌套函数调用中跳出，返回该函数的调用者之一（当前调用者或者调

用者的调用者，等等)，编码会更为简单。setjmp()和 longjmp()就提供了这一功能。

由于在 C 语言中，所有函数作用域的层级相同（即标准 C 语言不支持嵌套函数申明，尽管 gcc 将此功能作为其扩展功能），所以 goto 语句不能应用于函数间跳转。给定两个函数 X 和 Y，编译器无从知晓当调用 Y 时，X 函数的栈帧是否在栈上，所以也无法判断从 Y 函数跳转（goto）到 X 函数是否可行。支持嵌套函数声明的语言，比如 Pascal 语言，允许 goto 从一个嵌套函数跳转到其调用者，编译器得以根据函数的静态作用域来确定函数动态作用域的某些信息。因此，编译器若在词法解析时获悉函数 Y 嵌套于函数 X 之内¹，也必然能够推断当调用 Y 时，X 函数的栈帧一定已然在栈中存在（即动态作用域），并能为函数 Y 产生 goto 代码，从 Y 中跳转到 X 函数的某处。

```
#include <setjmp.h>

int setjmp(jmp_buf env);

                Returns 0 on initial call, nonzero on return via longjmp()

void longjmp(jmp_buf env, int val);
```

setjmp()调用为后续由 longjmp()调用执行的跳转确立了跳转目标。该目标正是程序发起 setjmp()调用的位置。从编程角度来看，调用 longjmp()函数后，看起来就和从第二次调用 setjmp()返回时完全一样。通过查看 setjmp()返回的整数值，可以区分 setjmp 调用是初始返回还是第二次“返回”。初始调用返回值为 0，后续“伪”返回的返回值为 longjmp()调用中 val 参数所指定的任意值。通过对 val 参数使用不同值，能够区分出程序中跳转至同一目标的不同起跳位置。

如果指定 longjmp()函数的 val 参数值为 0，而 longjmp 函数对此又不做检查，就会导致模拟 setjmp()时返回值为 0，如同初次调用 setjmp()函数返回时一样。出于这一原因，如果指定 val 参数值为 0，则 longjmp()调用实际会将其替换为 1。

这两个函数的入参 env 为成功实现跳转提供了黏合剂。setjmp()函数把当前进程环境的各种信息保存到 env 参数中。调用 longjmp()时必须指定相同的 env 变量，以此来执行“伪”返回。由于对 setjmp()函数和 longjmp()函数的调用分别位于不同函数（否则，使用简单的 goto 即可），所以应该将 env 参数定义为全局变量，或者将 env 作为函数入参来传递，后一种做法较为少见。

调用 setjmp()时，env 除了存储当前进程的其他信息外，还保存了程序计数寄存器（指向当前正在执行的机器语言指令）和栈指针寄存器（标记栈顶）的副本。这些信息能够使后续的 longjmp()调用完成两个关键步骤的操作。

- 将发起 longjmp()调用的函数与之前调用 setjmp()的函数之间的函数栈帧从栈上剥离。有时又将此过程称为“解开栈（unwinding the stack）”，这是通过将栈指针寄存器重置为 env 参数内的保存值来实现的。
- 重置程序计数寄存器，使程序得以从初始的 setjmp()调用位置继续执行。同样，此功能是通过 env 参数中的保存值（程序计数寄存器）来实现的。

程序示例

程序清单 6-5 展示了 setjmp()和 longjmp()函数的用法。该程序通过 setjmp()的初始调用建立

¹ 译者注：即上文所谓静态作用域。

了一个跳转目标，接下来的 switch（针对 setjmp()调用的返回值）用于检测是初次从 setjmp()调用返回还是在调用 longjmp()后返回。当 setjmp()调用返回值为 0 时，亦即对 setjmp()的初始调用完成后，将调用 f1()函数，f1()函数根据 argc 参数值（即命令行参数个数）来决定是立刻调用 longjmp()函数还是继续去调用 f2()函数。如果是调用 f2()函数，则 f2()函数将马上调用 longjmp()函数。两处对 longjmp()的调用都会使进程恢复到调用 setjmp()的位置。程序在两处调用中为 val 参数设定了不同值，以供 main()函数的 switch 语句区分发生跳转的函数，并打印相应信息。

在不带任何命令行参数的情况下运行程序清单 6-5 中的程序，结果如下所示：

```
$ ./longjmp
Calling f1() after initial setjmp()
We jumped back from f1()

指定命令行参数，会使程序跳转发生在函数 f2()中：
$ ./longjmp x
Calling f1() after initial setjmp()
We jumped back from f2()
```

程序清单 6-5：展示函数 setjmp()和 longjmp()的用法

```
proc/longjmp.c

#include <setjmp.h>
#include "tspi_hdr.h"

static jmp_buf env;

static void
f2(void)
{
    longjmp(env, 2);
}

static void
f1(int argc)
{
    if (argc == 1)
        longjmp(env, 1);
    f2();
}

int
main(int argc, char *argv[])
{
    switch (setjmp(env)) {
    case 0: /* This is the return after the initial setjmp() */
        printf("Calling f1() after initial setjmp()\n");
        f1(argc); /* Never returns... */
        break; /* ... but this is good form */

    case 1:
        printf("We jumped back from f1()\n");
        break;

    case 2:
        printf("We jumped back from f2()\n");
        break;
    }
}
```

```
    exit(EXIT_SUCCESS);
}
```

proc/longjmp.c

对 setjmp()函数的使用限制

SUSv3 和 C99 规定，对 setjmp()的调用只能在如下语境中使用。

- 构成选择或迭代语句中（if、switch、while 等）的整个控制表达式。
- 作为一元操作符!（not）的操作对象，其最终表达式构成了选择或迭代语句的整个控制表达式。
- 作为比较操作（==、!=、<等）的一部分，另一操作对象必须是一个整数常量表达式，且其最终表达式构成选择或迭代语句的整个控制表达式。
- 作为独立的函数调用，且没有嵌入到更大的表达式之中。

注意：C 语言赋值语句不在上述列表之列。以下形式的语句是不符合标准的：

```
s = setjmp(env);                /* WRONG! */
```

之所以规定这些限制，是因为作为常规函数的 setjmp()实现无法保证拥有足够信息来保存所有寄存器值和封闭表达式中用到的临时栈位置，以便于在 longjmp()调用后此类信息能得以正确恢复。因此，仅允许在足够简单且无需临时存储的表达式中调用 setjmp()。

滥用 longjmp()

如果将 env 缓冲区定义为全局变量，对所有函数可见（这也是通常用法），那么就可以执行如下操作序列。

1. 调用函数 x()，使用 setjmp()调用在全局变量 env 中建立一个跳转目标。
2. 从函数 x()中返回。
3. 调用函数 y()，使用 env 变量调用 longjmp()函数。

这是一个严重错误，因为 longjmp()调用不能跳转到一个已经返回的函数中。思考一下，在这种情况下，longjmp()函数会对栈打什么主意——尝试将栈解开，恢复到一个不存在的栈帧位置，这无疑将引起混乱。如果幸运的话，程序会一死（crash）了之。然而，取决于栈的状态，也可能会引起调用与返回间的死循环，而程序好像真地从一个当前并未执行的函数中返回了。（在多线程程序中有与之相类似的滥用，在线程某甲中调用 setjmp()函数，却在线程某乙中调用 longjmp()。）

SUSv3 规定，如果从嵌套的信号处理器（signal handler）（即信号某甲的处理器正在运行时，又发起对信号某乙处理器的调用）中调用 longjmp()函数，则该程序的行为未定义。

优化编译器的问题

优化编译器会重组程序的指令执行顺序，并在 CPU 寄存器中，而非 RAM 中存储某些变量。这种优化一般依赖于反映了程序词法结构的运行时（run-time）控制流程。由于 setjmp()和 longjmp()的跳转操作需在运行时才能得以确立和执行，并未在程序的词法结构中有所反映，故而编译器在进行优化时也无法将其考虑在内。此外，某些应用程序二进制接口（ABI）实现的语义要求 longjmp()函数恢复先前 setjmp()调用所保存的 CPU 寄存器副本。这意味着 longjmp()操作会致使经过优化的变量被赋以错误值。程序清单 6-6 中的程序行为就是其中一例。

程序清单 6-6: 编译器的优化和 longjmp()函数相互作用的示例

```
----- proc/setjmp_vars.c
#include <stdio.h>
#include <stdlib.h>
#include <setjmp.h>

static jmp_buf env;

static void
doJump(int nvar, int rvar, int vvar)
{
    printf("Inside doJump(): nvar=%d rvar=%d vvar=%d\n", nvar, rvar, vvar);
    longjmp(env, 1);
}

int
main(int argc, char *argv[])
{
    int nvar;
    register int rvar;          /* Allocated in register if possible */
    volatile int vvar;         /* See text */

    nvar = 111;
    rvar = 222;
    vvar = 333;

    if (setjmp(env) == 0) {    /* Code executed after setjmp() */
        nvar = 777;
        rvar = 888;
        vvar = 999;
        doJump(nvar, rvar, vvar);
    } else {                   /* Code executed after longjmp() */
        printf("After longjmp(): nvar=%d rvar=%d vvar=%d\n", nvar, rvar, vvar);
    }

    exit(EXIT_SUCCESS);
}
----- proc/setjmp_vars.c
```

以常规方式编译程序清单 6-6 中的程序，输出结果符合预期。

```
$ cc -o setjmp_vars setjmp_vars.c
$ ./setjmp_vars
Inside doJump(): nvar=777 rvar=888 vvar=999
After longjmp(): nvar=777 rvar=888 vvar=999
```

然而，若以优化方式编译该程序，结果就有些出乎预料了。

```
$ cc -O -o setjmp_vars setjmp_vars.c
$ ./setjmp_vars
Inside doJump(): nvar=777 rvar=888 vvar=999
After longjmp(): nvar=111 rvar=222 vvar=999
```

此处，在 longjmp()调用后，nvar 和 rvar 参数被重置为 setjmp()初次调用时的值。起因是优化器对代码的重组受到 longjmp()调用的干扰。作为候选优化对象的任一局部变量可能都难免会遇到这类问题，一般包含指针变量和 char、int、float、long 等任何简单类型的变量。

将变量声明为 `volatile`，是告诉优化器不要对其进行优化，从而避免了代码重组。在上面的程序输出中，无论编译优化与否，声明为 `volatile` 的变量 `vvar` 都得到了正确处理。

因为不同的优化器有着不同的优化方法，具备良好移植性的程序应在调用 `setjmp()` 的函数中，将上述类型的所有局部变量都声明为 `volatile`。

若在 GNU C 语言编译器中加入 `-Wextra`（产生额外的警告信息）选项，`setjmp_vars.c` 程序的编译结果将显示有帮助的警告信息如下：

```
$ cc -Wall -Wextra -O -o setjmp_vars setjmp_vars.c
setjmp_vars.c: In function `main':
setjmp_vars.c:17: warning: variable `nvar' might be clobbered by `longjmp' or
`vfork'
setjmp_vars.c:18: warning: variable `rvar' might be clobbered by `longjmp' or
`vfork'
```

无论优化与否，查看编译 `setjmp_vars.c` 程序所产生的汇编语言输出都是有益的。`cc -S` 命令产生一个以 `.s` 为扩展名的文件，内容为程序的汇编代码。

尽可能避免使用 `setjmp()` 函数和 `longjmp()` 函数

如果说 `goto` 语句会使程序难以阅读，那么非局部跳转会让事情的糟糕程度增加一个数量级，因为它能在程序中任意两个函数间传递控制。因此，应当慎用 `setjmp()` 函数和 `longjmp()` 函数。在设计和编码时花点心思来避免使用这两个函数，这通常是值得的。程序更具可读性，可能会更具可移植性。话虽如此，但在编写信号处理器时，这些函数偶尔还会派上用场——讨论信号时将重新论及这些函数的变体（参见 21.2.1 节中的 `sigsetjmp()` 函数和 `siglongjmp()` 函数）。

6.9 总结

每个进程都有一个唯一进程标识号（`process ID`），并保存有对其父进程号的记录。

进程的虚拟内存逻辑上被划分成许多段：文本段、（初始化和非初始化的）数据段、栈和堆。

栈由一系列帧组成，随函数调用而增，随函数返回而减。每个帧都包含有函数局部变量、函数实参以及单个函数调用的调用链接信息。

程序调用时，命令行参数通过 `argc` 和 `argv` 参数提供给 `main()` 函数。通常，`argv[0]` 包含调用程序的名称。

每个进程都会获得其父进程环境列表的一个副本，即一组“名称-值”键值对。全局变量 `environ` 和各种库函数允许进程访问和修改其环境列表中的变量。

`setjmp()` 函数和 `longjmp()` 函数提供了从函数某甲执行非局部跳转到函数某乙（栈解开）的方法。在调用这些函数时，为避免编译器优化所引发的问题，应使用 `volatile` 修饰符声明变量。非局部跳转会使程序难于阅读和维护，应尽量避免使用。

更多资料

[Tanenbaum, 2007]和[Vahalia, 1996]详细描述了虚拟内存管理。[Gorman, 2004]则详细描述了 Linux 内核内存管理算法和代码。

6.9 练习

- 6-1. 编译程序清单 6-1 中的程序 (`mem_segments.c`), 使用 `ls -l` 命令显示可执行文件的大小。虽然该程序包含一个大约 10MB 的数组, 但可执行文件大小要远小于此, 为什么?
- 6-2. 编写一个程序, 观察当使用 `longjmp()` 函数跳转到一个已经返回的函数时会发生什么?
- 6-3. 使用 `getenv()` 函数、`putenv()` 函数, 必要时可直接修改 `environ`, 来实现 `setenv()` 函数和 `unsetenv()` 函数。此处的 `unsetenv()` 函数应检查是否对环境变量进行了多次定义, 如果是多次定义则将移除对该变量的所有定义 (`glibc` 版本的 `unsetenv()` 函数实现了这一功能)。

第 7 章

内存分配

许多系统程序需要为动态数据结构（例如，链表和二叉树）分配额外内存，此类数据结构的大小由运行时所获取的信息决定。本章将介绍用于在堆或堆栈上分配内存的函数。

7.1 在堆上分配内存

进程可以通过增加堆的大小来分配内存，所谓堆是一段长度可变的连续虚拟内存，始于进程的未初始化数据段末尾，随着内存的分配和释放而增减（见图 6-1）。通常将堆的当前内存边界称为“program break”。

稍后将介绍 C 语言程序分配内存所惯用的 malloc 函数族，但首先还要从 malloc 函数族所基于的 brk()和 sbrk()开始谈起。

7.1.1 调整 program break: brk()和 sbrk()

改变堆的大小（即分配或释放内存），其实就像命令内核改变进程的 program break 位置一样简单。最初，program break 正好位于未初始化数据段末尾之后（如图 6-1 所示，与&end 位置相同）。

在 program break 的位置抬升后，程序可以访问新分配区域内的任何内存地址，而此时物理内存页尚未分配。内核会在进程首次试图访问这些虚拟内存地址时自动分配新的物理内存页。

传统的 UNIX 系统提供了两个操纵 program break 的系统调用：brk()和 sbrk()，在 Linux 中依然可用。虽然代码中很少直接使用这些系统调用，但了解它们有助于弄清内存分配的工作过程。

```
#include <unistd.h>

int brk(void *end_data_segment);
                                     Returns 0 on success, or -1 on error

void *sbrk(intptr_t increment);
                                     Returns previous program break on success, or (void *) -1 on error
```

系统调用 brk()会将 program break 设置为参数 end_data_segment 所指定的位置。由于虚拟

内存以页为单位进行分配，`end_data_segment` 实际会四舍五入到下一个内存页的边界处。

当试图将 `program break` 设置为一个低于其初始值（即低于 `&end`）的位置时，有可能会致无法预知的行为，例如，当程序试图访问的数据位于初始化或未初始化数据段中当前尚不存在的部分时，就会引发分段内存访问错误（`segmentation fault`）（`SIGSEGV` 信号，在 20.2 节描述）。`program break` 可以设定的精确上限取决于一系列因素，这包括进程中对数据段大小的资源限制（36.3 节中描述的 `RLIMIT_DATA`），以及内存映射、共享内存段、共享库的位置。

调用 `sbrk()` 将 `program break` 在原有地址上增加从参数 `increment` 传入的大小。（在 Linux 中，`sbrk()` 是在 `brk()` 基础上实现的一个库函数。）用于声明 `increment` 的 `intptr_t` 类型属于整数数据类型。若调用成功，`sbrk()` 返回前一个 `program break` 的地址。换言之，如果 `program break` 增加，那么返回值是指向这块新分配内存起始位置的指针。

调用 `sbrk(0)` 将返回 `program break` 的当前位置，对其不做改变。在意图跟踪堆的大小，或是监视内存分配函数包的行为时，可能会用到这一用法。

SUSv2 定义了 `brk()` 和 `sbrk()`，标记为 Legacy（传统）。但 SUSv3 删除了这些定义。

7.1.2 在堆上分配内存：`malloc()` 和 `free()`

一般情况下，C 程序使用 `malloc` 函数族在堆上分配和释放内存。较之 `brk()` 和 `sbrk()`，这些函数具备不少优点，如下所示。

- 属于 C 语言标准的一部分。
- 更易于在多线程程序中使用。
- 接口简单，允许分配小块内存。
- 允许随意释放内存块，它们被维护于一张空闲内存列表中，在后续内存分配调用时循环使用。

`malloc()` 函数在堆上分配参数 `size` 字节大小的内存，并返回指向新分配内存起始位置处的指针，其所分配的内存未经初始化。

```
#include <stdlib.h>

void *malloc(size_t size);

Returns pointer to allocated memory on success, or NULL on error
```

由于 `malloc()` 的返回类型为 `void*`，因而可以将其赋给任意类型的 C 指针。`malloc()` 返回内存块所采用的字节对齐方式，总是适宜于高效访问任何类型的 C 语言数据结构。在大多数硬件架构上，这实际意味着 `malloc` 是基于 8 字节或 16 字节边界来分配内存的。¹

SUSv3 规定：调用 `malloc(0)` 要么返回 `NULL`，要么是一小块可以（并且应该）用 `free()` 释放的内存。Linux 的 `malloc(0)` 行为遵循后者。

若无法分配内存（或许是因为已经抵达 `program break` 所能达到的地址上限），则 `malloc()` 返回 `NULL`，并设置 `errno` 以返回错误信息。虽然分配内存失败的可能性很小，但所有对 `malloc()` 以及后续提及的相关函数的调用都应对返回值进行错误检查。

¹ 译者注：遵作者邮件嘱改动原文，据称此为编译器提高内存访问效率的举措之一。

free()函数释放 ptr 参数所指向的内存块，该参数应该是之前由 malloc()，或者本章后续描述的其他堆内存分配函数之一所返回的地址。

```
#include <stdlib.h>

void free(void *ptr);
```

一般情况下，free()并不降低 program break 的位置，而是将这块内存填充到空闲内存列表中，供后续的 malloc()函数循环使用。这么做是出于以下几个原因。

- 被释放的内存块通常会位于堆的中间，而非堆的顶部，因而降低 program break 是不可能的。
- 它最大限度地减少了程序必须执行的 sbrk()调用次数。（正如 3.1 节指出的，系统调用的开销虽小，却也颇为可观。）
- 在大多数情况下，降低 program break 的位置不会对那些分配大量内存的程序有多少帮助，因为它们通常倾向于持有已分配内存或是反复释放和重新分配内存，而非释放所有内存后再持续运行一段时间。

如果传给 free()的是一个空指针，那么函数将什么都不做。（换句话说，给 free()传入一个空指针并不是错误代码。）

在调用 free()后对参数 ptr 的任何使用，例如将其再次传递给 free()，将产生错误，并可能导致不可预知的结果。

程序示例

程序清单 7-1 中的程序说明了 free()函数对 program break 的影响。该程序在分配了多块内存后，根据（可选）命令行参数来释放其中的部分或全部。

前两个命令行参数指定了分配内存块的数量和大小。第三个命令行参数指定了释放内存块的循环步长。如果是 1（这也是省略此参数时的默认值），那么程序将释放每块已分配的内存，如果为 2，那么每隔一块释放一块已分配内存，以此类推。第四个和第五个命令行参数指定需要释放的内存块范围。如果省略这两个参数，那么将（以第三个命令行参数所指定的步长）释放全部范围内的已分配内存。

程序清单 7-1：示范释放内存时 program break 的行为

```
----- memalloc/free_and_sbrk.c
#include "t1pi_hdr.h"

#define MAX_ALLOCS 1000000

int
main(int argc, char *argv[])
{
    char *ptr[MAX_ALLOCS];
    int freeStep, freeMin, freeMax, blockSize, numAllocs, j;

    printf("\n");

    if (argc < 3 || strcmp(argv[1], "--help") == 0)
        usageErr("%s num-allocs block-size [step [min [max]]]\n", argv[0]);

    numAllocs = getInt(argv[1], GN_GT_0, "num-allocs");
```

```

if (numAllocs > MAX_ALLOCS)
    cmdLineErr("num-allocs > %d\n", MAX_ALLOCS);

blockSize = getInt(argv[2], GN_GT_0 | GN_ANY_BASE, "block-size");

freeStep = (argc > 3) ? getInt(argv[3], GN_GT_0, "step") : 1;
freeMin = (argc > 4) ? getInt(argv[4], GN_GT_0, "min") : 1;
freeMax = (argc > 5) ? getInt(argv[5], GN_GT_0, "max") : numAllocs;

if (freeMax > numAllocs)
    cmdLineErr("free-max > num-allocs\n");

printf("Initial program break:          %10p\n", sbrk(0));

printf("Allocating %d*%d bytes\n", numAllocs, blockSize);
for (j = 0; j < numAllocs; j++) {
    ptr[j] = malloc(blockSize);
    if (ptr[j] == NULL)
        errExit("malloc");
}

printf("Program break is now:            %10p\n", sbrk(0));

printf("Freeing blocks from %d to %d in steps of %d\n",
       freeMin, freeMax, freeStep);
for (j = freeMin - 1; j < freeMax; j += freeStep)
    free(ptr[j]);

printf("After free(), program break is: %10p\n", sbrk(0));

exit(EXIT_SUCCESS);
}

```

memalloc/free_and_sbrk.c

用下面的命令行运行程序清单 7-1 的程序，将会分配 1000 个内存块，且每隔一个内存块释放一个内存块。

```
$ ./free_and_sbrk 1000 10240 2
```

输出结果显示，释放所有内存块后，program break 的位置仍然与分配所有内存块后的水平相当。

```

Initial program break:          0x804a6bc
Allocating 1000*10240 bytes
Program break is now:           0x8a13000
Freeing blocks from 1 to 1000 in steps of 2
After free(), program break is: 0x8a13000

```

下面的命令行要求除了最后一块内存块，释放所有已分配的内存块。再一次，program break 保持在了“高水位线”。

```

$ ./free_and_sbrk 1000 10240 1 1 999
Initial program break:          0x804a6bc
Allocating 1000*10240 bytes
Program break is now:           0x8a13000
Freeing blocks from 1 to 999 in steps of 1
After free(), program break is: 0x8a13000

```

但是，如果在堆顶部释放完整的一组连续内存块，会观察到 program break 从峰值上降下来，这表明 free() 使用了 sbrk() 来降低 program break。在这里，命令行释放了已分配内存的最后 500 个

内存块。

```
$ ./free_and_sbrk 1000 10240 1 500 1000
Initial program break:      0x804a6bc
Allocating 1000*10240 bytes
Program break is now:      0x8a13000
Freeing blocks from 500 to 1000 in steps of 1
After free(), program break is: 0x852b000
```

在这种情况下，`free()`函数的 `glibc` 实现会在释放内存时将相邻的空闲内存块合并为一整块更大的内存（这样做是为了避免在空闲内存列表中包含大量的小块内存碎片，这些碎片会因空间太小而难以满足后续的 `malloc()` 请求），因而也有能力识别出堆顶部的整个空闲区域。

仅当堆顶空闲内存“足够”大的时候，`free()`函数的 `glibc` 实现会调用 `sbrk()`来降低 `program break` 的地址，至于“足够”与否则取决于 `malloc` 函数包行为的控制参数（128 KB 为典型值）。这减少了必须对 `sbrk()`发起的调用次数（亦即对 `brk()`系统调用的调用次数）。

调用 `free()`还是不调用 `free()`

当进程终止时，其占用的所有内存都会返还给操作系统，这包括在堆内存中由 `malloc` 函数包内一系列函数所分配的内存。基于内存的这一自动释放机制，对于那些分配了内存并在进程终止前持续使用的程序而言，通常会省略对 `free()`的调用。这在程序中分配了多块内存的情况下可能会特别有用，因为加入多次对 `free()`的调用不但会消耗大量的 CPU 时间，而且可能会使代码趋于复杂。

虽然依靠终止进程来自动释放内存对大多数程序来说是可以接受的，但基于以下几个原因，最好能够在程序中显式释放所有的已分配内存。

- 显式调用 `free()`能使程序在未来修改时更具可读性和可维护性。
- 如果使用 `malloc` 调试库（如下所述）来查找程序的内存泄漏问题，那么会将任何未经显式释放处理的内存报告为内存泄漏。这会使发现真正内存泄漏的工作复杂化。

7.1.3 `malloc()`和 `free()`的实现

尽管 `malloc()`和 `free()`所提供的内存分配接口比之 `brk()`和 `sbrk()`要容易许多，但在使用时仍然容易犯下各种编程错误。理解 `malloc()`和 `free()`的实现，将使我们洞悉产生这些错误的原因以及如何才能避免此类错误。

`malloc()`的实现很简单。它首先会扫描之前由 `free()`所释放的空闲内存块列表，以求找到尺寸大于或等于要求的一块空闲内存。（取决于具体实现，采用的扫描策略会有所不同。例如，`first-fit` 或 `best-fit`。）如果这一内存块的尺寸正好与要求相当，就把它直接返回给调用者。如果是一块较大的内存，那么将对其进行分割，在将一块大小相当的内存返回给调用者的同时，把较小的那块空闲内存块保留在空闲列表中。

如果在空闲内存列表中根本找不到足够大的空闲内存块，那么 `malloc()`会调用 `sbrk()`以分配更多的内存。为减少对 `sbrk()`的调用次数，`malloc()`并未只是严格按所需字节数来分配内存，而是以更大幅度（以虚拟内存页大小的数倍）来增加 `program break`，并将超出部分置于空闲内存列表。

至于 `free()`函数的实现则更为有趣。当 `free()`将内存块置于空闲列表之上时，是如何知晓内存块大小的？这是通过一个小技巧来实现的。当 `malloc()`分配内存块时，会额外分配几个字节来存放记录这块内存大小的整数值。该整数位于内存块的起始处，而实际返回给调用者的内存地址恰好位于这一长度记录字节之后，如图 7-1 所示。

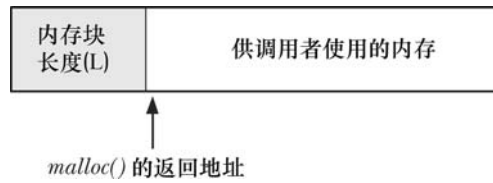


图 7-1: malloc()返回的内存块

当将内存块置于空闲内存列表（双向链表）时，free()会使用内存块本身的空间来存放链表指针，将自身添加到列表中，如图 7-2 所示。

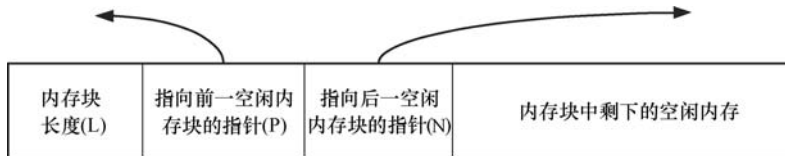


图 7-2: 空闲列表中的内存块

随着对内存不断地释放和重新分配，空闲列表中的空闲内存会和已分配的在用内存混杂在一起，如图 7-3 所示。

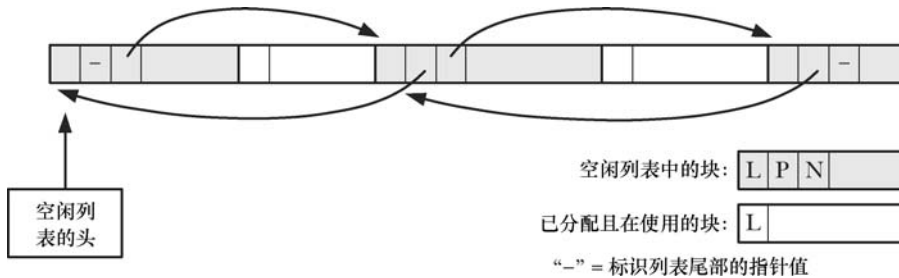


图 7-3: 包含有已分配内存和空闲内存列表的堆

应该认识到，C 语言允许程序创建指向堆中任意位置的指针，并修改其指向的数据，包括由 free()和 malloc()函数维护的内存块长度、指向前一空闲块和后一空闲块的指针。辅之以之前的描述，一旦推究起隐晦难解的编程缺陷来，这无疑形同掉进了火药桶。例如，假设经由一个错误指针，程序无意间增加了冠于一块已分配内存的长度值，并随即释放这块内存，free()因之会在空闲列表中记录下这块长度失真的内存。随后，malloc()也许会重新分配这块内存，从而导致如下场景：程序的两个指针分别指向两块它认为互不相干的已分配内存，但实际上这两块内存却相互重叠。至于其他的出错情况则数不胜数。

要避免这类错误，应该遵守以下规则。

- 分配一块内存后，应当小心谨慎，不要改变这块内存范围外的任何内容。错误的指针运算，或者循环更新内存块内容时出现的“off-by-one”（一字之偏）¹错误，都有可能导致这一情况。
- 释放同一块已分配内存超过一次是错误的。Linux 上的 glibc 库经常报出分段错误（SIGSEGV 信号）。这是好事，因为它提醒我们犯下了一个编程错误。然而，当两次

¹ 译者注：详见 http://en.wikipedia.org/wiki/Off-by-one_error。

释放同一块内存时，更常见的后果是导致不可预知的行为。

- 若非经由 `malloc` 函数包中函数所返回的指针，绝不能在调用 `free()` 函数时使用。
- 在编写需要长时间运行的程序（例如，`shell` 或网络守护进程）时，出于各种目的，如果需要反复分配内存，那么应当确保释放所有已使用完毕的内存。如若不然，堆将稳步增长，直至抵达可用虚拟内存的上限，在此之后分配内存的任何尝试都将以失败告终。这种情况被称之为“内存泄漏”。

malloc 调试的工具和库

如果不遵循上述准则，可能会在代码中引入既难以理解又难以重现的缺陷。而使用 `glibc` 提供的 `malloc` 调试工具或者任何一款 `malloc` 调试库，都会显著降低发现这些缺陷的难度，这也是设计它们的目的所在。

以下是 `glibc` 提供的 `malloc` 调试工具的部分功能。

- `mtrace()` 和 `muntrace()` 函数分别在程序中打开和关闭对内存分配调用进行跟踪的功能。这些函数要与环境变量 `MALLOC_TRACE` 搭配使用，该变量定义了写入跟踪信息的文件名。在被调用时，`mtrace()` 会检查是否定义了该文件，又是否可以打开文件并写入。如果一切正常，那么会在文件里跟踪和记录所有对 `malloc` 函数包中函数的调用。由于生成文件不易于理解，还提供有一个脚本（`mtrace`）用于分析文件，并生成易于理解的汇总报告。出于安全原因，设置用户 ID 和设置组 ID 的程序会忽略对 `mtrace()` 的调用。
- `mcheck()` 和 `mprobe()` 函数允许程序对已分配内存块进行一致性检查。例如，当程序试图在已分配内存之外进行写操作时，它们将捕获这个错误。这些函数提供的功能和下述 `malloc` 调试库有重叠之处。使用这些函数的程序，必须使用 `cc-lmcheck` 选项与 `mcheck` 库链接。
- `MALLOC_CHECK_` 环境变量（注意结尾处的下划线）提供了类似于 `mcheck()` 和 `mprobe()` 函数的功能。（两者之间的一个显著区别在于使用：`MALLOC_CHECK_` 无需对程序进行修改和重新编译。）通过为此变量设置不同的整数值，可以控制程序对内存分配错误的响应方式。可能的设置有：0，意即忽略错误；1，意即在标准错误输出（`stderr`）中打印诊断错误；2，意即调用 `abort()` 来终止程序。并非所有的内存分配和释放错误都是由 `MALLOC_CHECK_` 检测出的，它所发现的只是常见错误。但是，这种技术快速、易用，较之于 `malloc` 调试库具有较低的运行时开销。出于安全原因，设置用户 ID 和设置组 ID 的程序将忽略 `MALLOC_CHECK_` 设置。

关于以上所有功能更为详细的信息可以参考 `glibc` 手册。

而就 `malloc` 调试库而言，其提供了和标准 `malloc` 函数包相同的 API，但附加了捕获内存分配错误的功能。要使用调试库，需要在编译时链接调试库，而非标准 C 函数库的 `malloc` 函数包。由于调试库通常会降低运行速度，增加内存消耗，或是两者兼而有之，应当仅在调试时使用，而在正式发布产品时链接标准库的 `malloc` 包。这些库分别是：`Electric Fence` (<http://www.perens.com/FreeSoftware/>)、`dmalloc` (<http://dmalloc.com/>)、`Valgrind` (<http://valgrind.org/>)、`Insure++` (<http://www.parasoft.com/>)。

Valgrind 和 Insure++ 能够发现许多堆内存分配之外的其他类型错误。可以访问其各自网站，以获取详细信息。

控制和监测 malloc 函数包

glibc 手册介绍了一系列非标准函数，可用于监测和控制 malloc 包中函数的内存分配，其中包括如下几个函数。

- 函数 `malloc()` 能修改各项参数，以控制 `malloc()` 所采用的算法。例如，此类参数之一就指定了在调用 `sbrk()` 函数进行堆收缩之前，在空闲列表尾部必须保有的可释放内存空间的最小值。另一参数则规定了从堆中分配的内存块大小的上限，超出上限的内存块则使用 `mmap()` 系统调用（参见 49.7 节）来分配。
- `mallinfo()` 函数返回一个结构，其中包含由 `malloc()` 分配内存的各种统计数据。

众多 UNIX 实现提供各种版本的 `malloc()` 和 `mallinfo()`。然而，这些函数所提供的接口却随实现而不同，因而也无法移植。

7.1.4 在堆上分配内存的其他方法

除了 `malloc()`，C 函数库还提供了一系列在堆上分配内存的其他函数，在这里将逐一介绍。

用 `calloc()` 和 `realloc()` 分配内存

函数 `calloc()` 用于给一组相同对象分配内存。

```
#include <stdlib.h>

void *calloc(size_t numitems, size_t size);

Returns pointer to allocated memory on success, or NULL on error
```

参数 `numitems` 指定分配对象的数量，`size` 指定每个对象的大小。在分配了适当大小的内存块后，`calloc()` 返回指向这块内存起始处的指针（如果无法分配内存，则返回 `NULL`）。与 `malloc()` 不同，`calloc()` 会将已分配的内存初始化为 0。

下面是 `calloc()` 的一个使用范例：

```
struct { /* Some field definitions */ } myStruct;
struct myStruct *p;

p = calloc(1000, sizeof(struct myStruct));
if (p == NULL)
    errExit("calloc");
```

`realloc()` 函数用来调整（通常是增加）一块内存的大小，而此块内存应是之前由 `malloc` 包中函数所分配的。

```
#include <stdlib.h>

void *realloc(void *ptr, size_t size);

Returns pointer to allocated memory on success, or NULL on error
```

参数 `ptr` 是指向需要调整大小的内存块的指针。参数 `size` 指定所需调整大小的期望值。

如果成功，`realloc()` 返回指向大小调整后内存块的指针。与调用前的指针相比，二者指向的位置可能不同。如果发生错误，`realloc()` 返回 `NULL`，对 `ptr` 指针指向的内存块则原封不动（SUSv3 要求满足这一约定）。

若 `realloc()` 增加了已分配内存块的大小，则不会对额外分配的字节进行初始化。

使用 `calloc()` 或 `realloc()` 分配的内存应使用 `free()` 来释放。

调用 `realloc(ptr,0)` 等效于在 `free(ptr)` 之后调用 `malloc(0)`。若 `ptr` 为 `NULL`，则 `realloc(NULL, size)` 相当于调用 `malloc(size)`。

通常情况下，当增大已分配内存时，`realloc()` 会试图去合并空闲列表中紧随其后¹且大小满足要求的内存块。若原内存块位于堆的顶部，那么 `realloc()` 将对堆空间进行扩展。如果这块内存位于堆的中部，且紧邻其后的空闲内存空间大小不足，`realloc()` 会分配一块新内存，并将原有数据复制到新内存块中。最后这种情况最为常见，还会占用大量 CPU 资源。一般情况下，应尽量避免调用 `realloc()`。

既然 `realloc()` 可能会移动内存，对这块内存的后续引用就必须使用 `realloc()` 的返回指针。可以用 `realloc()` 来重新定位由变量 `ptr` 指向的内存块，代码如下：

```
nptr = realloc(ptr, newsize);
if (nptr == NULL) {
    /* Handle error */
} else {
    /* realloc() succeeded */
    ptr = nptr;
}
```

本例并没有把 `realloc()` 的返回值直接赋给 `ptr`，因为一旦调用 `realloc()` 失败，那么 `ptr` 会被置为 `NULL`，从而无法访问现有内存块。

由于 `realloc()` 可能会移动内存块，任何指向该内存块内部的指针在调用 `realloc()` 之后都可能不再可用。仅有一种内存块内的位置引用方法依然有效，即以指向此块内存起始处的指针再加上一个偏移量来进行定位，这将在 48.6 节中详细讨论。

分配对齐的内存: `memalign()` 和 `posix_memalign()`

设计函数 `memalign()` 和 `posix_memalign()` 的目的在于分配内存时，起始地址要与 2 的整数次幂边界对齐，该特征对于某些应用非常有用（例如程序清单 13-1）。

```
#include <malloc.h>

void *memalign(size_t boundary, size_t size);

Returns pointer to allocated memory on success, or NULL on error
```

函数 `memalign()` 分配 `size` 个字节的内存，起始地址是参数 `boundary` 的整数倍，而 `boundary` 必须是 2 的整数次幂。函数返回已分配内存的地址。

函数 `memalign()` 并非在所有 UNIX 实现上都存在。大多数提供 `memalign()` 的其他 UNIX 实现都要求引用 `<stdlib.h>` 而非 `<malloc.h>` 以获得函数声明。

SUSv3 并未纳入 `memalign()`，而是规范了一个类似函数，名为 `posix_memalign()`。该函数由标准委员会于近期创制，只是出现在了少数 UNIX 实现上。

```
#include <stdlib.h>

int posix_memalign(void **memptr, size_t alignment, size_t size);

Returns 0 on success, or a positive error number on error
```

¹ 译者注：参见图 7-3 堆中空闲内存块与已分配内存块的“杂居”状态，此处应指与 `ptr` 指向的已分配内存块地址相邻的空闲内存块。

函数 `posix_memalign()` 与 `memalign()` 存在以下两方面的不同。

- 已分配的内存地址通过参数 `memptr` 返回。
- 内存与 `alignment` 参数的整数倍对齐¹，`alignment` 必须是 `sizeof(void*)`（在大多数硬件架构上是 4 或 8 个字节）与 2 的整数次幂两者间的乘积。

还要注意该函数与众不同的返回值，出错时不是返回 -1，而是直接返回一个错误号（即通常在 `errno` 中返回的正整数）。

如果 `SizeOf(void*)` 为 4，就可以使用 `posix_memalign()` 分配 65536 字节的内存，并与 4096 字节的边界对齐，代码如下：

```
int s;
void *memptr;

s = posix_memalign(&memptr, 1024 * sizeof(void *), 65536);
if (s != 0)
    /* Handle error */
```

由 `memalign()` 或 `posix_memalign()` 分配的内存块应该调用 `free()` 来释放。

在一些 UNIX 实现中，无法通过调用 `free()` 来释放由 `memalign()` 分配的内存，因为此类 `memalign()` 在实现时使用 `malloc()` 来分配内存块，然后返回一个指针，指向该块内已对齐的适当地址。`glibc` 的 `memalign()` 则不受此限制。

7.2 在堆栈上分配内存： `alloca()`

和 `malloc` 函数包中的函数功能一样，`alloca()` 也可以动态分配内存，不过不是从堆上分配内存，而是通过增加栈帧的大小从堆栈上分配。根据定义，当前调用函数的栈帧位于堆栈的顶部，故而这种方法是可行的。因此，帧的上方存在扩展空间，只需修改堆栈指针值即可。

```
#include <alloca.h>

void *alloca(size_t size);
```

Returns pointer to allocated block of memory

参数 `size` 指定在堆栈上分配的字节数。函数 `alloca()` 将指向已分配内存块的指针作为其返回值。

不需要（实际上也绝不能）调用 `free()` 来释放由 `alloca()` 分配的内存。同样，也不可能调用 `realloc()` 来调整由 `alloca()` 分配的内存大小。

虽然 `alloca()` 不是 SUSv3 的一部分，但大多数 UNIX 实现都提供了此函数，因而也具备可移植性。

旧版本的 `glibc` 和其他一些 UNIX 实现（主要是 BSD 的衍生版本），要获取 `alloca()` 声明需引入 `<stdlib.h>` 而非 `<alloca.h>`。

若调用 `alloca()` 造成堆栈溢出，则程序的行为无法预知，特别是在没有收到一个 NULL 返回值通知错误的情况下。（事实上，在此情况下，可能会收到一个 SIGSEGV 信号。详情参见 21.3 节。）

请注意，不能在一个函数的参数列表中调用 `alloca()`，如下所示：

¹ 译者注：简而言之，该内存块的起始地址是 `alignment` 参数的整数倍。

```
func(x, alloca(size), z);          /* WRONG! */
```

这会使 `alloca()`分配的堆栈空间出现在当前函数参数的空间内（函数参数都位于栈帧内的固定位置）。相反，必须采用这样的代码：

```
void *y;

y = alloca(size);
func(x, y, z);
```

使用 `alloca()`来分配内存相对于 `malloc()`有一定优势。其中之一是，`alloca()`分配内存的速度要快于 `malloc()`，因为编译器将 `alloca()`作为内联代码处理，并通过直接调整堆栈指针来实现。此外，`alloca()`也不需要维护空闲内存块列表。

另一个优点在于，由 `alloca()`分配的内存随栈帧的移除而自动释放，亦即当调用 `alloca` 的函数返回之时。之所以如此，是因为函数返回时所执行的代码会重置栈指针寄存器，使其指向前一帧的末尾（即，假设堆栈向下增长，则指向恰好位于当前栈帧起始处之上的地址）。由于在函数的所有返回路径中都无需确保去释放所有的已分配内存，一些函数的编码也变得简单得多。

在信号处理程序中调用 `longjmp()`（6.8 节）或 `siglongjmp()`（21.2.1 节）以执行非局部跳转时，`alloca()`的作用尤其突出。此时，在“起跳”函数和“落地”函数之间的函数中，如果使用了 `malloc()`来分配内存，要想避免内存泄漏就极其困难，甚至是不可能的。与之相反，`alloca()`完全可以避免这一问题，因为堆栈是由这些调用展开的，所以已分配的内存会被自动释放¹。

7.3 总结

利用 `malloc` 函数族，进程可以动态分配和释放堆内存。在讨论这些函数的实现时，描述了程序对已分配内存处理失当的种种情况，还点出了一些有助于定位此类错误根源的调试工具。

函数 `alloca()`能够在堆栈上分配内存。该类内存会在调用 `alloca()`的函数返回时自动释放。

7.4 练习

- 7-1. 修改程序清单 7-1 中的程序（`free_and_sbrk.c`），在每次执行 `malloc()`后打印出 `program break` 的当前值。指定一个较小的内存分配尺寸来运行该程序。这将证明 `malloc()`不会在每次被调用时都调用 `sbrk()`来调整 `program break` 的位置，而是周期性地分配大块内存，并从中将小片内存返回给调用者。
- 7-2. （高级）实现 `malloc()`和 `free()`。

¹ 译者注：通过调整栈指针自然释放了栈中所分配的内存。

第 8 章

用户和组

每个用户都拥有一个唯一的用户名和一个与之相关的数值型用户标识符（UID）。用户可以隶属于一个或多个组。而每个组也都拥有唯一的一个名称和一个组标识符（GID）。

用户和组 ID 的主要用途有二：其一，确定各种系统资源的所有权；其二，对赋予进程访问上述资源的权限加以控制。比方说，每个文件都属于某个特定的用户和组，而每个进程也拥有相应的用户 ID 和组 ID 属性，这就决定了进程的所有者，以及进程访问文件时所拥有的权限（具体信息请参见第 9 章）。

本章首先会关注用于定义用户和组的系统文件，随后将描述用来从这些系统文件中获取信息的库函数。最后，将讨论用来加密和认证登录密码的 `crypt()` 函数。

8.1 密码文件：/etc/passwd

针对系统的每个用户账号，系统密码文件 `/etc/passwd` 会专列一行进行描述。每行都包含 7 个字段，之间用冒号分隔，如下所示：

```
mtk:x:1000:100:Michael Kerrisk:/home/mtk:/bin/bash
```

接下来，将按顺序介绍这 7 个字段。

- **登录名**：登录系统时，用户所必须输入的唯一名称。通常，也将其称为用户名。此外，也可将登录名视为人类可读的（符号）标识符，与数字用户标识符（稍后介绍）相对应。当使用诸如 `ls(1)` 这样的程序去显示文件的所有权时（比如，执行 `ls -l` 时），会显示出登录名，而非与文件关联的数值型用户 ID。
- **经过加密的密码**：该字段包含的是经过加密处理的密码，长度为 13 个字符，8.5 节会对此做深入讨论。如果密码字段中包含了任何其他字符串，特别是，当字符串长度超过 13 个字符时，将禁止此账户登录，原因是此类字符串不能代表一个经过加密的有效密码。不过，请注意，要是启用了 `shadow` 密码（这是常规做法），系统将会不解析该字段。这时，`/etc/passwd` 中的密码字段通常会包含字母“x”（当然，也可以是任何非空字符串），而经过加密处理的密码实际上却存储到 `shadow` 密码文件中（参见 8.2 节）。

若/etc/passwd 中密码字段为空，则该账户登录时无需密码（即便启用了 shadow 密码，也是如此）。

本章假定对密码的加密算法为 DES（数据加密标准），这也是一直为 UNIX 所广泛使用的密码加密算法。还可用其他加密算法（比如，MD5）来替代 DES，针对输入生成 128 位的消息摘要（hash 的一种）。在密码（或 shadow 密码）文件中，该消息摘要会以长度为 34 字符的字符串形式存储。

- 用户 ID (UID)：用户的数值型 ID。如果该字段的值为 0，那么相应账户即具有特权级权限。这种账号一般只有一个，其登录名为 root。在 Linux 2.2 或更早的版本中，用户 ID 为 16 位值，其范围为 0~65535。而 Linux 2.4 及其以后的版本则以 32 位值来存储用户 ID，因此能够支持更多的用户数。

在密码文件中，允许（但不常见）同一用户 ID 拥有多条记录，从而使得同一用户 ID 拥有多个登录名。如此一来，多个用户便能以不同密码（登录）去访问相同资源（比如，文件等）。此外，不同的登录名还可以关联一系列不同的组 ID。

- 组 ID (GID)：用户属组中首选属组的数值型 ID。关于用户与属组之间从属关系的进一步信息，会在系统组文件中加以定义。
- 注释：该字段存放关于用户的描述性文字。诸如 finger(1)之类的各种程序会显示此信息。
- 主目录：用户登录后所处的初始路径。会以该字段内容来设置 HOME 环境变量。
- 登录 shell：一旦用户登录，便交由该程序控制。通常，该程序为 shell 的一种（比如，bash），但也可以是其他任何程序。如果该字段为空，那么登录 shell 默认为/bin/sh（Bourne shell）。会以该字段值来设置 SHELL 环境变量。

在单机系统中，所有密码信息都存储在/etc/passwd 文件中¹。然而，如果使用了 NIS（网络信息系统）或 LDAP（轻型目录访问协议）在网络环境中分发密码，那么部分密码信息可能会由远端系统保存。只要访问密码信息的程序采用的是本章稍后描述的函数（getpwnam()、getpwuid()等），那么无论是使用 NIS 还是 LDAP，对应用程序来说都是透明的。类似论断同样适用于本章随后几节所讨论的 shadow 密码文件和组文件。

8.2 shadow 密码文件：/etc/shadow

很久以来，UNIX 一直在/etc/passwd 中维护所有的用户信息，这其中包括经过加密处理的密码。但这一举措也带来了安全问题。由于许多非特权级别系统工具需要读取密码文件中的其他信息，密码文件因而不得不对所有用户开放可读权限。这就为密码破解工具提供了可乘之机，它们会尝试对可能成为密码的大量词汇（比如，字典中的标准单词或人名）进行加密，然后再将结果与经过加密处理的密码进行比对。作为防范此类攻击的手段之一，shadow 密码文件/etc/shadow 应运而生。其理念是用户的所有非敏感信息存放于“人人可读”的密码文件中，而经过加密处理的密码则由 shadow 密码文件单独维护，仅供具有特权的程序读取。

shadow 密码文件包含有登录名（用来匹配密码文件中的相应记录）、经过加密的密码，以

¹ 译者注：此处有误，未考虑启用 shadow 密码的情况。

及其他若干与安全性相关的字段。`shadow(5)`手册页对这些字段作了详细描述。本章将着重关注经过加密的密码字段，将在 8.5 节介绍 `crypt()`库函数时做深入讨论。

SUSv3 并未对 `shadow` 密码作出规范，也并非所有的 UNIX 实现都提供这一特性，即使是都支持这一特性的各种实现，在关于 API 和文件位置上的细节也不尽相同。

8.3 组文件：/etc/group

出于各种管理方面的考虑，尤其是要控制对文件和其他系统资源的访问，对用户进行编组极具实用价值。

对用户所属各组信息的定义由两部分组成：一，密码文件中相应用户记录的组 ID 字段；二，组文件列出的用户所属各组。这种将信息分置于两个文件中的奇怪现状，自有其历史渊源。在早期 UNIX 实现中，一个用户同时只能从属于一个组。登录时，用户最初的属组关系由密码文件的组 ID 字段决定，在此之后，可使用 `newgrp(1)`命令去改变用户属组，但需要用户提供组密码（若该组处于密码的保护之下）。4.2BSD 引入了并发多属组（`multiple simultaneous group memberships`）的概念，POSIX.1-1990 随后对其进行了标准化。采用这种方案，组文件会列出每个用户所属的其他属组。（`groups(1)`命令会显示当前 shell 进程所属各组的信息，如果将一个或多个用户名作为其命令行参数，那么该命令将显示相应用户所属各组的信息。）

系统中的每个组在组文件 `/etc/group` 中都对应着一条记录。每条记录包含 4 个字段，之间以冒号分隔，如下所示：

```
users:x:100:  
jambit:x:106:claus,felli,frank,harti,markus,martin,mtk,paul
```

本节将依次介绍这 4 个字段。

- 组名：组的名称。与密码文件中的登录名相似，可以将其视为与数值型组标识符相对应的人类可读（符号）标识符。
- 经过加密处理的密码：组密码属于非强制特性，对应于该字段。随着多属组的出现，当今的 UNIX 系统已经很少使用组密码。不过，依然可以为组设置密码（特权用户可使用 `gpasswd` 命令来设置组密码）。如果用户并非某组的成员，那么在使用 `newgrp(1)` 启动新 shell 之前（新 shell 的属组包括该组），就需要用户提供此密码。如果启用了 `shadow` 密码，那么系统将不解析该字段（这时，该字段通常只包含字母 x，但也允许其内容为包括空字符串在内的任何字符串），而经过加密的密码实际上则存放于 `shadow` 组文件 `/etc/gshadow` 中，仅供具有特权的用户和程序访问。组密码的加密方式类似于用户密码（8.5 节）。
- 组 ID (GID)：该组的数值型 ID。正常情况下，对应于组 ID 号 0，只定义一个名为 `root` 的组（与 `/etc/passwd` 中用户 ID 为 0 的记录相近）。在 Linux 2.2 或更早的版本中，组 ID 为 16 位值，其范围为 0~65535；而自 Linux 2.4 以后的版本则以 32 位值来存储组 ID。
- 用户列表：属于该组的用户名列表，之间以逗号分隔。（这份列表包含的是用户名，而非用户 ID，原因在于如前所述，在密码文件的各条记录中，用户 ID 并不一定唯一。）为了证明用户 `avr` 是 `users`、`staff` 以及 `teach` 各组的成员，应能从密码文件中查看到如下记录：

```
avr:x:1001:100:Anthony Robins:/home/avr:/bin/bash
```

且在组文件中应有如下记录:

```
users:x:100:  
staff:x:101:mtk,avr,martinl  
teach:x:104:avr,rlb,alc
```

在密码文件记录的第 4 个字段中, 组 ID 为 100, 这说明 avr 是 users 组的成员之一。其他属组关系, 则见诸于组文件内包含 avr 的各条相关记录。

8.4 获取用户和组的信息

本节所要介绍的库函数, 其功能包括从密码文件、shadow 密码文件和组文件中获取单条记录, 以及扫描上述各个文件的所有记录。

从密码文件获取记录

函数 `getpwnam()` 和 `getpwuid()` 的作用是从密码文件中获取记录。

```
#include <pwd.h>  
  
struct passwd *getpwnam(const char *name);  
struct passwd *getpwuid(uid_t uid);  
  
Both return a pointer on success, or NULL on error;  
see main text for description of the “not found” case
```

为 `name` 提供一个登录名, `getpwnam()` 函数就会返回一个指针, 指向如下类型的结构, 其中包含了与密码记录相对应的信息:

```
struct passwd {  
    char *pw_name;        /* Login name (username) */  
    char *pw_passwd;     /* Encrypted password */  
    uid_t pw_uid;       /* User ID */  
    gid_t pw_gid;       /* Group ID */  
    char *pw_gecos;     /* Comment (user information) */  
    char *pw_dir;       /* Initial working (home) directory */  
    char *pw_shell;     /* Login shell */  
};
```

`passwd` 结构的 `pw_gecos` 和 `pw_passwd` 字段虽未在 SUSv3 中定义, 但获得了所有 UNIX 实现的支持。仅当未启用 shadow 密码的情况下, `pw_passwd` 字段才会包含有效信息。要确定是否启用了 shadow 密码, 最简单的编程方法是在成功调用 `getpwnam()` 之后, 紧接着调用 `getspnam()` (稍后介绍), 并观察后者是否能为同一用户名返回一条 shadow 密码记录。某些其他实现还会在该结构中定义额外的非标准字段。

`pw_gecos` 字段, 其命名源于早期的 UNIX 实现, 该字段所含信息原用于与运行 GECOS (通用电器综合操作系统) 的计算机进行通信。虽然这一用途早已过时, 但其名称却得以沿用至今, 只是将字段用途转而用于记录用户的相关信息。

函数 `getpwuid()` 的返回结果与 `getpwnam()` 完全一致, 但会使用提供给 `uid` 参数的数值型用户 ID 作为查询条件。

`getpwnam()` 和 `getpwuid()` 均会返回一个指针, 指向一个静态分配的结构。对此二者 (或是

下文描述的 `getpwent()` 函数) 的任何一次调用都会改写该数据结构。

由于 `getpwnam()` 和 `getpwuid()` 返回的指针指向由静态分配而成的内存, 故而二者都是不可重入的 (not reentrant)。实际上, 情况甚至要更加复杂, 因为返回的 `passwd` 结构还包含了指向其他信息 (比如, `pw_name`) 的指针, 而这些信息同样也是由静态分配而成的。21.1.2 节会解释可重入 (reentrancy) 概念。类似的论断同样适用于 `getgrnam()` 和 `getgrgid()` 函数 (稍后介绍)。

SUSv3 为该组函数定义了与之等价的一组可重入函数: `getpwnam_r()`、`getpwuid_r()`、`getgrnam_r()` 以及 `getgrgid_r()`。其参数包括 `passwd` (或 `group`) 结构, 以及一个缓冲区。这一缓冲区专门用来保存 `passwd(group)` 结构中各字段所指向的其他结构。可使用系统函数 `sysconf(_SC_GETPW_R_SIZE_MAX)` (若为与组相关的函数, 则使用 `sysconf(_SC_GETGR_R_SIZE_MAX)`), 来获得此缓冲区所需的字节数。以上函数的详细信息请查阅手册页。

SUSv3 规定, 如果在 `passwd` 文件中未发现匹配记录, 那么 `getpwnam()` 和 `getpwuid()` 将返回 `NULL`, 且不会改变 `errno`。这意味着, 可以使用如下代码, 对出错和 “未发现匹配记录” 这两种情况加以区分:

```
struct passwd *pwd;

errno = 0;
pwd = getpwnam(name);
if (pwd == NULL) {
    if (errno == 0)
        /* Not found */;
    else
        /* Error */;
}
```

然而, 不少 UNIX 实现在这点上并未遵守 SUSv3 规范。如果未能在 `passwd` 文件中发现一条匹配记录, 那么两个函数均会返回 `NULL`, 并将 `errno` 设置为非零值, 比如, `ENOENT` 或 `ESRCH`。针对这种情况, 2.7 版本之前的 `glibc` 会产生 `ENOENT` 错误, 而从 2.7 版本开始, `glibc` 开始遵守 SUSv3 规范。实现之间之所以存在上述差异, 部分原因是由于 `POSIX.1-1990` 不但不要求两个函数在出错时设置 `errno`, 而且还允许它们针对 “未发现匹配记录” 的情况去设置 `errno`。总而言之, 在使用这两个函数时, 若要区分上述这两种情况 (出错和 “未发现匹配记录”), 实际上将无法保证代码的可移植性。

从组文件获取记录

函数 `getgrnam()` 和 `getgrgid()` 的作用是从组文件中获取记录。

```
#include <grp.h>

struct group *getgrnam(const char *name);
struct group *getgrgid(gid_t gid);

Both return a pointer on success, or NULL on error;
see main text for description of the “not found” case
```

函数 `getgrnam()` 和 `getgrgid()` 分别通过组名和组 ID 来查找属组信息。两个函数都会返回一

个指针，指向如下类型结构：

```
struct group {
    char *gr_name;      /* Group name */
    char *gr_passwd;    /* Encrypted password (if not password shadowing) */
    gid_t gr_gid;      /* Group ID */
    char **gr_mem;     /* NULL-terminated array of pointers to names
                       of members listed in /etc/group */
};
```

SUSv3 并未就 `group` 结构中的 `gr_passwd` 字段做明确定义，但大多数 UNIX 实现都支持该字段。

与前述密码相关函数一样，对这两个函数的任何一次调用都会改写该结构的内容¹。

如果未能在 `group` 文件中发现匹配记录，那么这两个函数的行为变化与前述 `getpwnam()` 和 `getpwuid()` 函数相同²。

程序示例

对本节所述的函数来说，最常见的用法之一是在符号型用户名和组名与数值型 ID 之间进行相互转换。程序清单 8-1 以 `userNameFromId()`、`userIdFromName()`、`groupNameFromId()` 以及 `groupIdFromName()` 这 4 个函数的形式，演示了上述转换。为方便调用，`userIdFromName()` 和 `groupIdFromName()` 还允许 `name` 参数接受（纯）数值的字符串形式³。对于这种情况，会直接将字符串转换为数字返回给调用者。在本书后面的一些程序实例中，还会用到这几个函数。

程序清单 8-1：在用户名/组名和用户 ID/组 ID 之间互相转换的函数

```
----- users_groups/ugid_functions.c
#include <pwd.h>
#include <grp.h>
#include <ctype.h>
#include "ugid_functions.h" /* Declares functions defined here */

char *      /* Return name corresponding to 'uid', or NULL on error */
userNameFromId(uid_t uid)
{
    struct passwd *pwd;

    pwd = getpwuid(uid);
    return (pwd == NULL) ? NULL : pwd->pw_name;
}

uid_t      /* Return UID corresponding to 'name', or -1 on error */
userIdFromName(const char *name)
{
    struct passwd *pwd;
    uid_t u;
    char *endptr;
    if (name == NULL || *name == '\0') /* On NULL or empty string */
        return -1; /* return an error */
}
```

1 译者注：该结构也是由静态分配而成。

2 译者注：换言之，区分出错和“未发现匹配记录”情况的编程手法也与之类似。

3 译者注：例如“123”。

```

    u = strtol(name, &endptr, 10);    /* As a convenience to caller */
    if (*endptr == '\0')             /* allow a numeric string */
        return u;

    pwd = getpwnam(name);
    if (pwd == NULL)
        return -1;

    return pwd->pw_uid;
}

char *          /* Return name corresponding to 'gid', or NULL on error */
groupNameFromId(gid_t gid)
{
    struct group *grp;

    grp = getgrgid(gid);
    return (grp == NULL) ? NULL : grp->gr_name;
}

gid_t          /* Return GID corresponding to 'name', or -1 on error */
groupIdFromName(const char *name)
{
    struct group *grp;
    gid_t g;
    char *endptr;

    if (name == NULL || *name == '\0') /* On NULL or empty string */
        return -1;                    /* return an error */

    g = strtol(name, &endptr, 10);    /* As a convenience to caller */
    if (*endptr == '\0')             /* allow a numeric string */
        return g;

    grp = getgrnam(name);
    if (grp == NULL)
        return -1;

    return grp->gr_gid;
}

```

users_groups/ugid_functions.c

扫描密码文件和组文件中的所有记录

函数 `setpwent()`、`getpwent()`和 `endpwent()`的作用是按顺序扫描密码文件中的记录。

```

#include <pwd.h>

struct passwd *getpwent(void);

                Returns pointer on success, or NULL on end of stream or error

void setpwent(void);
void endpwent(void);

```

函数 `getpwent()`能够从密码文件中逐条返回记录，当不再有记录¹（或出错）时，该函数

¹ 译者注：即抵达流末端时。

返回 NULL。getpwent()一经调用，会自动打开密码文件。当密码文件处理完毕后，可调用 endpwent()将其关闭。

可使用以下代码遍历整个密码文件，并打印出登录名和用户 ID。

```
struct passwd *pwd;

while ((pwd = getpwent()) != NULL)
    printf("%-8s %5ld\n", pwd->pw_name, (long) pwd->pw_uid);

endpwent();
```

如果需要让后续的 getpwent()调用（也许是在程序的其他代码中，也许是在所调用的其他库函数中，该函数再次出现）再次打开密码文件并重启扫描过程，此处的 endpwent()调用就必不可少。此外，如果对该文件处理到中途时，还可以调用 setpwent()函数重返文件起始处。

函数 getgrent()、setgrent()和 endgrent()针对组文件执行类似的任务。由于这 3 个函数与前述的密码文件函数功能相似，故而其函数原型也就不再列出，详细信息请参考手册页。

从 shadow 密码文件中获取记录

下列函数的作用包括从 shadow 密码文件中获取个别记录，以及扫描该文件中的所有记录。

```
#include <shadow.h>

struct spwd *getspnam(const char *name);
                Returns pointer on success, or NULL on not found or error
struct spwd *getspent(void);
                Returns pointer on success, or NULL on end of stream or error
void setspent(void);
void endspent(void);
```

由于上述函数在操作上类似于相应的密码文件函数，故而此处对它们的介绍也就点到为止。（上述函数既未在 SUSv3 中明确定义，也未获得所有 UNIX 实现的支持。）

函数 getspnam()和 getspent()会返回指向 spwd 类型结构的指针。该结构的形式如下：

```
struct spwd {
    char *sp_namp;          /* Login name (username) */
    char *sp_pwdp;         /* Encrypted password */

    /* Remaining fields support "password aging", an optional
       feature that forces users to regularly change their
       passwords, so that even if an attacker manages to obtain
       a password, it will eventually cease to be usable. */

    long sp_lstchg;        /* Time of last password change
                           (days since 1 Jan 1970) */
    long sp_min;          /* Min. number of days between password changes */
    long sp_max;          /* Max. number of days before change required */
    long sp_warn;         /* Number of days beforehand that user is
                           warned of upcoming password expiration */
    long sp_inact;        /* Number of days after expiration that account
                           is considered inactive and locked */
```

```

long sp_expire;          /* Date when account expires
                        (days since 1 Jan 1970) */
unsigned long sp_flag; /* Reserved for future use */
};

```

在程序清单 8-2 中，将会演示对 `getspnam()` 的使用。

8.5 密码加密和用户认证

某些应用程序会要求用户对自身进行认证，通常会采取用户名（登录名）/密码的认证方式。出于这一目的，应用程序可能会维护其自有的用户名和密码数据库。然而，或许是由于势所必然，或许是为了方便起见，有时需要让用户输入标准的用户名/密码（定义于 `/etc/passwd` 和 `/etc/shadow` 之中）。（本节的剩余部分将假定系统启用了 `shadow` 密码，经过加密处理的密码也因此存储于 `/etc/shadow` 中。）需要登录到远程系统的网络应用程序，诸如 `ssh` 和 `ftp`，就是此类程序的典范，必须按标准的 `login` 程序那样，对用户名和密码加以验证。

由于安全方面的原因，UNIX 系统采用单向加密算法对密码进行加密，这意味着由密码的加密形式将无法还原出原始密码。因此，验证候选密码的唯一方法是使用同一算法对其进行加密，并将加密结果与存储于 `/etc/shadow` 中的密码进行匹配。加密算法封装于 `crypt()` 函数之中。

```

#define _XOPEN_SOURCE
#include <unistd.h>

char *crypt(const char *key, const char *salt);

Returns pointer to statically allocated string containing
encrypted password on success, or NULL on error

```

`crypt()` 算法会接受一个最长可达 8 字符的密钥（即密码），并施之以数据加密算法（DES）的一种变体。`salt` 参数指向一个两字符的字符串，用来扰动（改变）DES 算法，设计该技术，意在使得经过加密的密码更加难以破解。该函数会返回一个指针，指向长度为 13 个字符的字符串，该字符串为静态分配而成，内容即为经过加密处理的密码。

DES 的详细信息请参考 <http://www.itl.nist.gov/fipspubs/fip46-2.htm>。如前所述，除 DES 以外，也可以使用其他的加密算法。例如，使用 MD5 算法可以生成一个 34 字符的字符串，其首字符为美元符号（\$），这便于让 `crypt()` 将 DES 加密密码和 MD5 加密密码区分开来。

在关于密码加密的讨论中，本书对“加密”一词的使用相对宽松。确切说来，DES 会以给定的密码字符串作为加密密钥，编码得出固定位长的字符串，而 MD5 则是一种复杂的哈希函数。以上两种方法其实殊途同归，对输入密码的加密变换既不可逆又难以破解。

`salt` 参数和经过加密的密码，其组成成员均取自同一字符集合，范围在 `[a-zA-Z0-9/.]` 之间，共计 64 个字符。因此，两个字符的 `salt` 参数可使加密算法产生 4096（ 64×64 ）种不同变化。这意味着，预先对整部字典进行加密，再以其中的每个单词与经过加密处理的密码进行比对的做法并不可行，破解程序需要对照字典的 4096 种加密版本来检查密码。

由 `crypt()` 所返回的经过加密的密码中，头两个字符是对原始 `salt` 值的拷贝。也就是说，加密候选密码时，能够从已加密密码（存储于 `/etc/shadow` 内）中获取 `salt` 值。（加密新密码时，`passwd(1)` 这样的程序会生成一个随机 `salt` 值。）事实上，在 `salt` 字符串中，只有前两个字符对

crypt()函数有意义。因此，可以直接将已加密密码指定为 salt 参数。

要想在 Linux 中使用 crypt()，在编译程序时需开启 -lcrypt 选项，以便程序链接 crypt 库。

程序示例

程序清单 8-2 演示了如何使用 crypt()来验证用户。该程序首先读取用户名，然后会获取相应的密码记录以及（如开启了 shadow 密码功能）shadow 密码记录。若未能发现密码记录，或程序没有权限读取 shadow 密码文件（需要超级用户权限，或具有 shadow 组成员资格），该程序会打印一条错误消息并退出。接下来，该程序会使用 getpass()函数，读取用户密码。

```
#define _BSD_SOURCE
#include <unistd.h>

char *getpass(const char *prompt);

Returns pointer to statically allocated input password string
on success, or NULL on error
```

getpass()函数首先会屏蔽回显功能，并停止对终端特殊字符的处理（诸如中断字符，一般为 Control-C）。（第 62 章将论述如何更改这些终端设置。）然后，该函数会打印出 prompt 所指向的字符串，读取一行输入，返回以 NULL 结尾的输入字符串（剥离尾部的换行符）作为函数结果。（该字符串由静态分配而成，故而后续对 getpass()的调用会覆盖其原有内容。）返回结果之前，getpass()会将终端设置还原。

使用 getpass()读取密码之后，程序清单 8-2 所示程序会对密码进行验证——使用 crypt()加密密码，并将结果与 shadow 密码文件中经过加密的密码记录进行比对。若两者匹配，则显示用户 ID，如下所示：

```
$ su                                     Need privilege to read shadow password file
Password:
# ./check_password
Username: mtk
Password:                                We type in password, which is not echoed
Successfully authenticated: UID=1000
```

程序清单 8-2 中，以调用 sysconf(_SC_LOGIN_NAME_MAX)的返回值作为存放用户名字符串数组的大小，该调用获取了主机系统上用户名字符串的最大长度。11.2 节将介绍 sysconf()的使用。

程序清单 8-2：根据 shadow 密码文件验证用户

```
----- users_groups/check_password.c
#define _BSD_SOURCE    /* Get getpass() declaration from <unistd.h> */
#define _XOPEN_SOURCE /* Get crypt() declaration from <unistd.h> */
#include <unistd.h>
#include <limits.h>
#include <pwd.h>
#include <shadow.h>
#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
```

```

char *username, *password, *encrypted, *p;
struct passwd *pwd;
struct spwd *spwd;
Boolean authOk;
size_t len;
long lnmax;
lnmax = sysconf(_SC_LOGIN_NAME_MAX);
if (lnmax == -1) /* If limit is indeterminate */
    lnmax = 256; /* make a guess */

username = malloc(lnmax);
if (username == NULL)
    errExit("malloc");

printf("Username: ");
fflush(stdout);
if (fgets(username, lnmax, stdin) == NULL)
    exit(EXIT_FAILURE); /* Exit on EOF */

len = strlen(username);
if (username[len - 1] == '\n')
    username[len - 1] = '\0'; /* Remove trailing '\n' */

pwd = getpwnam(username);
if (pwd == NULL)
    fatal("couldn't get password record");
spwd = getsppnam(username);
if (spwd == NULL && errno == EACCES)
    fatal("no permission to read shadow password file");

if (spwd != NULL) /* If there is a shadow password record */
    pwd->pw_passwd = spwd->sp_pwdp; /* Use the shadow password */

password = getpass("Password: ");

/* Encrypt password and erase cleartext version immediately */

encrypted = crypt(password, pwd->pw_passwd);
for (p = password; *p != '\0'; )
    *p++ = '\0';

if (encrypted == NULL)
    errExit("crypt");

authOk = strcmp(encrypted, pwd->pw_passwd) == 0;
if (!authOk) {
    printf("Incorrect password\n");
    exit(EXIT_FAILURE);
}

printf("Successfully authenticated: UID=%ld\n", (long) pwd->pw_uid);

/* Now do authenticated work... */

exit(EXIT_SUCCESS);
}

```

users_groups/check_password.c

程序清单 8-2 展示了一个安全要点。读取密码的程序应立即加密密码，并尽快将密码的明文从内存中抹去。只有这样，才能基本杜绝如下事件的发生：恶意之徒借程序崩溃之机，读取内核转储文件以获取密码。

仍有可能采用其他方法曝光未经加密的密码。例如，如果包含密码的虚拟内存页执行了换出操作，那么特权级程序就能交换文件中读取密码。此外，拥有足够权限的进程可通过读取/dev/mem（虚拟设备之一，将计算机物理内存表示为有序字节流），来尝试发现密码。

SUSv2 将 `getpass()` 函数标记为 `LEGACY`，并特别指出该函数名容易产生误导，且其所提供的功能无论在何种情况下都极易于实现。SUSv3 摒弃了 `getpass()`，但在大多数 UNIX 实现中依然保留了对它的支持。

8.6 总结

每个用户都有一个唯一的用户名和一个与之对应的数值型用户 ID。用户可以隶属于一个或多个组，每个组都有一个唯一的名称和一个与之对应的数字标识符。这些标识符的主要用途在于确立各种系统资源（比如，文件）的所有权和访问这些资源的权限。

用户名和 ID 在 `/etc/passwd` 文件中加以定义，该文件也包含有关用户的其他信息。用户的属组则由 `/etc/passwd` 和 `/etc/group` 文件中的相关字段来定义。还有一个只能由特权级进程所读取的文件 `/etc/shadow`，其作用在于将敏感的密码信息与 `/etc/passwd` 中共用的用户信息分离开来。系统还提供有不同的库函数，用于从上述各个文件中获取信息。

`crypt()` 函数加密密码的方式与标准的 `login` 程序相同，这对需要认证用户的程序来说极为有用。

8.7 练习

- 8-1** 执行下列代码时，将会发现，尽管这两个用户在密码文件中对应不同的 ID，但该程序的输出还是会将同一个数字显示两次。请问为什么？

```
printf("%ld %ld\n", (long) (getpwnam("avr")->pw_uid),  
                (long) (getpwnam("tsr")->pw_uid));
```

- 8-2** 使用 `setpwent()`、`getpwent()` 和 `endpwent()` 来实现 `getpwnam()`。

第 9 章

进程凭证

每个进程都有一套用数字表示的用户 ID¹ (UID) 和组 ID(GID)。有时, 也将这些 ID 称之为进程凭证。具体如下所示。

- 实际用户 ID (real user ID) 和实际组 ID (real group ID)。
- 有效用户 ID (effective user ID) 和有效组 ID (effective group ID)。
- 保存的 set-user-ID (saved set-user-ID) 和保存的 set-group-ID (saved set-group-ID)。
- 文件系统用户 ID (file-system user ID) 和文件系统组 ID (file-system group ID) (Linux 专有)。
- 辅助组 ID。

本章将详细介绍这些进程 ID 的用途, 以及用于获取和修改此类 ID 的系统调用和库函数, 还将讨论特权级进程和非特权进程的概念, 并阐述了设置用户 ID 和设置组 ID 的使用机制, 采用该机制所创建的程序可以以特定用户或组的权限运行。

9.1 实际用户 ID 和实际组 ID

实际用户 ID 和实际组 ID 确定了进程所属的用户和组。作为登录过程的步骤之一, 登录 shell 从/etc/passwd 文件中读取相应用户密码记录的第三字段和第四字段, 置为其实际用户 ID 和实际组 ID (8.1 节)。当创建新进程 (比如, shell 执行一程序) 时, 将从其父进程中继承这些 ID。

9.2 有效用户 ID 和有效组 ID

在大多数 UNIX 实现 (Linux 实现略有差异, 具体参见 9.5 节的说明) 中, 当进程尝试执行各种操作 (即系统调用) 时, 将结合有效用户 ID、有效组 ID, 连同辅助组 ID 一起来确定

¹ 译者注: 即标识符。

授予进程的权限。例如，当进程访问诸如文件、System V 进程间通信 (IPC) 对象之类的系统资源时，此类 ID 会决定系统授予进程的权限，而这些资源的属主则另由与之相关的用户 ID 和组 ID 来决定。如 20.5 节所述，内核还会使用有效用户 ID 来决定一个进程是否能向另一个进程发送信号。

有效用户 ID 为 0 (root 的用户 ID) 的进程拥有超级用户的所有权限。这样的进程又称为特权级进程 (privileged process)。而某些系统调用只能由特权级进程执行。

第 39 章描述了 Linux 实现的能力 (capability) 方案，即把授予给超级用户的特权划分为若干不同单元，且能独立启用和禁用这些单元。

通常，有效用户 ID 及组 ID 与其相应的实际 ID 相等，但有两种方法能够致使二者不同。其一是使用 9.7 节中所讨论的系统调用，其二是执行 set-user-ID 和 set-group-ID 程序。

9.3 Set-User-ID 和 Set-Group-ID 程序

set-user-ID 程序会将进程的有效用户 ID 置为可执行文件的用户 ID (属主)，从而获得常规情况下并不具有的权限。set-group-ID 程序对进程有效组 ID 实现类似任务。(术语 set-user-ID 程序和 set-group-ID 程序有时也简称为 set-UID 程序和 set-GID 程序。)

与其他文件一样，可执行文件的用户 ID 和组 ID 决定了该文件的所有权。另外，可执行文件还拥有两个特别的权限位 set-user-ID 位和 set-group-ID 位。(实际上，任何文件都是如此，但此处只关注可执行文件的这两个权限位。) 可使用 chmod 命令来设置这些权限位。非特权用户能够对其拥有的文件进行设置，而特权级用户 (CAP_FOWNER) 能够对任何文件进行设置。例如：

```
$ su
Password:
# ls -l prog
-rwxr-xr-x  1 root    root      302585 Jun 26 15:05 prog
# chmod u+s prog
Turn on set-user-ID permission bit
# chmod g+s prog
Turn on set-group-ID permission bit
```

正如本例所示，也有可能对这两个权限位都进行设置，虽然这一做法并不常见。当使用 ls -l 命令查看文件权限时，如果为程序设置了 set-user-ID 权限位和 set-group-ID 权限位，那么通常用来表示文件可执行权限的 x 标识会被 s 标识所替换。

```
# ls -l prog
-rwsr-sr-x  1 root    root      302585 Jun 26 15:05 prog
```

当运行 set-user-ID 程序 (即通过调用 exec() 将 set-user-ID 程序载入进程的内存中) 时，内核会将进程的有效用户 ID 设置为可执行文件的用户 ID。set-group-ID 程序对进程有效组 ID 的操作与之类似。通过这种方法修改进程的有效用户 ID 或者组 ID，能够使进程 (换言之，执行该程序的用户) 获得常规情况下所不具有的权限。例如，如果一个可执行文件的属主为 root (超级用户)，且为此程序设置了 set-user-ID 权限位，那么当运行该程序时，进程会取得超级用户权限。

也可以利用程序的 set-user-ID 和 set-group-ID 机制，将进程的有效 ID 修改为 root 之外的其他用户。例如，为提供对一个受保护文件 (或其他系统资源) 的访问，采用如下方案就绰

绰有余：创建一个具有对该文件访问权限的专用用户（组）ID，然后再创建一个 `set-user-ID`（`set-group-ID`）程序，将进程有效用户（组）ID 变更为这个专用 ID。这样，无需拥有超级用户的所有权限，程序就能访问该文件。

有时会使用术语 `set-user-ID-root` 来表示 `root` 用户所拥有的 `set-user-ID` 程序，以示与由其他用户所拥有的 `set-user-ID` 程序有所区别，后者仅为进程提供其属主所具有的权限。

术语 `privileged`（特权级）有两种不同含义，其一是为早期定义而成的，有效用户 ID 为 0 的进程，拥有 `root` 用户的所有特权。然而，当 `set-user-ID` 程序的属主并非 `root` 用户时，进程也会获得 `set-user-ID` 程序属主的特权。各种情况下术语 `privileged` 的具体含义，可通过上下文来加以辨别。

出于 38.3 节所给出的理由，在 Linux 系统中，`set-user-ID` 和 `set-group-ID` 权限位对 shell 脚本无效。

Linux 系统中经常使用的 `set-user-ID` 程序包括：`passwd(1)`，用于更改用户密码；`mount(8)` 和 `umount(8)`，用于加载和卸载文件系统；`su(1)`，允许用户以另一用户的身份运行 shell。`set-group-ID` 程序的例子之一为 `wall(1)`，用来向 `tty` 组下辖的所有终端（通常情况下，所有终端都属于该组）写入一条消息。

8.5 节曾特别指出，程序清单 8-2 中的程序需要以 `root` 用户身份运行，以便获取对 `/etc/shadow` 文件的访问权限。欲使该程序可为任一用户执行，必须将其设置为 `set-user-ID-root` 程序，如下所示：

```
$ su
Password:
# chown root check_password           Make this program owned by root
# chmod u+s check_password           With the set-user-ID bit enabled

# ls -l check_password
-rwsr-xr-x  1 root  users   18150 Oct 28 10:49 check_password
# exit
$ whoami                               This is an unprivileged login
mtk
$ ./check_password                     But we can now access the shadow
Username: avr                           password file using this program
Password:
Successfully authenticated: UID=1001
```

`set-user-ID/set-group-ID` 技术集实用性与强大的功能于一身，但一旦设计欠佳也可能造成安全隐患。第 38 章总结了一整套良好的编程习惯，编写 `set-user-ID` 和 `set-group-ID` 程序时应多加参考。

9.4 保存 `set-user-ID` 和保存 `set-group-ID`

设计保存 `set-user-ID`（`saved set-user-ID`）和保存 `set-group-ID`（`saved set-group-ID`），意在将 `set-user-ID` 和 `set-group-ID` 程序结合使用。当执行程序时，将会（依次）发生如下事件（在诸多事件之中）。

1. 若可执行文件的 `set-user-ID`（`set-group-ID`）权限位已开启，则将进程的有效用户（组）ID 置为可执行文件的属主。若未设置 `set-user-ID`（`set-group-ID`）权限位，则进程的有效用户（组）ID 将保持不变。

2. 保存 set-user-ID 和保存 set-group-ID 的值由对应的有效 ID 复制而来。无论正在执行的文件是否设置了 set-user-ID 或 set-group-ID 权限位，这一复制都将进行。

举例说明上述操作的效果，假设某进程的实际用户 ID、有效用户 ID 和保存 set-user-ID 均为 1000，当其执行了 root 用户（用户 ID 为 0）拥有的 set-user-ID 程序后，进程的用户 ID 将发生如下变化：

```
real=1000 effective=0 saved=0
```

有不少系统调用，允许将 set-user-ID 程序的¹有效用户 ID 在实际用户 ID 和保存 set-user-ID 之间切换。针对 set-group-ID 程序对其进程有效组 ID 的修改，也有与之相类似的系统调用来支持。如此一来，对于与执行文件用户（组）ID 相关的任何权限，程序能够随时“收放自如”。（换言之，程序可以游走于两种状态之间：具备获取特权的潜力和以特权进行实际操作。）正如 38.2 节所述，只要 set-user-ID 程序和 set-group-ID 程序没有执行与特权级 ID（亦即实际 ID）相关的任何操作，就应将其置于非特权（即实际）ID 的身份之下，这是一种安全的编程手法。

有时也将保存 set-user-ID 和保存 set-group-ID 称之为保存用户 ID（saved user ID）和保存组 ID（saved group ID）。

保存设置 ID 由 System V 首创，后为 POSIX 所采用。4.4 之前的 BSD 版本不提供对此特性的支持。最初的 POSIX.1 标准将对这些 ID 的支持列为可选，但之后的版本（始于 1988 年诞生的 FIPS 151-1 标准）则强制要求提供这一特性。

9.5 文件系统用户 ID 和组 ID

在 Linux 系统中，要进行诸如打开文件、改变文件属主、修改文件权限之类的文件系统操作时，决定其操作权限的是文件系统用户 ID 和组 ID（结合辅助组 ID），而非有效用户 ID 和组 ID。（和其他 UNIX 实现一样，有效用户 ID 和组 ID 仍在使用，其用途在前面章节已有论述。）

通常，文件系统用户 ID 和组 ID 的值等同于相应的有效用户 ID 和组 ID（因而一般也等同于相应的实际用户 ID 和组 ID）。此外，只要有效用户或组 ID 发生了变化，无论是通过系统调用，还是通过执行 set-user-ID 或者 set-group-ID 程序，则相应的文件系统 ID 也将随之改变为同一值。由于文件系统 ID 对有效 ID 如此的“亦步亦趋”，这意味着在特权和权限检查方面，Linux 实际上跟其他 UNIX 实现非常类似。只有当使用 Linux 特有的两个系统调用（setfsuid() 和 setfsgid()）时，才可以刻意制造出文件系统 ID 与相应有效 ID 的不同，因而 Linux 也不同于其他的 UNIX 实现。

那么，Linux 为什么要提供文件系统 ID 呢？在何种情况下，需要使有效 ID 有别于文件系统 ID 呢？这主要是由于历史原因造成的。文件系统 ID 始见于 Linux 1.2 版本。在该版本的内核中，如果进程某甲的有效用户 ID 等同于进程某乙的实际用户 ID 或者有效用户 ID，那么发送者（某甲）就可以向目标进程（某乙）发送信号。这在当时影响到了不少程序，比如 Linux NFS（网络文件系统）服务器程序，在访问文件时就好像拥有着相应客户进程的有效 ID。然而，如果 NFS 服务器真地修改了自身的有效用户 ID，面对非特权用户进程的信号攻击，又将不堪一击。为了防范这一风险，文件系统用户 ID 和组 ID 应运而生。NFS 服务器将有效 ID 保持不变，而是通过修改文件系统 ID 伪装成另一用户，这样既达到了访问文件的目的，又避免了遭受信

¹ 译者注：进程。

号攻击。

自内核 2.0 起, Linux 开始在信号发送权限方面遵循 SUSv3 所强制规定的规则, 且这些规则不再涉及目标进程的有效用户 ID (参考 20.5 节)。因此, 从严格意义上来讲, 保留文件系统 ID 特性已无必要 (如今, 进程可以根据需要, 审慎而明智地利用本章稍后介绍的系统调用, 使以非特权值对有效用户 ID 的赋值来去自由, 以实现预期结果), 但为了与现有软件保持兼容, 这一功能得以保留了下来。

由于文件系统 ID 实属异类, 且一般都等同于相应的有效 ID, 本书后续部分在述及各种文件权限的检查, 以及设置新文件的属主时, 通常将根据进程有效 ID 来加以解释。即使是出于 Linux 系统的目的而真地使用了进程的文件系统 ID, 但在实践中, 这些标识的存在与否并不会带来显著差别。

9.6 辅助组 ID

辅助组 ID 用于标识进程所属的若干附加的组。新进程从其父进程处继承这些 ID, 登录 shell 从系统组文件中获取其辅助的组 ID。如前所述, 将这些 ID 与有效 ID 以及文件系统 ID 相结合, 就能决定对文件、System V IPC 对象和其他系统资源的访问权限。

9.7 获取和修改进程凭证

为了获取和变更本章已然论及的各种用户 ID 和组 ID, Linux 提供了一系列系统调用和库函数。SUSv3 仅对这些 API 中的部分做了规范, 余下部分中, 有一些在其他 UNIX 实现中得以广泛应用, 还有少量是 Linux 所特有的。在讨论每个 API 接口时, 将特别指出可移植性方面的问题。在本章结尾处, 表 9-1 总结了变更进程凭证的所有接口操作。

可以利用 Linux 系统特有的 `proc/PID/status` 文件, 通过对其中 `Uid`、`Gid` 和 `Groups` 各行信息的检查, 来获取任何进程的凭证, 这与下面即将介绍的系统调用有异曲同工之妙。`Uid` 和 `Gid` 各行, 按实际、有效、保存设置和文件系统 ID 的顺序来展示相应标识符。

在下列章节中所论及的特权级进程, 其定义是基于传统意义上的, 即进程的有效用户 ID 为 0。然而, 正如第 39 章所述, Linux 将超级用户权限划分成多种各不相同的能力 (capability)。在讨论修改用户 ID 和组 ID 的所有系统调用时, 将涉及其中的两种。

- `CAP_SETUID` 能力允许进程任意修改其用户 ID。
- `CAP_SETGID` 能力允许进程任意修改其组 ID。

9.7.1 获取和修改实际、有效和保存设置标识

下面段落将描述用于获取和修改实际、有效和保存设置 ID 的系统调用。能完成这些任务的系统调用有多个, 有时彼此间的功能还相互重叠, 这是由于各种系统调用分别源于不同的 UNIX 实现。

获取实际和有效 ID

系统调用 `getuid()` 和 `getgid()` 分别返回调用进程的实际用户 ID 和组 ID。而系统调用 `geteuid()` 和 `getegid()` 则对进程的有效 ID 实现类似功能。对这些系统函数的调用总会成功。

<code>#include <unistd.h></code>	
<code>uid_t getuid(void);</code>	Returns real user ID of calling process
<code>uid_t geteuid(void);</code>	Returns effective user ID of calling process
<code>gid_t getgid(void);</code>	Returns real group ID of calling process
<code>gid_t getegid(void);</code>	Returns effective group ID of calling process

修改有效 ID

`setuid()` 系统调用以给定的 `uid` 参数值来修改调用进程的有效用户 ID，也可能修改实际用户 ID 和保存 `set-user-ID`。系统调用 `setgid()` 则对相应组 ID 实现了类似功能。

<code>#include <unistd.h></code>	
<code>int setuid(uid_t uid);</code>	
<code>int setgid(gid_t gid);</code>	
	Both return 0 on success, or -1 on error

进程使用 `setuid()` 和 `setgid()` 系统调用能对其凭证做哪些修改呢？其规则取决于进程是否拥有特权（即有效用户 ID 为 0）。适用于 `setuid()` 系统调用的规则如下。

1. 当非特权进程调用 `setuid()` 时，仅能修改进程的有效用户 ID。而且，仅能将有效用户 ID 修改成相应的实际用户 ID 或保存 `set-user-ID`。（企图违反此约束将引发 `EPERM` 错误。）这意味着，对于非特权用户而言，仅当执行 `set-user-ID` 程序时，`setuid()` 系统调用才起作用，因为在执行普通程序时，进程的实际用户 ID、有效用户 ID 和保存 `set-user-ID` 三者之值均相等。在一些派生自 BSD 的实现中，非特权进程对 `setuid()` 或 `setgid()` 的调用，其语义有别于与其他 UNIX 实现：系统调用会修改实际、有效和保存设置 ID（将其改为当前的实际或有效 ID 值）。
2. 当特权进程以一个非 0 参数调用 `setuid()` 时，其实际用户 ID、有效用户 ID 和保存 `set-user-ID` 均被置为 `uid` 参数所指定的值。这一操作是单向的，一旦特权进程以此方式修改了其 ID，那么所有特权都将丢失，且之后也不能再使用 `setuid()` 调用将有效用户 ID 重置为 0。如果不希望发生这种情况，请使用稍后介绍的 `seteuid()` 或者 `setreuid()` 系统调用来替代 `setuid()`。

使用 `setgid()` 系统调用修改组 ID 的规则与之相类似，仅需要把 `setuid()` 替换为 `setgid()`，把用户替换为组。因之，规则 1 与前述完全一致，但在规则 2 中，由于对组 ID 的修改不会引起进程特权的丢失（拥有特权与否由有效用户 ID 决定），特权级程序可以使用 `setgid()` 对组 ID 进行任意修改。

对 `set-user-ID-root` 程序（即其有效用户 ID 的当前值为 0）而言，以不可逆方式放弃进程所有特权的首选方法是使用下面的系统调用（以实际用户 ID 值来设置有效用户 ID 和保存 `set-user-ID`）。

```
if (setuid(getuid()) == -1)
    errExit("setuid");
```

`set-user-ID` 程序的属主如果不是 `root` 用户，可使用 `setuid()` 将有效用户 ID 在实际用户 ID

和保存 set-user-ID 之间来回切换，其理由已在 9.4 节中予以阐述。然而，使用 seteuid() 来达成这个目的则更为可取，因为无论 set-user-ID 程序是否属于 root 用户，seteuid() 都能够实现同样的功能。

进程能够使用 seteuid() 来修改其有效用户 ID（改为参数 `eid` 所指定的值），还能使用 setegid() 来修改其有效组 ID（改为参数 `egid` 所指定的值）。

```
#include <unistd.h>
```

```
int seteuid(uid_t eid);  
int setegid(gid_t egid);
```

Both return 0 on success, or -1 on error

进程使用 seteuid() 和 setegid() 来修改其有效 ID 时，会遵循以下规则。

1. 非特权级进程仅能将其有效 ID 修改为相应的实际 ID 或者保存设置 ID。（换言之，对非特权级进程而言，除去前面讨论的 BSD 可移植性问题，seteuid() 和 setegid() 分别等效于 setuid() 和 setgid()。）
2. 特权级进程能够将其有效 ID 修改为任意值。若特权进程使用 seteuid() 将其有效用户 ID 修改为非 0 值，那么此进程将不再具有特权（但可以根据规则 1 来恢复特权）。

对于需要对特权“收放自如”的 set-user-ID 和 set-group-ID 程序，更推荐使用 seteuid()，示例如下：

```
eid = seteuid(); /* Save initial effective user ID (which  
                is same as saved set-user-ID) */  
if (seteuid(getuid()) == -1) /* Drop privileges */  
    errExit("seteuid");  
if (seteuid(eid) == -1) /* Regain privileges */  
    errExit("seteuid");
```

源于 BSD 系统的 seteuid() 和 setegid()，现已纳入 SUSv3 规范，并获得大多数 UNIX 系统实现的支持。

在 GNU C 语言函数库的早期版本中（glibc 2.0 及其之前的版本），将 seteuid(eid) 实现为 setreuid(-1, eid)。而在新版的 glibc 库中，则将 seteuid(eid) 实现为 setresuid(-1, eid, -1)。（稍后将给出对 setreuid()、setresuid() 及其类似函数的描述。）这两种实现都允许将 eid 参数值指定为当前有效用户 ID（即保持不变）。然而，SUSv3 并未对 seteuid() 的这个行为进行规范，并且其他一些 UNIX 实现对此也不支持。总的来说，这种潜在的差异在系统实现间并不明显，因为在通常情况下，有效用户 ID 要么与实际用户 ID 相同，要么与保存 set-user-ID 相同。（要想使有效用户 ID 与二者均不相同，在 Linux 系统中唯一的办法是采用非标准的 setresuid() 系统调用。）

在 glibc 库的所有版本（包括最新版本）中，是以 setegid(-1, egid) 来实现 setegid(egid) 的。如同 seteuid() 一样，这意味着能够将参数 egid 指定为当前有效组 ID，尽管 SUSv3 并未规范这一行为。还有一层含义是使用 setegid() 时，如果对有效组 ID 值的设置不同于当前的实际组 ID，那么还将改变保存 set-group-ID。（类似结论也适用于早期使用 setreuid() 来实现的 seteuid()。）同样，SUSv3 也不支持这一行为。

修改实际 ID 和有效 ID

`setreuid()`系统调用允许调用进程独立修改其实际和有效用户 ID。`setregid()`系统调用对实际和有效组 ID 实现了类似功能。

```
#include <unistd.h>
```

```
int setreuid(uid_t ruid, uid_t euid);  
int setregid(gid_t rgid, gid_t egid);
```

Both return 0 on success, or -1 on error

这两个系统调用的第一个参数都是新的实际 ID，第二个参数都是新的有效 ID。若只想修改其中的一个 ID，可以将另外一个参数指定为-1。

目前，最初派生自 BSD 的 `setreuid()`和 `setregid()`为 SUSv3 规范所接纳，并且获得了大多数 UNIX 系统的支持。

同本节介绍的其他系统调用一样，使用 `setreuid()`和 `setregid()`来作出变更也要遵循一定的规则。下面将从 `setreuid()`的视角来描述这些规则，除非另有说明，`setregid()`函数的规则也与之类似。

1. 非特权进程只能将其实际用户 ID 设置为当前实际用户 ID 值（即保持不变）或有效用户 ID 值，且只能将有效用户 ID 设置为当前实际用户 ID、有效用户 ID（即无变化）或保存 `set-user-ID`。

SUSv3 声称，对于非特权进程是否能使用 `setreuid()`将其实际用户 ID 修改为实际用户 ID、有效用户 ID 或者保存 `set-user-ID`的当前值，规范不做规定。至于真正能将实际用户 ID 修改成何值，这随 UNIX 实现的不同而不同。

SUSv3 对 `setregid()`的规定稍有不同，非特权进程能够将其实际组 ID 设置为保存 `set-group-ID`的当前值，或者将其有效组 ID 设置为实际组 ID 或保存 `set-group-ID`的当前值。但真正能对实际组 ID 做哪些修取决于具体的 UNIX 实现。

2. 特权级进程能够设置其实际用户 ID 和有效用户 ID 为任意值。
3. 不管进程拥有特权与否，只要如下条件之一成立，就能将保存 `set-user-ID` 设置成（新的）有效用户 ID。
 - a) `ruid` 不为-1（即设置实际用户 ID，即便是置为当前值）。
 - b) 对有效用户 ID 所设置的值不同于系统调用之前的实际用户 ID。

反过来说，如果进程使用 `setreuid()`仅将有效用户 ID 修改为实际用户 ID 的当前值，那么保存 `set-user-ID` 的值将保持不变，并且后续可调用 `setreuid()`（或 `seteuid()`）将有效用户 ID 恢复为保存 `set-user-ID` 的值。（`setreuid()`和 `setregid()`针对保存设置 ID 的这一效果，SUSv3 未做规定，但已被 SUSv4 纳入规范。）

规则 3 为 `set-user-ID` 程序提供了一个永久放弃特权的方法，使用如下调用：

```
setreuid(getuid(), getuid());
```

`set-user-ID-root` 进程若有意将用户凭证和组凭证改变为任意值，则应首先调用 `setregid()`，然后再调用 `setreuid()`。一旦调用顺序颠倒，那么调用 `setregid()`将会失败，因为调用 `setreuid()`

后，程序将不再具有特权。若使用 `setresuid()` 和 `setresgid()`（详见下述）来实现此功能，上述描述也同样适用。

直至 4.3BSD，BSD 发行版都不支持保存 `set-user-ID` 和保存 `set-group-ID`（如今已为 SUSv3 强制要求支持）。相反，在 BSD 中，`setreuid()` 和 `setregid()` 允许进程通过来回交换实际 ID 和有效 ID 来“收、放”特权。这一方式的不良副作用在于为了改变有效用户 ID 而改变实际用户 ID。

获取实际、有效和保存设置 ID

在大多数 UNIX 实现中，进程不能直接获取（或修改）其保存 `set-user-ID` 和保存 `set-group-ID` 的值。然而，Linux 提供了两个（非标准的）系统调用来实现此项功能：`getresuid()` 和 `getresgid()`。

```
#define _GNU_SOURCE
#include <unistd.h>

int getresuid(uid_t *ruid, uid_t *euid, uid_t *suid);
int getresgid(gid_t *rgid, gid_t *egid, gid_t *sgid);

Both return 0 on success, or -1 on error
```

`getresuid()` 系统调用将调用进程的当前实际用户 ID、有效用户 ID 和保存 `set-user-ID` 值返回至给定 3 个参数所指定的位置。`getresgid()` 系统调用针对相应的组 ID 实现了类似功能。

修改实际、有效和保存设置 ID

`setresuid()` 系统调用允许调用进程独立修改其 3 个用户 ID 的值。每个用户 ID 的新值由系统调用的 3 个参数给定。`setresgid()` 系统调用对相应的组 ID 实现了类似功能。

```
#define _GNU_SOURCE
#include <unistd.h>

int setresuid(uid_t ruid, uid_t euid, uid_t suid);
int setresgid(gid_t rgid, gid_t egid, gid_t sgid);

Both return 0 on success, or -1 on error
```

若不想同时修改这些 ID，则需将无意修改的 ID 参数值指定为 -1。例如，下列调用等同于 `setuid(x)` 调用：

```
setresuid(-1, x, -1);
```

关于 `setresuid()` 可做何种修改的规则（`setresgid()` 与之类似）如下所示。

1. 非特权进程能够将实际用户 ID、有效用户 ID 和保存 `set-user-ID` 中的任一 ID 设置为实际用户 ID、有效用户 ID 或保存 `set-user-ID` 之中的任一当前值。
2. 特权级进程能够对其实际用户 ID、有效用户 ID 和保存 `set-user-ID` 做任意设置。
3. 不管系统调用是否对其他 ID 做了任何改动，总是将文件系统用户 ID 设置为与有效用户 ID（可能是新值）相同。

`setresuid()` 和 `setresgid()` 调用具有 0/1 效应，即对 ID 的修改请求要么全都成功，要么全部失败。（这也适用于本章所述其他修改多个 ID 的系统调用。）

虽然 `setresuid()` 和 `setresgid()` 为修改进程凭证提供了最为直接的 API，但在应用程序中采用

这些调用会带来可移植性问题。SUSv3 规范并未包括这些调用，且其他 UNIX 实现对其也鲜有支持。

9.7.2 获取和修改文件系统 ID

前述所有修改进程有效用户 ID 或组 ID 的系统调用总是会修改相应的文件系统 ID。要想独立于有效 ID 而修改文件系统 ID，必须使用 Linux 特有的系统调用：`setfsuid()` 和 `setfsgid()`。

```
#include <sys/fsuid.h>

int setfsuid(uid_t fsuid);
                                Always returns the previous file-system user ID

int setfsgid(gid_t fsgid);
                                Always returns the previous file-system group ID
```

`setfsuid()` 系统调用将进程文件系统用户 ID 修改为参数 `fsuid` 所指定的值。`setfsgid()` 系统调用将文件系统组 ID 修改为参数 `fsgid` 所指定的值。

同样，此类变更也存在一些规则。`setfsgid()` 的规则类似于 `setfsuid()`，下面以 `setfsuid()` 为例。

1. 非特权进程能够将文件系统用户 ID 设置为实际用户 ID、有效用户 ID、文件系统用户 ID（即保持不变）或保存 `set-user-ID` 的当前值。
2. 特权级进程能够将文件系统用户 ID 设置为任意值。

这些系统调用的实现存在一些瑕疵。首先，没有相应的系统调用来获取当前的文件系统 ID。另外，这些系统调用根本不做错误检查。一旦非特权进程试图将文件系统 ID 设置为一个非法值，这一不轨企图也只是被静默地忽略掉。无论这些调用成功与否，其返回值都是之前相关文件系统的 ID。因此，这确实也是一种获得当前文件系统 ID 的方法，但却只能是在尝试修改这些值（不管是否成功）的同时进行。

在 Linux 系统中，使用 `setfsuid()` 和 `setfsgid()` 系统调用已不是必要的，若需要将应用程序移植到其他 UNIX 实现上，则应在设计时避免使用这两个调用。

9.7.3 获取和修改辅助组 ID

`getgroups()` 系统调用会将当前进程所属组的集合返回至由参数 `grouplist` 指向的数组中。

```
#include <unistd.h>

int getgroups(int gidsetsize, gid_t grouplist[]);
    Returns number of group IDs placed in grouplist on success, or -1 on error
```

像大多数 UNIX 实现一样，Linux 中的 `getgroups()` 仅返回调用进程的辅助组 ID。然而，SUSv3 规范还允许 UNIX 实现在返回的 `grouplist` 中包含调用进程的有效组 ID。

调用程序必须负责为 `grouplist` 数组分配存储空间，并在 `gidsetsize` 参数中指定其长度。若调用成功，`getgroups()` 会返回置于 `grouplist` 中的组 ID 数量。

若进程属组的数量超出 `gidsetsize`，则 `getgroups()` 将返回错误（错误号为 `EINVAL`）。为了避

免发生这种情况,可将 `grouplist` 数组的大小调整为常量 `NGROUPS_MAX+1`(考虑到可移植性,数组中可能包含了有效组 ID),该常量(定义于`<limits.h>`文件中)定义了进程属组的最大数量。因此,可声明 `grouplist` 如下:

```
gid_t grouplist[NGROUPS_MAX + 1];
```

在 Linux 内核版本 2.6.4 之前, `NGROUPS_MAX` 的值为 32。始于内核版本 2.6.4, `NGROUPS_MAX` 的值为 65536。

应用程序要在运行时获取 `NGROUPS_MAX` 的上限,还可使用如下方法。

- 调用 `sysconf(_SC_NGROUPS_MAX)`。(11.2 节解释了 `sysconf()` 的用法。)
- 从 Linux 特有的 `/proc/sys/kernel/ngroups_max` 只读文件中读取该限制。系统从内核 2.6.4 开始提供该文件。

除此之外,应用程序还能在调用 `getgroups()` 时将 `gidtsetsize` 参数指定为 0。这样一来, `grouplist` 数组未作修改,但调用的返回值却给出了进程属组的数量。

通过上述任意一种运行时技术所获取的 `NGROUPS_MAX` 值,可用于为后续的 `getgroups()` 调用动态分配 `grouplist` 数组。

特权级进程能够使用 `setgroups()` 和 `initgroups()` 来修改其辅助组 ID 集合。

```
#define _BSD_SOURCE
#include <grp.h>

int setgroups(size_t gidsetsize, const gid_t *grouplist);
int initgroups(const char *user, gid_t group);

Both return 0 on success, or -1 on error
```

`setgroups()` 系统调用用 `grouplist` 数组所指定的集合来替换调用进程的辅助组 ID。参数 `gidsetsize` 指定了置于参数 `grouplist` 数组中的组 ID 数量。

`initgroups()` 函数将扫描 `/etc/groups` 文件,为 `user` 创建属组列表,以此来初始化调用进程的辅助组 ID。另外,也会将参数 `group` 指定的组 ID 追加到进程辅助组 ID 的集合中。

`initgroups()` 函数的主要用户是创建登录会话的程序——例如 `login(1)`, 在用户调用登录 shell 之前,为进程设置各种属性。此类程序一般通过读取密码文件中用户记录的组属性来获取参数 `group` 的值。这稍微有点令人费解,因为密码文件中的组 ID 实际并非辅助组 ID,而是定义了登录 shell 初始的实际组 ID、有效组 ID 和保存 `set-group-ID`。尽管如此,这却是 `initgroups()` 函数的常用使用方式。

虽然未纳入 SUSv3, `setgroups()` 和 `initgroups()` 却获得了所有 UNIX 实现的支持。

9.7.4 修改进程凭证的系统调用总结

表 9-1 对修改进程凭证的各种系统调用及库函数的效果进行了总结。

图 9-1 提供了表 9-1 中信息的概括图示。本图内容是从修改用户 ID 的角度加以展示的,但修改组 ID 的规则与之类似。

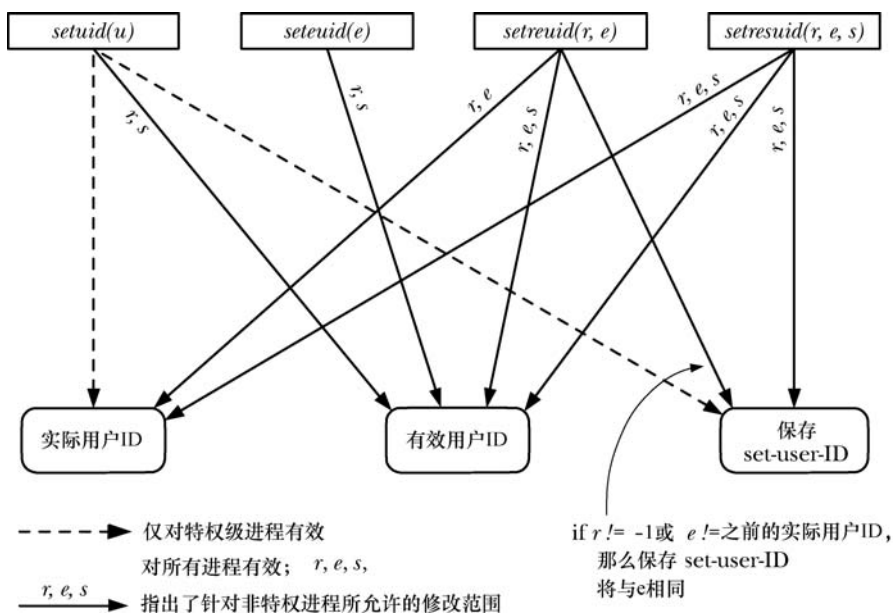


图 9-1: 凭证修改函数对进程用户 ID 的效果

表 9-1: 修改进程凭证的接口一览表

接 口	目的和效果应用于		可 移 植 性
	非特权进程	特权级进程	
setuid(u) setgid(g)	将有效 ID 修改为当前实际 ID 或保存设置 ID	将实际 ID、有效 ID 和保存设置 ID 修改为任何（一个）值	获得 SUSv3 规范的支持，但 BSD 的派生系统具有不同语义
seteuid(e) setegid(e)	将有效 ID 修改为当前的实际或保存设置 ID	修改有效 ID 为任意值	获得 SUSv3 规范支持
setreuid(r, e) setregid(r, e)	（独立）将实际 ID 修改为当前实际 ID 或有效 ID 值，将有效 ID 修改为当前实际 ID、有效 ID 或保存设置 ID	（独立）将实际 ID 和有效 ID 修改为任意值	获得 SUSv3 规范支持，但操作随系统实现不同而不同
setresuid(r, e, s) setresgid(r, e, s)	（独立）将实际 ID、有效 ID 和保存设置 ID 修改为当前实际 ID、有效 ID 或保存设置 ID	（独立）将实际 ID、有效 ID 和保存设置 ID 修改为任意值	未获 SUSv3 规范支持，并且鲜见于其他 UNIX 实现
setfsuid(u) setfsgid(u)	将文件系统 ID 修改为当前实际 ID、有效 ID、文件系统 ID 或者保存设置 ID	将文件系统 ID 修改为任意值	Linux 系统所特有
setgroups(n, 1)	非特权进程无法调用	设置辅助组 ID 为任意值	未见诸于 SUSv3 规范，但获得所有 UNIX 实现的支持

补充说明表 9-1 中的信息。

- glibc 库对 seteuid()(setresuid(-1, e, -1))和 setegid()(setregid(-1, e, -1))函数的实现方式允许将有效 ID 设置为有效 ID 的当前值，但 SUSv3 对此未作规范。此外，若将有效组 ID 设

置为当前实际组 ID 之外的值，那么 `setegid()` 的函数实现还会修改保存设置组 ID。（对于 `setegid()` 实现这一修改保存 `set-group-ID` 的行为，SUSv3 也未作规范。）

- 针对特权级进程和非特权进程调用 `setreuid()` 和 `setregid()` 的情况，若 `r` 的值不等于 -1，或者 `e` 的值有别于函数调用前的实际 ID，则将保存 `set-user-ID` 或保存 `set-group-ID` 设置为（新的）有效 ID。（`setreuid()` 和 `setregid()` 函数对保存设置 ID 的修改未获 SUSv3 支持。）
- 只要修改了有效用户（组）ID，就会将 Linux 特有的文件系统用户（组）ID 也修改为相同值。
- 不管有效用户 ID 是否改变，`setresuid()` 系统调用总是把文件系统用户 ID 修改为有效用户 ID，`setresgid()` 系统调用对文件系统组 ID 的效力与之类似。

9.7.5 示例：显示进程凭证

程序清单 9-1 中的程序使用前述系统调用和库函数来获取进程的所有用户 ID 和组 ID，并显示出来。

程序清单 9-1：显示进程的所有用户 ID 和组 ID

```
----- proccred/idshow.c
#define _GNU_SOURCE
#include <unistd.h>
#include <sys/fsuid.h>
#include <limits.h>
#include "ugid_functions.h" /* userNameFromId() & groupNameFromId() */
#include "tlpi_hdr.h"

#define SG_SIZE (NGROUPS_MAX + 1)

int
main(int argc, char *argv[])
{
    uid_t ruid, euid, suid, fsuid;
    gid_t rgid, egid, sgid, fsgid;
    gid_t suppGroups[SG_SIZE];
    int numGroups, j;
    char *p;

    if (getresuid(&ruid, &euid, &suid) == -1)
        errExit("getresuid");
    if (getresgid(&rgid, &egid, &sgid) == -1)
        errExit("getresgid");

    /* Attempts to change the file-system IDs are always ignored
       for unprivileged processes, but even so, the following
       calls return the current file-system IDs */

    fsuid = setfsuid(0);
    fsgid = setfsgid(0);

    printf("UID: ");
    p = userNameFromId(ruid);
    printf("real=%s (%ld); ", (p == NULL) ? "???" : p, (long) ruid);
    p = userNameFromId(euid);
    printf("eff=%s (%ld); ", (p == NULL) ? "???" : p, (long) euid);
    p = userNameFromId(suid);
    printf("saved=%s (%ld); ", (p == NULL) ? "???" : p, (long) suid);
```

```

p = userNameFromId(fsuid);
printf("fs=%s (%ld); ", (p == NULL) ? "???" : p, (long) fsuid);
printf("\n");

printf("GID: ");
p = groupNameFromId(rgid);
printf("real=%s (%ld); ", (p == NULL) ? "???" : p, (long) rgid);

p = groupNameFromId(egid);
printf("eff=%s (%ld); ", (p == NULL) ? "???" : p, (long) egid);
p = groupNameFromId(sgid);
printf("saved=%s (%ld); ", (p == NULL) ? "???" : p, (long) sgid);
p = groupNameFromId(fsgid);
printf("fs=%s (%ld); ", (p == NULL) ? "???" : p, (long) fsgid);
printf("\n");

numGroups = getgroups(SG_SIZE, suppGroups);
if (numGroups == -1)
    errExit("getgroups");

printf("Supplementary groups (%d): ", numGroups);
for (j = 0; j < numGroups; j++) {
    p = groupNameFromId(suppGroups[j]);
    printf("%s (%ld) ", (p == NULL) ? "???" : p, (long) suppGroups[j]);
}
printf("\n");

exit(EXIT_SUCCESS);
}

```

proccred/idshow.c

9.8 总结

每个进程都有一干与之相关的用户 ID 和组 ID（凭证）。实际 ID 定义了进程所属¹。在大多数的 UNIX 实现中，进程对诸如文件之类资源的访问，其许可权限由有效 ID 决定。然而，Linux 会使用文件系统 ID 来决定对文件的访问权限，而将有效 ID 用于检查其他权限。（因为文件系统 ID 一般等同于相应的有效 ID，所以 Linux 对文件权限的检查方式与其他 UNIX 实现相同。）进程辅助组 ID 则是出于权限检查目的而另行设立的进程属组集合。存在各种系统调用和库函数支持进程获取和修改其用户 ID 和组 ID。

set-user-ID 程序运行时，会将进程有效用户 ID 置为文件属主的用户 ID。运行某个特殊程序时，这种机制支持用户“假借”其他用户的身份和特权。相应的，set-group-ID 程序会修改运行该程序的进程的有效组 ID。保存 set-user-ID 和保存 set-group-ID 允许 set-user-ID 和 set-group-ID 程序临时性地放弃特权，并在之后恢复特权。

0 在用户 ID 中可谓卓尔不群。通常仅为一个名为 root 的账号所有。有效用户 ID 为 0 的进程属特权级进程。换言之，对于进程发起的各种系统调用，可免于接受通常所要历经的诸多权限检查（比如那些能够随意修改进程各种用户 ID 和组 ID 的调用）。

¹ 译者注：即属主和属组。

9.9 习题

- 9-1. 在下列每种情况中，假设进程用户 ID 的初始值分别为 `real`（实际）=1000、`effective`（有效）=0、`saved`（保存）=0、`file-system`（文件系统）=0。当执行这些调用后，用户 ID 的状态如何？

- a) `setuid(2000);`
- b) `setreuid(-1, 2000);`
- c) `seteuid(2000);`
- d) `setfsuid(2000);`
- e) `setresuid(-1, 2000, 3000);`

- 9-2. 拥有如下用户 ID 的进程享有特权吗？请予解释。

`real=0 effective=1000 saved=1000 file-system=1000`

- 9-3. 使用 `setgroups()` 及库函数从密码文件、组文件（参见 8.4 节）中获取信息，以实现 `initgroups()`。请注意，欲调用 `setgroups()`，进程必须享有特权。

`real=0 effective=1000 saved=1000 file-system=1000`

- 9-4. 假设某进程的所有用户标识均为 X，执行了用户 ID 为 Y 的 `set-user-ID` 程序，且 Y 为非 0 值，对进程凭证的设置如下：

`real=X effective=Y saved=Y`

（这里忽略了文件系统用户 ID，因为该 ID 随有效用户 ID 的变化而变化。）为执行如下操作，请分别列出对 `setuid()`、`seteuid()`、`setreuid()` 和 `setresuid()` 的调用。

a) 挂起和恢复 `set-user-ID` 身份（即将有效用户 ID 在实际用户 ID 和保存 `set-user-ID` 间切换）。

b) 永久放弃 `set-user-ID` 身份（即确保将有效用户 ID 和保存 `set-user-ID` 设置为实际用户 ID）。

（该练习还需要使用 `getuid()` 和 `geteuid()` 函数来获取进程的实际用户 ID 和有效用户 ID。）请注意，鉴于上述列出的某些系统调用，部分操作将无法实现。

- 9-5. 针对执行 `set-user-ID-root` 程序的进程，重复上述练习，进程凭证的初始值如下：

`real=X effective=0 saved=0`

第 10 章

时 间

程序可能会关注两种时间类型。

- 真实时间：度量这一时间的起点有二：一为某个标准点；二为进程生命周期内的某个固定时点（通常为程序启动）。前者为日历（calendar）时间，适用于需要对数据库记录或文件打上时间戳的程序；后者则称之为流逝（elapsed）时间或挂钟（wall clock）时间，主要针对需要周期性操作或定期从外部输入设备进行度量的程序。
- 进程时间：一个进程所使用的 CPU 时间总量，适用于对程序、算法性能的检查或优化。

大多数计算机体系结构都内置有硬件时钟，使内核得以计算真实时间和进程时间。本章将介绍系统调用对这两种时间的处理，以及在可读时间和机器时间之间互相转换的库函数。由于可读时间的表现形式与地理位置、语言和文化习俗有关，讨论这一话题自然引出对时区和地区的研究。

10.1 日历时间（Calendar Time）

无论地理位置如何，UNIX 系统内部对时间的表示方式均是以自 Epoch 以来的秒数来度量的，Epoch 亦即通用协调时间（UTC，以前也称为格林威治标准时间，或 GMT）的 1970 年 1 月 1 日早晨零点。这也是 UNIX 系统问世的大致日期。日历时间存储于类型为 `time_t` 的变量中，此类型是由 SUSv3 定义的整数类型。

在 32 位 Linux 系统，`time_t` 是一个有符号整数，可以表示的日期范围从 1901 年 12 月 13 日 20 时 45 分 52 秒至 2038 年 1 月 19 号 03:14:07。（SUSv3 未定义 `time_t` 值为负数时的含义。）因此，当前许多 32 位 UNIX 系统都面临一个 2038 年的理论问题，如果执行的计算工作涉及未来日期，那么在 2038 年之前就会与之遭遇。事实上，到了 2038 年，可能所有的 UNIX 系统都早已升级为 64 位或更多位数的系统，这一问题也许会随之而大为缓解。然而，32 位嵌入式系统，由于其寿命较之台式机硬件更长，故而仍然会受此问题的困扰。此外，对于依然以 32 位 `time_t` 格式保存时间的历史数据和应用程序，这个问题将依然存在。

系统调用 `gettimeofday()`，可于 `tv` 指向的缓冲区中返回日历时间。

```
#include <sys/time.h>

int gettimeofday(struct timeval *tv, struct timezone *tz);

Returns 0 on success, or -1 on error
```

参数 `tv` 是指向如下数据结构的一个指针：

```
struct timeval {
    time_t    tv_sec;        /* Seconds since 00:00:00, 1 Jan 1970 UTC */
    suseconds_t tv_usec;    /* Additional microseconds (long int) */
};
```

虽然 `tv_usec` 字段能提供微秒级精度，但其返回值的准确性则由依赖于构架的具体实现来决定。`tv_usec` 中的 `u` 源于与之形似的希腊字母 μ （读音“mu”），在公制系统中表示百万分之一。在现代 X86-32 系统上，`gettimeofday()` 的确可以提供微秒级的准确度（例如，Pentium 系统内置有时间戳计数寄存器，随每个 CPU 时钟周期而加一）。

`gettimeofday()` 的参数 `tz` 是个历史产物。早期的 UNIX 实现用其来获取系统的时区信息，**目前已遭废弃，应始终将其置为 NULL。**

如果提供了 `tz` 参数，那么将返回一个 `timezone` 的结构体，其内容为上次调用 `settimeofday()` 时传入的 `tz` 参数（已废弃）值。该结构包含两个字段 `tz_minuteswest` 和 `tz_dsttime`。`tz_minuteswest` 字段表示欲将本时区时间转换为 UTC 时间所必须增加的分钟数，如为负值，则表示此时区位于 UTC 以东（例如，如为欧洲中部时间，会提前 UTC 一小时，则将此字段设置为 -60）。`tz_dsttime` 字段内为一个常量，意在表示这个时区是否强制施行夏令时（DST）制。正由于夏令时制度无法用一个简单算法加以表达，故而 `tz` 参数已遭废弃。（Linux 从未支持过此参数。）详情请参考 `gettimeofday(2)` 手册页。

`time()` 系统调用返回自 Epoch 以来的秒数（和函数 `gettimeofday()` 所返回的 `tv` 参数中 `tv_sec` 字段的数值相同）。

```
#include <time.h>

time_t time(time_t *timep);

Returns number of seconds since the Epoch, or (time_t) -1 on error
```

如果 `timep` 参数不为 NULL，那么还会将自 Epoch 以来的秒数置于 `timep` 所指向的位置。

由于 `time()` 会以两种方式返回相同的值，而使用时唯一可能出错的地方是赋予 `timep` 参数一个无效地址（EFAULT），因此往往会简单地采用如下调用（不做错误检查）：

```
t = time(NULL);
```

之所以存在两个本质上目的相同的系统调用（`time()` 和 `gettimeofday()`），自有其历史原因。早期的 UNIX 实现提供了 `time()`。而 4.3BSD 又补充了更为精确的 `gettimeofday()` 系统调用。这时，再将 `time()` 作为系统调用就显得多余，可以将其实现为一个调用 `gettimeofday()` 的库函数。

10.2 时间转换函数

图 10-1 所示为用于在 `time_t` 值和其他时间格式之间相互转换的函数，其中包括打印输

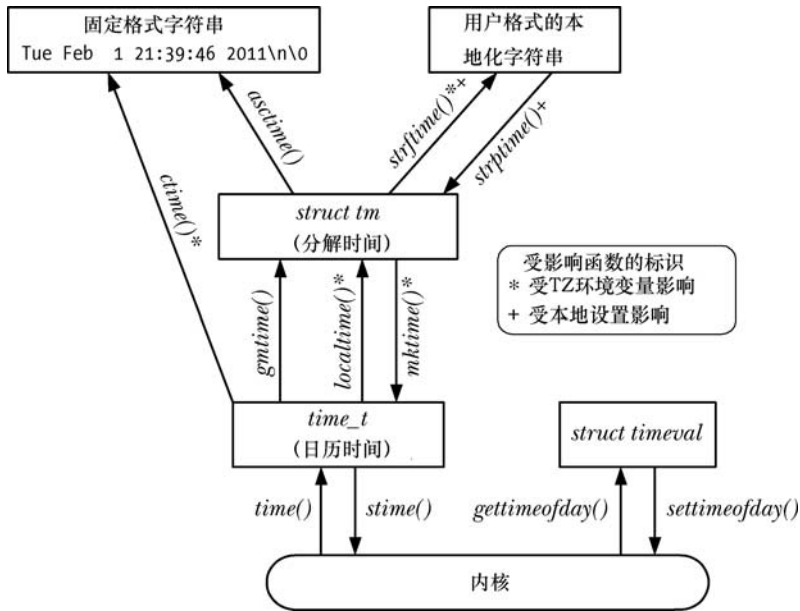


图 10-1: 获取和使用日历时间的函数

出。这些函数屏蔽了因时区、夏令时（DST）制和本地化等问题给转换所带来的种种复杂性。10.3 节将讨论时区（timezone），10.4 节将讨论地区（locale）。

10.2.1 将 `time_t` 转换为可打印格式

为了将 `time_t` 转换为可打印格式，`ctime()` 函数提供了一个简单方法。

```
#include <time.h>

char *ctime(const time_t *timep);

Returns pointer to statically allocated string terminated
by newline and \0 on success, or NULL on error
```

把一个指向 `time_t` 的指针作为 `timep` 参数传入函数 `ctime()`，将返回一个长达 26 字节的字符串，内含标准格式的日期和时间，如下例所示：

```
Wed Jun 8 14:22:34 2011
```

该字符串包含换行符和终止空字节各一。`ctime()` 函数在进行转换时，会自动对本地时区和 DST 设置加以考虑（10.3 节将解释这些设置的确定过程）。返回的字符串经由静态分配，下一次对 `ctime()` 的调用会将其覆盖。

SUSv3 规定，调用 `ctime()`、`gmtime()`、`localtime()` 或 `asctime()` 中的任一函数，都可能会覆盖由其他函数返回，且经静态分配的数据结构。换言之，这些函数可以共享返回的字符数组

和 `tm` 结构体，某些版本的 `glibc` 也正是这样实现的。如果有意在对这些函数的多次调用间维护返回的信息，那么必须将其保存在本地副本中。

`ctime_r()`是 `ctime()`的可重入版本。(21.1.2 节将解释重入。)该函数允许调用者额外指定一个指针参数，所指向的缓冲区（由调用者提供）用于返回时间字符串。本章所论及的其他函数的可重入版，其操作方式与之类似。

10.2.2 `time_t` 和分解时间之间的转换

函数 `gmtime()`和 `localtime()`可将一 `time_t` 值转换为一个所谓的分解时间（broken-down time）。分解时间被置于一个经由静态分配的结构中，其地址则作为函数结果返回。

```
#include <time.h>

struct tm *gmtime(const time_t *timep);
struct tm *localtime(const time_t *timep);

Both return a pointer to a statically allocated broken-down
time structure on success, or NULL on error
```

函数 `gmtime()`能够把日历时间转换为一个对应于 UTC 的分解时间。（字母 GM 源于格林威治标准时间）。相形之下，函数 `localtime()`需要考虑时区和夏令时设置，返回对应于系统本地时间的一个分解时间。

`gmtime_r()`和 `localtime_r()`分别是这些函数的可重入版。

在这些函数所返回的 `tm` 结构中，日期和时间被分解为多个独立字段，其形式如下：

```
struct tm {
    int tm_sec;        /* Seconds (0-60) */
    int tm_min;        /* Minutes (0-59) */
    int tm_hour;       /* Hours (0-23) */
    int tm_mday;       /* Day of the month (1-31) */
    int tm_mon;        /* Month (0-11) */
    int tm_year;       /* Year since 1900 */
    int tm_wday;       /* Day of the week (Sunday = 0)*/
    int tm_yday;       /* Day in the year (0-365; 1 Jan = 0)*/
    int tm_isdst;      /* Daylight saving time flag
                       > 0: DST is in effect;
                       = 0: DST is not effect;
                       < 0: DST information not available */
};
```

将字段 `tm_sec` 的上限设为 60（而非 59）以考虑闰秒，偶尔会用其将人类日历调整至精确的天文年（所谓的回归年）。

如果定义了 `_BSD_SOURCE` 功能测试宏，那么由 `glibc` 定义的 `tm` 结构还会包含两个额外字段，以描述关于所示时间的深入信息。第一个字段 `long int tm_gmtoff`，包含所示时间超出 UTC 以东的秒数。第二个字段 `const char* tm_zone`，是时区名称的缩写（例如，CEST 为欧洲中部夏令时间）。SUSv3 并未定义这些字段，它们只见诸于少数其他 UNIX 实现（主要为 BSD 衍生版本）。

函数 `mktime()` 将一个本地时区的分解时间翻译为 `time_t` 值，并将其作为函数结果返回。

调用者将分解时间置于一个 `tm` 结构，再以 `timeptr` 指针指向该结构。这一转换会忽略输入 `tm` 结构中的 `tm_wday` 和 `tm_yday` 字段。

```
#include <time.h>

time_t mktime(struct tm *timeptr);

Returns seconds since the Epoch corresponding to timeptr
on success, or (time_t) -1 on error
```

函数 `mktime()` 可能会修改 `timeptr` 所指向的结构体，至少会确保对 `tm_wday` 和 `tm_yday` 字段值的设置，会与其他输入字段的值能对应起来。

此外，`mktime()` 不要求 `tm` 结构体的其他字段受到前述范围的限制。任何一个字段的值超出范围，`mktime()` 都会将其调整回有效范围之内，并适当调整其他字段。所有这些调整，均发生于 `mktime()` 更新 `tm_wday` 和 `tm_yday` 字段并计算返回值 `time_t` 之前。

例如，如果输入字段 `tm_sec` 的值为 123，那么在返回时此字段的值将为 3，且 `tm_min` 字段值会在其之前值的基础上加 2。（如果这一改动造成 `tm_min` 溢出，那么将调整 `tm_min` 的值，并且递增 `tm_hour` 字段，以此类推。）这些调整甚至适用于字段负值。例如，指定 `tm_sec` 为 -1 即意味着前一分钟的第 59 秒。此功能允许以分解时间来计算日期和时间，故而非常有用。

`mktime()` 在进行转换时会对时区进行设置。此外，DST 设置的使用与否取决于输入字段 `tm_isdst` 的值。

- 若 `tm_isdst` 为 0，则将这一时间视为标准间（即，忽略夏令时，即使实际上每年的这一时刻处于夏令时阶段）。
- 若 `tm_isdst` 大于 0，则将这一时间视为夏令时（即，夏令时生效，即使每年的此时不处于夏令时阶段）。
- 若 `tm_isdst` 小于 0，则试图判定 DST 在每年的这一时间是否生效。这往往是众望所归的设置。

（无论 `tm_isdst` 的初始设置如何）在转换完成时，如果针对给定的时间，DST 生效，`mktime()` 会将 `tm_isdst` 字段置为正值，若 DST 未生效，则将 `tm_isdst` 置为 0。

10.2.3 分解时间和打印格式之间的转换

本节会介绍将分解时间和打印格式相互进行转换的函数。

从分解时间转换为打印格式

在参数 `tm` 中提供一个指向分解时间结构的指针，`asctime()` 则会返回一指针，指向经由静态分配的字符串，内含时间，格式则与 `ctime()` 相同。

```
#include <time.h>

char *asctime(const struct tm *timeptr);

Returns pointer to statically allocated string terminated by
newline and \0 on success, or NULL on error
```

相形于 `ctime()`，本地时区设置对 `asctime()` 没有影响，因为其所转换的是一个分解时间，

该时间通常要么已然通过 `localtime()` 作了本地化处理，要么早已经由 `gmtime()` 转换成了 UTC。如同 `ctime()` 一样，`asctime()` 也无法控制其所生成字符串的格式。

`asctime()` 的可重入版为 `asctime_r()`。

程序清单 10-1 演示 `asctime()` 以及直到本章结尾所述时间转换函数的用法。该程序获取当前的日历时间，随后使用各种时间转换函数并显示其结果。下例为冬季在德国慕尼黑运行此程序的结果，该地区处于欧洲中部时间这一时区，比 UTC 要早一小时。

```
$ date
Tue Dec 28 16:01:51 CET 2010
$ ./calendar_time
Seconds since the Epoch (1 Jan 1970): 1293548517 (about 40.991 years)
  gettimeofday() returned 1293548517 secs, 715616 microsecs
Broken down by gmtime():
  year=110 mon=11 mday=28 hour=15 min=1 sec=57 wday=2 yday=361 isdst=0
Broken down by localtime():
  year=110 mon=11 mday=28 hour=16 min=1 sec=57 wday=2 yday=361 isdst=0

asctime() formats the gmtime() value as: Tue Dec 28 15:01:57 2010
ctime() formats the time() value as: Tue Dec 28 16:01:57 2010
mktime() of gmtime() value: 1293544917 secs
mktime() of localtime() value: 1293548517 secs      3600 secs ahead of UTC
```

程序清单 10-1: 获取和转换日历时间

```
time/calendar_time.c

#include <locale.h>
#include <time.h>
#include <sys/time.h>
#include "tspi_hdr.h"
#define SECONDS_IN_TROPICAL_YEAR (365.24219 * 24 * 60 * 60)

int
main(int argc, char *argv[])
{
    time_t t;
    struct tm *gmp, *locp;
    struct tm gm, loc;
    struct timeval tv;

    t = time(NULL);
    printf("Seconds since the Epoch (1 Jan 1970): %ld", (long) t);
    printf(" (about %6.3f years)\n", t / SECONDS_IN_TROPICAL_YEAR);

    if (gettimeofday(&tv, NULL) == -1)
        errExit("gettimeofday");
    printf("  gettimeofday() returned %ld secs, %ld microsecs\n",
           (long) tv.tv_sec, (long) tv.tv_usec);

    gmp = gmtime(&t);
    if (gmp == NULL)
        errExit("gmtime");

    gm = *gmp;          /* Save local copy, since *gmp may be modified
                        by asctime() or gmtime() */
```

```

printf("Broken down by gmtime():\n");
printf(" year=%d mon=%d mday=%d hour=%d min=%d sec=%d ", gm.tm_year,
      gm.tm_mon, gm.tm_mday, gm.tm_hour, gm.tm_min, gm.tm_sec);
printf("wday=%d yday=%d isdst=%d\n", gm.tm_wday, gm.tm_yday, gm.tm_isdst);

locp = localtime(&t);
if (locp == NULL)
    errExit("localtime");

loc = *locp;          /* Save local copy */

printf("Broken down by localtime():\n");
printf(" year=%d mon=%d mday=%d hour=%d min=%d sec=%d ",
      loc.tm_year, loc.tm_mon, loc.tm_mday,
      loc.tm_hour, loc.tm_min, loc.tm_sec);
printf("wday=%d yday=%d isdst=%d\n\n",
      loc.tm_wday, loc.tm_yday, loc.tm_isdst);

printf("asctime() formats the gmtime() value as: %s", asctime(&gm));
printf("ctime() formats the time() value as: %s", ctime(&t));

printf("mktime() of gmtime() value: %ld secs\n", (long) mktime(&gm));
printf("mktime() of localtime() value: %ld secs\n", (long) mktime(&loc));

exit(EXIT_SUCCESS);
}

```

time/calendar_time.c

当把一个分解时间转换成打印格式时，函数 `strftime()` 可以提供更为精确的控制。令 `timeptr` 指向分解时间，`strftime()` 会将以 `null` 结尾、由日期和时间组成的相应字符串置于 `outstr` 所指向的缓冲区中。

```

#include <time.h>

size_t strftime(char *outstr, size_t maxsize, const char *format,
               const struct tm *timeptr);

```

Returns number of bytes placed in *outstr* (excluding terminating null byte) on success, or 0 on error

`outstr` 中返回的字符串按照 `format` 参数定义的格式做了格式化。`Maxsize` 参数指定 `outstr` 的最大长度。不同于 `ctime()` 和 `asctime()`，`strftime()` 不会在字符串的结尾包括换行符（除非 `format` 中定义有换行符）。

如果成功，`strftime()` 返回 `outstr` 所指缓冲区的字节长度，且不包括终止空字节。如果结果字符串的总长度，含终止空字节，超过了 `maxsize` 参数，那么 `strftime()` 会返回 0 以示出错，且此时无法确定 `outstr` 的内容。

`strftime()` 的 `format` 参数是一字符串，与赋予 `printf()` 的参数相类似。冠以百分号 (%) 的字符序列是对转换的定义，函数会将百分号后的说明符字符一一替换为日期和时间的组成部分。这是一套相当丰富的转换说明符，表 10-1 中所列的是其一个子集。（完整的列表可见诸于 `strftime(3)` 手册页。）除非特别注明，所有这些转换说明符都符合 SUSv3 标准。

%U 和 %W 说明符都生成一年中的周数。%U 的周数按以下方法计算。含有星期日的第一周编号为 1，此周的前一周编号为 0。如果星期天恰巧是当年的第一天，那么就没有第 0 周，

当年的最后一天则属于第 53 周。%W 的周数编号以同样的方式来计算，只不过计算对象是周一而非周日。

通常情况下，我们希望在本书的各种示范程序中显示当前时间。为此，本书提供了函数 currTime()，其返回一字符串，内含 strftime()按 format 参数格式化的当前时间。

```
#include "curr_time.h"

char *currTime(const char *format);

Returns pointer to statically allocated string, or NULL on error
```

程序清单 10-2 所示为 currTime()函数的实现。

表 10-1: strftime()的转换说明符选集

说 明 符	描 述	例 子
%%	百分号 (%) 字符	%
%a	星期几的缩写	Tue
%A	星期几的全称	Tuesday
%b, %h	月份名称的缩写	Feb
%B	月份全称	February
%c	日期和时间	Tue Feb 1 21:39:46 2011
%d	一个月的一天 (2 位数字, 01 至 31 天)	01
%D	美国日期格式 (与%m%d%y 相同)	02/01/11
%e	一个月中的一天 (2 个字符)	1
%F	ISO 日期格式 (与%Y-%m-%d 相同)	2011-02-01
%H	小时 (24 小时制, 2 位数)	21
%I	小时 (12 小时制, 2 位数)	09
%j	一年中的一天 (3 位数字, 从 001 到 366)	032
%m	十进制月 (2 位, 01 到 12)	02
%M	分 (2 位数)	39
%p	AM/PM	PM
%P	上午/下午 (GNU 扩展)	pm
%R	24 小时制的时间 (和%H:%M 格式相同)	21:39
%S	秒 (00 至 60)	46
%T	时间 (和%H:%M:%S 格式相同)	21: 39: 46
%u	星期几编号 (1 至 7, 星期一=1)	2
%U	以周日计算、一年中的周数 (00 到 53)	05
%w	星期几编号 (0 至 6, 星期日=0)	2

说 明 符	描 述	例 子
%W	以周一计算、一年中的周数 (00 到 53)	05
%x	日期 (本地化)	02/01/11
%X	时间 (本地化)	21: 39: 46
%y	2 位数字年份	11
%Y	4 位数字年份	2011
%Z	时区名称	CET

程序清单 10-2: 返回当前时间的字符串的函数

```

time/curr_time.c
#include <time.h>
#include "curr_time.h"          /* Declares function defined here */

#define BUF_SIZE 1000

/* Return a string containing the current time formatted according to
   the specification in 'format' (see strftime(3) for specifiers).
   If 'format' is NULL, we use "%c" as a specifier (which gives the
   date and time as for ctime(3), but without the trailing newline).
   Returns NULL on error. */
char *
currTime(const char *format)
{
    static char buf[BUF_SIZE]; /* Nonreentrant */
    time_t t;
    size_t s;
    struct tm *tm;

    t = time(NULL);
    tm = localtime(&t);
    if (tm == NULL)
        return NULL;

    s = strftime(buf, BUF_SIZE, (format != NULL) ? format : "%c", tm);

    return (s == 0) ? NULL : buf;
}
time/curr_time.c

```

将打印格式时间转换为分解时间

函数 `strptime()` 是 `strftime()` 的逆向函数，将包含日期和时间的字符串转换成一分解时间。

```

#define _XOPEN_SOURCE
#include <time.h>

char *strptime(const char *str, const char *format, struct tm *timeptr);
           Returns pointer to next unprocessed character in
           str on success, or NULL on error

```


函数 `strptime()` 按照参数 `format` 内的格式要求, 对由日期和时间组成的字符串 `str` 加以解析, 并将转换后的分解时间置于指针 `timeptr` 所指向的结构体中。

如果成功, `strptime()` 返回一指针, 指向 `str` 中下一个未经处理的字符。(如果字符串中还包含有需要应用程序处理的额外信息, 这一特性就能派上用场。) 如果无法匹配整个格式字符串, `strptime()` 返回 `NULL`, 以示出现错误。

`strptime()` 的格式规范类似于 `scanf(3)`, 包含以下类型的字符。

- 转换字符串冠以一个百分号 (%) 字符。
- 如包含空格字符, 则意味着其可匹配零个或多个空格。
- (% 之外的) 非空格字符必须和输入字符串中的相同字符严格匹配。

转换说明类似于之前为 `strptime()` 给出的内容 (表 10-1)。主要的区别在于, 此处的说明符更为通用。例如, 不拘于星期名称的全称或简称, `%a` 和 `%A` 都可接受, 而且 `%d` 和 `%e` 均可用于读取月中的个位数, 无论该数字前面是否有 0。此外, 不区分大小写, 例如, `May` 和 `MAY` 是相同的月份名称。使用字符串 `%%` 来匹配输入字符串中的百分号字符。 `strptime(3)` 手册页提供更多的细节。

`glibc` 在实现 `strptime()` 时, 并不修改 `tm` 结构体中那些未获 `format` 说明符初始化的字段。这也意味着可以根据多个字符串, 例如, 一个日期字符串和一个时间字符串, 发起多次 `strptime()` 调用, 来创建一个 `tm` 结构体。`SUSv3` 虽然允许这一行为, 但并不强制要求实现, 因此在其他 `UNIX` 实现上不能对其有所依赖。要保证应用的可移植性, 就必须确保, 要么 `str` 和 `format` 中所含输入信息足以设置最终 `tm` 结构的所有字段, 要么在调用 `strptime()` 之前对 `tm` 结构体已经做了适当的初始化处理。在大多数情况下, 用 `memset()` 把整个结构体置为 0 也就足够了, 但要留心, 在 `glibc` 和许多其他时间转换函数的实现中, `m_mday` 字段值为 0, 意为上月的最后一天。最后还要注意, `strptime()` 从不设置 `tm` 结构体的 `tm_isdst` 字段。

GNU C 库还提供有与 `strptime()` 功能类似的两个函数: `getdate()` (已由 `SUSv3` 规范, 且应用广泛) 及其可重入版 `getdate_r()` (`SUSv3` 中未定义, 仅获少数 `UNIX` 实现支持)。此处将不会介绍这些函数, 因为在指定用于扫描日期的格式时, 它们所采用的是外部文件 (由环境变量 `DATMSK` 定义), 这不但令其难以使用, 而且会在 `set-user-ID` 程序中造成安全漏洞。

程序清单 10-3 演示了 `strptime()` 和 `strftime()` 的用法。该程序从命令行参数中接受日期和时间, 然后用 `strptime()` 将其转换为一分解时间, 接着使用 `strftime()` 执行逆向转换并显示结果。该程序接收至多 3 个参数, 其中前两个为必需提供。第一个参数是包含日期和时间的字符串。第二个参数指定了 `strptime()` 在解析第一个参数时所采用的格式。可选的第三个参数是格式字符串, 用于 `strftime()` 的逆向转换。如果省略此参数, 将使用一个默认的格式字符串。(本程序中使用的 `setLocale()` 函数将在 10.4 节中加以介绍。) 以下 `shell` 会话日志显示了使用该程序的一些例子:

```
$ ./strptime "9:39:46pm 1 Feb 2011" "%I:%M:%S%p %d %b %Y"
calendar time (seconds since Epoch): 1296592786
strftime() yields: 21:39:46 Tuesday, 01 February 2011 CET
```

以下用法与之相似, 只不过这次为 `strftime()` 明确指定了格式:

```
$ ./strptime "9:39:46pm 1 Feb 2011" "%I:%M:%S%p %d %b %Y" "%F %T"
calendar time (seconds since Epoch): 1296592786
strftime() yields: 2011-02-01 21:39:46
```

```

#define _XOPEN_SOURCE
#include <time.h>
#include <locale.h>
#include "t1pi_hdr.h"

#define SBUF_SIZE 1000

int
main(int argc, char *argv[])
{
    struct tm tm;
    char sbuf[SBUF_SIZE];
    char *ofmt;

    if (argc < 3 || strcmp(argv[1], "--help") == 0)
        usageErr("%s input-date-time in-format [out-format]\n", argv[0]);

    if (setlocale(LC_ALL, "") == NULL)
        errExit("setlocale"); /* Use locale settings in conversions */

    memset(&tm, 0, sizeof(struct tm)); /* Initialize 'tm' */
    if (strptime(argv[1], argv[2], &tm) == NULL)
        fatal("strptime");

    tm.tm_isdst = -1; /* Not set by strptime(); tells mktime()
                       to determine if DST is in effect */
    printf("calendar time (seconds since Epoch): %ld\n", (long) mktime(&tm));

    ofmt = (argc > 3) ? argv[3] : "%H:%M:%S %A, %d %B %Y %Z";
    if (strftime(sbuf, SBUF_SIZE, ofmt, &tm) == 0)
        fatal("strftime returned 0");
    printf("strftime() yields: %s\n", sbuf);

    exit(EXIT_SUCCESS);
}

```

10.3 时区

不同的国家（有时甚至是同一国家内的不同地区）使用不同的时区和夏时制。对于要输入和输出时间的程序来说，必须对系统所处的时区和夏时制加以考虑。所幸的是，所有这些细节都已经由 C 语言函数库包办了。

时区定义

时区信息往往是既浩繁又多变的。出于这一原因，系统没有将其直接编码于程序或函数库中，而是以标准格式保存于文件中，并加以维护。

这些文件位于目录 `/usr/share/zoneinfo` 中。该目录下的每个文件都包含了一个特定国家或地区内时区制度的相关信息，且往往根据其所描述的时区来加以命名，诸如 EST（美国东部标准时间）、CET（欧洲中部时间）、UTC、Turkey 和 Iran。此外，可以利用子目录对相关时区进行有层次的分组。例如，Pacific 目录就可能包含文件 Auckland、Port_Moresby 和 Galapagos。在程

序中指定使用的时区，实际上是指定该目录下某一时区文件的相对路径名。

系统的本地时间由时区文件/etc/localtime 定义，通常链接到/usr/share/zoneinfo 下的一个文件。

时区文件的格式记述于 tzfile(5)手册页，其创建可通过 zic(8)（时区信息编译器，zone information compiler）工具来完成。zdump(8)命令可根据指定时区文件中的时区来显示当前时间。

为程序指定时区

为运行中的程序指定一个时区，需要将 TZ 环境变量设置为由一冒号(:)和时区名称组成的字符串，其中时区名称定义于/usr/share/zoneinfo 中。设置时区会自动影响到函数 ctime()、localtime()、mktime()和 strftime()。

为了获取当前的时区设置，上述函数都会调用 tzset(3)，对如下 3 个全局变量进行了初始化：

```
char *tzname[2]; /* Name of timezone and alternate (DST) timezone */
int daylight; /* Nonzero if there is an alternate (DST) timezone */
long timezone; /* Seconds difference between UTC and local
                standard time */
```

函数 tzset()会首先检查环境变量 TZ。如果尚未设置该变量，那么就采用/etc/localtime 中定义的默认时区来初始化时区。如果 TZ 环境变量的值为空，或无法与时区文件名相匹配，那么就使用 UTC。还可将 TZDIR 环境变量（非标准的 GNU 扩展）设置为搜寻时区信息的目录名称，以替代默认的/usr/share/zoneinfo 目录。

可以通过运行程序清单 10-4 中的程序来观察 TZ 变量的影响力。第一次运行输出的是相应系统的默认时区（欧洲中部时间，CET）。在第二次运行时，由于指定的时区为 New Zealand，其在每年此时已进入夏令时，时区要比 CET 提前 12 个小时。

```
$ ./show_time
ctime() of time() value is: Tue Feb  1 10:25:56 2011
asctime() of local time is: Tue Feb  1 10:25:56 2011
strftime() of local time is: Tuesday, 01 Feb 2011, 10:25:56 CET
$ TZ="Pacific/Auckland" ./show_time
ctime() of time() value is: Tue Feb  1 22:26:19 2011
asctime() of local time is: Tue Feb  1 22:26:19 2011
strftime() of local time is: Tuesday, 01 February 2011, 22:26:19 NZDT
```

程序清单 10-4：演示时区和地区的效果

```
----- time/show_time.c
#include <time.h>
#include <locale.h>
#include "tspi_hdr.h"

#define BUF_SIZE 200

int
main(int argc, char *argv[])
{
    time_t t;
    struct tm *loc;
    char buf[BUF_SIZE];

    if (setlocale(LC_ALL, "") == NULL)
```

```

    errExit("setlocale"); /* Use locale settings in conversions */

    t = time(NULL);

    printf("ctime() of time() value is: %s", ctime(&t));

    loc = localtime(&t);
    if (loc == NULL)
        errExit("localtime");

    printf("asctime() of local time is: %s", asctime(loc));

    if (strftime(buf, BUF_SIZE, "%A, %d %B %Y, %H:%M:%S %Z", loc) == 0)
        fatal("strftime returned 0");
    printf("strftime() of local time is: %s\n", buf);

    exit(EXIT_SUCCESS);
}

```

time/show_time.c

SUSv3 为设置 TZ 环境变量定义了两个通用方法。如前所述，可将 TZ 设置为由冒号外加字符串组成的字符序列，其中的字符串用以标识时区，并随系统实现的不同而不同，通常为时区描述文件的路径名。（在采用这种形式时，Linux 和其他一些 UNIX 实现允许将冒号省略，但 SUSv3 并未规范这一行为。为了保证代码的可移植性，应当始终包含冒号。）

设置 TZ 的另一种方法在 SUSv3 中有完整的定义。使用此方法，可以将如下形式的字符串赋给 TZ：

```
std offset [ dst [ offset ] [ , start-date [ /time ] , end-date [ /time ] ] ]
```

为了便于阅读，在上面这行字符串中加入了空格，但实际上任何空格都不应出现在 TZ 中。方括号 ([]) 用来表示可选项。

std 和 **dst** 部分是用以标识标准和 DST 时区名称的字符串。例如，CET 和 CEST 分别为欧洲中部时间和欧洲中部夏令时间。各种情况下的 **offset** 分别表示欲转换为 UTC，需要叠加在本地时间上的正、负调整值。最后四部分则提供了一个规则，描述何时从标准时间变更为夏令时。

可以多种格式指定 **date**，其中之一是 **Mm.n.d**，意即：**m**(1~12)月中，第 **n**(1~5，每月的最后 **d** 天总为第 5 周) 周，星期 **d** (0=星期一，6=星期天)。如果省略 **time**，则无论何种情况下均默认为 02:00:00（上午 2 点）。

以下将 TZ 定义为 **Central Europe**，该时区的标准时间比 UTC 提前 1 小时，且 DST 始于 3 月的最后一个星期日，直至 10 月的最后一个星期日结束，提前 UTC 2 小时。

```
TZ="CET-1:00:00CEST-2:00:00,M3.5.0,M10.5.0
```

此处省略了对 DST 转换时间的指定，因为默认其发生于 02:00:00。显然，较之于如下的 Linux 专有格式，上述形式的确缺乏可读性：

```
TZ=":Europe/Berlin"
```

10.4 地区 (Locale)

世界各地在使用数千种语言，其中在计算机系统上经常使用的占了相当比例。此外，在

显示诸如数字、货币金额、日期和时间之类的信息时，不同国家的习俗也不同。例如，大多数欧洲国家使用逗号，而非小数点来分隔实数的整数和小数部分，大多数国家日期的书写格式也与美国所采用的 MM/DD/YY 格式并不相同。SUSv3 对 locale 的定义为：用户环境中依赖于语言和文化习俗的一个子集。

理想情况下，意欲在多个地理区位运行的任何程序都应处理地区 (locales) 概念，以期以用户的语言和格式来显示和输入信息。这也构成了一个相当复杂的课题——国际化 (internationalization)。在理想情况下，程序只要一次经编写，则不论运行于何处，总会自动以正确方式来执行 I/O 操作，也就是说，完成本地化 (localization) 任务。尽管存在各种支持工具，程序国际化工作依然耗时不菲。诸如 glibc 之类的程序库也提供有工具，来帮助程序支持国际化。

经常将术语 internationalization 写为 I18N，意即：I 加上 18 个字母再加 N。这一形式既便于快速书写，又避免了单词本身在英语和美语间拼写方式不同的问题。

地区定义

和时区信息一样，地区信息同样是既浩繁且多变的。出于这一原因，与其要求各个程序和函数库来存储地区信息，还不如由系统按标准格式将地区信息存储于文件中，并加以维护。

地区信息维护于 /usr/share/local (在一些发行版本中为 /usr/lib/local) 之下的目录层次结构中。该目录下的每个子目录都包含一特定地区的信息。这些目录的命名约定如下：

```
language[_territory[.codeset]][@modifier]
```

language 是双字母的 ISO 语言代码。territory 是双字母的 ISO 国家代码。codeset 表示字符编码集。modifier 则提供了一种方法，用以区分多个地区目录下 language、territory 和 codeset 均相同的状况。de_DE.utf-8@euro 是完整地区目录名称的例子之一，代表地区如下：德语，德国，UTF-8 字符编码，并采用欧元作为货币单位。

正如命名格式中的方括号所示，可以将地区目录名称中的相应部分省略。通常情况下，命名只包括语言和国家。因此，en_US 是（说英语的）美国的地区目录，而 fr_CH 则是瑞士法语区的地区目录。

这里 CH 代表 Confoederatio Helvetica，在拉丁语（本地中性语言，locally language-neutral）中意即“瑞士”。由于有 4 门官方语言，瑞士在地区上类似于跨多个时区的国家。

当在程序中指定要使用的地区时，实际上是指定了 /usr/share/locale 下某个子目录的名称。如果程序指定地区不与任何子目录名称相匹配，那么 C 语言函数库将按如下顺序将各部分从指定地区 (locale) 中剥离，以寻求匹配：

1. codeset
2. normalized codeset
3. territory
4. modifier

标准化字符编码集 (normalized codeset) 是一个特定版本字符编码集的名称，剔除了所有非字母、非数字的字符，且将所有字母转换为小写，最终字符串前冠以 ISO 三个字符。标准化的目的，在于排除字符集名称中因大小写和标点符号（例如，额外的连字符）而发生的变化。

这里是剥离过程的一个例子，假设为一程序指定的地区为 fr_CH.utf-8，但并不存在以该

名称命名的地区目录，那么如果 fr_CH 目录存在，则与之匹配。如果 fr_CH 目录也不存在，那么将采用 fr 地区目录。万一 fr 目录也不存在，那么简而言之，setLocale()函数将会报错。

/user/share/locale/locale.alias 文件定义了为程序设定地区的替代方法。详见 locale.alias(5)手册页。

每个地区子目录中包括有标准的一套文件，指定了此地区的约定设置，如表 10-2 所示。关于本表中的信息，还要注意以下几点。

- 文件 LC_COLLATE 定义了一套规则，描述了如何在一字符集内对字符排序（例如 alphabetical “按字母顺序排列的”字符集顺序）。这些规则将决定函数 strcoll(3)和 strxfrm(3)的动作。即便是同属拉丁语系的语言，其遵循的排序规则也不相同。例如，一些欧洲语言有额外字母，在某些情况下排在字母 Z 之后。另外还有特殊情况，西班牙语的双字母序列 ll，排序时位于字母 l 之后。又比如德语的元音变音字符 ä，对应于 ae，并与该双字母排在相同位置。
- 目录 LC_MESSAGES 是程序显示信息迈向国际化的步骤之一。要实现更为全面的程序信息国际化，可以采用消息目录（参考 catopen(3)和 catgets(3)手册页）或是 GNU 的 gettext API（参见 <http://www.gnu.org/>）。

Glibc 的 2.2.2 版引入了一系列非标准的地区新类别。LC_ADDRESS 定义了特定于地区的邮政地址表示规则。LC_IDENTIFICATION 指定了识别地区的信息。LC_MEASUREMENT 定义了地区的度量系统（例如，公制/英制）。LC_NAME 定义了特定于地区的人名及头衔表示规则。LC_PAPER 定义了该地区的标准纸张尺寸（例如，美国信纸/其他大多数国家所使用的 A4 纸）。LC_TELEPHONE 则定义了特定于地区的国内及国际电话号码表示规则，以及国际长途国家代码和国际拨号前缀。

表 10-2: 特定于地区的子目录内容

文件名	目的
LC_CTYPE	该文件包含字符分类（参见 isalpha(3)手册页）以及大小写转换规则
LC_COLLATE	该文件包含针对一字符集的排序规则
LC_MONETARY	该文件包含对币值的格式化规则（见 localeconv(3)和<locale.h>）
LC_NUMERIC	该文件包含对币值以外数字的格式化规则（见 localeconv(3)和<locale.h>）
LC_TIME	该文件包含对日期和时间的格式化规则
LC_MESSAGES	该目录下所含文件，针对肯定和否定（是/否）响应，就格式及数值做了规定

系统中实际定义的地区可能会各有不同。除了必须定义一个名为 POSIX（与 C 同义，后者的存在是由于历史原因）的标准地区外，SUSv3 没有对此作出任何要求。POSIX 折射出 UNIX 系统的历史渊源。因之，系统建立于 ASCII 字符集之上，使用英文来描述日期，并以“yes/no”来响应。该地区的货币和数字格式则处于未定义状态。

locale 命令显示当前地区环境（本 shell 内）的相关信息。命令 locale -a 则将列出系统上定义的整套地区。

为程序设置地区

函数 `setlocale()` 既可设置也可查询程序的当前地区。

```
#include <locale.h>

char *setlocale(int category, const char *locale);

Returns pointer to a (usually statically allocated) string identifying
the new or current locale on success, or NULL on error
```

`category` 参数选择设置或查询地区的哪一部分，它仅能使用表 10-2 中列出的地区类别的常量名称。因此，它可以设置地区的时间显示格式是德国，而地区的货币符号是美元。或者，更常见的是，我们可以利用 `LC_ALL` 来指定我们要设置的地区的所有部分的值。

使用 `setLocale()` 设置地区有两种不同的方法。`locale` 参数可能是一个字符串，指定系统上已定义的一个地区（例如，`/usr /lib /locale` 中的子目录的名称），如 `de_DE` 或 `en_US`。另外，地区可能被指定为空字符串，这意味着从环境变量取得地区的设置。

```
setlocale(LC_ALL, "");
```

我们必须这样调用才能使程序使用环境变量中的地区。如果调用被省略，这些环境变量将不会对程序生效。

当运行程序调用了 `setLocale(LC_ALL, "")`，我们能够使用一系列环境变量来控制地区的各部分内容，环境变量的名称也是对应于表 10-2 中列出的类型：`LC_CTYPE`、`LC_COLLATE`、`LC_MONETARY`、`LC_NUMERIC`、`LC_TIME`、`LC_MESSAGES`。另外，我们可以使用 `LC_ALL` 或 `LANG` 环境变量指定整个地区的设置。如果设置了多个先前的环境变量，那么 `LC_ALL` 会覆盖所有其他的 `LC_*` 环境变量，同时 `LANG` 的优先级最低。因此，通常使用 `LANG` 为地区所有内容设置默认值，然后用单独的 `LC_*` 变量，设置地区的各个方面内容来覆盖默认值。

最后，`setLocale()` 返回一个指针指向标识这一类地区设置的字符串（通常是静态分配的）。如果我们仅需要查看地区的设置而不需要改变它，那么我们可以指定 `locale` 参数为 `NULL`。

地区设置影响众多 GNU/ Linux 实用程序，以及 `glibc` 的许多函数的功能。其中有函数 `strftime()` 和 `strptime()`（10.2.3 节），当我们在不同的地区运行程序清单 10-4，`strftime` 返回的结果如下：

```
$ LANG=de_DE ./show_time German locale
ctime() of time() value is: Tue Feb  1 12:23:39 2011
asctime() of local time is: Tue Feb  1 12:23:39 2011
strftime() of local time is: Dienstag, 01 Februar 2011, 12:23:39 CET
```

下一个运行演示 `LC_TIME` 比 `LANG` 的优先级高：

```
$ LANG=de_DE LC_TIME=it_IT ./show_time German and Italian locales
ctime() of time() value is: Tue Feb  1 12:24:03 2011
asctime() of local time is: Tue Feb  1 12:24:03 2011
strftime() of local time is: martedì, 01 febbraio 2011, 12:24:03 CET
```

而这个运行结果表明，LC_ALL 超过 LC_TIME 的优先级：

```
$ LC_ALL=fr_FR LC_TIME=en_US ./show_time          French and US locales
ctime() of time() value is: Tue Feb  1 12:25:38 2011
asctime() of local time is: Tue Feb  1 12:25:38 2011
strftime() of local time is: mardi, 01 février 2011, 12:25:38 CET
```

10.5 更新系统时钟

我们现在来看两个更新系统时钟的接口：settimeofday()和 adjtime()。这些接口都很少被应用程序使用，因为系统时间通常是由工具软件维护，如网络时间协议（Network Time Protocol）守护进程，并且它们需要调用者已被授权（CAP_SYS_TIME）。

系统调用 settimeofday()是 gettimeofday()的逆向操作（这是我们在 10.1 节中描述的）。它将 tv 指向 timeval 结构体里的秒数和微秒数，设置到系统的日历时间。

```
#define _BSD_SOURCE
#include <sys/time.h>

int settimeofday(const struct timeval *tv, const struct timezone *tz);

Returns 0 on success, or -1 on error
```

和函数 gettimeofday()一样，tz 参数已被废弃，**这个参数应该始终指定为 NULL**。

tv.tv_usec 字段的微秒精度并不意味着我们以微秒精度来设置系统时钟，因为时钟的精度可能会低于微秒。

虽然 SUSv3 没有定义 settimeofday()，但它在其他 UNIX 实现中被广泛使用。

Linux 还提供了 stime()系统调用来设置系统时钟。settimeofday()和 stime()之间的区别是，后者调用允许使用秒的精度来表示新的日历时间。和函数 time()与 gettimeofday()相同，stime()和 settimeofday()的并存是由历史原因造成的：拥有更高精确度的后一个函数，是由 4.3BSD 添加的。

settimeofday()调用所造成的那种系统时间的突然变化，可能会对依赖于系统时钟单调递增的应用造成有害的影响（例如，make(1)，数据库系统使用的时间戳或包含时间戳记的日志文件）。出于这个原因，当**对时间做微小调整时（几秒钟误差），通常是推荐使用库函数 adjtime()，它将系统时钟逐步调整到正确的时间。**

```
#define _BSD_SOURCE
#include <sys/time.h>

int adjtime(struct timeval *delta, struct timeval *olddelta);

Returns 0 on success, or -1 on error
```

delta 参数指向一个 timeval 结构体，指定需要改变时间的秒和微秒数。如果这个值是正数，那么每秒系统时间都会额外拨快一点点，直到增加完所需的时间。如果 delta 值为负时，时钟以类似的方式减慢。

Linux/x86-32 以每 2000 秒变化 1 秒（或每天 43.2 秒）的频率调整时钟。

在 `adjtime()` 函数执行的时间里，它可能无法完成时钟调整。在这种情况下，**剩余未经调整的时间存放在 `olddelta` 指向的 `timeval` 结构体内**。如果我们不关心这个值，我们可以指定 `olddelta` 为 `NULL`。相反，如果我们只关心当前未完成时间校正的信息，而并不想改变它，我们可以指定 `delta` 参数为 `NULL`。

虽然 SUSv3 未定义 `adjtime()`，可大多数 UNIX 实现提供了这个函数。

`adjtime()` 在 Linux 上，基于更通用和复杂的特定于 Linux 的系统调用 `adjtimex()` 来完成功能。这个系统调用也同时被网络时间协议 (NTP) 守护进程调用。如需进一步信息，请参阅 Linux 的源代码，Linux `adjtimex(2)` 帮助手册页和 NTP 规范 ([Mills, 1992])。

10.6 软件时钟 (jiffies)

在本书中所描述的时间相关的各种系统调用的精度是受限于系统软件时钟 (software clock) 的分辨率，它的度量单位被称为 `jiffies`。`jiffies` 的大小是定义在内核源代码的常量 `HZ`。这是内核按照 `round-robin` 的**分时调度算法 (35.1 节) 分配 CPU 进程的单位**。

在 2.4 或以上版本的 Linux/x86-32 内核中，软件时钟速度是 100 赫兹，也就是说，一个 `jiffy` 是 10 毫秒。

自 Linux 面世以来，由于 CPU 的速度已大大增加，Linux / x86-32 2.6.0 内核的软件时钟速度已经提高到 1000 赫兹。更高的软件时钟速率意味着定时器可以有更高的操作精度和时间可以拥有更高的测量精度。然而，这并非可以任意提高时钟频率，因为每个时钟中断会消耗少量的 CPU 时间，这部分时间 CPU 无法执行任何操作。

经过内核开发人员之间的的讨论，最终导致软件时钟频率成为一个可配置的内核的选项（包括处理器类型和特性，定时器的频率）。自 2.6.13 内核，时钟频率可以设置到 100、250（默认）或 1000 赫兹，对应的 `jiffy` 值分别为 10、4、1 毫秒。自内核 2.6.20，增加了一个频率：300 赫兹，它可以被两种常见的视频帧速率 25 帧每秒 (PAL) 和 30 帧每秒 (NTSC) 整除。

10.7 进程时间

进程时间是进程创建后使用的 CPU 时间数量。出于记录的目的，内核把 CPU 时间分成以下两部分。

- 用户 CPU 时间是在用户模式下执行所花费的时间数量。有时也称为虚拟时间 (virtual time)，这对于程序来说，是它已经得到 CPU 的时间。
- 系统 CPU 时间是在内核模式中执行所花费的时间数量。这是内核用于执行系统调用或代表程序执行的其他任务（例如，服务页错误）的时间。

有时候，进程时间是指处理过程中所消耗的总 CPU 时间。

当我们运行一个 `shell` 程序，我们可以使用的 `time(1)` 命令，同时获得这两个部分的时间值，以及运行程序所需的实际时间。

```
$ time ./myprog
real    0m4.84s
user    0m1.030s
sys     0m3.43s
```

系统调用 `times()`，检索进程时间信息，并把结果通过 `buf` 指向的结构体返回。

```
#include <sys/times.h>

clock_t times(struct tms *buf);

Returns number of clock ticks (sysconf(_SC_CLK_TCK)) since
“arbitrary” time in past on success, or (clock_t) -1 on error
```

`buf` 指向的 TMS 结构体有下列格式：

```
struct tms {
    clock_t tms_utime; /* User CPU time used by caller */
    clock_t tms_stime; /* System CPU time used by caller */
    clock_t tms_cutime; /* User CPU time of all (waited for) children */
    clock_t tms_cstime; /* System CPU time of all (waited for) children */
};
```

`tms` 结构体的前两个字段返回调用进程到目前为止使用的用户和系统组件的 CPU 时间。最后两个字段返回的信息是：父进程（比如，`times()`的调用者）执行了系统调用 `wait()`的所有已经终止的子进程使用的 CPU 时间。

数据类型 `clock_t` 是用时钟计时单元（clock tick）为单位度量时间的整型值，习惯用于计算 `tms` 结构体的 4 个字段。我们可以调用 `sysconf(_SC_CLK_TCK)`来获得每秒包含的时钟计时单元数，然后用这个数字除以 `clock_t` 转换为秒。（我们在 11.2 节叙述 `sysconf()`。）

在大多数 Linux 的硬件架构，`sysconf(_SC_CLK_TCK)`返回 100。与此对应的内核常量是 `USER_HZ`。然而 `USER_HZ` 在其他几个架构下可以被定义超过 100，如 Alpha 和 IA - 64。

如果成功，`times()`返回自过去的任意点流逝的以时钟计时单元为单位的（真实的）时间。SUSv3 特别未定义这点是什么，只是说，这将在是在调用进程的生命周期内的一个固定点。因此，这个返回值唯一的用法是通过计算一对 `times()`调用返回的值的差，来计算进程执行消耗的时间。然而，即使是这种用法，`times()`的返回值仍然不可靠的，因为它可能会溢出 `clock_t`的有效范围，这时 `times()`的返回值将再次从 0 开始计算（也就是说，一个稍后的 `times()`的调用返回的数值可能会低于一个更早的 `times()`调用）。可靠的测量经过时间的方法是使用函数 `gettimeofday()`（10.1 节所述）。

在 Linux 上，我们可以指定 `buf` 参数为 `NULL`。在这种情况下，`times()`只是简单地返回一个函数结果。然而，这是没有意义的。SUSv3 并未定义 `buf` 可以使用 `NULL`，因此许多其他 UNIX 实现需要这个参数必须为一个非 `NULL` 值。

函数 `clock()`提供了一个简单的接口用于取得进程时间。它返回一个值描述了调用进程使用的总的 CPU 时间（包括用户和系统）。

```
#include <time.h>

clock_t clock(void);

Returns total CPU time used by calling process measured in
CLOCKS_PER_SEC, or (clock_t) -1 on error
```

time()的返回值的计量单位是 CLOCKS_PER_SEC，所以我们必须除以这个值来获得进程所使用的 CPU 时间秒数。在 POSIX.1，CLOCKS_PER_SEC 是常量 10000，无论底层软件时钟（10.6 节）的分辨率是多少。clock()的精度最终仍然受限于软件时钟的分辨率。

虽然 clock()和 times()返回相同的数据类型 clock_t，这两个接口使用的测量单位却并不相同。这是历史原因造成了 clock_t 定义的冲突，一个是 POSIX.1 标准，而另一个是 C 编程语言标准。

即使 CLOCKS_PER_SEC 是常量 10000，SUSv3 注明，这个常量在不兼容 XSI (non-XSI-conformant)的系统上可以为整型变量，所以，我们不能简单地把它作为一个编译时常量（即，我们不能使用 #ifdef 预处理表达式）。它可能会被定义为一个长整数（即 1000000L），我们总是将这个常量转换为 long，因此我们可以简单地用 printf() 把它打印输出（见 3.6.2 节）。

SUSv3 描述 clock()应该返回“进程所使用的处理器时间”时有不同的解释。在一些 UNIX 的实现中，clock()返回的时间包含所有等待子进程使用的 CPU 时间。而在 Linux 上，它不包括。

示例程序

在程序清单 10-5 中的程序演示了如何使用本节中描述的功能。函数 displayProcessTimes() 首先打印由调用者提供的信息，然后使用 clock()和 times()来获得和显示进程时间。主程序首先调用函数 displayProcessTimes()，然后执行一个循环，通过重复调用 getppid()消耗一些 CPU 时间，再次调用 displayProcessTimes()来查看这个循环会消耗多少 CPU 时间。当我们使用这个程序调用 getppid()十万次，这就是我们看到的：

```
$ ./process_time 10000000
CLOCKS_PER_SEC=1000000 sysconf(_SC_CLK_TCK)=100

At program start:
  clock() returns: 0 clocks-per-sec (0.00 secs)
  times() yields: user CPU=0.00; system CPU: 0.00
After getppid() loop:
  clock() returns: 2960000 clocks-per-sec (2.96 secs)
  times() yields: user CPU=1.09; system CPU: 1.87
```

程序清单 10-5：获取进程 CPU 时间

```
time/process_time.c

#include <sys/times.h>
#include <time.h>
#include "tlpi_hdr.h"

static void          /* Display 'msg' and process times */
displayProcessTimes(const char *msg)
{
    struct tms t;
    clock_t clockTime;
    static long clockTicks = 0;

    if (msg != NULL)
        printf("%s", msg);
```

```

if (clockTicks == 0) {          /* Fetch clock ticks on first call */
    clockTicks = sysconf(_SC_CLK_TCK);
    if (clockTicks == -1)
        errExit("sysconf");
}

clockTime = clock();
if (clockTime == -1)
    errExit("clock");

printf("        clock() returns: %ld clocks-per-sec (%.2f secs)\n",
       (long) clockTime, (double) clockTime / CLOCKS_PER_SEC);

if (times(&t) == -1)
    errExit("times");
printf("        times() yields: user CPU=%.2f; system CPU: %.2f\n",
       (double) t.tms_utime / clockTicks,
       (double) t.tms_stime / clockTicks);
}

int
main(int argc, char *argv[])
{
    int numCalls, j;

    printf("CLOCKS_PER_SEC=%ld sysconf(_SC_CLK_TCK)=%ld\n\n",
           (long) CLOCKS_PER_SEC, sysconf(_SC_CLK_TCK));

    displayProcessTimes("At program start:\n");

    numCalls = (argc > 1) ? getInt(argv[1], GN_GT_0, "num-calls") : 100000000;
    for (j = 0; j < numCalls; j++)
        (void) getppid();

    displayProcessTimes("After getppid() loop:\n");

    exit(EXIT_SUCCESS);
}

```

time/process_time.c

10.8 总结

真实时间对应于时间定义的每一天。当真实时间通过一些标准点计算的时候，我们称它为日历时间。和经过的时间相对，它是度量一个进程生命周期中的一些点（通常是开始）。

进程时间是由一个进程使用的 CPU 时间量，并划分为用户时间和系统时间。

多种系统调用允许我们获取和设置系统时钟值（即日历时间，以秒为单位从 Epoch 计算），以及一系列的库函数能够完成从日历时间到其他时间格式之间的转换，包括分解时间和具有可读性字符串。描述这种转换把我们引入了地区和国际化的讨论。

使用和显示时间和日期是许多应用程序的一个重要组成部分，我们会在这本书后面的章节中经常使用到本节描述的功能。我们也会在第 23 章更多地介绍时间的度量。

进一步的信息

关于 Linux 内核如何度量时间的详细信息可以参考[Love, 2010]。

关于时区和国际化的进一步讨论，可以参考 GNU C 库手册（在线地址：<http://www.gnu.org/>）。SUSv3 文档也包括地区的详细描述。

10.9 练习

- 10-1.** 假设一个系统调用 `sysconf(_SC_CLK_TCK)` 返回的值是 100。假设 `times()` 返回 `clock_t` 的值是一个无符号的 32 位整数，需要多久这个值才能进入下一个从 0 开始的周期？对 `clock()` 返回的 `CLOCKS_PER_SEC` 值执行相同的计算。

第 11 章

系统限制和选项

但凡 UNIX 实现，无不对各种系统特性和资源加以限制，并提供（或者选择不提供）由各种标准所定义的选项，例如：

- 一个进程能同时拥有多少已打开的文件？
- 系统是否支持实时信号？
- `int` 类型变量可存储的最大值是多少？
- 一个程序的参数列表能有多大？
- 路径名的最大长度是多少？

尽管可以把假定的限制和选项硬性写入程序编码，但这将破坏程序的可移植性，因为限制和选项可能会有所不同。

- 在 UNIX 实现之间：虽然限制和选项在某个特定 UNIX 实现中可能是固定的，但在不同的 UNIX 实现之间，可能会有所不同。`int` 变量可存储的最大值就是此类限制的例子之一。
- 特定实现的运行环境：例如，可能重新配置了内核，改变了某个限制。又或者，在某个系统上编译的应用程序，却在另一个限制和选项有所不同的系统中运行。
- 从一个文件系统到另外一个文件系统：例如，传统的 System V 文件系统允许文件名长达 14 个字节，而传统的 BSD 文件系统和大多数“原生”Linux 文件系统则允许文件名高达 255 个字节。

因为系统限制和选项会影响应用程序的行为，所以可移植应用程序需要获取限制值，弄清系统对选项的支持情况。C 语言标准和 SUSv3 为此而提供了两种重要途径。

- 在编译程序时能够获得一些限制和选项。例如，`int` 类型的最大值取决于硬件结构和编译器的设计选择。此类限制可在头文件中记录。
- 另外一些限制和选项在程序运行时可能会有变化。对此，SUSv3 定义了 3 个函数 `sysconf()`、`pathconf()` 和 `fpathconf()`，供应用程序调用以检查系统实现的限制和选项。

SUSv3 规定有一系列限制，要求符合规范的实现必须支持，同时还规定了一套选项，特定系统可以有选择地对其中各个选项予以支持。本章介绍了部分限制和选项，其余则会在后续章节中适时加以描述。

11.1 系统限制

SUSv3 要求，针对其所规范的每个限制，所有实现都必须支持一个最小值。在大多数情况下，会将这些最小值定义为<limits.h>文件中的常量，其命名则冠以字符串 `_POSIX_`，而且（通常）还包含字符串 `_MAX`，因此，常量命名形如 `_POSIX_XXX_MAX`。

如果应用程序将本身限制在 SUSv3 对每个限制所要求的最小值之内，那么该程序对符合标准的所有实现都具有可移植性。然而，这一做法阻碍了应用程序去利用特定实现可提供的更高限制。因此，在特定系统上获取限制，通常更为可取的方法是使用<limits.h>文件、`sysconf()`或 `pathconf()`。

SUSv3 将其所定义的各类限制描述为最小值，但命名却使用了字符串 `_MAX`，这可能颇令人疑惑。换一种思路，将此类常量中的每一个都视为对某类资源或特性的上限，且标准要求这些上限都必须拥有一个确定的最小值，这种命名的用意也就不言自明了。

在某些情况下，会为某个限制提供最大值，并且在对这些值的命名中包含字符串 `_MIN`。对于这些常量，道理正好反过来；它们代表了对某些资源的下限，按照标准规定，在符合标准的实现中，该下限不能高于某个值。例如，限制 `FLT_MIN(1E-37)`为某个实现中所能表征的最小浮点数定义了最大值。所有满足标准的实现至少能够表征如此之小的浮点数。

每个限制都有一个名称，与上述最小值的名称相对应，但缺少了 `_POSIX_`前缀。某个实现可以在<limits.h>文件中以该名称定义一个常量，用以表示该实现的相应限制。若已然定义，则该限制值总是至少等同于前述最大值（即 `XXX_MAX >= _POSIX_XXX_MAX`）。

SUSv3 将其规定的限制归为 3 类：运行时恒定值、路径名变量值和运行时可增加值。在下列段落中将描述这些类别并提供一些例子。

运行时恒定值（可能不确定）

所谓运行时恒定值是指某一限制，若已然在<limits.h>文件中定义，则对于实现而言固定不变。然而该值可能是不确定的（因为该值可能依赖于可用的内存空间），因而在<limits.h>文件中会忽略对其定义。在这种情况下（即使在<limits.h>文件中已然定义了该限制），应用程序可以使用 `sysconf()`来获取运行时的值。

`MQ_PRIO_MAX` 限制就是运行时恒定值的例子之一。正如 52.5.1 节所述，针对 POSIX 消息队列中的消息，存在着优先级方面的限制。SUSv3 定义了值为 32 的常量 `_POSIX_MQ_PRIO_MAX`，将其作为符合规范的实现为该限制所必须提供的最小值。这意味着，所有符合规范的实现，其对消息优先级的支持至少应为从 0~31。一个 UNIX 实现可以为此限制设定更高值，并将该值在<limits.h>文件中以常量 `MQ_PRIO_MAX` 加以定义。例如，Linux 就将 `MQ_PRIO_MAX` 的值定义为 32768。也可以通过下列调用在运行时获取该值：

```
lim = sysconf(_SC_MQ_PRIO_MAX);
```

路径名变量值

所谓路径名变量值是指与路径名（文件、目录、终端等）相关的限制，每个限制可能是相对于某个系统实现的常量，也可能随文件系统的不同而不同。在限制可能因路径名而发生变化的情况下，应用程序可以使用 `pathconf()`或 `fpathconf()`来获取该值。

`NAME_MAX` 限制是路径名变量值的例子之一。此限制定义了在一个特定文件系统中文

件名的最大长度。SUSv3 定义了值为 14（老版本的 System V 文件系统限制）的常量 `_POSIX_NAME_MAX`，作为系统实现必须支持的最小限制值。系统实现可以定义一个高于此值的 `NAME_MAX` 限制，并/或向应用开放如下形式的调用，以获取特定文件系统的相关信息：

```
lim = pathconf(directory_path, _PC_NAME_MAX)
```

参数 `directory_path` 是目标文件系统上的目录路径名。

运行时可增加值

运行时可增加值是指某一限制，相对于特定实现其值固定，且运行此实现的所有系统至少都应支持这一最小值。然而，特定系统在运行时可能会增加该值，应用程序可以使用 `sysconf()` 来获得系统所支持的实际值。

运行时可增加值的例子之一是 `NGROUPS_MAX`，该限制定义了一进程可同时从属的辅助组 ID（9.6 节）的最大数量。SUSv3 定义了相应的最小值 `_POSIX_NGROUPS_MAX`，其值为 8。应用可在运行时通过调用 `sysconf(_SC_NGROUPS_MAX)` 来获取此限制值。

对选定 SUSv3 限制的总结

表 11-1 列举了与本书有关，由 SUSv3 所定义的部分限制（其他限制将在后续章节中加以介绍）。

表 11-1：选定的 SUSv3 限制

限制名称 (<code><limits.h></code>)	最小值	在 <code>sysconf()</code> / <code>pathconf()</code> 中的参数命名(<code><unistd.h></code>)	描 述
<code>ARG_MAX</code>	4096	<code>_SC_ARG_MAX</code>	提供给 <code>exec()</code> 的参数(<code>argv</code>)与环境变量(<code>environ</code>)所占存储空间之和的最大字节数（见 6.7 节和 27.2.3 节）
<code>none</code>	<code>none</code>	<code>_SC_CLK_TCK</code>	为 <code>times()</code> 提供的度量单位
<code>LOGIN_NAME_MAX</code>	9	<code>_SC_LOGIN_NAME_MAX</code>	登录名的最大长度（含终止空字符）
<code>OPEN_MAX</code>	20	<code>_SC_OPEN_MAX</code>	进程同时可打开的文件描述符的最大数量，比可用文件描述符的最大数量多 1 个（见 36.2 节）
<code>NGROUPS_MAX</code>	8	<code>_SC_NGROUPS_MAX</code>	进程所属辅助组 ID 数量的最大值（见 9.7.3 节）
<code>none</code>	1	<code>_SC_PAGESIZE</code>	一个虚拟内存页的大小 (<code>_SC_PAGE_SIZE</code> 与其同义)
<code>RTSIG_MAX</code>	8	<code>_SC_RTSIG_MAX</code>	单一实时信号的最大数量（见 22.8 节）
<code>SIGQUEUE_MAX</code>	32	<code>_SC_SIGQUEUE_MAX</code>	排队实时信号的最大数量（见 22.8 节）
<code>STREAM_MAX</code>	8	<code>_SC_STREAM_MAX</code>	同时可打开的 <code>stdio</code> 流的最大数量
<code>NAME_MAX</code>	14	<code>_PC_NAME_MAX</code>	排除终止空字符外，文件名称可达的最大字节长度
<code>PATH_MAX</code>	256	<code>_PC_PATH_MAX</code>	路径名称可达的最大字节长度，含尾部空字符
<code>PIPE_BUF</code>	512	<code>_PC_PIPE_BUF</code>	一次性（原子操作）写入管道或 FIFO 中的最大字节数（44.1 节）

表 11-1 中的第一列给出了限制的名称，可将其作为常量定义于<limits.h>文件中，用于表示特定实现下的限制。第二列是 SUSv3 为这些限制所定义的最小值（也定义于<limits.h>文件中）。在大多数情况下，会将每个限制的最小值定义为冠以字符串_POSIX_的常量。例如，常量_POSIX_RTSIG_MAX（其值为 8）为 SUSv3 实现对相应 RTSIG_MAX 常量的最低要求。第三列列出了为在运行期间获取实现的限制，调用 sysconf()或 pathconf()时应输入入参的常量名。冠以_SC_的常量用于 sysconf()，冠以_PC_的常量用于 pathconf()和 fpathconf()。

下列为对表 11-1 的补充信息，请予关注。

- getdtablesize()函数是确定进程文件描述符（OPEN_MAX）限制的备选方法，已遭弃用，该函数曾一度为 SUSv2 所定义（标记为 LEGACY），但 SUSv3 将其剔除。
- getpagesize()函数是确定系统页大小（_SC_PAGESIZE）的备选方法，已然废弃。该函数一度曾为 SUSv2 所定义（标记为 LEGACY），但 SUSv3 将其剔除。
- 定义于<stdio.h>文件中的常量 FOPEN_MAX，等同于常量 STREAM_MAX。
- NAME_MAX 不包含终止空字符，而 PATH_MAX 则包括。POSIX.1 标准在定义 PATH_MAX 时，对于是否包含终止空字符始终含糊不清，而上述差异则恰好弥补了这一缺陷。定义 PATH_MAX 中包含终止符也意味着，为路径名称分配了 PATH_MAX 个字节的应用程序依然符合标准。

从 shell 中获取限制和选项：getconf

在 shell 中，可以使用 getconf 命令获取特定 UNIX 系统中已然实现的限制和选项。该命令的格式一般如下：

```
$ getconf variable-name [ pathname ]
```

variable-name 标识用户意欲获取的限制，应是符合 SUSV3 标准的限制名称，例如：ARG_MAX 或 NAME_MAX。但凡限制与路径名相关，则还需要指定一个路径名，作为命令的第二个参数，如下第二个实例所示。

```
$ getconf ARG_MAX
131072
$ getconf NAME_MAX /boot
255
```

11.2 在运行时获取系统限制（和选项）

sysconf()函数允许应用程序在运行时获得系统限制值。

```
#include <unistd.h>

long sysconf(int name);
```

Returns value of limit specified by *name*,
or -1 if limit is indeterminate or an error occurred

参数 name 应为定义于<unistd.h>文件中的_SC_系列常量之一，其中部分在表 11-1 中已有所罗列。限制值将作为函数结果返回。

若无法确定某一限制，则 sysconf()返回-1。若调用 sysconf()函数时发生错误，也会返回-1。（唯一指定的错误是 EINVAL，表示 name 无效。）为区别上述两种情况，必须在调用函数

前将 `errno` 设置为 0，如果调用返回-1，且调用后 `errno` 值不为 0，那么调用 `sysconf()`函数时发生了错误。

由 `sysconf()`函数所返回的限制值类型总是（长）整型（`pathconf()`和 `fpathconf()`也是如此）。在对 `sysconf()`函数的原理描述中，SUSv3 特意指出，一度曾考虑将字符串作为可能的返回值，但由于实现和使用的复杂性而最终放弃了这一构想。

程序清单 11-1 所示为调用 `sysconf()`来展示各种系统限制。在某一 Linux 2.6.31/x86-32 系统上运行该程序，将产生如下结果：

```
$ ./t_sysconf
_SC_ARG_MAX:      2097152
_SC_LOGIN_NAME_MAX: 256
_SC_OPEN_MAX:     1024
_SC_NGROUPS_MAX:  65536
_SC_PAGESIZE:     4096
_SC_RTSIG_MAX:    32
```

程序清单 11-1：使用 `sysconf()`函数

```
----- syslim/t_sysconf.c
#include "tspi_hdr.h"

static void          /* Print 'msg' plus sysconf() value for 'name' */
sysconfPrint(const char *msg, int name)
{
    long lim;

    errno = 0;
    lim = sysconf(name);
    if (lim != -1) {      /* Call succeeded, limit determinate */
        printf("%s %ld\n", msg, lim);
    } else {
        if (errno == 0)  /* Call succeeded, limit indeterminate */
            printf("%s (indeterminate)\n", msg);
        else            /* Call failed */
            errExit("sysconf %s", msg);
    }
}

int
main(int argc, char *argv[])
{
    sysconfPrint("_SC_ARG_MAX:      ", _SC_ARG_MAX);
    sysconfPrint("_SC_LOGIN_NAME_MAX: ", _SC_LOGIN_NAME_MAX);
    sysconfPrint("_SC_OPEN_MAX:     ", _SC_OPEN_MAX);
    sysconfPrint("_SC_NGROUPS_MAX:  ", _SC_NGROUPS_MAX);
    sysconfPrint("_SC_PAGESIZE:     ", _SC_PAGESIZE);
    sysconfPrint("_SC_RTSIG_MAX:    ", _SC_RTSIG_MAX);
    exit(EXIT_SUCCESS);
}
----- syslim/t_sysconf.c
```

SUSv3 要求，针对特定限制，调用 `sysconf()`所获取的值在调用进程的生命周期内应保持不变。例如，就可以这样认定：针对 `_SC_PAGESIZE` 限制的返回值在进程运行期间不会改变。

在 Linux 系统中，对于上述要求，有一些（合理的）例外。进程能够使用 `setrlimit()`（见 36.2 节）修改进程的各种资源限制，这会波及由 `sysconf()` 所报告的限制值：`RLIMIT_NOFILE`，该限制确定进程能够打开的文件数量（`_SC_OPEN_MAX`）；`RLIMIT_NPROC`（实际并未纳入 SUSv3 中），即允许进程基于每用户所创建的子进程限额（`_SC_CHILD_MAX`）；`RLIMIT_STACK`，始于 Linux 2.6.23 版本，该限制确定了进程的命令行参数和环境变量所占存储空间的限额（`_SC_ARG_MAX`，具体参见 `execve(2)` 手册页）。

11.3 运行时获取与文件相关的限制（和选项）

`pathconf()` 和 `fpathconf()` 函数允许应用程序在运行时获取文件相关的限制值。

```
#include <unistd.h>

long pathconf(const char *pathname, int name);
long fpathconf(int fd, int name);

Both return value of limit specified by name,
or -1 if limit is indeterminate or an error occurred
```

`pathconf()` 和 `fpathconf()` 之间唯一的区别在于对文件或目录的指定方式。`pathconf()` 采用路径名方式来指定，而 `fpathconf()` 则使用（之前已经打开的）文件描述符。

参数 `name` 则是定义于 `<unistd.h>` 文件中的 `_PC_*` 系列常量之一，在表 11-1 中已经列举了其中的一部分。表 11-2 又针对表 11-1 中展示的 `_PC_*` 常量，提供了更深入的细节。

限制的值将作为函数结果返回。如要区分限制值不确定与发生错误的情况，应对方式与 `sysconf()` 相同。

有别于 `sysconf()` 函数，SUSv3 并不要求 `pathconf()` 和 `fpathconf()` 的返回值在进程的生命周期内保持恒定。这是因为，例如，在进程运行期间，可能会卸载一个文件系统，然后再以不同特性重新装载该文件系统。

表 11-2: `pathconf()` 函数中，选定 `_PC_*` 系列命名的详细说明

常 量	说 明
<code>_PC_NAME_MAX</code>	针对目录，返回该目录下文件命名的最大长度，对于其他文件类型，则未作规定
<code>_PC_PATH_MAX</code>	对于目录，返回该目录中相对路径名的最大长度，对于其他文件类型，则未作规定
<code>_PC_PIPE_BUF</code>	对于 FIFO 或者管道，返回一个应用于引用文件的值。对于目录，返回的值应用于在该目录下创建的一 FIFO。对于其他文件类型，则未作规定

程序清单 11-2 所示为针对由标准输入所指向的文件，使用 `fpathconf()` 函数获取各种限制。运行该程序时，若将 `ext2` 文件系统上的某一目录指定为标准输入，可产生如下结果：

```
$ ./t_fpathconf < .
_PC_NAME_MAX: 255
_PC_PATH_MAX: 4096
_PC_PIPE_BUF: 4096
```

```

_____  

syslim/t_fpathconf.c  

#include "tspi_hdr.h"

static void          /* Print 'msg' plus value of fpathconf(fd, name) */
fpathconfPrint(const char *msg, int fd, int name)
{
    long lim;

    errno = 0;
    lim = fpathconf(fd, name);
    if (lim != -1) {          /* Call succeeded, limit determinate */
        printf("%s %ld\n", msg, lim);
    } else {
        if (errno == 0)      /* Call succeeded, limit indeterminate */
            printf("%s (indeterminate)\n", msg);
        else                /* Call failed */
            errExit("fpathconf %s", msg);
    }
}

int
main(int argc, char *argv[])
{
    fpathconfPrint("_PC_NAME_MAX: ", STDIN_FILENO, _PC_NAME_MAX);
    fpathconfPrint("_PC_PATH_MAX: ", STDIN_FILENO, _PC_PATH_MAX);
    fpathconfPrint("_PC_PIPE_BUF: ", STDIN_FILENO, _PC_PIPE_BUF);
    exit(EXIT_SUCCESS);
}
_____  

syslim/t_fpathconf.c

```

11.4 不确定的限制

有时, 系统实现并未将一些系统限制定义为限制常量(比如: `PATH_MAX`), 并且 `sysconf()` 或 `pathconf()` 在返回相应限制(比如 `_PC_PATH_MAX`)时会将其归为不确定。对此, 可采用如下策略之一。

- 当编写一个可在多个 UNIX 实现间移植的应用程序时, 可选择使用 SUSv3 所规定的最低限制值。此类以 `_POSIX*_MAX` 形式命名的常量, 具体参见 11.1 节。此方法有时并不可行, 因为该限制之低已经超乎实际情况, 正如 `_POSIX_PATH_MAX` 和 `_POSIX_OPEN_MAX` 的情况。
- 在某些情况下, 切实可行的解决方法是省去对限制的检查, 取而代之以执行相关的系统调用或库函数。(类似观点也适用于 11.5 节中所描述的一些 SUSv3 选项。)如果调用失败, 且 `errno` 表明出错是由于超出了系统限制时, 那么可以根据需要调整应用的行为, 并再次尝试调用。例如, 对于可发送给进程的实时信号队列长度, 大多数 UNIX 实现都进行了强制限制。一旦达到限额, 试图进一步发送信号(使用 `sigqueue()` 函数)将以失败告终, 且会将错误号 `errno` 置为 `EAGAIN`。这时, 发送进程只需简单重试即可, 或许是在等待片刻之后。与之相类似, 试图打开一个文件时, 若文件名过长, 将会产生 `ENAMETOOLONG` 错误, 之后应用程序可以一个更加简短的命名进行重试。
- 自行编写程序或函数, 以推断或估算限制值。无论在哪一种情况下, 都会调用相关的

sysconf()或 pathconf(), 若限制不确定, 则函数将返回一合理估值。虽然有欠完美, 但这种解决方案往往在实践中是可行的。

- 也可以利用诸如 GNU Autoconf 之类的扩展工具, 该工具能够确定各种系统特性及限制存在与否、如何设置。Autoconf 程序可基于其收集到的信息而生成头文件, 并能在 C 程序中将其包含¹在内。关于 Autoconf 的更多信息, 请参考 <http://www.gnu.org/software/autoconf/>中的内容。

11.5 系统选项

除了对各种系统资源的限制加以规范外, SUSv3 还规定了 UNIX 实现可支持的各种选项。这包括对诸如实时信号、POSIX 共享内存、任务控制以及 POSIX 线程之类功能的支持。除少数特例外, 并未要求实现支持这些选项。相反, 对于实现在编译及运行时是否支持某一特定特性, SUSv3 允许实现自行给出建议。

通过在<unistd.h>文件中定义相应常量, 实现能够在编译时通告其对特定 SUSv3 选项的支持。此类常量的命名均会冠以前缀(比如 POSIX_ 或者 XOPEN_), 以标识其源于何种标准。

各个选项常量, 一经定义, 其值必为下列之一。

- 值为-1, 表示实现不支持该选项。此时, 系统实现无需定义与该选项有关的头文件、数据类型和函数接口。可以使用 #if 预处理程序指令, 通过条件编译来处理这种情况。
- 值为 0, 表示实现可能支持该选项。应用程序必须在运行时检查该选项是否获得支持。
- 值大于 0, 则表示实现支持该选项。实现定义了与该选项有关的所有头文件、数据类型和函数接口, 且其行为也符合规范要求。在很多情况下, SUSv3 要求这一正值为 200112L, 该常量对应于批准 SUSv3 标准的年、月。(SUSv4 中, 将类似功能的值设为 200809L。)

当定义常量为 0 时, 应用程序可使用 sysconf()和 pathconf() (或 f pathconf()) 在运行时检查选项是否获得实现的支持。传递给这些函数的入参 name, 其命名通常与编译时常量形式相同, 只是前缀为 SC_ 或 PC_ 所取代。系统实现必须至少提供头文件、常量以及实施运行时检查所必要的函数接口。

对于未定义的选项常量, 其含义到底等同于常量 0 (可能支持该选项) 还是-1 (不支持该选项), SUSv3 并未做出明确规定。随后, 标准委员会作出裁决, 规定这种情况应与定义为-1 的常量含义相同, 并且 SUSv4 对此明确作出了规定。

表 11-3 列举了 SUSv3 所规定的一些选项。表中第一列针对选项 (定义于<unistd.h>文件中) 给出了相关编译时常量的名称, 以及 sysconf() (_SC_*) 和 pathconf()(_PC_*)函数的相应入参 name 值。对于特定选项, 请注意以下几点。

- 某些选项在 SUSv3 中是必需的, 即编译时其常量值总应大于 0。历史上, 这些选项一度确实曾是可选项, 但如今已是时过境迁。“备注”栏会以字符“+”标识此类选项。(许多在 SUSv3 中的可选项在 SUSv4 中已经成为必选项。)

¹ 译者注: #include。

虽然这些选项在 SUSv3 中是必需的，但在安装一些 UNIX 系统时，若配置不当，系统依然会与规范不符。因此，对可移植的应用程序而言，不管标准是否对影响应用的选项作出了要求，总应检查系统是否支持该选项。

- 对于某些选项，其编译时常量必须为-1 以外的值。换言之，要么必须支持该选项，要么必须有方法可以检查出系统在运行时是否支持该选项。这些选项的“备注”栏以字符“*”标识这些选项。

表 11-3: 已选 SUSv3 选项

选项 (常量) 名 (sysconf() / pathconf() 入参 name 名)	描 述	备注
_POSIX_ASYNCHRONOUS_IO (_SC_ASYNCHRONOUS_IO)	异步 I/O	
_POSIX_CHOWN_RESTRICTED (_PC_CHOWN_RESTRICTED)	仅有特权级进程能够使用 chown() 和 fchown() 函数将文件的用户 ID 和组 ID 修改为任意值 (15.3.2 节)	*
_POSIX_JOB_CONTROL (_SC_JOB_CONTROL)	作业控制 (34.7 节)	+
_POSIX_MESSAGE_PASSING (_SC_MESSAGE_PASSING)	POSIX 消息队列 (第 52 章)	
_POSIX_PRIORITY_SCHEDULING (_SC_PRIORITY_SCHEDULING)	进程调度 (35.3 节)	
_POSIX_REALTIME_SIGNALS (_SC_REALTIME_SIGNALS)	实时信号扩展 (22.8 节)	+
_POSIX_SAVED_IDS(none)	进程拥有的保存 (saved)set-user-ID 和保存 (saved)set-group-ID (9.4 节)	
_POSIX_SEMAPHORES (_SC_SEMAPHORES)	POSIX 信号 (第 53 章)	
_POSIX_SHARED_MEMORY_OBJECTS (_SC_SHARED_MEMORY_OBJECTS)	POSIX 共享内存对象 (第 54 章)	
_POSIX_THREADS (_SC_THREADS)	POSIX 线程	
_XOPEN_UNIX (_SC_XOPEN_UNIX)	支持 XSI 扩展功能 (1.3.4 节)	

11.6 总结

对于系统实现必须支持的限制和可能支持的系统选项，SUSv3 都做了规范。

通常，不建议将对系统限制和选项的假设值硬性写入应用程序代码，因为这些值既可能随系统的不同而发生变化，也可能在同一个系统实现中因不同的运行期间或文件系统而不同。因此，SUSv3 规定了一干方法，借助于此，系统实现可发布其所支持的限制和选项。对于大

多数限制，SUSv3 规定了所有实现所必须支持的最小值。此外，每个实现还能在编译时(通过定义于<limits.h>或<unistd.h>文件中的常量)和/或运行时(通过调用 `sysconf()`、`pathconf()` 或 `fpathconf()`函数) 发布其特有的限制和选项。此类技术同样可应用于找出实现所支持的 SUSv3 选项。在一些情况下，无论使用上述何种方法，都不能获取某个特定限制的值。对于这些不确定的限制，必须采用特殊技术来确定应用程序所应遵循的限制。

更多信息

[Stevens & Rago, 2005]第 1 章和[Gallmeister, 1995] 第 2 章均涵盖了与本章相类似的知识，[Lewine, 1991]也提供了很多有用的(尽管稍有过时)背景知识。在 Linux 及 glibc 中与 POSIX 选项有关的一些详细信息，可见诸于 <http://people.redhat.com/drepper/posix-option-groups.html>。相关的 Linux 手册页如下：`sysconf(3)`、`pathconf(3)`、`feature_test_macros(7)`、`posixoptions(7)`和 `standards(7)`。

最佳的信息来源(虽然有时难以理解)还是 SUSv3 中的相关部分，特别是基本定义(XBD)的第 2 章以及针对<unistd.h>、<limits.h>、`sysconf()`和 `fpathconf()`的规格说明。[Josey, 2004]也为 SUSv3 的使用提供了指导。

11.7 练习

- 11-1.** 尝试在其他 UNIX 实现中运行程序清单 11-1 所列程序。
- 11-2.** 尝试在其他文件系统中运行程序清单 11-2 所列程序。

第 12 章

系统和进程信息

本章将研究一系列系统和进程信息的访问方法，重点讨论 `/proc` 文件系统。本章还描述了 `uname()` 系统调用，该调用用于获取各种系统标识。

12.1 `/proc` 文件系统

在较老的 UNIX 实现中，通常并无简单方法来获取（或修改）内核属性并回答如下问题：

- 系统中有多少进程正在运行，其属主是谁？
- 一个进程已经打开了什么文件？
- 目前锁定了什么文件，哪些进程持有这些锁？
- 系统正在使用什么套接字（socket）？

一些老版 UNIX 实现解决这一问题的方法是允许特权级程序深入内核内存中的数据结构。然而，这会带来各种问题。特别是，这要求对内核数据结构具有专业知识，并且这些结构可能因内核版本的演进而发生改变，故而需要加以重写。

为了提供更为简便的方法来访问内核信息，许多现代 UNIX 实现提供了一个 `/proc` 虚拟文件系统。该文件系统驻留于 `/proc` 目录中，包含了各种用于展示内核信息的文件，并且允许进程通过常规文件 I/O 系统调用来方便地读取，有时还可以修改这些信息。之所以将 `/proc` 文件系统称为虚拟，是因为其包含的文件和子目录并未存储于磁盘上，而是由内核在进程访问此类信息时动态创建而成。

本节展示了 `/proc` 文件系统的概况。后续各章将视各自主题来描述特定的 `/proc` 文件。虽然许多 UNIX 实现提供了 `/proc` 文件系统，但 SUSv3 并未对其进行规范，本书所述细节是 Linux 专有的。

12.1.1 获取与进程有关的信息：`/proc/PID`

对于系统中每个进程，内核都提供了相应的目录，命名为 `/proc/PID`，其中 `PID` 是进程的 ID。在此目录中的各种文件和子目录包含了进程的相关信息。例如，通过查看 `/proc/1` 目录下

的文件，可以获取 init 进程的信息，该进程的 ID 总是为 1。

每个/proc/PID 目录中都存在一个命名为 status 的文件，提供了有关该进程的一系列信息。

```
$ cat /proc/1/status
Name:      init                Name of command run by this process
State:    S (sleeping)        State of this process
Tgid:     1                    Thread group ID (traditional PID, getpid())
Pid:      1                    Actually, thread ID (gettid())
PPid:     0                    Parent process ID
TracerPid: 0                    PID of tracing process (0 if not traced)
Uid:      0      0      0      0    Real, effective, saved set, and FS UIDs
Gid:      0      0      0      0    Real, effective, saved set, and FS GIDs
FDSize:   256                  # of file descriptor slots currently allocated
Groups:                                Supplementary group IDs
VmPeak:   852 kB                Peak virtual memory size
VmSize:   724 kB                Current virtual memory size
VmLck:    0 kB                  Locked memory
VmHWM:    288 kB                Peak resident set size
VmRSS:    288 kB                Current resident set size
VmData:   148 kB                Data segment size
VmStk:    88 kB                 Stack size
VmExe:    484 kB                Text (executable code) size
VmLib:    0 kB                  Shared library code size
VmPTE:    12 kB                Size of page table (since 2.6.10)
Threads:  1                     # of threads in this thread's thread group
SigQ:     0/3067                Current/max. queued signals (since 2.6.12)
SigPnd:   0000000000000000      Signals pending for thread
ShdPnd:   0000000000000000      Signals pending for process (since 2.6)
SigBlk:   0000000000000000      Blocked signals
SigIgn:   ffffffff5770d8fc      Ignored signals
SigCgt:   00000000280b2603      Caught signals
CapInh:   0000000000000000      Inheritable capabilities
CapPrm:   00000000ffffffff       Permitted capabilities
CapEff:   00000000fffffeff       Effective capabilities
CapBnd:   00000000ffffffff       Capability bounding set (since 2.6.26)
Cpus_allowed: 1                  CPUs allowed, mask (since 2.6.24)
Cpus_allowed_list: 0             Same as above, list format (since 2.6.26)
Mems_allowed: 1                  Memory nodes allowed, mask (since 2.6.24)
Mems_allowed_list: 0             Same as above, list format (since 2.6.26)
voluntary_ctxt_switches: 6998    Voluntary context switches (since 2.6.23)
nonvoluntary_ctxt_switches: 107  Involuntary context switches (since 2.6.23)
Stack usage: 8 kB                Stack usage high-water mark (since 2.6.32)
```

上面的输出来自于内核 2.6.32。正如伴随文件输出的“始于”说明所示，该文件格式随着时间的推移而不断演进，在不同内核版本中增加了新字段（极少情况下，也会移除字段）。（除了注释中 Linux 2.6 所带来的改变之外，Linux 2.4 增加了 Tgid、TracerPid、FDSize 和 Threads 字段。）

该文件内容随着时间而改变，这一事实揭示出关于/proc 文件使用的要点所在。当这些文件由多个条目组成时，对其解析应当谨慎从事，在这种情况下，应查找包含特殊字符串（如，PPid）的匹配行记录，而非按照（逻辑）行号来处理文件。

表 12-1 列举了在每个/proc/PID 目录中的部分其他文件。

表 12-1: 每个/proc/PID 目录下的文件节选

文 件	描述 (进程属性)
cmdline	以\0 分隔的命令行参数
cwd	指向当前工作目录的符号链接
Environ	NAME=value 键值对环境列表, 以\0 分隔
exe	指向正在执行文件的符号链接
fd	文件目录, 包含了指向由进程打开文件的符号链接
maps	内存映射
mem	进程虚拟内存 (在 I/O 操作之前必须调用 lseek()移至有效偏移量)
mounts	进程的安装点
root	指向根目录的符号链接
status	各种信息 (比如, 进程 ID、凭证、内存使用量、信号)
task	为进程中的每个线程均包含一个子目录 (始自 Linux 2.6)

/proc/PID/fd 目录

/proc/PID/fd 目录为进程打开的每个文件描述符都包含了一个符号链接, 每个符号链接的名称都与描述符的数值相匹配。例如, /proc/1968/1 是 ID 为 1968 的进程中指向标准输出的符号链接, 更多信息参见 5.11 节。

为方便起见, 任何进程都可使用符号链接/proc/self 来访问其自己的/proc/PID 目录。

线程: /proc/PID/task 目录

Linux 2.4 增加了线程组概念, 正式支持 POSIX 线程模型。因为线程组中的一些属性对于线程而言是唯一的, 所以 Linux 2.4 在/proc/PID 目录下增加了一个 task 子目录。针对进程中的每个线程, 内核提供了以/proc/PID/task/TID 命名的子目录, 其中 TID 是该线程的线程 ID。(此值等同于在线程中调用 gettid()函数的返回值。)

每个/proc/PID/task/TID 子目录中都有一套类似于/proc/PID 目录内容的文件和目录。因为线程共享了多个属性, 所以这些文件中的许多信息对进程中各个线程而言都是相同的。然而, 这些文件也显示了每个线程的独特信息, 故而是合理的。例如, 在线程组的/proc/PID/task/TID/status 文件中, 存在那种对每个线程而言, 内容都有可能不同的字段, State、Pid、SigPnd、SigBlk、CapInh、CapPrm、CapEff 和 CapBnd 就在此列。

12.1.2 /proc 目录下的系统信息

/proc 目录下的各种文件和子目录提供了对系统级信息的访问。图 12-1 展示了其中的部分。

图 12-1 中的许多文件在本书的其他章节进行描述。表 12-2 总结了图 12-1 所示/proc 子目录的一般用途。

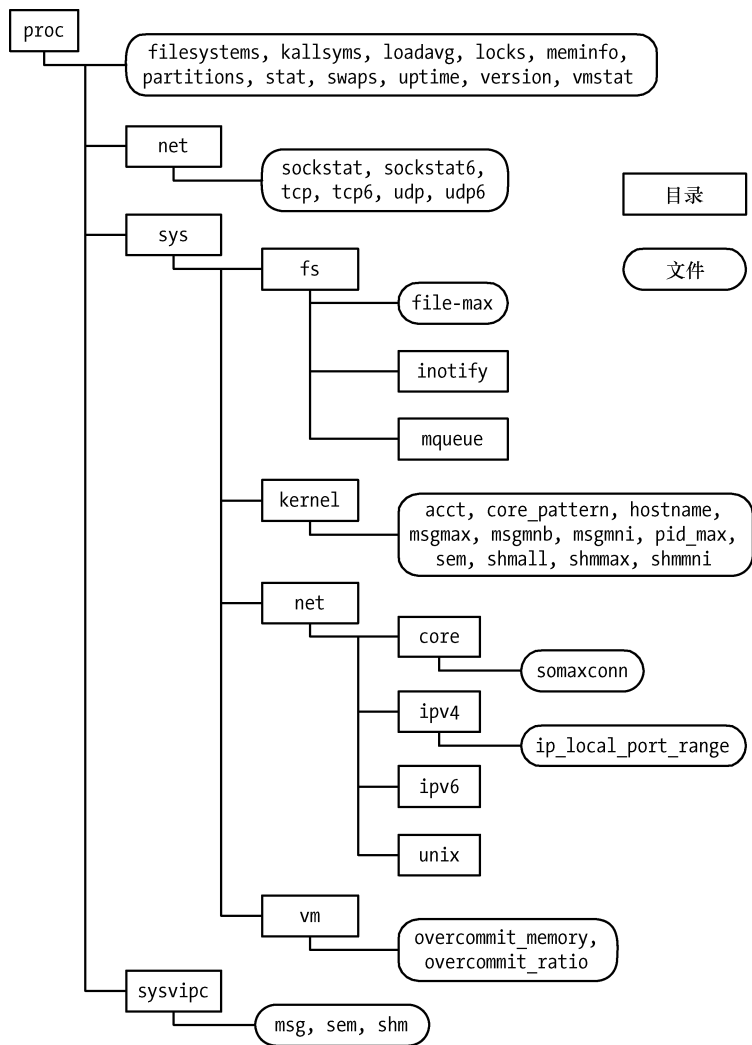


图 12-1: /proc 目录下文件和子目录的节选

表 12-2: 节选/proc 子目录的用途

目 录	目录中文件表达的信息
/proc	各种系统信息
/proc/net	有关网络和套接字的状态信息
/proc/sys/fs	文件系统相关设置
/proc/sys/kernel	各种常规的内核设置
/proc/sys/net	网络和套接字的设置
/proc/sys/vm	内存管理设置
/proc/sysvipc	有关 System V IPC 对象的信息

12.1.3 访问/proc 文件

通常使用 shell 脚本来访问/proc 目录下的文件（使用诸如 Python 或者 Perl 之类的脚本语

言，很容易解析大多数/proc 目录下包含有多个值的文件)。例如，使用如下 shell 命令，就可以修改和查看/proc 目录下的文件内容：

```
# echo 100000 > /proc/sys/kernel/pid_max
# cat /proc/sys/kernel/pid_max
100000
```

也可以从程序中使用常规 I/O 系统调用来访问/proc 目录下的文件。但在访问这些文件时，有如下一些限制。

- /proc 目录下的一些文件是只读的，即这些文件仅用于显示内核信息，但无法对其进行修改。/proc/PID 目录下的大多数文件就属于此类型。
- /proc 目录下的一些文件仅能由文件拥有者（或特权级进程）读取。例如，/proc/PID 目录下的所有文件都属于拥有相应进程的用户，而且即使是对文件的属主，其中的部分文件（如：proc/PID/environ 文件）也仅仅授予了读权限。
- 除了/proc/PID 子目录中的文件，/proc 目录的其他文件大多属于 root 用户，并且也只有 root 用户能够修改那些可修改的文件。

访问/proc/PID 目录中的文件

/proc/PID 目录内容变化不定。每个目录随着含有相应进程 ID 的进程创建而生，又随进程的终止而灭。这意味着要确定特定/proc/PID 目录的存在，就需要干净利落地处理如下可能性：当打开此目录下的文件时，进程已经终止，并且也已经删除了相应的/proc/PID 目录。

示例程序

程序清单 12-1 展示了如何读取和修改一个/proc 目录下的文件。该程序读取并显示了/proc/sys/kernel/pid_max 文件的内容。若提供了命令行参数，则程序将使用此参数对文件进行更新。该文件（Linux 2.6 的新增文件）规定了进程 ID 的上限（见 6.2 节）。此处是运用该程序的一个例子：

```
$ su                                     Privilege is required to update pid_max file
Password:
# ./procfs_pidmax 10000
Old value: 32768
/proc/sys/kernel/pid_max now contains 10000
```

程序清单 12-1：访问/proc/sys/kernel/pid_max 文件

```
----- sysinfo/procfs_pidmax.c
#include <fcntl.h>
#include "tlpi_hdr.h"

#define MAX_LINE 100

int
main(int argc, char *argv[])
{
    int fd;
    char line[MAX_LINE];
    ssize_t n;

    fd = open("/proc/sys/kernel/pid_max", (argc > 1) ? O_RDWR : O_RDONLY);
    if (fd == -1)
```

```

    errExit("open");

n = read(fd, line, MAX_LINE);
if (n == -1)
    errExit("read");

if (argc > 1)
    printf("Old value: ");
printf("%.s", (int) n, line);

if (argc > 1) {
    if (write(fd, argv[1], strlen(argv[1])) != strlen(argv[1]))
        fatal("write() failed");

    system("echo /proc/sys/kernel/pid_max now contains "
           "`cat /proc/sys/kernel/pid_max`");
}

exit(EXIT_SUCCESS);
}

```

sysinfo/procfs_pidmax.c

12.2 系统标识: uname()

uname()系统调用返回了一系列关于主机系统的标识信息，存储于 `utsbuf` 所指向的结构中。

```

#include <sys/utsname.h>

int uname(struct utsname *utsbuf);

```

Returns 0 on success, or -1 on error

`utsbuf` 参数是一个指向 `utsname` 结构的指针，其定义如下：

```

#define _UTSNAME_LENGTH 65

struct utsname {
    char sysname[_UTSNAME_LENGTH];    /* Implementation name */
    char nodename[_UTSNAME_LENGTH];   /* Node name on network */
    char release[_UTSNAME_LENGTH];    /* Implementation release level */
    char version[_UTSNAME_LENGTH];    /* Release version level */
    char machine[_UTSNAME_LENGTH];    /* Hardware on which system
                                     is running */
#ifdef _GNU_SOURCE                    /* Following is Linux-specific */
    char domainname[_UTSNAME_LENGTH]; /* NIS domain name of host */
#endif
};

```

SUSv3 规范了 `uname()`，但对 `utsname` 结构中各种字段的长度未加定义，仅要求字符串以空字节终止。在 Linux 中，这些字段长度均为 65 个字节，其中包括了空字节终止符所占用的空间。而在一些 UNIX 实现中，这些字段更短，但在其他操作系统（如 Solaris）中，这些字段的长度长达 257 个字节。

`utsname` 结构中的 `sysname`、`release`、`version` 和 `machine` 字段由内核自动设置。

在 Linux 中, /proc/sys/kernel 目录下的 3 个文件提供了与 utsname 结构的 sysname、release 和 version 字段返回值相同的信息, 这些只读文件分别为 ostype、osrelease 和 version。另外一个文件/proc/version, 也包含了这些信息, 并且还包含了有关内核编译的步骤信息(即执行编译的用户名、用于编译的主机名, 以及使用的 gcc 版本)。

nodename 字段的返回值由 sethostname()系统调用设置(详情请参考此系统调用的手册页)。通常, 该值类似于系统 DNS 域名中的前缀主机名。

domainname 字段的返回值由 setdomainname()系统调用设置(详情请参考此系统调用的手册页)。该值是主机的网络信息服务(NIS)域名(与主机域名不同)。

gethostname()系统调用(是 sethostname()函数的反向操作)用于获取系统主机名, 也可利用 hostname(1)命令和 Linux 特有的/proc/hostname 文件来查看和设置系统主机名。

getdomainname()系统调用(setdomainname()函数的反向操作)用于获取 NIS 域名, 也可利用 domainname(1)命令和 Linux 特有的/proc/domainname 文件来查看和设置 NIS 域名。

sethostname()和 setdomainname()系统调用在应用程序中鲜有使用。通常, 会在系统启动时运行启动脚本来确立主机名和 NIS 域名。

程序清单 12-2 中程序展示了 uname()的返回信息。下面是运行该程序可能看到的输出信息:

```
$ ./t_uname
Node name:  tekapo
System name: Linux
Release:    2.6.30-default
Version:    #3 SMP Fri Jul 17 10:25:00 CEST 2009
Machine:    i686
Domain name:
```

程序清单 12-2: 使用 uname()

```
----- sysinfo/t_uname.c
#define _GNU_SOURCE
#include <sys/utsname.h>
#include "t_lpi_hdr.h"

int
main(int argc, char *argv[])
{
    struct utsname uts;

    if (uname(&uts) == -1)
        errExit("uname");

    printf("Node name:  %s\n", uts.nodename);
    printf("System name: %s\n", uts.sysname);
    printf("Release:    %s\n", uts.release);
    printf("Version:    %s\n", uts.version);
    printf("Machine:    %s\n", uts.machine);
#ifdef _GNU_SOURCE
    printf("Domain name: %s\n", uts.domainname);
#endif
    exit(EXIT_SUCCESS);
}
----- sysinfo/t_uname.c
```

12.3 总结

/proc 文件系统向应用程序暴露了一系列内核信息。每个/proc/PID 子目录都包含有许多文件和子目录，是进程 ID 为 PID 的进程提供的相关信息。/proc 目录下的其他许多文件和目录，则暴露了应用程序可以读取，有时还可以修改的系统级信息。

使用 `uname()` 系统调用，能够获取 UNIX 的实现信息以及应用程序所运行的机器类型。

进阶信息

关于/proc 文件系统的深入信息可见诸于 `proc (5)` 手册页、内核源文件 `Documentation/filesystems/proc.txt` 以及 `Documentation/sysctl` 目录下的各种文件。

12.4 练习

- 12-1. 编写一个程序，以用户名作为命令行参数，列表显示该用户下所有正在运行的进程 ID 和命令名。（程序清单 8-1 中的 `userIdFromName()` 函数对本题程序的编写可能会有所帮助。）通过分析系统中/proc/PID/status 文件的 `Name:` 和 `Uid:` 各行信息，可以实现此功能。遍历系统的所有/proc/PID 目录需要使用 `readdir(3)` 函数，18.8 节对其进行了描述。程序必须能够正确处理如下可能性：在确定目录存在与程序尝试打开相应/proc/PID/status 文件之间，/proc/PID 目录消失了。
- 12-2. 编写一个程序绘制树状结构，展示系统中所有进程的父子关系，根节点为 `init` 进程。对每个进程而言，程序应该显示进程 ID 和所执行的行命令。程序输出类似于 `ps tree(1)` 的输出结果，但也无需像后者那样复杂。每个进程的父进程可通过对/proc/PID/status 系统文件中 `PPid:` 行的分析获得。但是需要小心处理如下可能性：在扫描所有/proc/PID 目录的过程中，进程的父进程（以及父进程的/proc/PID 目录）消失了。
- 12-3. 编写一个程序，列表展示打开同一特定路径名文件的所有进程。可以通过分析所有/proc/PID/fd/*符号链接的内容来实现此功能。这需要利用 `readdir(3)` 函数来嵌套循环，扫描所有/proc/PID 目录以及每个/proc/PID 目录下所有/proc/PID/fd 的条目内容。读取/proc/PID/fd/n 符号链接的内容，需要使用 `readlink()`，18.5 节对其进行了描述。

第 13 章

文件 I/O 缓冲

出于速度和效率考虑，系统 I/O 调用（即内核）和标准 C 语言库 I/O 函数（即 `stdio` 函数）在操作磁盘文件时会对数据进行缓冲。本章描述了这两种类型的缓冲，并讨论了其对应用程序性能的影响。本章还讨论了可以屏蔽或影响缓冲的各种技术，以及直接 I/O 技术——在某些需要绕过内核缓冲的场景中非常有用。

13.1 文件 I/O 的内核缓冲：缓冲区高速缓存

`read()`和 `write()`系统调用在操作磁盘文件时不会直接发起磁盘访问，而是仅仅在用户空间缓冲区与内核缓冲区高速缓存（`kernel buffer cache`）之间复制数据。例如，如下调用将 3 个字节的数据从用户空间内存传递到内核空间的缓冲区中：

```
write(fd, "abc", 3);
```

`write()`随即返回。在后续某个时刻，内核会将其缓冲区中的数据写入（刷新至）磁盘。（因此，可以说系统调用与磁盘操作并不同步。）如果在此期间，另一进程试图读取该文件的这几个字节，那么内核将自动从缓冲区高速缓存中提供这些数据，而不是从文件中（读取过期的内容）。

与此同理，对输入而言，内核从磁盘中读取数据并存储到内核缓冲区中。`read()`调用将从该缓冲区中读取数据，直至把缓冲区中的数据取完，这时，内核会将文件的下一段内容读入缓冲区高速缓存。（这里的描述有所简化。对于序列化的文件访问，内核通常会尝试执行预读，以确保在需要之前就将文件的下一数据块读入缓冲区高速缓存中。更多关于预读的内容请参考 13.5 节。）

采用这一设计，意在使 `read()`和 `write()`调用的操作更为快速，因为它们不需要等待（缓慢的）磁盘操作。同时，这一设计也极为高效，因为这减少了内核必须执行的磁盘传输次数。

Linux 内核对缓冲区高速缓存的大小没有固定上限。内核会分配尽可能多的缓冲区高速缓存页，而仅受限于两个因素：可用的物理内存总量，以及出于其他目的对物理内存的需求（例如，需要将正在运行进程的文本和数据页保留在物理内存中）。若可用内存不足，则内核会将

一些修改过的缓冲区高速缓存页内容刷新到磁盘，并释放其供系统重用。

更确切地说，从内核 2.4 开始，Linux 不再维护一个单独的缓冲区高速缓存。相反，会将文件 I/O 缓冲区置于页面高速缓存中，其中还含有诸如内存映射文件的页面。然而，正文的讨论采用了“缓冲区高速缓存 (buffer cache)”这一术语，因为这是 UNIX 实现中历史悠久的通称。

缓冲区大小对 I/O 系统调用性能的影响

无论是让磁盘写 1000 次，每次写入一个字节，还是一次写入 1000 个字节，内核访问磁盘的字节数都是相同的。然而，我们更属意于后者，因为它只需要一次系统调用，而前者则需要调用 1000 次。尽管比磁盘操作要快许多，但系统调用所耗费的时间总量也相当可观：内核必须捕获调用，检查系统调用参数的有效性，在用户空间和内核空间之间传输数据（详情参见 3.1 节）。

为 BUF_SIZE (BUF_SIZE 指定了每次调用 read()和 write() 时所传输的字节数) 设定不同的大小来运行程序清单 4-1，可以观察到不同大小的缓冲区对执行文件 I/O 所产生的影响。表 13-1 所示为在 Linux ext2 文件系统上复制大小为 100MB 的文件，该程序在使用不同 BUF_SIZE 值时所需要的时间。有关本表中的信息，需要注意以下几点。

- 总用时和总 CPU 时间这两列含义很明显。而用户 CPU 和系统 CPU 两列是将总 CPU 用时分解为在用户模式下执行代码所需的时间和执行内核代码所需的时间（比如，系统调用）。
- 表中测试结果得自于 2.6.30 普通 (vanilla) 内核下，块大小为 4096 字节的 ext2 文件系统。

所谓普通内核 (vanilla kernel)，意指未打补丁的主线 (mainline) 内核。与之形成鲜明对比的是大多数发行商所提供的内核，常包含各种补丁来修复错误和添加新功能。

- 每行显示的结果为在给定缓冲区大小下运行 20 次的均值。在这些测试以及本章后续提及的其他测试里，在程序每次的执行间隔中，会卸载并再次重新装配文件系统，以确保文件系统的缓冲区高速缓存为空。计时则由 shell 命令 time 完成。

表 13-1: 复制 100MB 大小的文件所需时间

BUF_SIZE	时间 (秒)			
	总用时	总 CPU 用时	用户 CPU 用时	系统 CPU 用时
1	107.43	107.32	8.20	99.12
2	54.16	53.89	4.13	49.76
4	31.72	30.96	2.30	28.66
8	15.59	14.34	1.08	13.26
16	7.50	7.14	0.51	6.63
32	3.76	3.68	0.26	3.41
64	2.19	2.04	0.13	1.91

续表

BUF_SIZE	时间 (秒)			
	总用时	总 CPU 用时	用户 CPU 用时	系统 CPU 用时
128	2.16	1.59	0.11	1.48
256	2.06	1.75	0.10	1.65
512	2.06	1.03	0.05	0.98
1024	2.05	0.65	0.02	0.63
4096	2.05	0.38	0.01	0.38
16384	2.05	0.34	0.00	0.33
65536	2.06	0.32	0.00	0.32

因为采用不同的缓冲区大小时，数据的传输总量（因此招致磁盘操作的数量）是相同的，表 13-1 所示为发起 read()和 write()调用的开销。缓冲区大小为 1 字节时，需要调用 read()和 write()1 亿次，缓冲区大小为 4096 个字节时，需要调用 read()和 write() 24000 次左右，几乎达到最优性能。设置再超过这个值，对性能的提升就不显著了，这是因为与在用户空间和内核空间之间复制数据以及执行实际磁盘 I/O 所花费的时间相比，read()和 write() 系统调用的成本就显得微不足道了。

从表 13-1 的最后一行中可以粗略估算出在用户空间与内核空间之间传输数据以及执行文件 I/O 的总耗时。因为此时系统调用的次数相对较少，所以它们所花费的时间相对于总耗时和 CPU 时间可以忽略不计。据此可认为，系统 CPU 时间主要是测量用户空间与内核空间之间数据传输所消耗的时间。而总耗时则是对与磁盘传输数据所需时间的估算。（正如下面将提到的，时间主要花在了对磁盘的读操作上。）

总之，如果与文件发生大量的数据传输，通过采用大块空间缓冲数据，以及执行更少的系统调用，可以极大地提高 I/O 性能。

表 13-1 度量了一系列因素：执行 read()和 write()系统调用所需的时间、内核空间和用户空间缓冲区之间传输数据所需的时间、内核缓冲区与磁盘之间传输数据所需的时间。再进一步考虑一下最后一个要素，显然，将输入文件的内容传输到缓冲区高速缓存是不可避免的。然而，当数据从用户空间传输到内核空间后，write()调用立即返回。由于测试系统上的 RAM 大小(4 GB) 远超欲复制文件的大小 (100MB)，据此推断，当程序完成时，输出文件实际尚未写入磁盘。因此，再进一步做个实验，运行一个程序，使用不同大小的缓冲区，以 write() 随意向文件中写入一些数据。运行结果如表 13-2 所示。

同样，表 13-2 中数据是来自内核 2.6.30，以及块大小为 4096 字节的 ext2 文件系统，并且每行显示为运行了 20 次后的均值。本节并未列出测试程序 (filebuff/ write_bytes.c) 的代码清单，但可从随本书一起发行的源码中获取。

表 13-2: 写一个 100MB 大小的文件所需要的时间

BUF_SIZE	时间 (秒)			
	总用时	总 CPU 用时	用户 CPU 用时	系统 CPU 用时
1	72.13	72.11	5.00	67.11
2	36.19	36.17	2.47	33.70

续表

BUF_SIZE	时间 (秒)			
	总用时	总 CPU 用时	用户 CPU 用时	系统 CPU 用时
4	20.01	19.99	1.26	18.73
8	9.35	9.32	0.62	8.70
16	4.70	4.68	0.31	4.37
32	2.39	2.39	0.16	2.23
64	1.24	1.24	0.07	1.16
128	0.67	0.67	0.04	0.63
256	0.38	0.38	0.02	0.36
512	0.24	0.24	0.01	0.23
1024	0.17	0.17	0.01	0.16
4096	0.11	0.11	0.00	0.11
16384	0.10	0.10	0.00	0.10
65536	0.09	0.09	0.00	0.09

表 13-2 显示为使用不同大小的缓冲区调用 `write()` 从用户空间向内核缓冲区高速缓存传输数据所花费的成本。缓冲区越大，与表 13-1 中数据的差异就越明显。例如，对于一个 65536 字节大小的缓冲区，在表 13-1 中总耗时为 2.06 秒，而表 13-2 中仅为 0.09 秒。这是因为在后者的情况下并未执行实际的磁盘 I/O 操作。换言之，表 13-1 中采用大缓冲区时的耗时绝大部分花在了对磁盘的读取上。

正如 13.3 节所述，若强制在数据传输到磁盘前阻塞输出操作，则调用 `write()` 所需的时间会显著上升。

最后，值得注意的是，表 13-2（以及表 13-3）中信息仅仅代表了对文件系统评价基准的形式之一，还不完善。此外，文件系统不同，结果可能也会有所不同。对文件系统的度量还有各种其他标准，比如多用户、高负载下的性能表现，创建和删除文件的速度，在一个大型目录下搜索一个文件所需的时间，存储小文件所需的空间，或者在遭遇系统崩溃时对文件完整性的维护。只要 I/O 或是其他文件系统操作的性能至关重要，那么在目标平台上针对特定应用的测试基准就不可替代。

13.2 stdio 库的缓冲

当操作磁盘文件时，缓冲大块数据以减少系统调用，C 语言函数库的 I/O 函数（比如，`fprintf()`、`fscanf()`、`fgets()`、`fputs()`、`fputc()`、`fgetc()`）正是这么做的。因此，使用 `stdio` 库可以使编程者免于自行处理对数据的缓冲，无论是调用 `write()` 来输出，还是调用 `read()` 来输入。

设置一个 `stdio` 流的缓冲模式

调用 `setvbuf()` 函数，可以控制 `stdio` 库使用缓冲的形式。

```
#include <stdio.h>

int setvbuf(FILE *stream, char *buf, int mode, size_t size);

Returns 0 on success, or nonzero on error
```

参数 `stream` 标识将要修改哪个文件流的缓冲。打开流后，必须在调用任何其他 `stdio` 函数之前先调用 `setvbuf()`。`setvbuf()`调用将影响后续在指定流上进行的所有 `stdio` 操作。

不要将 `stdio` 库所使用的流与 System V 系统的 STREAMS 机制相混淆，Linux 的主线内核中并未实现 System V 系统的 STREAMS。

参数 `buf` 和 `size` 则针对参数 `stream` 要使用的缓冲区，指定这些参数有如下两种方式。

- 如果参数 `buf` 不为 `NULL`，那么其指向 `size` 大小的内存块以作为 `stream` 的缓冲区。因为 `stdio` 库将要使用 `buf` 指向的缓冲区，所以应该以动态或静态在堆中为该缓冲区分配一块空间(使用 `malloc()`或类似函数)，而不应是分配在栈上的函数本地变量。否则，函数返回时将销毁其栈帧，从而导致混乱。
- 若 `buf` 为 `NULL`，那么 `stdio` 库会为 `stream` 自动分配一个缓冲区（除非选择非缓冲的 I/O，如下所述）。SUSv3 允许，但不强制要求库实现使用 `size` 来确定其缓冲区的大小。glibc 实现会在该场景下忽略 `size` 参数。

参数 `mode` 指定了缓冲类型，并具有下列值之一。

`_IONBF`

不对 I/O 进行缓冲。每个 `stdio` 库函数将立即调用 `write()`或者 `read()`，并且忽略 `buf` 和 `size` 参数，可以分别指定两个参数为 `NULL` 和 `0`。`stderr` 默认属于这一类型，从而保证错误能立即输出。

`_IOLBF`

采用行缓冲 I/O。指代终端设备的流默认属于这一类型。对于输出流，在输出一个换行符（除非缓冲区已经填满）前将缓冲数据。对于输入流，每次读取一行数据。

`_IOFBF`

采用全缓冲 I/O。单次读、写数据（通过 `read()`或 `write()`系统调用）的大小与缓冲区相同。指代磁盘的流默认采用此模式。

下面的代码演示了 `setvbuf()`函数的用法：

```
#define BUF_SIZE 1024
static char buf[BUF_SIZE];

if (setvbuf(stdout, buf, _IOFBF, BUF_SIZE) != 0)
    errExit("setvbuf");
```

注意：`setvbuf()`出错时返回非 0 值（而不一定是-1）。

`setbuf()`函数构建于 `setvbuf()`之上，执行了类似任务。

```
#include <stdio.h>

void setbuf(FILE *stream, char *buf);
```

`setbuf(fp,buf)`调用除了不返回函数结果外，就相当于：

```
setvbuf(fp, buf, (buf != NULL) ? _IOFBF: _IONBF, BUFSIZ);
```

要么将参数 `buf` 指定为 `NULL` 以表示无缓冲，要么指向由调用者分配的 `BUFSIZ` 个字节的缓冲区。（`BUFSIZ` 定义于 `<stdio.h>` 头文件中。glibc 库实现将此常量定义为一个典型值 8192。）

`setbuffer()` 函数类似于 `setbuf()` 函数，但允许调用者指定 `buf` 缓冲区大小。

```
#define _BSD_SOURCE
#include <stdio.h>

void setbuffer(FILE *stream, char *buf, size_t size);
```

对 `setbuffer(fp, buf, size)` 的调用相当于如下调用：

```
setvbuf(fp, buf, (buf != NULL) ? _IOFBF : _IONBF, size);
```

SUSv3 并未对 `setbuffer()` 函数加以定义，但大多数 UNIX 实现均支持它。

刷新 stdio 缓冲区

无论当前采用何种缓冲区模式，在任何时候，都可以使用 `fflush()` 库函数强制将 `stdio` 输出流中的数据（即通过 `write()` 刷新到内核缓冲区中。此函数会刷新指定 `stream` 的输出缓冲区。

```
#include <stdio.h>

int fflush(FILE *stream);

Returns 0 on success, EOF on error
```

若参数 `stream` 为 `NULL`，则 `fflush()` 将刷新所有的 `stdio` 缓冲区。

也能将 `fflush()` 函数应用于输入流，这将丢弃业已缓冲的输入数据。（当程序下一次尝试从流中读取数据时，将重新装满缓冲区。）

当关闭相应流时，将自动刷新其 `stdio` 缓冲区。

在包括 glibc 库在内的许多 C 函数库实现中，若 `stdin` 和 `stdout` 指向一终端，那么无论何时从 `stdin` 中读取输入时，都将隐含调用一次 `fflush(stdout)` 函数。这将刷新写入 `stdout` 的任何提示，但不包括终止换行符（比如，`printf("Date: ")`）。然而，SUSv3 和 C99 并未规定这一行为，也并非所有的 C 语言函数库都实现了这一行为。要保证程序的可移植性，应用应使用显式的 `fflush(stdout)` 调用来确保显示这些提示。

若打开一个流同时用于输入和输出，则 C99 标准中提出了两项要求。首先，一个输出操作不能紧跟一个输入操作，必须在二者之间调用 `fflush()` 函数或是一个文件定位函数（`fseek()`、`fsetpos()` 或者 `rewind()`）。其次，一个输入操作不能紧跟一个输出操作，必须在二者之间调用一个文件定位函数，除非输入操作遭遇文件结尾。

13.3 控制文件 I/O 的内核缓冲

强制刷新内核缓冲区到输出文件是可能的。这有时很有必要，例如，当应用程序（诸如数据库的日志进程）要确保在继续操作前将输出真正写入磁盘（或者至少写入磁盘的硬件高

速缓存中)。

在描述用于控制内核缓冲的系统调用之前，有必要先熟悉一下 SUSv3 中的相关定义。

同步 I/O 数据完整性和同步 I/O 文件完整性

SUSv3 将同步 I/O 完成¹定义为：某一 I/O 操作，要么已成功完成到磁盘的数据传递，要么被诊断为不成功。

SUSv3 定义了两种不同类型的同步 I/O 完成，二者之间的区别涉及用于描述文件的元数据（关于数据的数据），亦即内核针对文件而存储的数据。14.4 节在描述文件 i-node 时将详细讨论文件的元数据，但就目前而言，了解文件元数据包含了些什么，诸如文件属主、属组、文件权限、文件大小、文件（硬）链接数量，表明文件最近访问、修改以及元数据发生变化的时间戳，指向文件数据块的指针，就足够了。

SUSv3 定义的第一种同步 I/O 完成类型是 **synchronized I/O data integrity completion**²，旨在确保针对文件的一次更新传递了足够的信息（到磁盘），以便于之后对数据的获取。

- 就读操作而言，这意味着被请求的文件数据已经（从磁盘）传递给进程。若存在任何影响到所请求数据的挂起写操作，那么在执行读操作之前，会将这些数据传递到磁盘。
- 就写操作而言，这意味着写请求所指定的数据已传递（至磁盘）完毕，且用于获取数据的所有文件元数据也已传递（至磁盘）完毕。此处的要点在于要获取文件数据，并非需要传递所有经过修改的文件元数据属性。发生修改的文件元数据中需要传递的属性之一是文件大小（如果写操作确实扩展了文件）。相形之下，如果是文件时间戳发生了变化，就无需在下次获取数据前将其传递到磁盘。

Synchronized I/O file integrity completion 是 SUSv3 定义的另一种同步 I/O 完成，也是上述 **synchronized I/O data integrity completion** 的超集。该 I/O 完成模式的区别在于在对文件的一次更新过程中，要将所有发生更新的文件元数据都传递到磁盘上，即使有些在后续对文件数据的读操作中并不需要。

用于控制文件 I/O 内核缓冲的系统调用

`fsync()` 系统调用将使缓冲数据和与打开文件描述符 `fd` 相关的所有元数据都刷新到磁盘上。调用 `fsync()` 会强制使文件处于 **Synchronized I/O file integrity completion** 状态。

```
#include <unistd.h>

int fsync(int fd);

Returns 0 on success, or -1 on error
```

仅在对磁盘设备（或者至少是其高速缓存）的传递完成后，`fsync()` 调用才会返回。

`fdatasync()` 系统调用的运作类似于 `fsync()`，只是强制文件处于 **synchronized I/O data integrity completion** 的状态。

1 译者注：synchronized I/O completion。

2 译者注：忍无可忍，保留原文。

```
#include <unistd.h>

int fdatsync(int fd);
```

Returns 0 on success, or -1 on error

`fdatsync()`可能会减少对磁盘操作的次数，由 `fsync()`调用请求的两次变为一次。例如，若修改了文件数据，而文件大小不变，那么调用 `fdatsync()`只强制进行了数据更新。（前面已然述及，针对 `synchronized I/O data completion` 状态，如果是诸如最近修改时间戳之类的元数据属性发生了变化，那么是无需传递到磁盘的。）相比之下，`fsync()`调用会强制将元数据传递到磁盘上。

对某些应用而言，以这种方式来减少磁盘 I/O 操作的次数是很有用的，比如对性能要求极高，而对某些元数据（比如时间戳）的准确性要求不高的应用。当应用程序同时进行多处文件更新时，二者存在相当大的性能差异，因为文件数据和元数据通常驻留在磁盘的不同区域，更新这些数据需要反复在整个磁盘上执行寻道操作。

Linux 2.2 以及更早版本的内核将 `fdatsync()`实现为对 `fsync()`的调用，因而性能也未获得提升。

始于内核 2.6.17, Linux 提供了非标准的系统调用 `sync_file_range()`, 当刷新文件数据时, 该调用提供比 `fdatsync()`调用更为精准的控制。调用者能够指定待刷新的文件区域, 并且还能指定标志, 以控制该系统调用在遭遇写磁盘时是否阻塞。更详细的信息请参阅 `sync_file_range(2)`手册页。

`sync()`系统调用会使包含更新文件信息的所有内核缓冲区（即数据块、指针块、元数据等）刷新到磁盘上。

```
#include <unistd.h>

void sync(void);
```

在 Linux 实现中，`sync()`调用仅在所有数据已传递到磁盘上（或者至少高速缓存）时返回。然而，SUSv3 却允许 `sync()`实现只是简单调度一下 I/O 传递，在动作未完成之前即可返回。

若内容发生变化的内核缓冲区在 30 秒内未经显式方式同步到磁盘上，则一条长期运行的内核线程会确保将其刷新到磁盘上。这一做法是为了规避缓冲区与相关磁盘文件内容长期处于不一致状态（以至于在系统崩溃时发生数据丢失）的问题。在 Linux 2.6 版本中，该任务由 `pdflush` 内核线程执行。（在 Linux 2.4 版本中，则由 `kupdated` 内核线程执行。）

文件 `/proc/sys/vm/dirty_expire_centisecs` 规定了在 `pdflush` 刷新之前脏缓冲区必须达到的“年龄”（以 1%秒为单位）。位于同一目录下的其他文件则控制了 `pdflush` 操作的其他方面。

使所有写入同步：O_SYNC

调用 `open()`函数时如指定 `O_SYNC` 标志，则会使所有后续输出同步（`synchronous`）。

```
fd = open(pathname, O_WRONLY | O_SYNC);
```

调用 `open()`后，每个 `write()`调用会自动将文件数据和元数据刷新到磁盘上（即，按照 `Synchronized I/O file integrity completion` 的要求执行写操作）。

早期 BSD 系统曾使用 `O_FSYNC` 标志来提供 `O_SYNC` 标志的功能。在 `glibc` 库中，将 `O_FSYNC` 定义为与 `O_SYNC` 标志同义。

O_SYNC 对性能的影响

采用 `O_SYNC` 标志（或者频繁调用 `fsync()`、`fdatasync()` 或 `sync()`）对性能的影响极大。表 13-3 所示为采用不同缓冲区大小，在有、无 `O_SYNC` 标识的情况下将一百万字节写入一个（位于 `ext2` 文件系统上的）新创建文件所需要的时间。运行（随本书一同发行源码中的 `filebuff/write_bytes.c` 程序）结果取自于 `vanilla 2.6.30` 内核以及块大小为 4096 字节的 `ext2` 文件系统。每行数据均为在给定缓冲区大小下运行 20 次的平均值。

从表中可以看出，`O_SYNC` 标志使运行总用时大为增加——在缓冲区为 1 字节的情况下，运行时间相差 1000 多倍。还要注意，以 `O_SYNC` 标志执行写操作时运行总用时和 CPU 时间之间的巨大差异。这是因为系统在将每个缓冲区中数据向磁盘传递时会把程序阻塞起来。

表 13-3 所示的结果中还略去了使用 `O_SYNC` 时影响性能的一个深层次因素。现代磁盘驱动器均内置大型高速缓存，而默认情况下，使用 `O_SYNC` 只是将数据传递到该缓存中。如果禁用磁盘上的高速缓存（使用命令 `hdparm -W0`），那么 `O_SYNC` 对性能的影响将变得更为极端。在缓冲区大小为 1 字节的情况下，运行总用时从 1030 秒攀升到 16000 秒左右。而当缓冲区大小为 4096 字节时，运行总用时也会从 0.34 秒上升到 4 秒。

表 13-3: `O_SYNC` 标志对写入 1MB 速度的影响

BUF_SIZE	所需时间（秒）			
	无 O_SYNC		有 O_SYNC	
	运行总用时	总 CPU 时间	运行总用时	总 CPU 时间
1	0.73	0.73	1030	98.8
16	0.05	0.05	65.0	0.40
256	0.02	0.02	4.07	0.03
4096	0.01	0.01	0.34	0.03

总之，如果需要强制刷新内核缓冲区，那么在设计应用程序时就应考虑是否可以使用大尺寸的 `write()` 缓冲区，或者在调用 `fsync()` 或 `fdatasync()` 时谨慎行事，而不是在打开文件时就使用 `O_SYNC` 标志。

O_DSYNC 和 O_RSYNC 标志

`SUSv3` 规定了两个与同步 I/O 有关的、更为细化的打开文件状态标志：`O_DSYNC` 和 `O_RSYNC`。

`O_DSYNC` 标志要求写操作按照 `synchronized I/O data integrity completion` 来执行（类似于 `fdatasync()`）。与之相映成趣的是 `O_SYNC` 标志，遵从 `synchronized I/O file integrity completion`（类似于 `fsync()` 函数）。

`O_RSYNC` 标志是与 `O_SYNC` 标志或 `O_DSYNC` 标志配合一起使用的，将这些标志对写操作的作用结合到读操作中。如果在打开文件时同时指定 `O_RSYNC` 和 `O_DSYNC` 标志，那么就意味着会遵照 `synchronized I/O data integrity completion` 的要求来完成所有后续读操作（即，在执行读操作之前，像执行 `O_DSYNC` 标志一样完成所有待处理的写操作）。而在打开

文件时指定 `O_RSYNC` 和 `O_SYNC` 标志，则意味着会遵照 `synchronized I/O file integrity completion` 的要求来完成所有后续读操作（即，在执行读操作之前，像执行 `O_SYNC` 标志一样完成所有待处理的写操作）。

2.6.33 版本之前的 Linux 内核并未实现 `O_DSYNC` 和 `O_RSYNC` 标志。glibc 头文件当时只是将这些常量定义为 `O_SYNC` 标志。（以上描述实际上不适用于 `O_RSYNC` 标志，因为 `O_SYNC` 与读操作无关。）

始于 2.6.33 版本，Linux 内核实现了 `O_DSYNC` 标志的功能，而 `O_RSYNC` 标志的功能则可望在未来的版本中添加。

在 2.6.33 版本之前，Linux 内核并未完全实现 `O_SYNC` 的语义，而是将其实现为 `O_DSYNC` 标志。针对基于较老内核而构建的应用程序，为保证其行为的一致性，与老版 GNU C 函数库链接的应用程序会继续为 `O_SYNC` 标志提供 `O_DSYNC` 标志的语义，即使在 Linux 2.6.33 及更高版本的运行环境下。

13.4 I/O 缓冲小结

图 13-1 概括了 `stdio` 函数库和内核所采用的缓冲（针对输出文件），以及对各种缓冲类型的控制机制。从图中自上而下，首先是通过 `stdio` 库将用户数据传递到 `stdio` 缓冲区，该缓冲区位于用户态内存区。当缓冲区填满时，`stdio` 库会调用 `write()` 系统调用，将数据传递到内核高速缓冲区（位于内核态内存区）。最终，内核发起磁盘操作，将数据传递到磁盘。

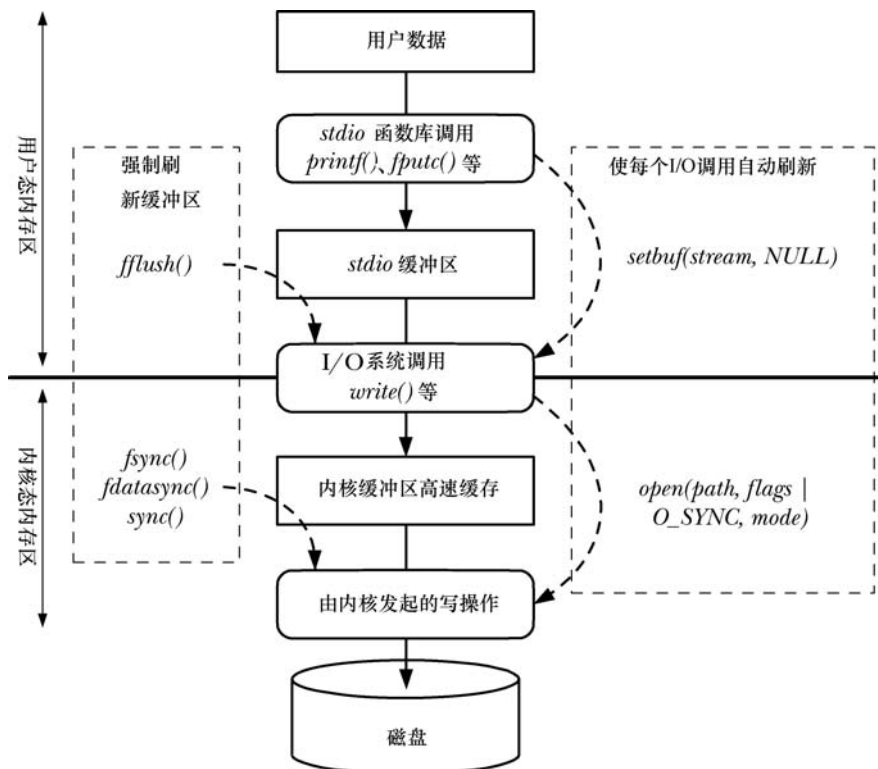


图 13-1: I/O 缓冲小结

图 13-1 左侧所示为可于任何时刻显式强制刷新各类缓冲区的调用。图右侧所示为促使刷新自动化的调用：一是通过禁用 `stdio` 库的缓冲，二是在文件输出类的系统调用中启用同步，从而使每个 `write()`调用立刻刷新到磁盘。

13.5 就 I/O 模式向内核提出建议

`posix_fadvise()`系统调用允许进程就自身访问文件数据时可能采取的模式通知内核。

```
#define _XOPEN_SOURCE 600
#include <fcntl.h>

int posix_fadvise(int fd, off_t offset, off_t len, int advice);

Returns 0 on success, or a positive error number on error
```

内核可以（但不必非要）根据 `posix_fadvise()`所提供的信息来优化对缓冲区高速缓存的使用，进而提高进程和整个系统的性能。调用 `posix_fadvise()`对程序语义并无影响。

参数 `fd` 所指为一文件描述符，调用期望通知内核进程对 `fd` 指代文件的访问模式。参数 `offset` 和 `len` 确定了建议所适用的文件区域。`offset` 指定了区域起始的偏移量，`len` 指定了区域的大小（以字节数为单位）。`len` 为 0 表示从 `offset` 开始，直至文件结尾。（在内核 2.6.6 版本之前，`len` 为 0 就表示长度为 0 个字节。）

参数 `advice` 表示进程期望对文件采取的访问模式。具体为下列参数之一：

POSIX_FADV_NORMAL

进程对访问模式并无特别建议。如果没有建议，这就是默认行为。在 Linux 中，该操作将文件预读窗口大小置为默认值（128KB）。

POSIX_FADV_SEQUENTIAL

进程预计会从低偏移量到高偏移量顺序读取数据。在 Linux 中，该操作将文件预读窗口大小置为默认值的两倍。

POSIX_FADV_RANDOM

进程预计以随机顺序访问数据。在 Linux 中，该选项会禁用文件预读。

POSIX_FADV_WILLNEED

进程预计会在不久的将来访问指定的文件区域。内核将由 `offset` 和 `len` 指定区域的文件数据预先填充到缓冲区高速缓存中。后续对该文件的 `read()`调用将不会阻塞磁盘 I/O，只需从缓冲区高速缓存中抓取数据即可。对于从文件读取的数据在缓冲区高速缓存中能保留多长时间，内核并无保证。如果其他进程或内核的活动对内存存在强劲需求，那么最终会重用到这些页面。换言之，如果内存压力高，程序员就应该确保 `posix_fadvise()`调用和后续 `read()`调用间的总运行时长较短。（Linux 特有的系统调用 `readahead()`提供了与 `POSIX_FADV_WILLNEED` 操作等效的功能。）

POSIX_FADV_DONTNEED

进程预计在不久的将来将不会访问指定的文件区域。这一操作给内核的建议是释放相关的高速缓存页面（如果存在的话）。在 Linux 中，该操作将分两步执行。首先，如果底层设备目前没有挤满一系列排队的写操作请求，那么内核会对指定区域中已修改的页面进行刷新。之后，

内核会尝试释放该区域的高速缓存页面。仅当该区域中已修改的页面在第一步中成功写入底层设备时，第二步才可能操作成功，也就是说，在该设备的写入操作请求没有发生拥塞的情况下。因为应用程序无法控制设备的拥塞（congestion），所以要确保释放高速缓存页面，变通的方法之一是在 `POSIX_FADV_DONTNEED` 操作之前对指定的参数 `fd` 调用 `sync()` 或 `fdatasync()`。

POSIX_FADV_NOREUSE

进程预计会一次性地访问指定文件区域，不再复用。这等于提示内核对指定区域访问一次后即可释放页面。在 Linux 中，该操作目前不起作用。

对 `posix_fadvise()` 的规范是 SUSv3 中的新增内容，并非所有 UNIX 实现都支持该接口。Linux 内核从 2.6 版本开始提供 `posix_fadvise()`。

13.6 绕过缓冲区高速缓存：直接 I/O

始于内核 2.4，Linux 允许应用程序在执行磁盘 I/O 时绕过缓冲区高速缓存，从用户空间直接将数据传递到文件或磁盘设备。有时也称此为直接 I/O（direct I/O）或者裸 I/O（raw I/O）。

此处的描述细节为 Linux 所特有，SUSv3 并未对其进行规范。尽管如此，大多数 UNIX 实现均对设备和文件提供了某种形式的直接 I/O 访问。

有时会将直接 I/O 误认为获取快速 I/O 性能的一种手段。然而，对于大多数应用而言，使用直接 I/O 可能会大大降低性能。这是因为为了提高 I/O 性能，内核针对缓冲区高速缓存做了不少优化，其中包括：按顺序预读取，在成簇（clusters）磁盘块上执行 I/O，允许访问同一文件的多个进程共享高速缓存的缓冲区。应用如使用了直接 I/O 将无法受益于这些优化举措。直接 I/O 只适用于有特定 I/O 需求的应用。例如数据库系统，其高速缓存和 I/O 优化机制均自成一體，无需内核消耗 CPU 时间和内存去完成相同任务。

可针对一个单独文件或块设备（比如，一块磁盘）执行直接 I/O。要做到这点，需要在调用 `open()` 打开文件或设备时指定 `O_DIRECT` 标志。

`O_DIRECT` 标志自内核 2.4.10 开始有效，并非所有 Linux 文件系统和内核版本都支持该标志。绝大多数原生（native）文件系统都支持 `O_DIRECT`，但许多非 UNIX 文件系统（比如 VFAT）则不支持。对于所关注的文件系统，有必要进行相关测试（若文件系统不支持 `O_DIRECT`，则 `open()` 将失败并返回错误号 `EINVAL`）或是阅读内核源码，以此来加以验证。

若一进程以 `O_DIRECT` 标志打开某文件，而另一进程以普通方式（即使用了高速缓存缓冲区）打开同一文件，则由直接 I/O 所读写的数据与缓冲区高速缓存中内容之间不存在一致性。应尽量避免这一场景。

`raw(8)` 手册页描述了一个获取对磁盘设备进行原始访问的老技术（现在已过时）。

直接 I/O 的对齐限制

因为直接 I/O（针对磁盘设备和文件）涉及对磁盘的直接访问，所以在执行 I/O 时，必须遵守一些限制。

- 用于传递数据的缓冲区，其内存边界必须对齐为块大小的整数倍。
- 数据传输的开始点，亦即文件和设备的偏移量，必须是块大小的整数倍。
- 待传递数据的长度必须是块大小的整数倍。

不遵守上述任一限制均将导致 EINVAL 错误。在上述列表中，块大小（block size）指设备的物理块大小（通常为 512 字节）。

当执行直接 I/O 时，Linux 2.4 比 Linux 2.6 限制更为严格：对齐、长度及偏移量必须是底层文件系统逻辑块大小的整数倍。（典型文件系统的逻辑块大小为 1024、2048 或 4096 字节。）

示例程序

程序清单 13-1 提供了一个使用 O_DIRECT 标志打开一个文件读取数据的简单例子。该程序可指定多达 4 个命令行参数，依次为要读取的文件、要从文件中读取的字节数、读之前在文件中定位（seek）的偏移量和传递给 read() 的数据缓冲区对齐。最后两个为可选参数，默认值分别为 0 字节和 4096 字节。下面是运行该程序的一些示例：

```
$ ./direct_read /test/x 512          Read 512 bytes at offset 0
Read 512 bytes                      Succeeds
$ ./direct_read /test/x 256
ERROR [EINVAL Invalid argument] read Length is not a multiple of 512
$ ./direct_read /test/x 512 1
ERROR [EINVAL Invalid argument] read Offset is not a multiple of 512
$ ./direct_read /test/x 4096 8192 512
Read 4096 bytes                     Succeeds
$ ./direct_read /test/x 4096 512 256
ERROR [EINVAL Invalid argument] read Alignment is not a multiple of 512
```

程序清单 13-1 中程序使用 memalign() 函数来分配一块内存，其内存块与第一个参数的整数倍对齐。7.1.4 节对 memalign() 函数有所描述。

程序清单 13-1：使用 O_DIRECT 跳过缓冲区高速缓存

```
----- filebuff/direct_read.c
#define _GNU_SOURCE /* Obtain O_DIRECT definition from <fcntl.h> */
#include <fcntl.h>
#include <malloc.h>
#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    int fd;
    ssize_t numRead;
    size_t length, alignment;
    off_t offset;
    void *buf;

    if (argc < 3 || strcmp(argv[1], "--help") == 0)
        usageErr("%s file length [offset [alignment]]\n", argv[0]);
    length = getLong(argv[2], GN_ANY_BASE, "length");
    offset = (argc > 3) ? getLong(argv[3], GN_ANY_BASE, "offset") : 0;
    alignment = (argc > 4) ? getLong(argv[4], GN_ANY_BASE, "alignment") : 4096;

    fd = open(argv[1], O_RDONLY | O_DIRECT);
    if (fd == -1)
        errExit("open");
```

```

/* memalign() allocates a block of memory aligned on an address that
   is a multiple of its first argument. The following expression
   ensures that 'buf' is aligned on a non-power-of-two multiple of
   'alignment'. We do this to ensure that if, for example, we ask
   for a 256-byte aligned buffer, then we don't accidentally get
   a buffer that is also aligned on a 512-byte boundary.

   The '(char *)' cast is needed to allow pointer arithmetic (which
   is not possible on the 'void *' returned by memalign()). */

buf = (char *) memalign(alignment * 2, length + alignment) + alignment;
if (buf == NULL)
    errExit("memalign");

if (lseek(fd, offset, SEEK_SET) == -1)
    errExit("lseek");

numRead = read(fd, buf, length);
if (numRead == -1)
    errExit("read");
printf("Read %ld bytes\n", (long) numRead);

exit(EXIT_SUCCESS);
}

```

filebuff/direct_read.c

13.7 混合使用库函数和系统调用进行文件 I/O

在同一文件上执行 I/O 操作时，还可以将系统调用和标准 C 语言库函数混合使用。fileno() 和 fdopen() 函数有助于完成这一工作。

```

#include <stdio.h>

int fileno(FILE *stream);
           Returns file descriptor on success, or -1 on error

FILE *fdopen(int fd, const char *mode);
           Returns (new) file pointer on success, or NULL on error

```

给定一个（文件）流，fileno() 函数将返回相应的文件描述符（即 stdio 库在该流上已经打开的文件描述符）。随即可以在诸如 read()、write()、dup() 和 fcntl() 之类的 I/O 系统调用中正常使用该文件描述符。

fdopen() 函数与 fileno() 函数的功能相反。给定一个文件描述符，该函数将创建了一个使用该描述符进行文件 I/O 的相应流。mode 参数与 fopen() 函数中 mode 参数含义相同。例如，r 为读，w 为写，a 为追加。若该参数与文件描述符 fd 的访问模式不一致，则对 fdopen() 的调用将失败。

fdopen() 函数对非常规文件描述符特别有用。正如后续章节将提及的，创建套接字和管道的系统调用总是返回文件描述符。为了在这些文件类型上使用 stdio 库函数，必须使用 fdopen() 函数来创建相应文件流。

当使用 stdio 库函数，并结合系统 I/O 调用来实现对磁盘文件的 I/O 操作时，必须将缓冲

问题牢记于心。I/O 系统调用会直接将数据传递到内核缓冲区高速缓存，而 `stdio` 库函数会等到用户空间的流缓冲区填满，再调用 `write()` 将其传递到内核缓冲区高速缓存。请考虑如下向标准输出写入的代码：

```
printf("To man the world is twofold, ");
write(STDOUT_FILENO, "in accordance with his twofold attitude.\n", 41);
```

通常情况下，`printf()` 函数的输出往往在 `write()` 函数的输出之后出现。因此，代码产生如下输出：

```
in accordance with his twofold attitude.
To man the world is twofold,
```

将 I/O 系统调用和 `stdio` 函数混合使用时，使用 `fflush()` 来规避这一问题，是明智之举。也可以使用 `setvbuf()` 或 `setbuf()` 使缓冲区失效，但这样做可能会影响应用的 I/O 性能，因为每个输出操作将引起一次 `write()` 系统调用。

要将 I/O 系统调用和 `stdio` 函数混合使用，SUSv3 针对此类应用的要求有所规范。详情参见系统接口卷（System Interfaces (XSH)）“通用信息”一章中“文件描述符和标准 I/O 流”一节。

13.8 总结

输入输出数据的缓冲由内核和 `stdio` 库完成。有时可能希望阻止缓冲，但这需要了解其对应用程序性能的影响。可以使用各种系统调用和库函数来控制内核和 `stdio` 缓冲，并执行一次性的缓冲区刷新。

进程使用 `posix_fadvise()` 函数，可就进程对特定文件可能采取的数据访问模式向内核提出建议。内核可籍此来优化对缓冲区高速缓存的应用，进而提高 I/O 性能。

在 Linux 环境下，`open()` 所特有的 `O_DIRECT` 标识允许特定应用跳过缓冲区高速缓存。

在对同一个文件执行 I/O 操作时，`fileno()` 和 `fdopen()` 有助于系统调用和标准 C 语言库函数的混合使用。给定一个流，`fileno()` 将返回相应的文件描述符，`fdopen()` 则反其道而行之，针对指定的打开文件描述符创建一个新的流。

补充信息

[Bach, 1986] 描述了 System V 中缓冲区高速缓存的实现和优势。[Goodheart & Cox, 1994] 和 [Vahalia, 1996] 也描述了 System V 缓冲区高速缓存的基本原理和实现。更多关于 Linux 环境下的相关信息参见 [Bovet & Cesati, 2005] 和 [Love, 2010]。

13.9 练习

- 13-1.** 使用 shell 内嵌的 `time` 命令，测算程序清单 4-1(copy.c) 在当前环境下的用时。
- 使用不同的文件和缓冲区大小进行试验。编译应用程序时使用 `-DBUF_SIZE=nbytes` 选项可设置缓冲区大小。
 - 对 `open()` 的系统调用加入 `O_SYNC` 标识，针对不同大小的缓冲区，速度存在多大差异？

c) 在一系列文件系统（比如，ext3、XFS、Btrfs 和 JFS）中执行这些计时测试。结果相似吗？当缓冲区大小从小变大时，用时趋势相同吗？

13-2. 测定 filebuff/write_bytes.c（随本书发行版提供源码）程序在不同的缓冲区大小以及文件系统下的用时。

13-3. 如下语句的执行效果是什么？

```
fflush(fp);  
fsync(fileno(fp));
```

13-4. 试解释取决于将标准输出重定向到终端还是磁盘文件，为什么如下代码的输出结果不同。

```
printf("If I had more time, \n");  
write(STDOUT_FILENO, "I would have written you a shorter letter.\n", 43);
```

13-5. tail [-n num] file 命令打印名为 file 文件的最后 num 行（默认为 10 行）。使用 I/O 系统调用（lseek()、read()、write()等）来实现该命令。牢记本章所描述的缓冲问题，力求实现的高效性。

第 14 章

系统编程概念

本书第 4 章、第 5 章及第 13 章介绍了文件 I/O，尤其侧重于常规（磁盘）文件。本章和后续几章则会深入探讨与文件相关的一系列主题。

- 本章会介绍文件系统。
- 第 15 章将会讨论与文件相关的各种属性，其中包括时间戳、所有权以及权限。
- 第 16 章和第 17 章则会关注 Linux 2.6 的两个新特性：扩展属性和访问控制列表(ACL)。扩展属性可将任意元数据与一文件进行关联，而 ACL 则是对传统 UNIX 文件权限模型的扩展。
- 第 18 章将讨论目录和链接。

文件系统是对文件和目录的组织集合，本章的绝大多数内容都与文件系统相关。本章会解释一系列与文件系统有关的概念，举例时将采用传统的 Linux ext2 文件系统。此外，本章还会简要介绍一些 Linux 支持的日志文件系统。

在本章结尾，将会讨论用于挂载（mount）和卸载（unmount）文件系统的系统调用，以及用来获取已挂载文件系统信息的库函数。

14.1 设备专用文件（设备文件）

本章会经常提到磁盘设备，因此这里先简要介绍一下设备文件的概念。

设备专用文件与系统的某个设备相对应。在内核中，每种设备类型都有与之相对应的设备驱动程序，用来处理设备的所有 I/O 请求。设备驱动程序属内核代码单元，可执行一系列操作，（通常）与相关硬件的输入/输出动作相对应。由设备驱动程序提供的 API 是固定的，包含的操作对应于系统调用 `open()`、`close()`、`read()`、`write()`、`mmap()` 以及 `ioctl()`。每个设备驱动程序所提供的接口一致，这隐藏了每个设备在操作方面的差异，从而满足了 I/O 操作的通用性（请参见 4.2 节）。

某些设备是实际存在的，比如鼠标、磁盘和磁带设备。而另一些设备则是虚拟的，亦即并不存在相应硬件，但内核会（通过设备驱动程序）提供一种抽象设备，其所携带的 API 与

真实设备一般无异。

可将设备划分为以下两种类型。

- 字符型设备基于每个字符来处理数据。终端和键盘都属于字符型设备。
- 块设备则每次处理一块数据。块的大小取决于设备类型，但通常为 512 字节的倍数。

磁盘和磁带设备都属于块设备。

与其他类型的文件一样，设备文件总会出现于文件系统中，通常位于 `/dev` 目录下。超级用户可使用 `mknod` 命令创建设备文件，特权级程序（`CAP_MKNOD`）执行 `mknod()` 系统调用亦可完成相同任务。

本书不会对 `mknod()`（make file-system i-node 创建，文件系统 i 节点）系统调用做详细介绍，因为该系统调用的用法一目了然，而且如今仅用于创建设备文件，一般应用程序鲜有问津。当然，也可以使用 `mknod()` 创建 FIFO（参见 44.7 节），但最好使用 `mkfifo()` 函数来完成该任务。早先，某些 UNIX 实现会使用 `mknod()` 来创建目录，但如今已为 `mkdir()` 系统调用所取代。然而，还有一些 UNIX 实现（Linux 不在此列），为了保持向后兼容性，仍然在 `mknod()` 中保留了这一能力，详情请见 `mknod(2)` 手册页。

在 Linux 的早期版本中，`/dev` 包含了系统中所有可能设备的条目，即使某些设备实际并未与系统连接。这意味着 `/dev` 会包含数以千计的未用设备项，从而导致了两个缺点：其一，对于需要扫描该目录内容的应用而言，降低了程序的执行速度；其二，根据该目录下的内容无法发现系统中实际存在哪些设备。Linux 2.6 运用 `udev` 程序解决了上述问题。该程序所依赖的 `sysfs` 文件系统，是装载于 `/sys` 下的伪文件系统，将设备和其他内核对象的相关信息导出至用户空间。

[Kroah-Hartman, 2003] 一书简要介绍了 `udev`，并概述了该程序较之于 `devfs` 的优势，后者是 Linux 2.4 内核对此类问题的解决方案。与 `sysfs` 文件系统有关的内容可见诸于 Linux 2.6 内核源码文件 `Documentation/filesystems/sysfs.txt` 和 [Mochel, 2005] 一书。

设备 ID

每个设备文件都有主、辅 ID 号各一。主 ID 号标识一般的设备等级，内核会使用主 ID 号查找与该类设备相应的驱动程序。辅 ID 号能够在一般等级中唯一标识特定设备。命令 `ls -l` 可显示出设备文件的主、辅 ID。

设备文件的 i 节点中记录了设备文件的主、辅 ID（本章第 4 节将介绍 i 节点）。每个设备驱动程序都会将自己与特定主设备号的关联关系向内核注册，藉此建立设备专用文件和设备驱动程序之间的关系。内核是不会使用设备文件名来查找驱动程序的。

在 Linux 2.4 以及更早的版本中，系统的设备总数受限于这一事实：设备的主、辅 ID 只能由 8 位数来表示。加之主设备 ID 固定不变，且为统一分配（由 Linux 命名和编号机构分配，请见 <http://www.lanana.org>），使得上述问题更为严重。Linux 2.6 采用了更多位数来存放主、辅 ID（分别为 12 位和 20 位），从而缓解了这一问题。

14.2 磁盘和分区

常规文件和目录通常都存放在硬盘设备里。（其他设备也能存放文件和目录，比如，

CD-ROM、flash 内存卡以及虚拟磁盘等，但这里主要关注的是硬盘设备。) 下面几节会介绍磁盘的组织方式，以及如何对其分区。

磁盘驱动器

硬盘驱动器是一种机械装置，由一个或多个高速旋转（每分钟旋转数以千计）的盘片组成。通过在磁盘上快速移动的读/写磁头，便可获取/修改磁盘表面的磁性编码信息。磁盘表面信息物理上存储于称为磁道（track）的一组同心圆上。磁道自身又被划分为若干扇区，每个扇区则包含一系列物理块。物理块的容量一般为 512 字节（或 512 的倍数），代表了驱动器可读/写的最小信息单元。

尽管现代磁盘速度很快，但读写磁盘信息耗时依然不菲。首先，磁头要移动到相应磁道（寻道时间）；然后，在相应扇区旋转到磁头下之前，驱动器必须一直等待（旋转延迟）；最后，还要从所请求的块上传输数据（传输时间）。执行上述操作所耗费的时间总量通常以毫秒为单位。相形之下，同样的时间可供现代 CPU 执行数百万条指令。

磁盘分区

可将每块磁盘划分为一个或多个（不重叠的）分区。内核则将每个分区视为位于 `/dev` 路径下的单独设备。

系统管理员可使用 `fdisk` 命令来决定磁盘分区的编号、大小和类型。命令 `fdisk -l` 会列出磁盘上的所有分区。`Linux` 专有文件 `/proc/partitions` 记录了系统中每个磁盘分区的主辅设备编号、大小和名称。

磁盘分区可容纳任何类型的信息，但通常只会包含以下之一。

- **文件系统**：用来存放常规文件，请参阅本章第 3 节。
- **数据区域**：可做为裸设备对其进行访问，请参阅 13.6 节（一些数据库管理系统会使用该技术）。
- **交换区域**：供内核的内存管理之用。

可通过 `mkswap(8)` 命令来创建交换区域。特权级进程 (`CAP_SYS_ADMIN`) 可利用 `swapon()` 系统调用向内核报告将磁盘分区用作交换区域。`swapoff()` 系统调用则会执行反向功能——告之内核，停止将磁盘分区用作交换区域。尽管 `SUSv3` 并未对上述系统调用进行规范，但它们却获得了许多 `UNIX` 实现的支持。其他信息请参考 `swapon(2)`、`swapon(8)` 手册页。

可使用 `Linux` 专有文件 `/proc/swaps` 来查看系统中当前已激活交换区域的信息。其中包括每个交换区域的大小，以及在用交换区域的个数。

14.3 文件系统

文件系统是对常规文件和目录的组织集合。用于创建文件系统的命令是 `mkfs`。

`Linux` 的强项之一便是支持种类繁多的文件系统，如下所示。

- 传统的 `ext2` 文件系统。
- 各种原生（native）`UNIX` 文件系统，比如，`Minix`、`System V` 以及 `BSD` 文件系统。

- 微软的 FAT、FAT32 以及 NTFS 文件系统。
- ISO 9660 CD-ROM 文件系统。
- Apple Macintosh 的 HFS。
- 一系列网络文件系统，包括广为使用的 SUN NFS（Linux 对 NFS 的实现信息请参见 <http://nfs.sourceforge.net/>）、IBM 和微软的 SMB、Novell NCP 以及 Carnegie Mellon 大学开发的 Coda 文件系统。
- 一系列日志文件系统，包括 ext3、ext4、Reiserfs、JFS、XFS 以及 Btrfs。

从 Linux 的专有文件/proc/filesystems 中可以查看当前为内核所知的文件系统类型。

Linux 2.6.14 中，添加了 FUSE（用户空间文件系统）工具。采用这一机制，可为内核添加挂钩（hook），以便以用户空间程序来完整实现文件系统，而无需对内核进行修补或重新编译。详细信息请见 <http://fuse.sourceforge.net/>。

ext2 文件系统

多年来，ext2（扩展文件系统二世）是 Linux 上使用最为广泛的文件系统，也是原始 Linux 文件系统——ext 的继任者。近来，随着各种日志文件系统的兴起，对 ext2 的使用也日趋减少。有时，在介绍通用文件系统概念时，以一款特定的文件系统实现为例会容易一些，出于这一目的，本章将以 ext2 为例来介绍文件系统。

ext2 文件系统由 Remy Card 编写。ext2 的源码篇幅不大（约 5000 行 C 语言代码），是其他几种文件系统实现的原型。ext2 文件系统的主页为 <http://e2fsprogs.sourceforge.net/ext2.html>。该站点上有一篇概括 ext2 实现的优秀论文。此外，David Rusling 所著的在线书籍“*The Linux kernel*”（可从 <http://www.tldp.org/> 下载）对 ext2 也有描述。

文件系统结构

在文件系统中，用来分配空间的基本单位是逻辑块，亦即文件系统所在磁盘设备上若干连续的物理块。例如，在 ext2 文件系统上，逻辑块的大小为 1024、2048 或 4096 字节。（使用 mkfs(8) 命令创建文件系统时，可指定逻辑块的大小作为命令行参数。）

特权级程序（CAP_SYS_RAWIO）可利用 ioctl() 的 FIBMAP 操作，来判定文件指定逻辑块的物理位置。该调用的第三个参数是整型值，同时用于返回结果。调用之前，应将参数设置为逻辑块编号（第一个逻辑块编号为 0）；调用之后，其中返回的为存储该逻辑块的起始物理块编号。

图 14-1 所示为磁盘分区和文件系统之间的关系，以及一般文件系统的组成。



图 14-1：磁盘分区和文件系统布局

文件系统由以下几部分组成。

- 引导块：总是作为文件系统的首块。引导块不为文件系统所用，只是包含用来引导操作系统的信息。操作系统虽然只需一个引导块，但所有文件系统都设有引导块（其中的绝大多数都未使用）。
- 超级块：紧随引导块之后的一个独立块，包含与文件系统有关的参数信息，其中包括：
 - i 节点表容量；
 - 文件系统中逻辑块的大小；
 - 以逻辑块计，文件系统的大小；

驻留于同一物理设备上的不同文件系统，其类型、大小以及参数设置（比如，块大小）都可以有所不同。这也是将一块磁盘划分为多个分区的原因之一。

- i 节点表：文件系统中的每个文件或目录在 i 节点表中都对对应着唯一一条记录。这条记录登记了关乎文件的各种信息。下一节会深入讨论 i 节点。有时也将 i 节点表称为 i-list。
- 数据块：文件系统的大部分空间都用于存放数据，以构成驻留于文件系统之上的文件和目录。

就 ext2 文件系统而言，情况要比正文中的描述稍微复杂一点。在起始的引导块之后，ext2 文件系统被划分为一系列大小相等的块组（block group）。每个块组都包含了一份超级块的拷贝、与块组有关的参数信息，以及该块组的 i 节点表和数据块。ext2 文件系统会尽量在同一块组内存储一个文件的所有块，以期在对文件线性访问时缩短寻道时间。更多详情，请参考 Linux 源码 Documentation/filesystems/ext2.txt、dumpe2fs 程序的源代码（作为 e2fsprogs 软件包的一部分发布），以及[Bovet & Cesati, 2005]。

14.4 i 节点

针对驻留于文件系统上的每个文件，文件系统的 i 节点表会包含一个 i 节点（索引节点的简称）。对 i 节点的标识，采用的是 i 节点表中的顺延位置，以数字表示。文件的 i 节点号（或简称为 i 号）是 `ls -li` 命令所显示的第一列。i 节点所维护的信息如下所示。

- 文件类型（比如，常规文件、目录、符号链接，以及字符设备等）。
- 文件属主（亦称用户 ID 或 UID）。
- 文件属组（亦称为组 ID 或 GID）。
- 3 类用户的访问权限：属主（有时也称为用户）、属组以及其他用户（属主和属组用户之外的用户）。详情请见 15.4 节。
- 3 个时间戳：对文件的最后访问时间（`ls -lu` 所显示的时间）、对文件的最后修改时间（也是 `ls -l` 所默认显示的时间），以及文件状态的最后改变时间（`ls -lc` 所显示的最后改变 i 节点信息的时间）。值得注意的是，与其他 UNIX 实现一样，大多数 Linux 文件系统不会记录文件的创建时间。
- 指向文件的硬链接数量。
- 文件的大小，以字节为单位。
- 实际分配给文件的块数量，以 512 字节块为单位。这一数字可能不会简单等同于文件的字节大小，因为考虑文件中包含空洞（请参见 4.7 节）的情形，分配给文件的块数可能会低于根据文件正常大小（以字节为单位）所计算出的块数。

- 指向文件数据块的指针。

ext2 中的 i 节点和数据块指针

类似于大多数 UNIX 文件系统，ext2 文件系统在存储文件时，数据块不一定连续，甚至不一定按顺序存放（尽管 ext2 会尝试将数据块彼此靠近存储）。为了定位文件数据块，内核在 i 节点内维护有一组指针。图 14-2 所示为在 ext2 文件系统上完成上述任务的情况。

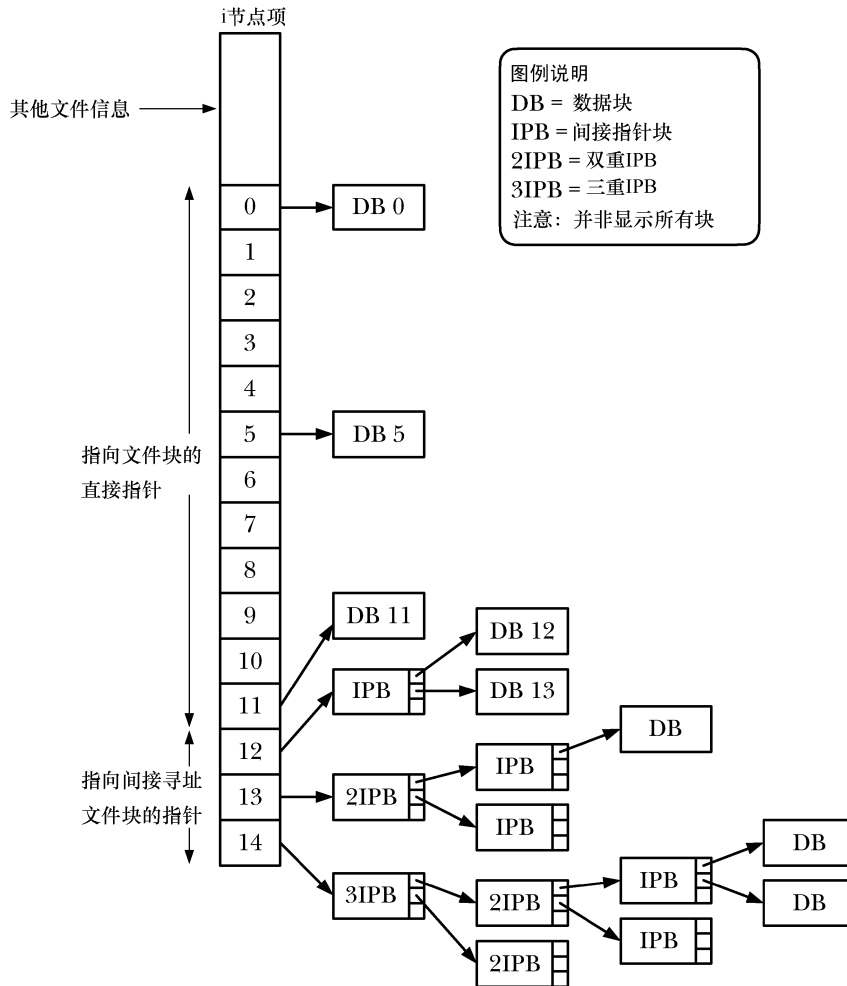


图 14-2: ext2 文件系统中文件的文件块结构

无需连续存储文件块，使得文件系统对磁盘空间的利用更为高效。特别是，还能降低空闲磁盘空间的碎片化程度，即因众多不连续空闲磁盘碎片（因其空间太小而无法使用）而导致的磁盘空间浪费。换言之，对空闲磁盘空间的高效利用，是以已分配磁盘空间中文件的碎片化为代价的。

在 ext2 中，每个 i 节点包含 15 个指针。其中的前 12 个指针（图 14-2 中编号为 0~11 的指针）指向文件前 12 个块在文件系统中的位置。接下来，是一个指向指针块的指针，提供了文件的第 13 个以及后续数据块的位置。指针块中指针的数量取决于文件系统中块的大小。每

个指针需占用 4 字节，因此指针的数量可能在 256（块容量为 1024 字节）~1024（块容量为 4096 字节）之间。这样就考虑了大型文件的情况。即便是对于巨型文件，第 14 个指针（图中编号为 13）是一个双重间接指针——指向指针块，其块中指针进而指向指针块，此块中指针最终才指向文件的数据块。只要有体量巨大的文件，就会随之产生更深一层的递进：图中 i 节点的最后一个指针属于三重间接指针。

这一貌似复杂的系统，其设计意图是为了满足多重需求。首先，该系统在维持 i 节点结构大小固定的同时，支持任意大小的文件。其次，文件系统既可以以不连续方式来存储文件块，又可通过 `lseek()` 随机访问文件，而内核只需计算所要遵循的指针。最后，对于在大多数系统中占绝对多数的小文件而言，这种设计满足了对文件数据块的快速访问：通过 i 节点的直接指针访问，一击必中。

试举一例，笔者对一个包含约 150 000 个文件的系统进行了度量。其中 30% 多的文件大小在 1000 字节以下，80% 的文件占用了 10 000 字节或者更少的空间。假定块的大小为 1024 字节，只要使用 12 个直接指针便能引用大小为 10000 字节及以下的文件，可访问总计 12288 字节的块。若块大小为 4096 字节，则该上限可达 49152 字节（系统中 95% 的文件大小都处于该容量限制之下）。

上述设计同样考虑了巨型文件的处理，对于大小为 4096 字节的块而言，理论上，文件大小可略高于 $1024 \times 1024 \times 1024 \times 4096$ 字节，或 4TB（4096 GB）。（之所以说“略高于”，是因为指针指向块的方式可以为直接、间接或双重间接。与三重间接指针所指向的范围相比，多出来的那些空间实在是微不足道。）

该设计的另一优点在于文件可以有黑洞（如 4.7 节所述）。文件系统只需将 i 节点和间接指针块中的相应指针打上标记（值 0），表明这些指针并未指向实际的磁盘块即可，而无需为文件黑洞分配空字节数据块。

14.5 虚拟文件系统（VFS）

Linux 所支持的各种文件系统，其实现细节均不相同。举例来说，这些差异包括文件块的分配方式，以及目录的组织方式。如果每个与文件打交道的程序都需要理解各种文件系统的具体细节，那么编写与各类文件系统交互的程序将近乎于不可能完成的任务。虚拟文件系统（VFS，有时也称为虚拟文件交换）是一种内核特性，通过为文件系统操作创建抽象层来解决上述问题（参见图 14-3）。VFS 背后的原理其实很直白。

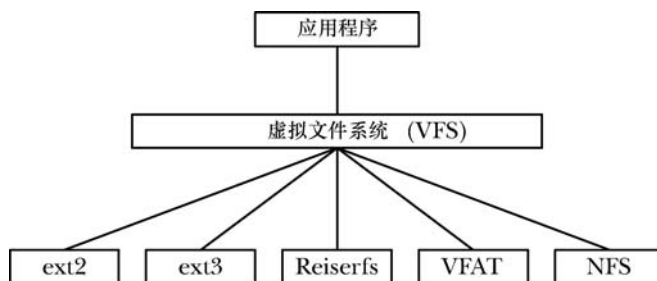


图 14-3: 虚拟文件系统

- VFS 针对文件系统定义了一套通用接口。所有与文件交互的程序都会按照这一接口来进行操作。
- 每种文件系统都会提供 VFS 接口的实现。

这样一来，程序只需理解 VFS 接口，而无需过问具体文件系统的实现细节。

VFS 接口的操作与涉及文件系统和目录的所有常规系统调用相对应，这些系统调用有 `open()`、`read()`、`write()`、`lseek()`、`close()`、`truncate()`、`stat()`、`mount()`、`umount()`、`mmap()`、`mkdir()`、`link()`、`unlink()`、`symlink()` 以及 `rename()`。

VFS 的抽象层建模精确仿照传统的 UNIX 文件系统模型。当然，还有一些文件系统，尤其是非 UNIX 文件系统，并不支持所有的 VFS 操作。（比如，微软的 VFAT 就不支持使用 `symlink()` 创建的符号链接概念。）对于这种情况，底层文件系统会将错误代码传回 VFS 层，表明不支持相应操作，而 VFS 随之会将错误代码传递给应用程序。

14.6 日志文件系统

`ext2` 文件系统是传统 UNIX 文件系统的优秀典范，自然也受制于其短板：系统崩溃之后，为确保文件系统的完整性，重启时必须对文件系统的一致性进行检查（`fsck`）。由于系统每次崩溃时，对文件的更新可能只完成了一部分，而文件系统元数据（目录项、`i` 节点信息以及文件数据块指针）也将处于不一致状态，一旦这一问题得不到修复，那么文件系统会遭到进一步破坏，因此上述举措实属必要。如有可能，就必须进行修复，否则，将会丢弃那些无法获取的信息（可能会包含文件数据）。

问题在于，一致性检查需要遍历整个文件系统。如果文件系统较小，只需几秒或几分钟便可完成。而在大型文件系统中，上述操作可能会历时数小时，这对于需要保持高可用性的系统来说（比如，网络服务器），情况就非常严重。

采用日志文件系统，则无需在系统崩溃后对文件进行漫长的一致性检查。在实际更新元数据之前，日志文件系统会将这些更新操作记录于专用的磁盘日志文件中。对元数据更新的记录是按其相关性分组（以事务的方式记录）进行的。在事务处理过程中，一旦系统崩溃，系统重启时便可利用日志重做（`redo`）任何不完整的更新，同时为文件系统恢复一致性状态。（借用数据库的说法，日志文件系统能够确保总是将文件元数据事务作为一个完整单元来提交。）系统崩溃之后，即便是超大型的日志文件系统，通常也会在几秒之内复原，因而对于有高可用性需求的系统极具吸引力。

日志文件系统最为昭著的臭名在于增加了文件更新的时间，当然，良好的设计可以降低这方面的开销。

某些日志文件系统只会确保文件元数据的一致性。由于不记录文件数据，因此一旦系统崩溃，可能会造成数据丢失。`ext3`、`ext4` 和 `Reiserfs` 文件系统提供了记录数据更新的选项，但若记录的东西过多，则会降低文件 I/O 的性能。

以下列出了 Linux 所支持的日志文件系统。

- `Reiserfs` 是首个被集成进内核（版本号为 2.4.1）的日志文件系统。`Reiserfs` 提供了一种名为 `tail packing`（或 `tail merging`）的特性：可将小文件（以及较大文件的最后一块）与文件元数据装入相同的磁盘块。而许多系统都拥有（或由应用程序创建了）众多小文

件，因此这会节省大量的磁盘空间。

- ext3 文件系统，源于一个旨在以最小改动为 ext2 追加日志功能的项目。从 ext2 升级到 ext3 非常简单（无需备份和恢复操作），还支持反向降级。内核版本 2.4.15 集成了 ext3。
- JFS 由 IBM 开发，内核版本 2.4.20 对其进行了集成。
- XFS (<http://oss.sgi.com/projects/xfs/>)最初是由 SGI (Silicon Graphics) 于 20 世纪 90 年代初期开发，所针对的是自己的私有 UNIX 实现: Irix。2001 年, XFS 被移植到了 Linux 平台，并成为自由软件项目。2.4.24 内核对其进行了集成。

配置内核时，可在“File systems”菜单下激活对不同文件系统支持的内核设置选项。

写作本书之际，还有两种提供了日志功能，且支持多种其他高级特性的文件系统尚在开发之中。

- ext4 文件系统 (<http://ext4.wiki.kernel.org/>) 是 ext3 文件系统的“接班人”。Linux 2.6.19 将其首个实现并入，内核的后续版本中又陆续添加了各种特性。ext4 的规划（或已实现的）特性包括 extents（预留连续存储块）、旨在降低文件碎片化的其他分配特性、在线文件系统的磁盘碎片整理、更为快捷的文件系统检查以及对纳秒级时间戳的支持。
- Btrfs (B-树 FS，一般读作“butter FS”，<http://btrfs.wiki.kernel.org/>) 是一种自下而上进行设计的新型文件系统，意在提供一系列现代化特性，其中包括 extents、可写快照（等价于对元数据和数据的日志功能）、对数据和元数据的校验和、在线文件系统检查、在线文件系统的磁盘碎片整理、高效利用空间的小文件打包存放和可检索目录。内核版本 2.6.29 中集成了该文件系统。

14.7 单根目录层级和挂载点

与其他 UNIX 系统一样，Linux 上所有文件系统上的文件都位于单根目录树下，树根就是根目录“/”。其他的文件系统都挂载在根目录之下，被视为整个目录层级的子树（subtree）。超级用户可使用如下命令来挂载文件系统：

```
$ mount device directory
```

这条命令会将名为 device 的文件系统挂接到目录层级中由 directory 所指定的目录，即文件系统的挂载点。可使用 unmount 命令卸载文件系统，然后在另一个挂载点再次挂载文件系统，从而改变文件系统的挂载点。

自 Linux 版本 2.4.19 以后，情况变得更为复杂。如今，内核支持针对每个进程的挂载命名空间（mount namespace）。这意味着每个进程都可能拥有属于自己的一组文件系统挂载点，因此进程视角下的单根目录层级彼此会有所不同。本书将在 28.2.1 节介绍 CLONE_NEWNS 标记时，对上述内容做深入讨论。

不带任何参数来执行 mount 命令，可以列出当前已挂载的文件系统，如下例所示（与实际输出相比，略有删减）：


```

$ mount
/dev/sda6 on / type ext4 (rw)
proc on /proc type proc (rw)
sysfs on /sys type sysfs (rw)
devpts on /dev/pts type devpts (rw,mode=0620,gid=5)
/dev/sda8 on /home type ext3 (rw,acl,user_xattr)
/dev/sda1 on /windows/C type vfat (rw,noexec,nosuid,nodev)
/dev/sda9 on /home/mtk/test type reiserfs (rw)

```

图 14-4 所示的部分目录及文件结构就出自于执行上述 mount 命令的系统。该图演示了将安装点映射到目录层级的方法。

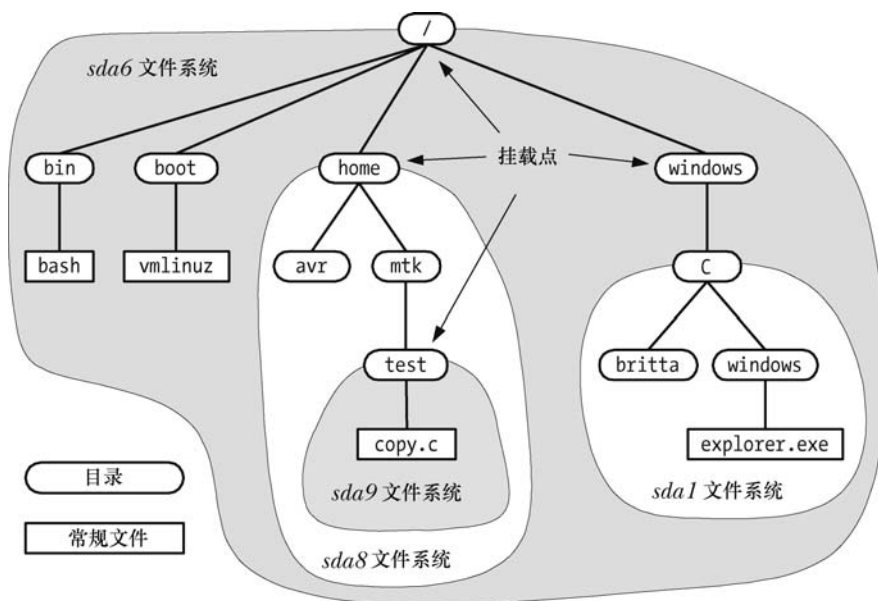


图 14-4: 演示文件系统挂载点的目录层级示例

14.8 文件系统的挂载和卸载

系统调用 `mount()` 和 `umount()` 运行特权级进程 (`CAP_SYS_ADMIN`) 以挂载或卸载文件系统。大多数 UNIX 实现都提供了这两个系统调用。不过，SUSv3 并未对其进行规范，因此其操作也随 UNIX 实现和文件系统的不同而不同。

在讨论这两个系统调用之前，需要先了解以下 3 个文件，其中包含了当前已挂载或可挂载的文件系统信息。

- 通过 Linux 特有的虚拟文件 `/proc/mounts`，可查看当前已挂载文件系统的列表。`/proc/mounts` 是内核数据结构的接口，因此总是包含已挂载文件系统的精确信息。

随着引入了前述的每进程挂载命名空间特性，如今，每个进程都拥有一个 `/proc/PID/mounts` 文件，其中会列出组成进程挂载空间的挂载点，而 `/proc/mounts` 只是指向 `/proc/self/mounts` 的符号链接。

- `mount(8)` 和 `umount(8)` 命令会自动维护 `/etc/mntab` 文件，该文件所包含的信息与 `/proc/mounts` 的内容相类似，只是略微详细一些。特别是，`etc/mntab` 包含了传递给

mount(8)的文件系统专有选项，这并未在/proc/mounts 中出现。但是，因为系统调用 mount()和 umount()并不更新/etc/mtab，如果某些挂载或卸载了设备的应用程序没有更新该文件，那么/etc/mtab 可能会变得不准确。

- /etc/fstab（由系统管理员手工维护）包含了对系统支持的所有文件系统的描述，该文件可供 mount(8)、umount(8)以及 fsck(8)所用。

/proc/mounts、/etc/mtab 和/etc/fstab 的格式相同，请参考 fstab(5)手册页。以下示例摘自 /proc/mounts 中的一条记录（一行）：

```
/dev/sda9 /boot ext3 rw 0 0
```

这条记录包含了 6 个字段。

1. 已挂载设备名。
2. 设备的挂载点。
3. 文件系统类型。
4. 挂载标志。上例的 rw 表示以可读写方式挂载文件系统。
5. 一个数字，dump(8)会使用其来控制对文件系统的备份操作。只有/etc/fstab 文件才会用到该字段和第 6 个字段，在/proc/mounts 和/etc/mtab 中，该字段总是为 0。
6. 一个数字，在系统引导时，用于控制 fsck(8)对文件系统的检查顺序。
getfsent(3)和 getmntent(3)手册页记录了用于从上述文件中读取记录的函数。

14.8.1 挂载文件系统：mount()

mount()系统调用将由 source 指定设备所包含的文件系统，挂载到由 target 指定的目录下。

```
#include <sys/mount.h>

int mount(const char *source, const char *target, const char *fstype,
          unsigned long mountflags, const void *data);

Returns 0 on success, or -1 on error
```

头两个参数分别命名为 source 和 target，其原因在于，除了将磁盘文件系统挂载到一目录下之外，mount()还可以执行其他任务。

参数 fstype 是一字符串，用来标识设备所含文件系统的类型，比如，ext4 或 btrfs。

参数 mountflags 为一位掩码，通过对表 14-1 中所示的 0 个或多个标志进行或（OR）操作而得出，稍后将做详细介绍。

表 14-1: 供 mount()使用的 mountflags 值

标 记	用 途
MS_BIND	建立绑定挂载（始于 Linux 2.4）
MS_DIRSYNC	同步更新路径（始于 Linux 2.6）
MS_MANDLOCK	允许强制锁定文件
MS_MOVE	以原子操作将挂载点移到新位置
MS_NOATIME	不更新文件的最后访问时间
MS_NODEV	不允许访问设备

标 记	用 途
MS_NODIRATIME	不更新目录的最后访问时间
MS_NOEXEC	不允许程序执行
MS_NOSUID	禁用 set-user-ID 和 set-group-ID 程序
MS_RDONLY	以只读方式挂载；不能修改或创建文件
MS_REC	递归挂载（始于 Linux 2.6.20）
MS_RELATIME	只有当最后访问时间早于最后修改时间或最后状态变更时间时，才对前者进行更新（始于 Linux 2.4.11）
MS_REMOUNT	使用新的 mountflags 和 data 重新挂载
MS_STRICTATIME	总是更新最后访问时间（始于 Linux 2.6.30）
MS_SYNCHRONOUS	使得所有文件和目录同步更新

mount()的最后一个参数 data 是一个指向信息缓冲区的指针，对其信息的解释则取决于文件系统。就大多数文件系统而言，该参数是一字符串，包含了以逗号分隔的选项设置。在 mount(8)手册页中，有这些选项的完整列表。若未见之于 mount(8)手册页，请查找相关文件系统的文档。

mountflags 参数是标志的位掩码，用来修改 mount()操作。在 mountflags 中，可以指定 0 到多个如下标志：

MS_BIND（始于 Linux 2.4）

用来建立绑定挂载。14.9.4 节将描述这一特性。如果指定了该标志，那么 mount()会忽略 fstype、data 参数，以及 mountflags 中除 MS_REC 之外的标志（见后续描述）。

MS_DIRSYNC（始于 Linux 2.6）

用来同步更新路径。该标志的效果类似于 open()的 O_SYNC 标志（参见 13.3 节），但只针对路径。后面介绍的 MS_SYNCHRONOUS 提供了 MS_DIRSYNC 功能的超集，可同时同步更新文件和目录。采用 MS_DIRSYNC 标志的应用程序在确保同步更新目录（比如，open(pathname, O_CREAT)、rename()、link()、unlink()、symlink()以及 mkdir()）的同时，还无需消耗同步更新文件所带来的成本。FS_DIRSYNC_FL 标志的用途与之相近，其区别在于可将 MS_DIRSYNC 应用于单个目录。此外，在 Linux 上，针对指代目录的文件描述符调用 fsync()，可对目标目录进行更新。（SUSv3 并未对 fsync()的这一 Linux 专有行为加以规范。）

MS_MANDLOCK

允许对该文件系统中的文件强行锁定记录。第 55 章将描述记录锁定。

MS_MOVE

将由 source 指定的现有挂载点移到由 target 指定的新位置，整个动作为一原子操作，不可分割。这与 mount(8)命令的--move 选项相对应。实际上，这等同于卸载子树，并将其重新装载到另一位置，只是卸载子树的时点并无意义¹。source 参数为一字符串，其内容应与前一个 mount()调用中的 target 相同。一旦使用了这一标志，那么 mount()将忽略 fstype、data 参数，

¹ 译者注：因为是原子操作。

以及 `mountflags` 中的其他标志。

MS_NOATIME

针对该文件系统中的文件，不更新其最后访问时间。与下面将要介绍的 `MS_NODIRATIME` 标志一样，使用该标志意在消除额外的磁盘访问，避免在每次访问文件时都去更新文件 `i` 节点。对某些应用程序来说，维护这一时间戳意义不大，而放弃这一做法还能显著提升性能。`MS_NOATIME` 标志与 `FS_NOATIME_FL` 标志（见 15.5 节）目的相似，区别在于可将 `FS_NOATIME_FL` 标志应用于单个文件。此外，Linux 还可以运用 `open()` 的 `O_NOATIME` 标志，来提供类似功能，可以针对打开的单个文件来选择这一行为（参见 4.3.1 节）。

MS_NODEV

不允许访问此文件系统上的块设备和字符设备。设计这一特性的目的是为了保障系统安全，规避如下情况：假设用户插入了可移动磁盘，而磁盘中又包含了可随意访问系统的设备专有文件。

MS_NODIRATIME

不更新此文件系统中目录的最后访问时间（该标志提供了 `MS_NOATIME` 标志的部分功能，`MS_NOATIME` 标志不会对所有文件类型的最后访问时间进行更新）。

MS_NOEXEC

不允许在此文件系统上执行程序（或脚本）。该标志用于文件系统包含非 Linux 可执行文件的场景。

MS_NOSUID

禁用此文件系统上的 `set-user-ID` 和 `set-group-ID` 程序。这属于安全特性，意在防止用户从可移动磁盘上运行 `set-user-ID` 和 `set-group-ID` 程序。

MS_RDONLY

以只读方式挂载文件系统，在此文件系统上既不能创建文件，也不能修改现有文件。

MS_REC (始于 Linux 2.4.11)

该标志与其他标志（比如，`MS_BIND`）结合使用，以递归方式将挂载动作施之于子树下的所有挂载。

MS_RELATIME (始于 Linux 2.6.20)

在此文件系统中，只有当文件最后访问时间戳的当前值¹小于或等于最后一次修改或状态变更的时间戳时，才对其进行更新。这不但吸取了 `MS_NOATIME` 性能上的一些优点，而且还可应用于如下场景：程序能了解到，自上次更新以来，有无读取过文件。自 Linux 2.6.30 以来，系统会默认提供 `MS_RELATIME` 的行为（除非明确指定了 `MS_NOATIME` 标志），要获取传统行为，必须使用 `MS_STRICTATIME` 标志。此外，只要文件最后访问时间戳距今超过 24 小时，即便其大于最近修改和状态改变时间戳，系统仍会更新该文件的最后访问时间戳。（该标志对于监控目录的系统程序来说极为有用，可以了解最近有无对文件进行过访问。）

MS_REMOUNT

针对已挂载的文件系统，改变其 `mountflag`（装备标记）和 `data`（数据）。（例如，令只读

¹ 译者注：即上次更新的时间。

文件系统可写。)使用该标志时, source 和 target 参数应该与最初用于 mount()系统调用的参数相同, 而对 fstype 参数则予以忽略。使用该标志可以避免对磁盘进行卸载和重新挂载, 在某些场合中, 这是不可能做到的。比方说, 如果有进程打开了文件系统上的文件, 或进程的当前工作目录位于文件系统之内(对 root 文件系统来说, 情况总是如此), 就无法卸载相应的文件系统。使用 MS_REMOUNT 的另一场景是 tmpfs (基于内存的) 文件系统, 一旦卸载了这一文件系统, 其内容便会丢失。并非所有的 mountflag 都是可修改的, 具体信息请参考 mount(2) 手册页。

MS_STRICTATIME (始于 Linux 2.6.30)

只要访问此文件系统上的文件, 就总是更新文件的最后访问时间戳。Linux 2.6.30 之前, 这是系统的默认行为。只要定义了 MS_STRICTATIME, 即使在 mountflag 中定义了 MS_NOATIME 和 MS_RELATIME, 也会将其忽略。

MS_SYNCHRONOUS

对文件系统上的所有文件和目录保持同步更新。(对文件来说, 就如同总是以 O_SYNC 标志调用 open()来打开文件一样。)

从内核 2.6.15 起, 为支持共享子树 (shared subtree) 的概念, Linux 提供了 4 个新的挂载标志, 分别是 MS_PRIVATE、MS_SHARED、MS_SLAVE、MS_UNBINDABLE。这些标志可与 MS_REC 结合使用, 从而将其效果传播至一挂载子树 (mount subtree) 下的所有子挂载 (submount)。设计共享子树的目的, 是为了支持某些高级文件系统特性, 比如, 每进程挂载命名空间 (请参见 28.2.1 一节对 CLONE_NEWNS 的描述), 以及用户空间的文件系统 (FUSE) 工具。共享子树机制允许以一种受控方式在挂载命名空间之间传播文件系统的挂载。关于共享子树的详细信息, 可查阅内核源码文件 Documentation/filesystems/sharedsubtree.txt 和[Viro & Pai, 2006]一书。

程序示例

程序清单 14-1 提供了对 mount(2)系统调用的命令行级接口。其实, 也是 mount(8)命令的简化版。以下 shell 会话日志演示了对该程序的运用。首先创建一个目录作为挂载点, 并挂载文件系统。

```
$ su Need privilege to mount a file system
Password:
# mkdir /testfs
# ./t_mount -t ext2 -o bsdgroups /dev/sda12 /testfs
# cat /proc/mounts | grep sda12 Verify the setup
/dev/sda12 /testfs ext3 rw 0 0 Doesn't show bsdgroups
# grep sda12 /etc/mtab
```

这里可以发现, 上面的 grep 命令并未产生任何输出, 因为该程序并未更新/etc/mtab。现继续以只读方式重新挂载文件系统:

```
# ./t_mount -f Rr /dev/sda12 /testfs
# cat /proc/mounts | grep sda12 Verify change
/dev/sda12 /testfs ext3 ro 0 0
```

从/proc/mounts 输出的字符串“ro”表明这是一次只读方式的挂载。
最后，再将挂载点移动至目录层级内的新位置：

```
# mkdir /demo
# ./t_mount -f m /testfs /demo
# cat /proc/mounts | grep sda12          Verify change
/dev/sda12 /demo ext3 ro 0
```

程序清单 14-1：使用 mount()

filesys/t_mount.c

```
#include <sys/mount.h>
#include "tspi_hdr.h"

static void
usageError(const char *progName, const char *msg)
{
    if (msg != NULL)
        fprintf(stderr, "%s", msg);

    fprintf(stderr, "Usage: %s [options] source target\n\n", progName);
    fprintf(stderr, "Available options:\n");
#define fpe(str) fprintf(stderr, "    " str) /* Shorter! */
    fpe("-t fstype      [e.g., 'ext2' or 'reiserfs']\n");
    fpe("-o data        [file system-dependent options,\n");
    fpe("                e.g., 'bsdgroups' for ext2]\n");
    fpe("-f mountflags    can include any of:\n");
#define fpe2(str) fprintf(stderr, "    " str)
    fpe2("b - MS_BIND      create a bind mount\n");
    fpe2("d - MS_DIRSYNC     synchronous directory updates\n");
    fpe2("l - MS_MANDLOCK     permit mandatory locking\n");
    fpe2("m - MS_MOVE        atomically move subtree\n");
    fpe2("A - MS_NOATIME     don't update atime (last access time)\n");
    fpe2("V - MS_NODEV      don't permit device access\n");
    fpe2("D - MS_NODIRATIME don't update atime on directories\n");
    fpe2("E - MS_NOEXEC     don't allow executables\n");
    fpe2("S - MS_NOSUID     disable set-user/group-ID programs\n");
    fpe2("r - MS_RDONLY     read-only mount\n");
    fpe2("c - MS_REC        recursive mount\n");
    fpe2("R - MS_REMOUNT    remount\n");
    fpe2("s - MS_SYNCHRONOUS make writes synchronous\n");
    exit(EXIT_FAILURE);
}

int
main(int argc, char *argv[])
{
    unsigned long flags;
    char *data, *fstype;
    int j, opt;

    flags = 0;
    data = NULL;
    fstype = NULL;

    while ((opt = getopt(argc, argv, "o:t:f:")) != -1) {
        switch (opt) {
```

```

    case 'o':
        data = optarg;
        break;

    case 't':
        fstype = optarg;
        break;

    case 'f':
        for (j = 0; j < strlen(optarg); j++) {
            switch (optarg[j]) {
                case 'b': flags |= MS_BIND;           break;
                case 'd': flags |= MS_DIRSYNC;       break;
                case 'l': flags |= MS_MANDLOCK;       break;
                case 'm': flags |= MS_MOVE;          break;
                case 'A': flags |= MS_NOATIME;        break;
                case 'V': flags |= MS_NODEV;         break;
                case 'D': flags |= MS_NODIRATIME;     break;
                case 'E': flags |= MS_NOEXEC;        break;
                case 'S': flags |= MS_NOSUID;        break;
                case 'r': flags |= MS_RDONLY;        break;
                case 'c': flags |= MS_REC;           break;
                case 'R': flags |= MS_REMOUNT;       break;
                case 's': flags |= MS_SYNCHRONOUS;    break;
                default: usageError(argv[0], NULL);
            }
        }
        break;

    default:
        usageError(argv[0], NULL);
}

if (argc != optind + 2)
    usageError(argv[0], "Wrong number of arguments\n");

if (mount(argv[optind], argv[optind + 1], fstype, flags, data) == -1)
    errExit("mount");

exit(EXIT_SUCCESS);
}

```

filesys/t_mount.c

14.8.2 卸载文件系统：umount()和 umount2()

umount()系统调用用于卸载已挂载的文件系统。

```
#include <sys/mount.h>
```

```
int umount(const char *target);
```

Returns 0 on success, or -1 on error

target 参数指定待卸载文件系统的挂载点。

对于内核版本为 2.2 及更早的 Linux 系统来说，存在两种方法来标识文件系统：其一，通过挂载点；其二，通过包含文件系统的设备名。自内核版本 2.4 之后，Linux 不再允许使用第二种方法，其原因是如今可以在多个位置挂载单个文件系统，以第二种方式为 `target` 指定文件系统就会混淆不清。本书的 14.9.1 节将详细介绍这一点。

无法卸载正在使用中的（`busy`）文件系统，意即这一文件系统上有文件被打开，或者进程的当前工作目录驻留在此文件系统下。针对使用中的文件系统调用 `umount()`，系统会返回 `EBUSY` 错误。

系统调用 `umount2()` 是 `umount()` 的扩展版。通过 `flags` 参数，`umount2()` 可对卸载操作施以更精密的控制。

```
#include <sys/mount.h>

int umount2(const char *target, int flags);

Returns 0 on success, or -1 on error
```

这一标志位掩码参数由下列 0 个或多个值相或（`OR`）而成。

`MNT_DETACH` (始于 Linux 2.4.11)

执行 `lazy` 卸载。对挂载点加以标记，一方面允许已使用了挂载点的进程得以继续使用，同时禁止任何其他进程对该挂载点发起新的访问。当所有进程不再使用访问点时，系统会卸载相应的文件系统。

`MNT_EXPIRE` (始于 Linux 2.6.8)

将挂载点标记为到期（`expired`）。若首次调用 `umount2()` 时指定了该标志，且挂载点处于空闲状态，则该调用将以失败告终，并返回 `EAGAIN` 错误，同时将挂载点标记为到期。（如果挂载点处于在用状态，那么调用也将失败，并返回 `EBUSY` 错误，但不会将挂载点标记为到期。）只要无任何后续进程发起对挂载点的访问，该挂载点便会一直保持到期状态。再度调用 `umount2()` 时，如指定 `MNT_EXPIRE` 标志，将卸载到期的挂载点。这就提供了一种机制，以卸载在某段时间内未用的文件系统。该标志不能与 `MNT_DETACH` 或 `MNT_FORCE` 标志一并使用。

`MNT_FORCE`

即便文件系统（只对 `NFS` 挂载有效）处于在用状态，依然将其强行卸载。采用这一选项可能会造成数据丢失。

`UMOUNT_NOFOLLOW` (始于 Linux 2.6.34)

若 `target` 为符号链接，则不对其进行解引用。该标志专为某些 `set-user-ID-root` 程序而设计——此类程序允许非特权级用户执行卸载操作，意在避免安全性问题的发生（例如，若 `target` 为符号链接，且被改变以指向另外的位置）。

14.9 高级挂载特性

本节会介绍挂载文件系统时可采用的若干高级特性。这里会使用 `mount(8)` 命令来演示其

中的大多数特性，在程序中调用 `mount(2)`也能起到相同效果。

14.9.1 在多个挂载点挂载文件系统

内核版本 2.4 之前，一个文件系统只能挂载于单个挂载点。从内核版本 2.4 开始，可以将一个文件系统挂载于文件系统内的多个位置。由于每个挂载点下的目录子树内容都相同，在一个挂载点下对目录子树所做的改变，同样可见于其他挂载点，如下列 shell 会话所示：

```
$ su Privilege is required to use mount(8)
Password:
# mkdir /testfs Create two directories for mount points
# mkdir /demo
# mount /dev/sda12 /testfs Mount file system at one mount point
# mount /dev/sda12 /demo Mount file system at second mount point
# mount | grep sda12 Verify the setup
/dev/sda12 on /testfs type ext3 (rw)
/dev/sda12 on /demo type ext3 (rw)
# touch /testfs/myfile Make a change via first mount point
# ls /demo View files at second mount point
lost+found myfile
```

如 `ls` 命令的输出所示，在挂载点一(/testfs)下对目录子树所做的改变，在挂载点二(/demo)下完全可见。

14.9.4 节在介绍绑定挂载时，将举例说明多点挂载文件系统的用处所在。

正因为可在多点挂载一个设备，在 Linux 2.4 及其后续版本中，`umount()`系统调用不再将设备作为其入参。

14.9.2 多次挂载同一挂载点

在内核版本 2.4 之前，一个挂载点只能使用一次。从内核 2.4 开始，Linux 允许针对同一挂载点执行多次挂载。每次新挂载都会隐藏之前可见于挂载点下的目录子树。卸载最后一次挂载时，挂载点下上次挂载的内容会再次显示，请参考以下 shell 会话：

```
$ su Privilege is required to use mount(8)
Password:
# mount /dev/sda12 /testfs Create first mount on /testfs
# touch /testfs/myfile Make a file in this subtree
# mount /dev/sda13 /testfs Stack a second mount on /testfs
# mount | grep testfs Verify the setup
/dev/sda12 on /testfs type ext3 (rw)
/dev/sda13 on /testfs type reiserfs (rw)
# touch /testfs/newfile Create a file in this subtree
# ls /testfs View files in this subtree
newfile
# umount /testfs Pop a mount from the stack
# mount | grep testfs
/dev/sda12 on /testfs type ext3 (rw) Now only one mount on /testfs
# ls /testfs Previous mount is now visible
lost+found myfile
```

在现有且在用的挂载点上执行新的挂载操作是此类堆叠挂载的用法之一。持有打开文件描述符的进程、建立 `chroot` 监禁区 (`jail`) 的进程，以及工作目录位于老挂载点之下的进程将继续在旧有挂载下运行，而针对挂载点发起新访问的进程将使用新挂载。结合 `MNT_DETACH`

标志下的 `umount` 操作，则无需将文件系统置为单用户模式，即可为其提供平滑迁移。14.10 节在讨论 `tmpfs` 时会举例说明堆叠挂载的另一用法。

14.9.3 基于每次挂载的挂载标志

在内核 2.4 版本以前，文件和挂载点之间是一一对应的关系。由于从 Linux 2.4 开始，这一特征不再适用，故而 14.8.1 节所述的某些 `mountflag` 标志值可以基于每次挂载来设置。这包括 `MS_NOATIME` (始于 Linux 2.6.16)、`MS_NODEV`、`MS_NODIRATIME` (始于 Linux 2.6.16)、`MS_NOEXEC`、`MS_NOSUID`、`MS_RDONLY` (始于 Linux 2.6.26)，以及 `MS_RELATIME`。以下 shell 会话演示了使用 `MS_NOEXEC` 标志的效果：

```
$ su
Password:
# mount /dev/sda12 /testfs
# mount -o noexec /dev/sda12 /demo
# cat /proc/mounts | grep sda12
/dev/sda12 /testfs ext3 rw 0 0
/dev/sda12 /demo ext3 rw,noexec 0 0
# cp /bin/echo /testfs
# /testfs/echo "Art is something which is well done"
Art is something which is well done
# /demo/echo "Art is something which is well done"
bash: /demo/echo: Permission denied
```

14.9.4 绑定挂载

始于内核版本 2.4，Linux 支持了创建绑定挂载。绑定挂载（由使用 `MS_BIND` 标志的 `mount()` 调用来创建）是指在文件系统目录层级的另一处挂载目录或文件。这将导致文件或目录在两处同时可见。绑定挂载有些类似于硬链接，但存在两个方面的差异。

- 绑定挂载可以跨越多个文件系统挂载点，甚至不拘于 `chroot` 监禁区（`jail`）。
- 可针对目录执行绑定挂载。

可使用 `mount(8)` 的 `bind` 选项，在 shell 中创建绑定挂载，如下面几个例子所示。

第一个例子在另一处绑定挂载了一个目录，并展示了在一处目录中所创建的文件，对另一处目录同样可见。

```
$ su                                     Privilege is required to use mount(8)
Password:
# pwd
/testfs
# mkdir d1                               Create directory to be bound at another location
# touch d1/x                             Create file in the directory
# mkdir d2                               Create mount point to which d1 will be bound
# mount --bind d1 d2                     Create bind mount: d1 visible via d2
# ls d2                                  Verify that we can see contents of d1 via d2
x
# touch d2/y                             Create second file in directory d2
# ls d1                                  Verify that this change is visible via d1
x y
```

第二个例子会在另一处绑定挂载文件，并展示了在一处挂载下，可以看见另一处挂载下对文件所做的改变。

```

# cat > f1 Create file to be bound to another location
Chance is always powerful. Let your hook be always cast.
Type Control-D
# touch f2 This is the new mount point
# mount --bind f1 f2 Bind f1 as f2
# mount | egrep '(d1|f1)' See how mount points look
/testfs/d1 on /testfs/d2 type none (rw,bind)
/testfs/f1 on /testfs/f2 type none (rw,bind)
# cat >> f2 Change f2
In the pool where you least expect it, will be a fish.
# cat f1 The change is visible via original file f1
Chance is always powerful. Let your hook be always cast.
In the pool where you least expect it, will be a fish.
# rm f2 Can't do this because it is a mount point
rm: cannot unlink `f2': Device or resource busy
# umount f2 So unmount
# rm f2 Now we can remove f2

```

绑定挂载的应用场景之一是创建 chroot 监禁区 (jail) (参见 18.12 节)。在监禁区下，无需将各种标准目录 (诸如 /lib) 复制过来，为这些路径创建绑定挂载 (可能是以只读方式) 即可轻而易举地解决问题。

14.9.5 递归绑定挂载

默认情况下，如果使用 MS_BIND 为某个目录创建了绑定挂载，那么只会将该目录挂载到新位置。假设源目录下还存在子挂载 (submount)，则不会将这些子挂载复制到挂载 target 之下。Linux 2.4.11 添加了 MS_REC 标志，若与 MS_BIND 相或 (OR) 并作为标志参数的一部分传入 mount()，则会将子挂载复制到挂载目标下，此之谓递归绑定挂载。采用 mount(8) 命令所提供的 --rbind 选项，可在 shell 中完成相同任务，参见如下 shell 会话。

首先创建了一个目录树(src1)，并将其挂载在 top 之下。top 目录树下 (top/sub)，包括了一个子挂载(src2)。

```

$ su
Password:
# mkdir top This is our top-level mount point
# mkdir src1 We'll mount this under top
# touch src1/aaa
# mount --bind src1 top Create a normal bind mount
# mkdir top/sub Create directory for a submount under top
# mkdir src2 We'll mount this under top/sub
# touch src2/bbb
# mount --bind src2 top/sub Create a normal bind mount
# find top Verify contents under top mount tree
top
top/aaa
top/sub This is the submount
top/sub/bbb

```

现在以 top 作为源目录，另行创建绑定挂载(dir1)。由于属于非递归操作，新挂载不会复制子挂载。

```

# mkdir dir1
# mount --bind top dir1 Here we use a normal bind mount
# find dir1
dir1
dir1/aaa
dir1/sub

```

输出中并未发现 `dir1/sub/bbb`，这表明并未复制子挂载 `top/sub`。

再以 `top` 作为源目录来创建递归绑定挂载。

```
# mkdir dir2
# mount --rbind top dir2
# find dir2
dir2
dir2/aaa
dir2/sub
dir2/sub/bbb
```

从输出中可以发现 `dir2/sub/bbb`，这表明已然复制了子挂载 `top/sub`。

14.10 虚拟内存文件系统：tmpfs

到目前为止，本章已论及的所有文件系统均驻留在磁盘之上。然而，Linux 同样支持驻留于内存中的虚拟文件系统。对应用程序来说，此类文件系统看起来与任何其他文件系统别无二致——可施以相同操作（`open()`、`read()`、`write()`、`link()`、`mkdir()`等）。不过，二者之间还是存在一个重要差别：由于不涉及磁盘访问，虚拟文件系统的文件操作速度极快。

在 Linux 上，已经开发出了林林总总基于内存的文件系统。迄今为止，其中最为复杂的则非 tmpfs 文件系统莫属，该系统在 Linux 2.4 中首度出现。较之于其他基于内存的文件系统，其独特之处在于它属于虚拟内存文件系统。这意味着，该文件系统不但使用 RAM，而且在 RAM 耗尽的情况下，还会利用交换空间。（虽然此处描述的 tmpfs 文件系统为 Linux 所专有，但大多数 UNIX 实现都提供某种形式的基于内存的文件系统。）

tmpfs 文件系统是一个 Linux 内核的可选组件，通过 `CONFIG_TMPFS` 选项加以配置。

要创建 tmpfs 文件系统，请使用如下形式的命令：

```
# mount -t tmpfs source target
```

其中“source”可以是任意名称，其唯一的意义是在 `/proc/mounts` 中“抛头露面”，并通过 `mount` 和 `df` 命令显示出来。与往常一样，`target` 是该文件系统的挂载点。请注意，无需使用 `mkfs` 预先创建一个文件系统，内核会将此视为 `mount()` 系统调用的一部分自动加以执行。

作为使用 tmpfs 的例子之一，可采用堆叠挂载（无需顾忌 `/tmp` 目录目前是否处于在用状态），创建一 tmpfs 文件系统并将其挂载至 `/tmp`，如下所示：

```
# mount -t tmpfs newtmp /tmp
# cat /proc/mounts | grep tmp
newtmp /tmp tmpfs rw 0 0
```

有时，会使用如上命令（或 `/etc/fstab` 中的等价条目）来改善应用程序（比如，编译器）的性能，此类应用程序因创建临时性文件而频繁使用 `/tmp` 目录。

默认情况下，允许将 tmpfs 文件系统的大小提高至 RAM 容量的一半，但在创建文件系统或之后重新挂载时，可使用 `mount` 的 `size=nbytes` 选项为该文件系统的大小设置不同的上限值。（tmpfs 文件系统仅会根据其当前所持有的文件来消耗内存和交换空间。）

一旦卸载 tmpfs 文件系统，或者遭遇系统崩溃，那么该文件系统中的所有数据都将丢失，“tmpfs”正是得名于此。

除了用于用户应用程序以外，tmpfs 文件系统还有以下两个特殊用途。

- 由内核内部挂载的隐形 tmpfs 文件系统, 用于实现 System V 共享内存 (第 48 章) 和共享匿名内存映射 (第 49 章)。
- 挂载于 /dev/shm 的 tmpfs 文件系统, 为 glibc 用以实现 POSIX 共享内存和 POSIX 信号量。

14.11 获得与文件系统有关的信息: statvfs()

statvfs()和 fstatvfs()库函数能够获得与已挂载文件系统有关的信息。

```
#include <sys/statvfs.h>

int statvfs(const char *pathname, struct statvfs *statvfsbuf);
int fstatvfs(int fd, struct statvfs *statvfsbuf);

Both return 0 on success, or -1 on error
```

两者之间唯一的区别在于其标识文件系统的方式。statvfs()需使用 `pathname` 来指定文件系统中任一文件的名称。而 fstatvfs()则需使用打开文件描述符 `fd`, 来指代文件系统中的任一文件。二者均返回一个 `statvfs` 结构, 属于由 `statvfsbuf` 所指向的缓冲区, 其中包含了关乎文件系统的信息。`statvfs` 结构的形式如下:

```
struct statvfs {
    unsigned long f_bsize;      /* File-system block size (in bytes) */
    unsigned long f_frsize;    /* Fundamental file-system block size
                               (in bytes) */
    fsblkcnt_t   f_blocks;     /* Total number of blocks in file
                               system (in units of 'f_frsize') */
    fsblkcnt_t   f_bfree;     /* Total number of free blocks */
    fsblkcnt_t   f_bavail;    /* Number of free blocks available to
                               unprivileged process */
    fsfilcnt_t   f_files;     /* Total number of i-nodes */
    fsfilcnt_t   f_ffree;     /* Total number of free i-nodes */
    fsfilcnt_t   f_favail;    /* Number of i-nodes available to unprivileged
                               process (set to 'f_ffree' on Linux) */
    unsigned long f_fsid;     /* File-system ID */
    unsigned long f_flag;     /* Mount flags */
    unsigned long f_namemax;  /* Maximum length of filenames on
                               this file system */
};
```

上述注释已然清晰地描述出 `statvfs` 结构中大多数字段的用途。对其中一些字段, 这里还要深入交代几句。

- `fsblkcnt_t` 和 `fsfilcnt_t` 数据类型是由 SUSv3 所定义的整型。
- 对绝大多数 Linux 文件系统而言, `f_bsize` 和 `f_frsize` 的取值是相同的。然而, 某些文件系统支持块片段的概念, 在无需使用完整数据块的情况下, 可在文件尾部分配较小的存储单元, 从而避免因分配完整块而导致的空间浪费。在此类文件系统上, `f_frsize` 和 `f_bsize` 分别为块片段和整个块的大小。(据 [McKusick et al., 1984] 所述, UNIX 文件系统的块片段概念首现于 20 世纪 80 年代初期的 4.2BSD 快速文件系统。)
- 许多原生 UNIX 和 Linux 文件系统, 都支持为超级用户预留一部分文件系统块, 如此一来, 即便在文件系统空间耗尽的情况下, 超级用户仍可以登录系统解决故障。如果文件系统中确有预留块, 那么 `statvfs` 结构中 `f_bfree` 和 `f_bavail` 字段间的差值则

为预留块数。

- `f_flag` 字段是一个位掩码标志，用于挂载文件系统。也就是说，该字段所包含的信息类似于传入 `mount(2)` 的 `mountflags` 参数。然而，该字段所使用的标志位在命名时均冠以 `ST_`，这不同于 `mountflags` 中冠以 `MS_` 的命名手法。SUSv3 仅规范了 `ST_RDONLY` 和 `ST_NOSUID` 常量，而 `glibc` 实现则支持与 `MS_` 系列(参见 `mount()` 中对 `mountflags` 参数的描述)相对应的全系列常量。
- 某些 UNIX 实现会使用 `f_fsid` 字段来返回文件系统的唯一标识符，比方说，根据文件系统所驻留设备的标识符来取值。对大多数 UNIX 实现来说，该字段为 0。

SUSv3 规范了 `statvfs()` 和 `fstatvfs()`。对于 Linux (其他几种 UNIX 实现也一样)，二者均位于与其颇为相似的 `statfs()` 和 `fstatfs()` 系统调用之上。(有些 UNIX 实现只提供 `statfs()` 系统调用，而不提供 `statvfs()`。) 以下列出函数与系统调用间的主要区别 (除去字段命名差异以外)。

- `statvfs()` 和 `fstatvfs()` 函数均返回 `f_flag` 字段，内含关于文件系统的挂载标志信息。(`glibc` 实现通过扫描 `/proc/mounts` 或 `/etc/mntab` 来获取上述信息。)
- `statfs()` 和 `fstatfs()` 系统调用返回 `f_type` 字段，内含文件系统类型 (比如，返回值为 `0xef53` 则表示文件系统类型为 `ext2`)。

随本书发布源码的 `filesystem` 子目录中包含了 `t_statvfs.c` 和 `t_statfs.c` 文件，用来演示对 `statvfs()` 和 `statfs()` 的运用。

14.12 总结

设备都由 `/dev` 下的文件来表示。每个设备都有相应的设备驱动程序，用以执行一套标准的操作，与之对应的系统调用包括 `open()`、`read()`、`write()` 和 `close()`。设备既可以是实际存在的，也可以是虚拟的，这分别表明了硬件设备的存在与否。无论如何，内核都会提供一种设备驱动程序，并实现与真实设备相同的 API。

可将硬盘划分为一个或多个分区，**每个分区都可包含一个文件系统**。文件系统是对常规文件和目录的组织集合。Linux 实现的文件系统多种多样，其中包括传统的 `ext2` 文件系统。`ext2` 文件系统在概念上类似于早期的 UNIX 文件系统，由引导块、超级块、`i` 节点表和包含文件数据块的数据区域组成。每个文件在文件系统 `i` 节点表中都有一条对应记录，记录了与文件相关的各种信息，其中包括文件类型、大小、链接数、所有权、权限、时间戳，以及指向文件数据块的指针。

Linux 还提供了若干日志文件系统，其中包括 `Reiserfs`、`ext3`、`ext4`、`XFS`、`JFS` 以及 `Btrfs`。在实际更新文件之前，日志文件系统会记录元数据更新 (还可选地记录数据更新和文件系统更新)。这也意味着，一旦系统崩溃，系统可以重放 (`replay`) 日志文件，并迅速将文件系统恢复到一致状态。日志文件系统的最大优点在于系统崩溃后，无需像常规 UNIX 文件系统那样对文件系统进行漫长的一致性检查。

Linux 系统上的所有文件系统都被挂载于单根目录树之下，其树根为目录 `/`。目录树中挂载文件系统的位置被称为文件系统挂载点。

特权级进程可使用 `mount()` 和 `umount()` 系统调用来挂载、卸载文件系统。可使用 `statvfs()` 来获取与已挂载文件系统有关的信息。

进阶阅读

与设备和设备驱动程序有关的详细信息请参阅[Bovet & Cesati, 2005]和[Corbet et al., 2005], 尤其是后者。内核源码文件 `Documentation/devices.txt` 中, 也能找到一些与设备相关的有用信息。

以下几本著作都提供了关于文件系统的深度信息。[Tanenbaum, 2007]对文件系统的结构和实现做了一般性介绍。[Bach, 1986]介绍了 UNIX 文件系统的实现, 主要针对 System V。[Vahalia, 1996]和[Goodheart & Cox, 1994]也描述了 System V 文件系统。[Love, 2010]和[Bovet & Cesati, 2005]则讨论了 Linux VFS 的实现。

在内核源码子目录 `Documentation/filesystems` 下, 可以找到关于各种文件系统的文档。针对 Linux 所支持的大多数文件系统实现, 不少 WEB 站点也有论述。

14.13 练习

- 14-1.** 编写一程序, 试对在单目录下创建和删除大量 1 字节文件所需的时间进行度量。该程序应以 `xNNNNNN` 命名格式来创建文件, 其中 `NNNNNN` 为随机的 6 位数字。文件的创建顺序与生成文件名相同, 为随机方式, 删除文件则按数字升序操作 (删除与创建的顺序不同)。文件的数量 (FN) 和文件所在目录应由命令行指定。针对不同的 NF 值 (比如, 在 1000 和 20000 之间取值) 和不同的文件系统 (比如 `ext2`、`ext3` 和 `XFS`) 来测量时间。随着 NF 的递增, 每个文件系统下耗时的变化模式如何? 不同文件系统之间, 情况又是如何呢? 如果按数字升序来创建文件 (`x000001`、`x000001`、`x000002` 等), 然后以相同顺序加以删除, 结果会改变吗? 如果会, 原因何在? 此外, 上述结果会随文件系统类型的不同而改变吗?

第 15 章

文件属性

本章将探讨文件的各种属性（文件元数据）。首先介绍的是系统调用 `stat()`，可利用其返回一个包含多种文件属性（包括文件时间戳、文件所有权以及文件权限）的结构。然后，将描述用来改变文件属性的各种系统调用。（对文件权限的探讨会在第 17 章继续进行，讲述访问控制列表。）本章将在结尾处讨论 `i` 节点标志（也称为 `ext2` 扩展文件属性），可利用其控制内核对文件处理的方方面面。

15.1 获取文件信息：stat()

利用系统调用 `stat()`、`lstat()` 以及 `fstat()`，可获取与文件有关的信息，其中大部分提取自文件 `i` 节点。

```
#include <sys/stat.h>
```

```
int stat(const char *pathname, struct stat *statbuf);  
int lstat(const char *pathname, struct stat *statbuf);  
int fstat(int fd, struct stat *statbuf);
```

All return 0 on success, or -1 on error

以上 3 个系统调用之间仅有的区别在于对文件的描述方式不同。

- `stat()` 会返回所命名文件的相关信息。
- `lstat()` 与 `stat()` 类似，区别在于如果文件属于符号链接，那么所返回的信息针对的是符号链接自身（而非符号链接所指向的文件）。
- `fstat()` 则会返回由某个打开文件描述符所指代文件的相关信息。

系统调用 `stat()` 和 `lstat()` 无需对其所操作的文件本身拥有任何权限，但针对指定 `pathname` 的父目录要有执行（搜索）权限。而只要供之以有效的文件描述符，`fstat()` 系统调用总是成功。

上述所有系统调用都会在缓冲区中返回一个由 `statbuf` 指向的 `stat` 结构，其格式如下：

```
struct stat {
    dev_t    st_dev;           /* IDs of device on which file resides */
    ino_t    st_ino;          /* I-node number of file */
    mode_t   st_mode;         /* File type and permissions */
    nlink_t  st_nlink;        /* Number of (hard) links to file */
    uid_t    st_uid;          /* User ID of file owner */
    gid_t    st_gid;          /* Group ID of file owner */
    dev_t    st_rdev;         /* IDs for device special files */
    off_t    st_size;         /* Total file size (bytes) */
    blksize_t st_blksize;     /* Optimal block size for I/O (bytes) */
    blkcnt_t st_blocks;       /* Number of (512B) blocks allocated */
    time_t   st_atime;        /* Time of last file access */
    time_t   st_mtime;        /* Time of last file modification */
    time_t   st_ctime;        /* Time of last status change */
};
```

在 SUSv3 中，明确定义了供 `stat` 结构各字段使用的不同数据类型。更多与这些数据类型有关的信息，请参考 3.6.2 节。

根据 SUSv3，将 `lstat()` 应用于符号链接时，只需在 `st_size` 字段和描述文件类型的 `st_mode` 字段（稍后介绍）返回有效信息，并不要求所返回的其他字段（比如，`time` 字段）信息有效。如此一来，出于效率的原因，系统实现可选择维护此类字段。说透彻一点，早期 UNIX 标准的意图在于把符号链接要么实现为 `i` 节点，要么实现为目录中的一条记录。如果是后一种实现方式，在实现中顾及 `stat` 结构的所有字段是不现实的。（在所有当代主流的 UNIX 实现中，符号链接都是以 `i` 节点的方式来实现的。）在 Linux 上，将 `lstat()` 应用于符号链接时，会返回所有 `stat` 字段的信息。

接下来，会对 `stat` 结构的某些字段做重点介绍。最后，还会给出展示完整 `stat` 结构的程序示例。

设备 ID 和 `i` 节点号

`st_dev` 字段标识文件所驻留的设备。`st_ino` 字段则包含了文件的 `i` 节点号。利用以上两者，可在所有文件系统中唯一标识某个文件。`dev_t` 类型记录了设备的主、辅 ID（见 14.1 节）。

如果是针对设备的 `i` 节点，那么 `st_rdev` 字段则包含设备的主、辅 ID。

利用宏 `major()` 和 `minor()`，可提取 `dev_t` 值的主、辅 ID。获取对两个宏声明的头文件则随 UNIX 实现而各异。在 Linux 系统上，若定义了 `_BSD_SOURCE` 宏，则两个宏定义于 `<sys/types.h>` 中。

由 `major()` 和 `minor()` 所返回的整形值大小也随 UNIX 实现的不同而各不相同。为保证可移植性，打印时应总是将返回值强制转换为 `long`（见 3.6.2 节）。

文件所有权

`st_uid` 和 `st_gid` 字段分别标识文件的属主（用户 ID）和属组（组 ID）。

链接数

`st_nlink` 字段包含了指向文件的（硬）链接数。本书第 18 章将详细介绍链接。

文件类型及权限

`st_mode` 字段内含有位掩码，起标识文件类型和指定文件权限的双重作用。图 15-1 所示

为该字段所含各位的布局情况。

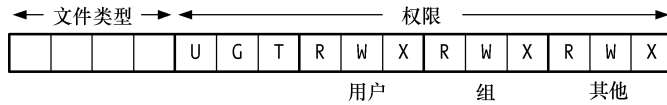


图 15-1: st_mode 位掩码的布局

与常量 S_IFMT 相与 (&), 可从该字段中析取文件类型。(Linux 使用了 st_mode 字段中的 4 位来标识文件类型位。但由于 SUSv3 并未对文件类型位的表示方式做出任何规定, 故而其具体细节随各实现而异。) 将计算结果与一系列常量进行比较, 即可确定文件类型, 如下所示:

```
if ((statbuf.st_mode & S_IFMT) == S_IFREG)
    printf("regular file\n");
```

鉴于上述操作属于常见操作, 因此可利用标准宏将其简化为:

```
if (S_ISREG(statbuf.st_mode))
    printf("regular file\n");
```

表 15-1 所列为全套文件类型宏 (定义于 <sys/stat.h>)。这些宏均由 SUSv3 定义, 并为 Linux 所支持。一些其他的 UNIX 实现还定义了别的文件类型 (比如, 用于 Solaris door files 的 S_IFDOOR)。因为调用 stat() 时会循符号链接而直抵实际文件, 所以只有在调用 lstat() 时才有可能返回类型 S_IFLNK。

想从 <sys/stat.h> 中获得 S_IFSOCK 和 S_ISSOCK() 的定义, 必须定义 _BSD_SOURCE 特性测试宏, 或是将 _XOPEN_SOURCE 定义为不小于 500 的值。(具体规则随 glibc 版本而异。在某些情况下, 需将 _XOPEN_SOURCE 的值定义为不小于 600。)

最初的 POSIX.1 标准并未定义表 15-1 中第一列所列的常量, 尽管其中的大部分已为多数 UNIX 实现所支持。而 SUSv3 则把这些常量纳入规范。

表 15-1: 针对 stat 结构中的 st_mode 来检查文件类型的宏

常 量	测 试 宏	文 件 类 型
S_IFREG	S_ISREG()	常规文件
S_IFDIR	S_ISDIR()	目录
S_IFCHR	S_ISCHR()	字符设备
S_IFBLK	S_ISBLK()	块设备
S_IFIFO	S_ISFIFO()	FIFO 或管道
S_IFSOCK	S_ISSOCK()	套接字
S_IFLNK	S_ISLNK()	符号链接

st_mode 字段的低 12 位定义了文件权限, 会在 15.4 节介绍。目前, 只要知道其中最低 9 位分别用来表示文件属主、属组以及其他用户的读、写、执行权限。

文件大小、已分配块以及最优 I/O 块大小

对于常规文件, st_size 字段表示文件的字节数。对于符号链接, 则表示链接所指路径名的长度, 以字节为单位。对于共享内存对象 (见第 54 章), 该字段则表示对象的大小。

st_blocks 字段表示分配给文件的总块数, 块大小为 512 字节, 其中包括了为指针块所分配的

空间（参见图 14-2）。之所以选择 512 字节大小的块作为度量单位，有其历史原因——对于 UNIX 所实现的任何文件系统而言，最小的块大小即为 512 字节。更为现代的 UNIX 文件系统则使用更大尺寸的逻辑块。例如，对于 ext2 文件系统，取决于其逻辑块大小为 1024、2048 还是 4096 字节，`st_blocks` 的取值将总是 2、4、8 的倍数。

SUSv3 并未定义度量 `st_blocks` 时所使用的单位，故而 UNIX 实现可以不使用 512 字节作为其单位。大多数 UNIX 实现使用 512 字节作为 `st_blocks` 字段的单位，但 HP-UX 11 所使用的单位则视文件系统而定（有时为 1024 字节）。

`st_blocks` 字段记录了实际分配给文件的磁盘块数量。如果文件内含空洞（见 4.7 节），该值将小于从相应文件字节数字段（`st_size`）的值。（执行显示磁盘使用情况的 `du -k file` 命令，便可获悉分配给文件的实际空间，单位为 KB。亦即，得自对文件 `st_blocks` 值，而非 `st_size` 值的计算结果。）

`st_blksize` 字段的命名多少有些令人费解。其所指并非底层文件系统的块大小，而是针对文件系统上文件进行 I/O 操作时的最优块大小（以字节为单位）。若 I/O 所采用的块大小小于该值，则被视为低效（参阅 13.1 节）。一般而言，`st_blksize` 的返回值为 4096。

文件时间戳

`st_atime`、`st_mtime` 和 `st_ctime` 字段，分别记录了对文件的上次访问时间、上次修改时间，以及文件状态发生改变的上次时间。这 3 个字段的类型均属 `time_t`，是标准的 UNIX 时间格式，记录了自 Epoch 以来的秒数。15.2 节对此有深入描述。

程序示例

程序清单 15-1 所列程序使用 `stat()` 去获取文件（文件名由该程序的命令行提供）的相关信息。若以 `-l` 选项执行命令，程序会改用 `lstat()`，以获取与符号链接（而非该链接所指代的文件）有关的信息。该程序会将返回 `stat` 结构的所有字段一一打印出来。（至于程序中将 `st_size` 和 `st_blocks` 字段强制转换为 `long long` 类型的原因，请参考 5.10 节。）该程序所调用的 `filePermStr()` 函数源码见之于程序清单 15-4。

以下为对该程序的执行情况。

```
$ echo 'All operating systems provide services for programs they run' > apue
$ chmod g+s apue          Turn on set-group-ID bit; affects last status change time
$ cat apue                Affects last file access time
All operating systems provide services for programs they run
$ ./t_stat apue
File type:                regular file
Device containing i-node: major=3  minor=11
I-node number:           234363
Mode:                    102644 (rw-r--r--)
  special bits set:      set-GID
Number of (hard) links:   1
Ownership:               UID=1000  GID=100
File size:                61 bytes
Optimal I/O block size:  4096 bytes
512B blocks allocated:   8
Last file access:        Mon Jun  8 09:40:07 2011
Last file modification:  Mon Jun  8 09:39:25 2011
Last status change:     Mon Jun  8 09:39:51 2011
```

```
#define _BSD_SOURCE /* Get major() and minor() from <sys/types.h> */
#include <sys/types.h>
#include <sys/stat.h>
#include <time.h>
#include "file_perms.h"
#include "tlpi_hdr.h"

static void
displayStatInfo(const struct stat *sb)
{
    printf("File type:          ");

    switch (sb->st_mode & S_IFMT) {
        case S_IFREG: printf("regular file\n");          break;
        case S_IFDIR: printf("directory\n");            break;
        case S_IFCHR: printf("character device\n");      break;
        case S_IFBLK: printf("block device\n");          break;
        case S_IFLNK: printf("symbolic (soft) link\n");  break;
        case S_IFIFO: printf("FIFO or pipe\n");          break;
        case S_IFSOCK: printf("socket\n");               break;
        default:      printf("unknown file type?\n");    break;
    }

    printf("Device containing i-node: major=%ld  minor=%ld\n",
           (long) major(sb->st_dev), (long) minor(sb->st_dev));

    printf("I-node number:          %ld\n", (long) sb->st_ino);

    printf("Mode:                  %lo (%s)\n",
           (unsigned long) sb->st_mode, filePermStr(sb->st_mode, 0));

    if (sb->st_mode & (S_ISUID | S_ISGID | S_ISVTX))
        printf("    special bits set:   %s%s%s\n",
               (sb->st_mode & S_ISUID) ? "set-UID " : "",
               (sb->st_mode & S_ISGID) ? "set-GID " : "",
               (sb->st_mode & S_ISVTX) ? "sticky " : "");

    printf("Number of (hard) links:  %ld\n", (long) sb->st_nlink);

    printf("Ownership:              UID=%ld  GID=%ld\n",
           (long) sb->st_uid, (long) sb->st_gid);

    if (S_ISCHR(sb->st_mode) || S_ISBLK(sb->st_mode))
        printf("Device number (st_rdev): major=%ld; minor=%ld\n",
               (long) major(sb->st_rdev), (long) minor(sb->st_rdev));

    printf("File size:                %lld bytes\n", (long long) sb->st_size);
    printf("Optimal I/O block size:  %ld bytes\n", (long) sb->st_blksize);
    printf("%512B blocks allocated: %lld\n", (long long) sb->st_blocks);
    printf("Last file access:        %s", ctime(&sb->st_atime));
    printf("Last file modification: %s", ctime(&sb->st_mtime));
    printf("Last status change:     %s", ctime(&sb->st_ctime));
}
}
```

```

int
main(int argc, char *argv[])
{
    struct stat sb;
    Boolean statLink;          /* True if "-l" specified (i.e., use lstat) */
    int fname;                /* Location of filename argument in argv[] */

    statLink = (argc > 1) && strcmp(argv[1], "-l") == 0;
                                /* Simple parsing for "-l" */
    fname = statLink ? 2 : 1;

    if (fname >= argc || (argc > 1 && strcmp(argv[1], "--help") == 0))
        usageErr("%s [-l] file\n"
                "      -l = use lstat() instead of stat()\n", argv[0]);

    if (statLink) {
        if (lstat(argv[fname], &sb) == -1)
            errExit("lstat");
    } else {
        if (stat(argv[fname], &sb) == -1)
            errExit("stat");
    }

    displayStatInfo(&sb);

    exit(EXIT_SUCCESS);
}

```

files/t_stat.c

15.2 文件时间戳

stat 结构的 st_atime、st_mtime 和 st_ctime 字段所含为文件时间戳，分别记录了对文件的上次访问时间、上次修改时间，以及文件状态（即文件 i 节点内信息）上次发生变更的时间。对时间戳的记录形式为自 1970 年 1 月 1 日（参见 10.1 节）以来所历经的秒数。

大多数原生 Linux 和 UNIX 文件系统都支持上述所有的时间戳字段，但某些非 UNIX 文件系统则未必如此。

表 15-2 总结了本书所述各种系统调用及库函数所改变的相应时间戳字段（有时则是指父目录的类似字段）。本表标题中的 a、m 和 c 分别表示 st_atime、st_mtime 和 st_ctime 字段。在大多数情况下，系统调用会将相关时间戳置为当前时间。但 utime() 及类似调用（将在 15.2.1 和 15.2.2 节讨论）则不在此列，可利用这些系统调用显式将对文件的上次访问时间和上次修改时间设定为任意值。

表 15-2: 各种函数对文件时间戳的影响

函 数	文件或目录			父目录			注 释
	a	m	c	a	m	c	
chmod()			•				与 fchmod() 相同
chown()			•				与 lchown() 和 fchown() 相同
exec()	•						
link()			•		•	•	影响第二个参数的父目录

函 数	文件或目录			父目录			注 释
	a	m	c	a	m	c	
mkdir()	•	•	•		•	•	
mkfifo()	•	•	•		•	•	
mknod()	•	•	•		•	•	
mmap()	•	•	•				仅当对具有 MAP_SHARED 属性的映射进行更新时, 才会改变 st_mtime 和 st_ctime
msync()		•	•				仅当修改文件时, 才会改变时间戳
open(), creat()	•	•	•		•	•	新建文件时
open(), creat()		•	•				截断现有文件时
pipe()	•	•	•				
read()	•						与 readv()、pread()和 preadv()相同
readdir()	•						readdir()可缓冲目录条目, 仅当读取目录时, 才会更新各时间戳
removexattr()			•				与 fremovexattr()和 removexattr()相同
rename()			•		•	•	同时影响文件(更名前后)的父目录。SUSv3 并未强制要求改变文件的 st_ctime, 只是顺带指出有一些实现是照此办理的
rmdir()					•	•	与 remove(directory)相同
sendfile()	•						改变了输入文件的时间戳
setxattr()			•				与 fsetxattr()和 lsetxattr()相同
symlink()	•	•	•		•	•	设置链接(而非目标文件)的时间戳
truncate()		•	•				与 ftruncate()相同, 仅当文件大小改变时, 才会改变时间戳
unlink()			•		•	•	与 remove(file)相同。若之前的链接计数大于 1, 会改变文件的 st_ctime
utime()	•	•	•				与 utimes()、futimes()、futimens()、lutimes() 和 utimensat()相同
write()		•	•				与 writev()、pwrite()和 pwritev()相同

本书 14.8.1 节和 15.5 节分别介绍了可阻止对文件上次访问时间进行更新的 `mount(2)` 选项¹, 以及作用于单个文件的标志。4.3.1 节中介绍的 `open()` `O_NOATIME` 标志也能起到类似作用。由于采用此标志可降低对磁盘的操作次数, 提升某些应用的文件访问性能, 故而颇具实用价值。

虽然大多数 UNIX 系统都不会记录文件的创建时间, 但最新的 BSD 系统会使用名为 `st_birthtime` 的 `stat` 字段来记录这一时间。

¹ 译者注: 对整个文件系统起作用。

纳秒时间戳

对于 `stat` 结构所含的 3 个时间戳字段，Linux 从 2.6 版本将其精度提升至纳秒级。纳秒级分辨率将提高某些程序的精度，因为此类程序需要根据文件时间戳的先后顺序来作决定（比如，`make(1)`）。

SUSv3 并未强制要求 `stat` 结构对纳秒级时间戳的支持，但 SUSv4 对此则有明文规定。

并非所有文件系统都支持纳秒级精度的时间戳。JFS、XFS、`ext4`，以及 `Btrfs` 文件系统都支持，但 `ext2`、`ext3` 以及 `Reiserfs` 文件系统则不然。

`glibc` API（自版本 2.3 起）将每个时间戳字段都定义为 `timespec` 结构（本节稍后在介绍 `utimensat()` 时将会讲解该结构），此结构是以秒和纳秒为单位来表示时间的一种组件。使用恰当的宏定义，该组件的秒级部分可见诸于传统字段（`st_atime`、`st_mtime`，以及 `st_ctime`）中。而对其纳秒级部分的访问则会采用如下手法：通过诸如 `st_atim.tv_nsec` 之类的字段名来获取文件上次访问时间的纳秒级部分。

15.2.1 使用 `utime()` 和 `utimes()` 来改变文件时间戳

使用 `utime()` 或与之相关的系统调用集之一，可显式改变存储于文件 `i` 节点中的文件上次访问时间戳和上次修改时间戳。解压文件时，`tar(1)` 和 `unzip(1)` 之类的程序会使用这些系统调用来重置文件的时间戳。

```
#include <utime.h>

int utime(const char *pathname, const struct utimbuf *buf);
                Returns 0 on success, or -1 on error
```

参数 `pathname` 用来标识欲修改时间的文件。若该参数为符号链接，则会进一步解除引用。参数 `buf` 既可为 `NULL`，也可为指向 `utimbuf` 结构的指针。

```
struct utimbuf {
    time_t actime;    /* Access time */
    time_t modtime;  /* Modification time */
};
```

该结构中的字段记录了自 Epoch（见 10.1 节）以来的秒数。

`utime()` 的运作方式则视以下两种不同情况而定。

- 如果 `buf` 为 `NULL`，那么会将文件的上次访问和修改时间同时置为当前时间。这时，进程要么具有特权级别（`CAP_FOWNER` 或 `CAP_DAC_OVERRIDE`），要么其有效用户 ID 与该文件的用户 ID（属主）相匹配，且对文件有写权限（逻辑上，对文件拥有写权限的进程在调用其他系统调用时，可能会于无意间改变这些时间戳）。（准确地说，如 9.5 节所述，在 Linux 系统中，用来与文件用户 ID 做比对的是进程的文件系统用户 ID，而非其有效用户 ID。）
- 若将 `buf` 指定为指向 `utimbuf` 结构的指针，则会使用该结构的相应字段去更新文件的上次访问和修改时间。此时，要么调用程序具有特权级别（`CAP_FOWNER`），要么进程的有效用户 ID 必需匹配文件的用户 ID（仅对文件拥有写权限是不够的）。

为更改文件时间戳中的一项，可以先利用 `stat()` 来获取两个时间，并使用其中之一来初始化 `utimbuf` 结构，然后再将另一时间置为期望值。下列代码演示了这一操作，将文件的上次修

改时间改为与上次访问时间相同。

```
struct stat sb;
struct utimbuf utb;

if (stat(pathname, &sb) == -1)
    errExit("stat");
utb.actime = sb.st_atime;      /* Leave access time unchanged */
utb.modtime = sb.st_atime;
if (utime(pathname, &utb) == -1)
    errExit("utime");
```

只要调用 `utime()` 成功，总会将文件的上次状态更改时间置为当前时间。

Linux 还提供了源于 BSD 的 `utimes()` 系统调用，其功用类似于 `utime()`。

```
#include <sys/time.h>

int utimes(const char *pathname, const struct timeval tv[2]);

Returns 0 on success, or -1 on error
```

`utime()` 与 `utimes()` 之间最显著的差别在于后者可以以微秒级精度来指定时间值（`timeval` 结构请见 10.1 节）。Linux 2.6 为文件时间戳提供了纳秒级的精度支持，在这里也部分得以体现。新的文件访问时间在 `tv[0]` 中指定，新的文件修改时间在 `tv[1]` 中指定。

`utimes()` 的使用例子请参考随本书一同发行的源码中的 `files/t_utimes.c` 文件。

`futimes()` 和 `lutimes()` 库函数的功能与 `utimes()` 大同小异。前两者与后者之间的差异在于，用来指定要更改时间戳文件的参数不同。

```
#include <sys/time.h>

int futimes(int fd, const struct timeval tv[2]);
int lutimes(const char *pathname, const struct timeval tv[2]);

Both return 0 on success, or -1 on error
```

调用 `futimes()` 时，使用打开文件描述符 `fd` 来指定文件。

调用 `lutimes()` 时，使用路径名来指定文件，有别于调用 `utimes()` 的是：对于 `lutimes()`，若路径名指向一符号链接，则调用不会对该链接进行解引用，而是更改链接自身的时间戳。

`glibc` 自 2.3 版本开始支持 `futimes()` 函数，自 2.6 版本开始支持 `lutimes()` 函数。

15.2.2 使用 `utimensat()` 和 `futimens()` 改变文件时间戳

`utimensat()` 系统调用（内核自 2.6.22 版本开始支持）和 `futimens()` 库函数（`glibc` 自版本 2.6 开始支持）为设置对文件的上次访问和修改时间戳提供了扩展功能。以下对这两个编程接口的优点列举一二。

- 可按纳秒级精度设置时间戳。相对于提供微秒级精度的 `utimes()`，这是重大改进。
- 可独立设置某一时间戳（一次只设置其一）。如前所述，要使用旧编程接口去改变时间戳之一，需要首先调用 `stat()` 获取另一时间戳的值。然后再将获取值与打算变更的时间戳一同指定。（若另一进程在这两步之间执行了更新时间戳的操作，将会导致竞争状态。）
- 可独立将任一时间戳置为当前时间。要使用旧编程接口将一个时间戳改为当前时间，

需要调用 `stat()` 去获取那些保持不变的时间戳的设置情况，并调用 `gettimeofday()` 以获得当前时间。

在 SUSv3 中并未定义以上两个接口，但 SUSv4 将其纳入规范。

`utimensat()` 系统调用会把由 `pathname` 指定文件的时间戳更新为由数组 `times` 指定的值。

```
#define _XOPEN_SOURCE 700      /* Or define _POSIX_C_SOURCE >= 200809 */
#include <sys/stat.h>

int utimensat(int dirfd, const char *pathname,
              const struct timespec times[2], int flags);

Returns 0 on success, or -1 on error
```

若将 `times` 指定为 `NULL`，则会将以上两个文件时间戳都更新为当前时间。若 `times` 值为非 `NULL`，则会针对指定文件在 `times[0]` 中放置新的上次访问时间，在 `times[1]` 中放置新的上次修改时间。数组 `times` 所含的每一元素都是如下格式的一个结构：

```
struct timespec {
    time_t tv_sec;      /* Seconds ('time_t' is an integer type) */
    long tv_nsec;      /* Nanoseconds */
};
```

结构所含的字段分别指定自 Epoch (10.1 节) 以来的秒数和纳秒数。

若有意将时间戳之一置为当前时间，则可将相应的 `tv_nsec` 字段指定为特殊值 `UTIME_NOW`。若希望某一时间戳保持不变，则需把相应的 `tv_nsec` 字段指定为特殊值 `UTIME_OMIT`。无论是上述哪一种情况，都将忽略相应 `tv_sec` 字段中的值。

可以将 `dirfd` 参数指定为 `AT_FDCWD`，此时对 `pathname` 参数的解读与 `utimes()` 相类似。或者，也可以将其指定为指代目录的文件描述符，18.11 节将描述这一用法的目的所在。

`flags` 参数可以为 0，或者 `AT_SYMLINK_NOFOLLOW`，意即当 `pathname` 为符号链接时，不会对其解引用（也就是说，改变的是符号链接自身的时间戳）。相形之下，`utimes()` 总是对符号链接进行解引用。

以下代码片段在将对文件的上次访问时间置为当前时间的同时，上次修改时间则保持不变。

```
struct timespec times[2];

times[0].tv_sec = 0;
times[0].tv_nsec = UTIME_NOW;
times[1].tv_sec = 0;
times[1].tv_nsec = UTIME_OMIT;
if (utimensat(AT_FDCWD, "myfile", times, 0) == -1)
    errExit("utimensat");
```

利用 `utimensat()` (和 `futimens()`) 改变时间戳时所遵循的权限规则与旧有 API 函数相类似，`utimensat(2)` 手册页对此有详细讨论。

使用 `futimens()` 库函数可更新打开文件描述符 `fd` 所指代文件的各个文件时间戳。

```
#include _GNU_SOURCE
#include <sys/stat.h>

int futimens(int fd, const struct timespec times[2]);

Returns 0 on success, or -1 on error
```

其中，`times` 参数的使用方法与 `utimensat()` 相同。

15.3 文件属主

每个文件都有一个与之关联的用户 ID (UID) 和组 ID (GID)，籍此可以判定文件的属主和属组。

15.3.1 新建文件的属主

文件创建时，其用户 ID “取自” 进程的有效用户 ID。而新建文件的组 ID 则“取自” 进程的有效组 ID (等同于 System V 系统的默认行为)，或父目录的组 ID (BSD 系统的行为)。当为项目创建目录时，需要该目录下的所有文件隶属于某一特定组，并且可为该组所有成员所访问。这时，采用后一种行为就非常实用。新建文件的组 ID 在这两者间如何取舍是由多种因素决定的，新文件所在文件系统的类型就是其中之一。这里先介绍一下 `ext2` 和某些其他类型文件系统所遵循的规则。

为求精确，本节所使用的术语有效用户 ID 或组 ID，实际是指文件系统用户 ID 或组 ID (见 9.5 节)。

装配 `ext2` 文件系统时，`mount` 命令要么带有 `-o grpuid` 的选项 (或等效的 `-o bsdgroups` 选项)，要么带有 `-o nogrpuid` 选项 (或等效的 `-o sysvgroups` 选项)。(若两者均未指定，`mount` 命令的默认选项为 `-o nogrpuid`。) 若指定了 `-o grpuid` 选项，那么新建文件总是继承其父目录的组 ID。若指定了 `-o nogrpuid` 选项，那么在默认情况下，新建文件的组 ID 则“取自” 进程的有效组 ID。不过，如果已将目录的 `set-group-ID` 位置位 (通过 `chmod g+s` 命令)，那么文件的组 ID 又将从其父目录处继承。表 15-3 对上述规则做了总结。

正如 18.6 节所述，一旦将某一目录的 `set-group-ID` 位置位后，该目录下所有子目录的 `set-group-ID` 位也将被置位。如此一来，正文中所描述的 `set-group-ID` 行为会遍布整个目录树。

表 15-3: 确定新建文件组所有权的规则

文件系统装配选项	有无设置父目录的 Set-group-ID 位	新建文件的组所有权取自何处
<code>-o grpuid</code> , <code>-o bsdgroups</code>	忽略	父目录组 ID
<code>-o nogrpuid</code> , <code>-o sysvgroups</code> (默认)	无	父目录组 ID
	有	父目录组 ID

写作本书之际，支持 `grpuid` 和 `nogrpuid` 装配选项的文件系统仅限于 `ext2`、`ext3`、`ext4` 以及自 Linux 2.6.14 出现的 XFS。其他类型的文件系统则遵循 `nogrpuid` 规则。

15.3.2 改变文件属主：`chown()`、`fchown()`和 `lchown()`

系统调用 `chown()`、`lchown()`和 `fchown()` 可用来改变文件的属主 (用户 ID) 和属组 (组 ID)。

```

#include <unistd.h>

int chown(const char *pathname, uid_t owner, gid_t group);

#define _XOPEN_SOURCE 500    /* Or: #define _BSD_SOURCE */
#include <unistd.h>

int lchown(const char *pathname, uid_t owner, gid_t group);
int fchown(int fd, uid_t owner, gid_t group);

All return 0 on success, or -1 on error

```

以上 3 个系统调用之间的区别类似于 `stat()` 系统调用一族。

- `chown()` 改变由 `pathname` 参数命名文件的所有权。
- `lchown()` 用途与 `chown()` 相同，不同之处在于若参数 `pathname` 为一符号链接，则将会改变链接文件本身的所有权，而与该链接所指代的文件无干。
- `fchown()` 也会改变文件的所有权，只是文件由打开文件描述符 `fd` 所引用。

参数 `owner` 和 `group` 分别为文件指定新的用户 ID 和组 ID。若只打算改变其中之一，只需将另一参数置为 -1，即可令与之相关的 ID 保持不变。

Linux 2.2 之前，`chown()` 不对符号链接进行解引用。从 Linux 2.2 开始，`chown()` 的语义发生了变化，并且添加了新系统调用 `lchown()`，以提供老系统调用 `chown()` 的行为。

只有特权级进程 (CAP_CHOWN) 才能使用 `chown()` 改变文件的用户 ID。非特权级进程可使用 `chown()` 将自己所拥有文件的组 ID 改为其所从属的任一属组的 ID，前提是进程的有效用户 ID 与文件的用户 ID 相匹配。特权级进程则可将文件的组 ID 修改为任意值。

当超级用户改变可执行文件的属主和属组时，是否应当屏蔽 `set-user-ID` 和 `set-group-ID` 位？SUSv3 对此未置可否。Linux 2.0 确实会屏蔽以上各位，但某些 2.2 版本（不超过 2.2.12）的早期内核则不然。其后的 2.2 内核又回归了 2.0 内核的行为——造成变化的无论是超级用户还是其他用户，系统在处理时都将一视同仁。（但若以 `root` 登录后执行 `chown(1)` 命令来改变文件的所有权，则 `chown` 命令会在调用 `chown(2)` 之后利用系统调用 `chmod()` 来重新激活 `set-user-ID` 和 `set-group-ID` 位。）

如果文件组的属主或属组发生了改变，那么 `set-user-ID` 和 `set-group-ID` 权限位也会随之关闭。这一安全举措是为了防止如下行为：普通用户若能打开某一可执行文件的 `set-user-ID`（或 `set-group-ID`）位，然后再设法令其为某些特权级用户（或组）所拥有，就能在执行该文件时获得特权用户身份。

改变文件的属主和属组时，如果已然屏蔽了属组的可执行权限位，或者要改变的是目录的所有权时，那么将不会屏蔽 `set-group-ID` 权限位。在上述两种情况下，`set-group-ID` 位的用途并非是去创建一个启用了 `set-group-ID` 的程序，因此将该位屏蔽并不可取。`set-group-ID` 的其他用途如下所示。

- 若屏蔽了属组的可执行权限位，则可利用 `set-group-ID` 权限位来启用强制文件锁定（请参阅 55.4 节）。
- 当作用于目录时，可利用 `set-group-ID` 位来控制在该目录下创建文件的所有权（见 15.3.1 节）。

程序清单 15-2 演示了 `chown()` 的用法，该程序允许用户改变任意数量文件（由命令行参数指定）的属主和属组。（该程序使用程序清单 8-1 中的 `userIdFromName()` 和 `groupIdFromName()` 函数，将用户名和组名转换为相应的数字 ID。）

程序清单 15-2：改变文件的属主和属组

```
files/t_chown.c
#include <pwd.h>
#include <grp.h>
#include "ugid_functions.h"          /* Declarations of userIdFromName()
                                   and groupIdFromName() */
#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    uid_t uid;
    gid_t gid;
    int j;
    Boolean errFnd;

    if (argc < 3 || strcmp(argv[1], "--help") == 0)
        usageErr("%s owner group [file...]\n"
                 "          owner or group can be '-', "
                 "meaning leave unchanged\n", argv[0]);

    if (strcmp(argv[1], "-") == 0) {          /* "-" ==> don't change owner */
        uid = -1;
    } else {                                  /* Turn user name into UID */
        uid = userIdFromName(argv[1]);
        if (uid == -1)
            fatal("No such user (%s)", argv[1]);
    }

    if (strcmp(argv[2], "-") == 0) {          /* "-" ==> don't change group */
        gid = -1;
    } else {                                  /* Turn group name into GID */
        gid = groupIdFromName(argv[2]);
        if (gid == -1)
            fatal("No group user (%s)", argv[1]);
    }

    /* Change ownership of all files named in remaining arguments */

    errFnd = FALSE;
    for (j = 3; j < argc; j++) {
        if (chown(argv[j], uid, gid) == -1) {
            errMsg("chown: %s", argv[j]);
            errFnd = TRUE;
        }
    }

    exit(errFnd ? EXIT_FAILURE : EXIT_SUCCESS);
}
files/t_chown.c
```

15.4 文件权限

本节将介绍应用于文件和目录的权限方案。尽管此处所讨论的权限主要是针对普通文件和目录，但其规则可适用于所有文件类型，包括设备文件、FIFO 以及 UNIX 域套接字等。此外，System V 和 POSIX 进程间通信对象（共享内存、信号量和消息队列）也具有权限掩码，而适用于此类对象的权限规则也与文件的权限规则相类似。

15.4.1 普通文件的权限

如 15.1 节所述，stat 结构中 st_mod 字段的低 12 位定义了文件权限。其中的前 3 位为专用位，分别是 set-user-ID 位、set-group-ID 位和 sticky 位（在图 15-1 中分别被标注为 U、G、T 位），将在 15.4.5 节中详细介绍。其余 9 位则构成了定义权限的掩码，分别授予访问文件的各类用户。文件权限掩码分为 3 类。

- Owner（亦称为 user）：授予文件属主的权限。

chmod(1)之类的命令使用术语 user 的缩写 u 来指代该类权限。

- Group：授予文件属组成员用户的权限。
- Other：授予其他用户的权限。

可为每一类用户授予的权限如下所示。

- Read：可阅读文件的内容。
- Write：可更改文件的内容。
- Execute：可以执行文件（亦即，文件是程序或脚本）。要执行脚本文件（比如，一个 bash 脚本），需同时具备读权限和执行权限。

执行 ls-l 命令，可查看文件的权限和所有权，如下所示：

```
$ ls -l myscript.sh
-rwxr-x--- 1 mtk users 1667 Jan 15 09:22 myscript.sh
```

在以上输出中，将文件权限显示为“rwxr-x---”（该字符串起始处的连接号“-”表明该文件属于普通文件）。在解释该字符串时，需将其一剖为三，以 3 个字符为一组，分别表示读、写、可执行权限具备与否。第一组字符用来表示文件属主的权限，在本例中，则是读、写、执行权限俱全。第二组字符用来表示属组权限，对于本例，组内用户具有读和可执行权限，但不具有写权限。最后一组字符用来表示其他用户的权限，本例中的其他用户没有任何权限。

头文件<sys/stat.h>定义了可与 stat 结构中 st_mode 相与（&）的常量，用于检查特定权限位置位与否。（<fcntl.h>为 open()系统调用提供了原型，在程序中包含该头文件也可定义这些常量。）表 15-4 列出了这些常量。

表 15-4: 用来表示文件权限位的常量

常 量	其 他 值	权 限 位
S_ISUID	04000	Set-user-ID
S_ISGID	02000	Set-group-ID
S_ISVTX	01000	Sticky

续表

常 量	其 他 值	权 限 位
S_IRUSR	0400	User-read
S_IWUSR	0200	User-write
S_IXUSR	0100	User-execute
S_IRGRP	040	Group-read
S_IWGRP	020	Group-write
S_IXGRP	010	Group-execute
S_IROTH	04	Other-read
S_IWOTH	02	Other-write
S_IXOTH	01	Other-execute

除表 15-4 所列常量以外，还分别将各类（属主、属组及其他）权限掩码定义为常量：S_IRWXU (0700)、S_IRWXG (070)和 S_IRWXO (07)。

程序清单 15-3 声明的函数 filePermStr(), 会针对给定的文件权限掩码返回一个静态分配的字符串，以 ls(1)所采用的风格来表示该掩码。

程序清单 15-3: file_perms.c 文件的头文件

```
----- files/file_perms.h
#ifndef FILE_PERMS_H
#define FILE_PERMS_H

#include <sys/types.h>

#define FP_SPECIAL 1          /* Include set-user-ID, set-group-ID, and sticky
                             bit information in returned string */

char *filePermStr(mode_t perm, int flags);

#endif
----- files/file_perms.h
```

如果在 filePermStr()的 flag 参数中设置了 FP_SPECIAL 标志，那么返回的字符串将包括 set-user-ID、set-group-ID，以及 sticky 位的设置信息，其表现形式同样会沿袭 ls(1)的风格。

程序清单 15-4 展示了 filePermStr()函数的实现。程序清单 15-1 中的程序调用了该函数。

程序清单 15-4: 将文件权限掩码转换为字符串

```
----- files/file_perms.c
#include <sys/stat.h>
#include <stdio.h>
#include "file_perms.h"          /* Interface for this implementation */

#define STR_SIZE sizeof("rwxrwxrwx")

char *          /* Return ls(1)-style string for file permissions mask */
filePermStr(mode_t perm, int flags)
{
    static char str[STR_SIZE];
```

```

snprintf(str, STR_SIZE, "%c%c%c%c%c%c%c%c",
    (perm & S_IRUSR) ? 'r' : '-', (perm & S_IWUSR) ? 'w' : '-',
    (perm & S_IXUSR) ?
        (((perm & S_ISUID) && (flags & FP_SPECIAL)) ? 's' : 'x') :
        (((perm & S_ISUID) && (flags & FP_SPECIAL)) ? 'S' : '-'),
    (perm & S_IRGRP) ? 'r' : '-', (perm & S_IWGRP) ? 'w' : '-',
    (perm & S_IXGRP) ?
        (((perm & S_ISGID) && (flags & FP_SPECIAL)) ? 's' : 'x') :
        (((perm & S_ISGID) && (flags & FP_SPECIAL)) ? 'S' : '-'),
    (perm & S_IROTH) ? 'r' : '-', (perm & S_IWOTH) ? 'w' : '-',
    (perm & S_IXOTH) ?
        (((perm & S_ISVTX) && (flags & FP_SPECIAL)) ? 't' : 'x') :
        (((perm & S_ISVTX) && (flags & FP_SPECIAL)) ? 'T' : '-'));

    return str;
}

```

files/file_perms.c

15.4.2 目录权限

目录与文件拥有相同的权限方案，只是对 3 种权限的含义另有所指。

- 读权限：可列出（比如，通过 `ls` 命令）目录之下的内容（即目录下的文件名）。

在实验验证对目录读权限位的操作时，应当了解有些 Linux 发行版对 `ls` 做了别名处理，命令所携带的一些选项（比如，`-F`）需要访问目录中文件的 `i` 节点信息，而这又需要拥有对目录的执行权限。为确保使用的是 `ls` 命令本身，执行时要给出命令的完整路径名（`/bin/ls`）。

- 写权限：可在目录内创建、删除文件。注意，要删除文件，对文件本身无需有任何权限。
- 可执行权限：可访问目录中的文件。因此，有时也将对目录的执行权限称为 `search`（搜索）权限。

访问文件时，需要拥有对路径名所列所有目录的执行权限。例如，想读取文件 `/home/mtk/x`，则需拥有对目录 `/`、`/home` 以及 `/home/mtk` 的执行权限（还要有对文件 `x` 自身的读权限）。若当前的工作目录为 `/home/mtk/sub1`，访问相对路径名 `../sub2/x` 时，需握有 `/home/mtk` 和 `/home/mtk/sub2` 这两个目录的可执行权限（不必有对 `/` 或 `/home` 的执行权限）。

拥有对目录的读权限，用户只是能查看目录中的文件列表。要想访问目录内文件的内容或是这些文件的 `i` 节点信息，还需握有对目录的执行权限。

反之，若拥有对目录的可执行权限，而无读权限，只要知道目录内文件的名称，仍可对其进行访问，但不能列出目录下的内容（即目录所含的其他文件名）。在控制对公共目录内容的访问时，这是一种常用技术，简单而且实用。

要想在目录中添加或删除文件，需要同时拥有对该目录的执行和写权限。

15.4.3 权限检查算法

只要在访问文件或目录的系统调用中指定了路径名称，内核就会检查相应文件的权限。如果赋予系统调用的路径名还包含目录前缀时，那么内核除去会检查对文件本身所需的权限以外，还会检查前缀所含每个目录的可执行权限。内核会使用进程的有效用户 ID、有效组 ID 以及辅助组 ID，来执行权限检查。（准确说来，Linux 内核会使用文件系统用户 ID 和组 ID，而非相应的有效用户 ID 和组 ID，来进行文件权限检查，这一点 9.5 节已经提及。）

一旦调用 `open()` 打开了文件，针对返回描述符的后续系统调用（比如，`read()`、`write()`、`fstat()`、`fcntl()`，以及 `mmap()`）将不再进行任何权限检查。

检查文件权限时，内核所遵循的规则如下。

1. 对于特权级进程，授予其所有访问权限。
2. 若进程的有效用户 ID 与文件的用户 ID（属主）相同，内核会根据文件的属主权限，授予进程相应的访问权限。比方说，若文件权限掩码中的属主读权限（owner-read permission）位被置位，则授予进程读权限。否则，则拒绝进程对文件的读取操作。
3. 若进程的有效组 ID 或任一附属组 ID 与文件的组 ID（属组）相匹配，内核会根据文件的属组权限，授予进程对文件的相应访问权限。
4. 若以上三点皆不满足，内核会根据文件的 other(其他)权限，授予进程相应权限。

其实，内核代码在实现上述检查规则时，在构造上也颇具匠心。只有当进程通过其他测试未能获得所需要的权限时，才去检查进程是否属于特权级进程。这就省去了对 ASU 进程记账标志的设置，该标志用于标记进程是否曾利用过超级用户特权（见 28.1 节）。

内核会依次执行针对属主、属组以及其他用户的权限检查，只要匹配上述检查规则之一，便会停止检查。这样得出的结果可能会在意料之外，比方说，若组权限超过了属主权限，那么文件属主所拥有的权限要低于组成员的权限，如下例所示：

```
$ echo 'Hello world' > a.txt
$ ls -l a.txt
-rw-r--r--  1 mtk   users   12 Jun 18 12:26 a.txt
$ chmod u-rw a.txt           Remove read and write permission from owner
$ ls -l a.txt
----r--r--  1 mtk   users   12 Jun 18 12:26 a.txt
$ cat a.txt
cat: a.txt: Permission denied  Owner can no longer read file
$ su avr                       Become someone else...
Password:
$ groups                         who is in the group owning the file...
users staff teach cs
$ cat a.txt                       and thus can read the file
Hello world
```

若为文件的其他用户分配的权限大于文件属主或属组，上述论述也同样适用。

由于文件的权限及所有权信息都维护于文件的 `i` 节点之内，故而也为指向同一 `i` 节点的所有文件名（链接）所共享。

Linux 2.6 支持访问控制列表，从而可以以每用户或每组为基础来定义文件权限。若文件与一 ACL 挂钩，内核则会在上述算法的基础上略作改动。本书第 17 章将会介绍 ACL。

检查特权级别进程的权限

上文曾提及，若进程为特权级进程，则内核在检查权限时将授予进程所有的访问权限。这一论述成立，其实还要加个限制条件。对于非目录文件，仅当该文件的 3 种权限类型（至少）之一具有可执行权限时，Linux 才会将该权限赋予一特权级进程。而在其他一些 UNIX 的实行中，即使文件的任何权限类型都不具有可执行权限，特权级进程还是能执行该文件。而当访问目录时，特权级进程总是拥有可执行（搜索）权限。

就两种 Linux 进程能力：CAP_DAC_READ_SEARCH 和 CAP_DAC_OVERRIDE（参见 39.2 节）而言，有必要修改之前对特权级进程的描述。具备 CAP_DAC_READ_SEARCH 能力的进程对任何类型的文件都拥有读权限，对于目录则总是具有可执行和写权限（即总能访问目录中的文件，并能读取目录中的文件列表）。具备 CAP_DAC_OVERRIDE 能力的进程对任何类型的文件都拥有读、写权限，对于目录或是文件在权限分类中的至少一类具有可执行权限的情况下，则该进程对其还拥有可执行权限。

15.4.4 检查对文件的访问权限：access()

如上节所述，当进程访问文件时，系统会以其 effective(有效)用户 ID、effective(有效)组 ID 以及附属组 ID 来确定权限。当然，对于程序（比如，set-user-ID 或 set-group-ID 程序）来说，根据进程的 real（真实）用户 ID 和组 ID 来检查对文件的访问权限，也并非没有可能。

系统调用 access()就是根据进程的真实用户 ID 和组 ID（以及附属组 ID），去检查对 pathname 参数所指定文件的访问权限。

```
#include <unistd.h>

int access(const char *pathname, int mode);

Returns 0 if all permissions are granted, otherwise -1
```

若 pathname 为符号链接，access()将对其解引用。

参数 mode 是由表 15-5 中常量相或 (|) 而成的位掩码。若由 pathname 所指定的文件具备 mode 参数包含的所有权限，access()将返回 0；只要有一项权限未得到满足（或者有错误发生），access()则返回-1。

表 15-5: access()的 mode 常量

常 量	描 述
F_OK	有这个文件吗
R_OK	对该文件有读权限吗
W_OK	对该文件有写权限吗
X_OK	对该文件有执行权限吗

由于对某一文件调用 access()与对同一文件的后续操作之间存在时间差，因此（不论间隔多么短暂）执行后续操作时，也无法保证在对文件的后续操作时由 access()所返回的信息依然正确。在某些应用程序设计中，上述情形可能会导致安全漏洞。

比方说，假设有一 set-user-ID-root 程序，使用 access()来检查程序的真实用户 id 是否可以访问某文件，如果可以访问，就对其执行（open()或 exec()之类的)操作。

问题是，若输入 access()的路径名为符号链接，而恶意用户可抢在第二步检查之前设法更改该链接，使其指向另一文件，则最终会导致 set-user-ID-root 去操作真实用户 ID 并无权限的文件。（这也是对 38.6 节所述检查时间与调用时间之间竞争条件的例证。）正因如此，建议杜绝使用 access()(参见[Borisov, 2005])。对于前文所举示例，可以暂时更改 set-user-ID 进程的

有效（或文件系统）用户 ID 来实施（`open()`或 `exec()`之类的）文件操作，并通过对返回值和 `errno` 的检查来判断，操作失败是否应归咎于权限问题。

GNU C 库提供了一个功能相似的非标准函数 `uidaccess()`（及其同义函数 `eaccess()`），该函数使用进程的有效用户 ID 来检查对文件的访问权限。

15.4.5 Set-User-ID、Set-Group-ID 和 Sticky 位

除了 9 位用来表明属主、属组和其他用户的权限之外，文件权限掩码还另设有 3 个附加位，分别为 `set-user-ID` (bit 04000)、`set-group-ID` (bit 02000)和 `sticky` (bit 01000)位。9.3 节讨论了创建特权级程序时对 `set-user-ID` 和 `set-group-ID` 权限位的使用。`set-group-ID` 位还有两种其他用途：对于在以 `nogroupid` 选项装配的目录下所新建的文件，控制其群组从属关系；可用于强制锁定文件。以上两种用途分别在 15.3.1 节和 55.4 节有所介绍。本节将重点讨论 `sticky` 位的用途。

在老的 UNIX 实现中，提供 `sticky` 位的目的在于让常用程序的运行速度更快。若对某程序文件设置了 `sticky` 位，则首次执行程序时，系统会将其文本¹拷贝保存于交换区中，即“粘”（`stick`）在交换区内，故而能提高后续执行的加载速度。现代 UNIX 实现对内存的管理更为精准，故而也将权限位的这一用法束之高阁。

表 15-4 所示 `Sticky` 权限位的常量名称——`S_ISVTX` 源于对 `sticky` 位的别称：`saved-text` 位。

在现代 UNIX 实现（包括 Linux）中，`sticky` 权限位所起的作用全然不同于老的 UNIX 实现。作用于目录时，`sticky` 权限位起限制删除位的作用。为目录设置该位，则表明仅当非特权进程具有对目录的写权限，且为文件或目录的属主时，才能对目录下的文件进行删除（`unlink()`、`rmdir()`）和重命名（`rename()`）操作。（具有 `CAP_FOWNER` 能力的进程可省去对属主的检查。）可藉此机制来创建为多个用户共享的一个目录，各个用户可在其下创建或删除属于自己的文件，但不能删除隶属于其他用户的文件。为 `/tmp` 目录设置 `sticky` 权限位，原因正在于此。

可通过 `chmod` 命令（`chmod +t file`）或 `chmod()` 系统调用来设置文件的 `sticky` 权限位。若对某文件设置了 `sticky` 权限位，则当执行 `ls-l` 命令显示该文件时，会在其他用户执行权限字段上看到字母 `T`，其大小写则要取决于是否对文件开启了其他用户执行权限位，如下所示：

```
$ touch tfile
$ ls -l tfile
-rw-r--r--  1 mtk  users    0 Jun 23 14:44 tfile
$ chmod +t tfile
$ ls -l tfile
-rw-r--r-T  1 mtk  users    0 Jun 23 14:44 tfile
$ chmod o+x tfile
$ ls -l tfile
-rw-r--r-t  1 mtk  users    0 Jun 23 14:44 tfile
```

15.4.6 进程的文件模式创建掩码：umask()

本节将针对新建文件或目录的权限设置展开深入讨论。对于新建文件，内核会使用 `open()` 或 `creat()` 中 `mode` 参数所指定的权限。对于新建目录，则会根据 `mkdir()` 的 `mode` 参数来设置权

¹ 译者注：段。

限。然而，文件模式创建掩码（简称为 `umask`）会对这些设置进行修改。`umask` 是一种进程属性，当进程新建文件或目录时，该属性用于指明应屏蔽哪些权限位。

进程的 `umask` 通常继承自其父 `shell`，其结果往往正如人们所期望的那样：用户可以使用 `shell` 的内置命令 `umask` 来改变 `shell` 进程的 `umask`，从而控制在 `shell` 下运行程序的 `umask`。

大多数 `shell` 的初始化文件会将 `umask` 默认置为八进制值 `022` (`---w--w-`)。其含义为对于同组或其他用户，应总是屏蔽写权限。因此，假定 `open()` 调用中的 `mode` 参数为 `0666`（即令所有用户享有读、写权限，通常如此），那么对新建文件来说，其属主拥有读、写权限，所有其他用户只具有读权限（针对文件执行 `ls-l` 命令，会显示“`rw-r--r--`”）。同理，假定将 `mkdir()` 的 `mode` 参数指定为 `0777`（即所有用户享有所有权限），那么对于新建目录来说，其属主享有所有权限，同组和其他用户则只拥有读取和执行权限（即 `rxrx-r-x`）。

系统调用 `umask()` 将进程的 `umask` 改变为 `mask` 参数所指定的值。

```
#include <sys/stat.h>

mode_t umask(mode_t mask);

Always successfully returns the previous process umask
```

可以以八进制数或是表 15-4 中所列常量相或 (`|`) 来指定 `mask` 参数。

对 `umask()` 的调用总会成功，并返回进程的前一 `umask`。

程序清单 15-5 演示了 `umask()` 与 `open()` 和 `mkdir()` 的相互配合。运行该程序的结果如下：

```
$ ./t_umask
Requested file perms: rw-rw----           This is what we asked for
Process umask:      ----wx-wx           This is what we are denied
Actual file perms:  rw-r-----         So this is what we end up with

Requested dir. perms: rwxrwxrwx
Process umask:      ----wx-wx
Actual dir. perms:  rwxr--r--
```

程序清单 15-5 使用 `mkdir()` 和 `rmdir()` 系统调用来创建和删除目录，使用 `unlink()` 系统调用来删除文件。以上系统调用将在第 18 章再做讲解。

程序清单 15-5：使用 `umask()`

```
files/t_umask.c

#include <sys/stat.h>
#include <fcntl.h>
#include "file_perms.h"
#include "tspi_hdr.h"

#define MYFILE "myfile"
#define MYDIR  "mydir"
#define FILE_PERMS (S_IRUSR | S_IWUSR | S_IRGRP | S_IWGRP)
#define DIR_PERMS (S_IRWXU | S_IRWXG | S_IRWXO)
#define UMASK_SETTING (S_IWGRP | S_IXGRP | S_IWOTH | S_IXOTH)

int
main(int argc, char *argv[])
{
```

```

int fd;
struct stat sb;
mode_t u;

umask(UMASK_SETTING);

fd = open(MYFILE, O_RDWR | O_CREAT | O_EXCL, FILE_PERMS);
if (fd == -1)
    errExit("open-%s", MYFILE);
if (mkdir(MYDIR, DIR_PERMS) == -1)
    errExit("mkdir-%s", MYDIR);

u = umask(0);          /* Retrieves (and clears) umask value */

if (stat(MYFILE, &sb) == -1)
    errExit("stat-%s", MYFILE);
printf("Requested file perms: %s\n", filePermStr(FILE_PERMS, 0));
printf("Process umask:      %s\n", filePermStr(u, 0));
printf("Actual file perms:  %s\n\n", filePermStr(sb.st_mode, 0));
if (stat(MYDIR, &sb) == -1)
    errExit("stat-%s", MYDIR);
printf("Requested dir. perms: %s\n", filePermStr(DIR_PERMS, 0));
printf("Process umask:      %s\n", filePermStr(u, 0));
printf("Actual dir. perms:   %s\n", filePermStr(sb.st_mode, 0));

if (unlink(MYFILE) == -1)
    errMsg("unlink-%s", MYFILE);
if (rmdir(MYDIR) == -1)
    errMsg("rmdir-%s", MYDIR);
exit(EXIT_SUCCESS);
}

```

files/t_umask.c

15.4.7 更改文件权限：chmod()和fchmod()

可利用系统调用 `chmod()`和 `fchmod()`去修改文件权限。

```

#include <sys/stat.h>

int chmod(const char *pathname, mode_t mode);

#define _XOPEN_SOURCE 500    /* Or: #define _BSD_SOURCE */
#include <sys/stat.h>

int fchmod(int fd, mode_t mode);

```

Both return 0 on success, or -1 on error

系统调用 `chmod()`更改由 `pathname` 参数所指定文件的权限。若该参数所指为符号链接，调用 `chmod()`会改变符号链接所指代文件的访问权限，而非对符号链接自身的访问权限。（符号链接自创建起，其所有权限便为所有用户共享，且这些权限也不得更改。对符号链接解引用时，将忽略所有这些权限。）

系统调用 `fchmod()`更改由打开文件描述符 `fd` 所指代文件的权限。

参数 `mode` 用于描述文件的新权限，可以采用八进制数字形式，亦或是由表 15-4 所列权限位相或 (`|`) 而成的掩码。要想更改文件权限，进程要么具有特权级别 (`CAP_FOWNER`)，

要么其有效用户 ID 于文件的用户 ID (属主) 相匹配。(准确说来, 对于 Linux 系统上的非特权级进程, 需与文件用户 ID 相匹配的是进程的文件系统用户 ID, 而非其有效用户 ID, 如 9.5 节所述。)

要将文件权限设为使所有用户仅具有读权限, 需执行如下系统调用:

```
if (chmod("myfile", S_IRUSR | S_IRGRP | S_IROTH) == -1)
    errExit("chmod");
/* Or equivalently: chmod("myfile", 0444); */
```

要修改文件的特定权限位, 需先调用 stat() 来获取文件的现有权限, 调整想修改的权限位, 然后使用 chmod() 去更新权限。

```
struct stat sb;
mode_t mode;

if (stat("myfile", &sb) == -1)
    errExit("stat");
mode = (sb.st_mode | S_IWUSR) & ~S_IROTH;
/* owner-write on, other-read off, remaining bits unchanged */
if (chmod("myfile", mode) == -1)
    errExit("chmod");
```

执行以上代码, 等价于执行如下 shell 命令:

```
$ chmod u+w,o-r myfile
```

15.3.1 节曾提及, 若一目录驻留于以 `-o bsdgroups` 选项装配的 ext2 文件系统之上, 或是驻留于以 `-o sysvgroups` 选项装配的 ext2 文件系统上, 并且开启了该目录的 `set-group-ID` 权限位, 那么在该目录下新建的文件会继承其父目录 (而非文件创建进程的有效组 ID) 的组所有权。可能会出现这样一种情况, 即文件的组 ID 与创建文件进程的任一组 ID 都不匹配。正因如此, 当非特权级 (不具备 CAP_FSETID 能力的) 进程调用 chmod() (或 fchmod()) 时, 若文件的组 ID 不等于进程的有效组 ID 或是任一辅助组 ID, 内核则总是清除文件的 `set-group-ID` 权限位。这一安全措施意在防止用户为其不隶属的组创建 `set-group-ID` 程序。以下 shell 命令演示了上述安全措施所堵住的安全漏洞。

```
$ mount | grep test           Hmmm, /test is mounted with -o bsdgroups
/dev/sda9 on /test type ext3 (rw,bsdgroups)
$ ls -ld /test                Directory has GID root, writable by anyone
drwxrwxrwx  3 root   root   4096 Jun 30 20:11 /test
$ id                           I'm an ordinary user, not part of root group
uid=1000(mtk) gid=100(users) groups=100(users),101(staff),104(teach)
$ cd /test
$ cp ~/myprog .               Copy some mischievous program here
$ ls -l myprog                 Hey! It's in the root group!
-rwxr-xr-x  1 mtk   root   19684 Jun 30 20:43 myprog
$ chmod g+s myprog             Can I make it set-group-ID to root?
$ ls -l myprog                 Hmm, no...
-rwxr-xr-x  1 mtk   root   19684 Jun 30 20:43 myprog
```

15.5 i 节点标志 (ext2 扩展文件属性)

某些 Linux 文件系统允许为文件和目录设置各种各样的 i-node flags(I 节点标志)。该特性是一种非标准的 Linux 扩展功能。

现代 BSD 支持类似于 I 节点标志的特性，使用 `chflags(1)`和 `chflags(2)`加以设置。

`ext2` 是首个支持 i 节点标志的 Linux 文件系统，有时人们也将这些标志称为 `ext2` 扩展文件属性。随后，其他文件系统，诸如 `Btrfs`、`ext3`、`ext4`、`Reiserfs`（自 Linux 2.4.19 起）、`XFS`（自 Linux 2.4.25 和 2.6 起）以及 `JFS`（自 Linux 2.6.17 起），也纷纷加入对 i 节点标志的支持。

各种文件系统对 i 节点标志的支持范围略有不同。要在 `Reiserfs` 文件系统上使用 I 节点标志，需在装配文件系统时带上 `-o attrs` 选项。

在 shell 中，可通过执行 `chattr` 和 `lsattr` 命令来设置和查看 i 节点标志，如下例所示：

```
$ lsattr myfile
----- myfile
$ chattr +ai myfile           Turn on Append Only and Immutable flags
$ lsattr myfile
----ia-- myfile
```

在程序中，可利用 `ioctl()`系统调用来获取并修改 i 节点标志，本节稍后会加以详述。

对普通文件或目录均可设置 i 节点标志。大多数 i 节点标志是供普通文件使用的，也有少部分兼供（或专供）目录使用。表 15-6 对于所支持的 i 节点标志作了总结，展示了程序调用 `ioctl()`时所使用的相应标志名称（定义于 `<linux/fs.h>`中），以及配合 `chattr` 命令使用的选项字母。

表 15-6: i 节点标记

常 量	Chattr 选项	用 途
<code>FS_APPEND_FL</code>	a	仅能在尾部追加（需要特权）
<code>FS_COMPR_FL</code>	c	启用文件压缩（未实现）
<code>FS_DIRSYNC_FL</code>	D	目录更新同步（自 Linux 2.6 起）
<code>FS_IMMUTABLE_FL</code>	i	不可变更（需要特权）
<code>FS_JOURNAL_DATA_FL</code>	j	针对数据启用日志功能（需要特权）
<code>FS_NOATIME_FL</code>	A	不更新文件的上次访问时间
<code>FS_NODUMP_FL</code>	d	不转储
<code>FS_NOTAIL_FL</code>	t	禁用尾部打包
<code>FS_SECRM_FL</code>	s	安全删除（未实现）
<code>FS_SYNC_FL</code>	S	文件（和目录）同步更新
<code>FS_TOPDIR_FL</code>	T	以 Orlov 策略来处理顶层目录（自 Linux 2.6 起）
<code>FS_UNRM_FL</code>	u	可恢复已删除的文件（未实现）

Linux 2.6.19 之前，`<linux/fs.h>`尚未定义表 15-6 所列的 FS 系列常量。相反，针对各种文件系统有一套专门的头文件，将相同值定义为各文件系统所专有的常量名。因此，同一值在 `ext2` 文件系统中被定义为 `<linux/ext2_fs.h>`内的 `EXT2_APPEND_FL`，在 `Reiserfs` 文件系统中被定义为 `<linux/reiser_fs.h>`内的 `REISERFS_APPEND_FL`，其他文件系统以此类推。由于每种文件系统的头文件将相同的值定义为相应常量，因此在不提供 `<linux/fs.h>` 定义的老系统中，可以包含上述任一类型的头文件，并使用各文件系统所专有的名称。

FL 系列变量及其含义如下所示

FS_APPEND_FL

仅当指定 `O_APPEND` 标志时，方能打开文件并写入。（因而迫使所有对文件的更新都追加到文件尾部。）例如，可以用该标志来写入日志文件。只有特权级进程（具备 `CAP_LINUX_IMMUTABLE` 能力）方可设置该标志。

FS_COMPR_FL

文件内容经压缩后存储于磁盘之上。在主流的纯 Linux 文件系统上，`FS_COMPR_FL` 不属于标配特性。（有软件包针对 `ext2` 和 `ext3` 文件系统实现了该特性。）考虑到磁盘存储的低廉成本，以及压缩和解压缩所要耗费的 CPU 开销，再加之一旦将文件压缩起来，对其内容的随机访问也不会像原来那么随心所欲（通过 `lseek()`），故而许多应用都会对此避之不及。

FS_DIRSYNC_FL（自 Linux 2.6 以后）

使得对目录的更新（例如：`open(pathname, O_CREAT)`、`link()`、`unlink()`、`mkdir()`）同步发生。这类似于 13.3 节所述的文件同步更新机制。同样，目录同步更新也存在性能问题。这一设置可以只应用于目录。（14.8.1 节所述的 `MS_DIRSYNC` 装配标志提供了类似功能，但是是针对每个装配而言的。）

FS_IMMUTABLE_FL

将文件设置为不可更改，既不能更新文件数据（`write()`和 `truncate()`），也不能改变文件元数据（即 `chmod()`、`chown()`、`unlink()`、`link()`、`rename()`、`rmdir()`、`utime()`、`setxattr()` 和 `removexattr()`）只有特权级进程（具备 `CAP_LINUX_IMMUTABLE` 能力的进程）可为文件设置这一标志。该标志一旦设定，即便是特权级进程也无法改变文件的内容或元数据。

FS_JOURNAL_DATA_FL

对数据启用日志功能。只有 `ext3` 和 `ext4` 文件系统才支持该标志。这些文件系统提供 3 种层次的日志记录：`journal`（日志）、`ordered`（排序），以及 `writeback`（回写）。所有模式都会记录对文件元数据的更新，而 `journal`（日志）模式额外还记录了对文件数据的变更。而在以排序或回写模式运行日志功能的文件系统上，特权级（具有 `CAP_SYS_RESOURCE` 能力的）进程可以为单个文件设置此标志，从而启用对该文件数据更新的日志功能。（`mount(8)`手册页描述了排序与回写两模式之间的区别。）

FS_NOATIME_FL

访问文件时不更新文件的上次访问时间。这省去了每次访问文件时对 I 节点的更新，故而改进了 I/O 性能。（参见 14.8.1 节介绍 `MS_NOATIME` 标志的内容。）

FS_NODUMP_FL

在使用 `dump(8)`备份系统时跳过具有此标志的文件。正如 `dump(8)`手册页所载，该标志有效与否取决于此命令的 `-h` 选项。

FS_NOTAIL_FL

禁用尾部打包。只有 `Reiserfs` 文件系统才支持该标志。此标志屏蔽了 `Reiserfs` 的尾部打包特性，即尝试将小文件（或是较大文件的最后一段）与其元数据置于同一磁盘块中。

装配 Reiserfs 文件系统时，mount 如带有 `-o noatime` 选项将对整个文件系统禁用尾部打包。

FS_SECRM_FL

安全删除文件。该特性尚未实现，其用意在于删除文件时能够万无一失，将被删除文件的数据覆盖掉，以免磁盘扫描程序能够读取并重建该文件。（要做到对数据真正的安全删除其实颇为复杂。要想稳妥地“抹去”先前记录的数据，需要在磁盘介质上执行多次写入操作，详见[Gutmann, 1996]。）

FS_SYNC_FL

令对文件的更新保持同步。当应用于文件时，该标志将致使对文件的写入操作同步完成（就好像对该文件执行的所有 `open()` 调用都引用了 `O_SYNC` 标志一样）。当应用于目录时，该标志的作用等同于前述的同步目录更新标志。

FS_TOPDIR_FL（自 Linux 2.6 起）

这标志着将在 Orlov 块分配策略的指导下对某一目录进行特殊处理。Orlov 策略的灵感来自于 BSD 系统，是对 ext2 文件系统块分配策略的一种改良，试图增大相关文件（例如：同一目录下的各个文件）在磁盘中比邻而居的几率，进而缩短磁盘的寻道时间。详情请见[Corbet, 2002]和[Kumar, et al. 2008]。只有 EXT2 及其升级版本 EXT3、EXT4 文件系统支持 `FS_TOPDIR_FL`。

FS_UNRM_FL

允许该文件在遭删除后能得以恢复。由于可在内核之外实现文件的恢复机制，因此该特性尚未实现。

一般而言，如果针对某一目录设置了 `i` 节点标志，那么新建于其下的文件和子目录会自动将其继承。不过也有例外。

- `FS_DIRSYNC_FL` (`chattr +D`)标志只能应用于目录，故而也只能为新建于该目录下的子目录所继承。
- 当将 `FS_IMMUTABLE_FL` (`chattr +i`)标志应用于目录时，不会有创建于该目录下的文件或子目录继承此标志，因为该标志会阻止在此目录中添加任何新的条目。

在程序中可以分别调用 `ioctl()` 的 `FS_IOC_GETFLAGS` 和 `FS_IOC_SETFLAGS` 操作，来获取和修改 `i` 节点标志（这两个常量定义于 `<linux/fs.h>`）。以下代码演示了如何为打开文件描述符 `fd` 所指代的文件设置 `FS_NOATIME_FL` 标志。

```
int attr;

if (ioctl(fd, FS_IOC_GETFLAGS, &attr) == -1)    /* Fetch current flags */
    errExit("ioctl");
attr |= FS_NOATIME_FL;
if (ioctl(fd, FS_IOC_SETFLAGS, &attr) == -1)    /* Update flags */
    errExit("ioctl");
```

想改变文件的 `i` 节点标志，至少要满足下列两种条件之一：其一，进程的有效用户 ID 需匹配文件的用户 ID（属主）；其二，进程享有特权级别（具备 `CAP_FOWNER` 能力）。严格说来，对于 Linux 上运行的非特权进程，与文件的用户 ID 相匹配的是其文件系统用户 ID，而非有效用户 ID（详见 9.5 节）。

15.6 总结

`stat()`系统调用可获取某一文件的相关信息（元数据），其中大部分取自文件的 `i` 节点，这些信息包括文件的所有权、文件权限以及文件时间戳。

程序可调用 `utime()`、`utimes()`或类似编程接口，去更改文件的上次访问时间及上次修改时间。

每个文件都有一个与之相关的用户 ID（属主）和组 ID，以及一组权限位。为了限制用户对文件的访问权限，把用户划分为 3 类：文件属主（亦称用户）、属组¹以及其他用户。可把 3 种权限授予上述 3 类用户，分别是读、写、可执行权限。目录也与之相同，但权限位的含义则略有不同。可利用系统调用 `chown()`和 `chmod()`来更改文件的所有权及权限。系统调用 `umask()`则用来设置权限的位掩码，当进程新建文件时，会按位掩码来关闭相应权限位。

文件和目录还用到了 3 个额外的权限位。可将 `set-user-ID` 和 `set-group-ID` 权限位应用于程序文件，在进程的执行过程中假借另一有效用户或组 `id`（亦即属于该程序文件）的身份从而获得特权。在以 `nogroupid` (`sysvgroup`)选项装配的文件系统上，对驻留于其上的目录，可通过设置 `set-group-ID` 权限位来控制如下行为：该目录下新建文件的组 ID 是继承进程的有效组 ID，还是父目录的组 ID。当将 `sticky` 权限位应用于目录时，其作用相当于限制删除标志。

`I` 节点标记控制着文件和目录的各种行为。尽管发源于 `ext2`，但如今已得到了几种其他文件系统的支持。

15.7 练习

- 15-1.** 15.4 节中描述了针对各种文件系统操作所需的权限。请使用 `shell` 命令或编写程序来回答或验证以下说法。
- 将文件属主的所有权限“剥夺”后，即使“本组”和“其他”用户仍有访问权，属主也无法访问文件。
 - 在一个可读但无可执行权限的目录下，可列出其中的文件名，但无论文件本身的权限如何，也不能访问其内容。
 - 要创建一个新文件，打开一个文件进行读操作，打开一个文件进行写操作，以及删除一个文件，父目录和文件本身分别需要具备何种权限？对文件执行重命名操作时，源及目标目录分别需要具备何种权限？若重命名操作的目标文件已存在，该文件需要具备何种权限？为目录设置 `sticky` 位(`chmod +t`)，将如何影响重命名和删除操作？
- 15-2.** 你认为系统调用 `stat()`会改变文件 3 个时间戳中的任意之一吗？请解释原因。
- 15-3.** 在运行 Linux 2.6 的系统上修改程序清单 15-1(`t_stat.c`)，令其可以纳秒级精度来显示文件时间戳。
- 15-4.** 系统调用 `access()`会利用进程的实际用户和组 ID 来检查权限。请编写相应函数，根据进程的有效用户和组 ID 来进行权限检查。

¹ 译者注：同组用户。

15-5. 如 15.4.6 节所述, `umask()`总会在设置进程 `umask` 的同时返回老 `umask` 的拷贝。请问, 如何在不改变进程当前 `umask` 的同时获取到其拷贝?

15-6. 命令 `chmod a+rX file` 的作用是对所有各类用户授予读权限, 并且, 当 `file` 是目录, 或者 `file` 的任一用户类型具有可执行权限时, 将向所有各类用户授予可执行权限, 如下例所示:

```
$ ls -ld dir file prog
dr----- 2 mtk users  48 May  4 12:28 dir
-r----- 1 mtk users 19794 May  4 12:22 file
-r-x----- 1 mtk users 19336 May  4 12:21 prog
$ chmod a+rX dir file prog
$ ls -ld dir file prog
dr-xr-xr-x 2 mtk users  48 May  4 12:28 dir
-r--r--r-- 1 mtk users 19794 May  4 12:22 file
-r-xr-xr-x 1 mtk users 19336 May  4 12:21 prog
```

使用 `stat()`和 `chmod()`编写一程序, 令其等效于执行 `chmod a+rX` 命令。

15-7. 编写 `chattr(1)`命令的简化版来修改文件的 `i` 节点标志。参阅 `chattr(1)` 手册页以掌握 `chattr` 命令行接口的细节。(无需实现 `-R`、`-V`、`-v` 选项。)

第 16 章

扩展属性

本章将介绍文件的扩展属性（EA），即以名称-值对形式将任意元数据与文件 *i* 节点关联起来的技术。Linux 自版本 2.6 起，开始支持 EA。

16.1 概述

EA 可用于实现访问列表（第 17 章）和文件能力（第 39 章）。但就设计而论，其能力绝不仅限于此。例如，还可利用 EA 去记录文件的版本号、与文件的 MIME 类型/字符集有关的信息，或是指向图符的指针。

SUSv3 并未对 EA 加以规范。但少数其他 UNIX 实现却提供了类似的特性，其中知名的有现代 BSD（详见 `extattr(2)`）系列和 Solaris 9 及其后续版本（详见 `fsattr(5)`）。

EA 需要有底层文件系统来提供支撑，Btrfs、ext2、ext3、ext4、JFS、Reiserfs 以及 XFS 等文件系统都支持 EA。

各类文件系统对 EA 的支持都属可选项，受内核配置选项中的“File systems”菜单控制。Reiserfs 文件系统自 Linux 2.6.7 开始支持 EA。

EA 命名空间

EA 的命名格式为 `namespace.name`。其中 `namespace` 用来把 EA 从功能上划分为截然不同的几大类，而 `name` 则用来在既定命名空间内唯一标识某个 EA。

可供 `namespace` 使用的值有 4 个：`user`、`trusted`、`system` 以及 `security`。这 4 类 EA 的用途如下所示。

- `user` EA 将在文件权限检查的制约下由非特权级进程操控。欲获取 `user` EA 值，需要有文件的读权限；欲改变 `user` EA 值，则需要写权限。（若无所需权限，将会导致 `EACCES` 错误。）在 `ext2`、`ext3`、`ext4` 或 `Reiserfs` 文件系统上，如欲将 `user` EA 与一文件关联，在装配底层文件系统时需带有 `user_xattr` 选项。

```
$ mount -o user_xattr device directory
```

- trusted EA 也可由用户进程“驱使”，这点与 user EA 相似。而区别则在于，要操纵 trusted EA，进程必须具有特权（CAP_SYS_ADMIN）。
- system EA 供内核使用，将系统对象与一文件关联。目前仅支持访问控制列表（第 17 章）。
- security EA 的作用有二：其一，用来存储服务于操作系统安全模块的文件安全标签；其二，将可执行文件与能力关联起来（39.9.2 节）。而发明 security EA 的初衷是为了支持安全强化版的 Linux(SELinux, <http://www.nsa.gov/research/selinux/>)。

一个 i 节点可以拥有多个相关 EA，其所从属的命名空间可以相同，也可不同。在各命名空间内的 EA 名均自成一体。在 user 和 trusted 命名空间内，EA 名可以为任意字符串。而在 system 命名空间内，只有经内核明确认可的（例如，用于访问控制列表的）命名方可使用。

JFS 支持另一种命名空间——os2，其他文件系统均未实现。提供这一命名空间是为了支持传统的 OS/2 文件系统 EA。进程无需特权，就能创建 OS2 EA。

通过 shell 创建并查看 EA

在 shell 中，可执行 setfattr(1)和 getfattr(1)命令来设置和查看文件的 EA。

```
$ touch tfile
$ setfattr -n user.x -v "The past is not dead." tfile
$ setfattr -n user.y -v "In fact, it's not even past." tfile
$ getfattr -n user.x tfile          Retrieve value of a single EA
# file: tfile                     Informational message from getfattr
user.x="The past is not dead."    The getfattr command prints a blank
                                  line after each file's attributes
$ getfattr -d tfile              Dump values of all user EAs
# file: tfile
user.x="The past is not dead."
user.y="In fact, it's not even past."

$ setfattr -n user.x tfile      Change value of EA to be an empty string
$ getfattr -d tfile
# file: tfile
user.x
user.y="In fact, it's not even past."

$ setfattr -x user.y tfile      Remove an EA
$ getfattr -d tfile
# file: tfile
user.x
```

以上 shell 会话所展示的要点之一是，EA 值可以为空字符串，这不同于未定义的 EA 值。（由 shell 会话结尾处的例子可知，user.x 的值为空字符串，user.y 的值为未定义。）

默认情况下，getfattr 只会列出 user EA 值。还可利用 -m 选项来指定一正则表达式，来筛选想要显示的 EA 名：

```
$ getfattr -m 'pattern' file
```

pattern 的默认值为 “^user\.”。可执行如下命令，列出一个文件的所有 EA 值。

```
$ getfattr -m - file
```

16.2 扩展属性的实现细节

本节是对上一节内容的延伸，描述 EA 实现的部分细节。

对 user 扩展属性的限制

user EA 只能施之于文件或目录，之所以将其他文件类型排除在外，原因如下。

- 对于符号链接，会对所有用户开启所有权限，且不容更改。（如 18.2 节所述，符号链接的权限在 Linux 上毫无意义。）这意味着，无法利用权限来阻止任意用户将 user EA 置于符号链接之上。要想解决这个问题，就得防止所有用户针对符号链接创建 user EA。
- 对于设备文件、套接字以及 FIFO 而言，授予用户权限，意在对其针对底层对象所执行的 I/O 操作加以控制。如欲操控这些权限，转而求取对创建 user EA 的控制，则二者间会产生冲突。

此外，若某一目录启用了粘性位（sticky 位）（15.4.5 节），且为其他用户所拥有，则非特权进程不能将一 user EA 置于该目录之上。惟其如此，才能防止任一用户将 EA 附着于诸如/tmp 之类的目录，由于其可写权限对所有用户开放（从而导致任意用户均可操纵此目录的 EA），而设置粘性位，意在防止用户删除该目录下为其他用户所拥有的文件。

EA 在实现方面的限制

Linux VFS 针对所有文件系统上的 EA 均施以如下限制。

- EA 名称的长度不能超过 255 个字节。
- EA 值的容量为 64KB。

此外，某些文件系统对可与文件挂钩的 EA 数量及其大小还有更为严格的限制。

- 在 ext2、ext3 以及 ext4 文件系统中，与一文件关联的所有 EA 命名和 EA 值的总字节数不会超过单个逻辑磁盘块（14.3 节）的大小：1024 字节、2048 字节或 4096 字节。
- 在 JFS 上，为某一文件所使用的所有 EA 名和 EA 值的总字节数上限为 128KB。

16.3 操控扩展属性的系统调用

本节将会介绍用来更新、获取以及删除 EA 的系统调用。

创建和修改 EA

系统调用 `setxattr()`、`lsetxattr()` 以及 `fsetxattr()` 用来设置文件的 EA 值之一。

```
#include <sys/xattr.h>

int setxattr(const char *pathname, const char *name, const void *value,
             size_t size, int flags);
int lsetxattr(const char *pathname, const char *name, const void *value,
             size_t size, int flags);
int fsetxattr(int fd, const char *name, const void *value,
             size_t size, int flags);

All return 0 on success, or -1 on error
```

这 3 个系统调用之间的区别类似于 `stat()`、`lstat()` 以及 `fstat()` (15.1 节) 三者间的差异。

- `setxattr()` 通过 `pathname` 来标识文件，若文件名为符号链接，则对其解引用。
- `lsetxattr()` 通过 `pathname` 来标识文件，但不会对符号链接解引用。
- `fsetxattr()` 则通过打开文件描述符 `fd` 来标识文件。

以上 3 者之间的差异同样适用于本节下面将要介绍的其他各组系统调用。

参数 `name` 是一个以空字符结尾的字符串，定义了 EA 的名称。参数 `value` 是一个指向缓冲区的指针，包含了为 EA 定义的新值。参数 `size` 则指明了缓冲区大小。

默认情况下，若具有给定名称 (`name`) 的 EA 不存在，上述系统调用会创建一个新 EA。若 EA 已经存在，则将替换 EA 值。可利用参数 `flags` 将这一行为控制得更为精准。将该参数指定为 0，以获得默认行为，或者可将其指定为如下常量之一。

XATTR_CREATE

若具有给定名称 (`name`) 的 EA 已经存在，则失败。

XATTR_REPLACE

若具有给定名称 (`name`) 的 EA 不存在，则失败。

下例使用 `setxattr()` 创建了一个 user EA:

```
char *value;

value = "The past is not dead.";

if (setxattr(pathname, "user.x", value, strlen(value), 0) == -1)
    errExit("setxattr");
```

获取 EA 值

可利用系统调用 `getxattr()`、`lgetxattr()` 以及 `fgetxattr()` 来获取 EA 值。

```
#include <sys/xattr.h>

ssize_t getxattr(const char *pathname, const char *name, void *value,
                size_t size);
ssize_t lgetxattr(const char *pathname, const char *name, void *value,
                 size_t size);
ssize_t fgetxattr(int fd, const char *name, void *value,
                 size_t size);

All return (nonnegative) size of EA value on success, or -1 on error
```

参数 `name` 是一个以空字符结尾的字符串，用来标识欲取值的 EA。返回的 EA 值保存于参数 `value` 所指向的缓冲区中。该缓冲区必须由调用者分配，其大小应在 `size` 中指定。若调用成功，上述系统调用会返回复制到 `value` 所指缓冲区中的字节数。

若文件不含名为 “`name`” 的属性，上述系统调用则会失败，并会返回错误 `ENODATA`。若 `size` 值过小，上述系统调用也会失败，并返回错误 `ERANGE`。

可把 `size` 指定为 0，对于这种情况，将忽略 `value` 值，但系统调用仍将返回 EA 值的大小。可利用这一机制来确定后续系统调用在实际获取 EA 值时所需的 `value` 缓冲区大小。但是应当注意，这并不能保证后续在通过系统调用获取 EA 值时，上述返回值就足够大。系统调用期间，

另一进程可能为文件的这一属性分配了较大的值，或是将其完全删除。

删除 EA

系统调用 `removexattr()`、`lremovexattr()`以及 `fremovexattr()`用来删除文件的 EA。

```
#include <sys/xattr.h>

int removexattr(const char *pathname, const char *name);
int lremovexattr(const char *pathname, const char *name);
int fremovexattr(int fd, const char *name);

All return 0 on success, or -1 on error
```

`name` 所含以空字符结尾的字符串，用于标识打算删除的 EA。若试图删除不存在的 EA，调用将失败，并会返回错误 `ENODATA`。

获取与文件相关联的所有 EA 的名称

执行系统调用 `listxattr()`、`llistxattr()`以及 `flistxattr()`，所返回的列表会包含与某文件关联的所有 EA 的名称。

```
#include <sys/xattr.h>

ssize_t listxattr(const char *pathname, char *list, size_t size);
ssize_t llistxattr(const char *pathname, char *list, size_t size);
ssize_t flistxattr(int fd, char *list, size_t size);

All return number of bytes copied into list on success, or -1 on error
```

调用将 EA 的名称列表以一系列以空字符结尾的字符串形式置于 `list` 所指向的缓冲区中。缓冲区的大小由 `size` 指定。一旦成功，上述系统调用会返回复制到 `list` 中的字节数。

与 `getxattr()` 一样，也可将 `size` 指定为 0，系统调用将忽略 `list`，并返回后续调用实际获取 EA 名称列表（假定该列表尚未改变）时所需的缓冲区大小。

想获取与某文件相关联的 EA 名列表，只需对文件拥有“访问”权限（亦即对 `pathname` 下的所有路径均拥有执行权限），对文件本身则无需任何权限。

出于安全考虑，`list` 中返回的 EA 名称可能不包含调用进程无权访问的属性名。比方说，在非特权进程中调用 `listxattr()` 时，大多数文件系统都会略去 `trusted` 属性。请注意上一句中的“可能”二字，这表明文件系统实现并非一定要如此。因而，使用 `list` 中返回的 EA 名去调用 `getxattr()`，是有可能失败的，因为进程并不具有获得该 EA 值所需的特权。（同样，当另一进程在 `listxattr()` 和 `getxattr()` 调用之间将该属性删除，也会发生类似错误。）

程序示例

程序清单 16-1 所示程序将获取并显示命令行所列文件的所有 EA 名和 EA 值。该程序使用 `listxattr()`，去获取与每个文件相关联的所有 EA 名称，随后循环调用 `getxattr()`，为每个名称获取相应的值。默认以纯文本方式显示属性值。若带有 `-x` 选项，那么属性值将以十六进制字符串形式显示。以下 shell 会话记录展示了该程序的使用。

```

$ setfattr -n user.x -v "The past is not dead." tfile
$ setfattr -n user.y -v "In fact, it's not even past." tfile
$ ./xattr_view tfile
tfile:
    name=user.x; value=The past is not dead.
    name=user.y; value=In fact, it's not even past.

```

程序清单 16-1: 显示文件的扩展属性

xattr/xattr_view.c

```

#include <sys/xattr.h>
#include "tspi_hdr.h"

#define XATTR_SIZE 10000

static void
usageError(char *progName)
{
    fprintf(stderr, "Usage: %s [-x] file...\n", progName);
    exit(EXIT_FAILURE);
}

int
main(int argc, char *argv[])
{
    char list[XATTR_SIZE], value[XATTR_SIZE];
    ssize_t listLen, valueLen;
    int ns, j, k, opt;
    Boolean hexDisplay;

    hexDisplay = 0;
    while ((opt = getopt(argc, argv, "x")) != -1) {
        switch (opt) {
            case 'x': hexDisplay = 1;          break;
            case '?': usageError(argv[0]);
        }
    }

    if (optind >= argc + 2)
        usageError(argv[0]);
    for (j = optind; j < argc; j++) {
        listLen = listxattr(argv[j], list, XATTR_SIZE);
        if (listLen == -1)
            errExit("listxattr");

        printf("%s:\n", argv[j]);

        /* Loop through all EA names, displaying name + value */

        for (ns = 0; ns < listLen; ns += strlen(&list[ns]) + 1) {
            printf("    name=%s; ", &list[ns]);

            valueLen = getxattr(argv[j], &list[ns], value, XATTR_SIZE);
            if (valueLen == -1) {
                printf("couldn't get value");
            } else if (!hexDisplay) {
                printf("value=%.*s", (int) valueLen, value);
            } else {

```



```
        printf("value=");
        for (k = 0; k < valueLen; k++)
            printf("%02x ", (unsigned int) value[k]);
    }

    printf("\n");
}

printf("\n");
}

exit(EXIT_SUCCESS);
}
```

xattr/xattr_view.c

16.4 总结

自 2.6 版本以来，Linux 开始支持扩展属性，允许以名称-值对的形式将任意元数据与文件关联起来。

16.5 练习

- 16-1.** 编写一程序，可创建或修改文件的 user EA（亦即，setfattr(1)的简化版）。应将文件名、EA 名以及 EA 值以命令行参数形式提供给该程序。

第 17 章

访问控制列表

15.4 节已经介绍了传统 UNIX（及 Linux）对文件权限的规划方案。对于诸多应用来说，这一方案已经能够满足要求。但还有些应用，需要在为特定用户和组授权时进行更为精密的控制。为满足这一需求，许多 UNIX 系统对传统的 UNIX 文件权限模型进行了名为访问控制列表（ACL）的扩展。利用 ACL，可以在任意数量的用户和组之中，为单个用户或组指定文件权限。自版本 2.6 起，Linux 内核开始支持 ACL。

各文件系统对 ACL 的支持属于可选项，由“File systems”菜单下的内核配置选项控制。Reiserfs 文件系统自内核 2.6.7 起开始支持 ACL。

要想在 ext2、ext3、ext4 或 reiserfs 文件系统上创建 ACL，装配相应的文件系统时需要带 `mount -o acl` 选项。

针对 UNIX 系统，从未正式出台过 ACL 的相关标准。人们曾以 POSIX.1e 和 POSIX.2c 标准草案的形式对此进行过尝试，二者分别定义了服务于 ACL 的 API 和 shell 命令（以及诸如能力之类的其他特性）。最终，这一标准化进程以失败告终，标准草案也随之撤销。不过，诸多 UNIX（包括 Linux）对 ACL 的实现还是遵循上述标准草案（通常是根据最终稿，即第 17 号草案）。但由于在各种 ACL 实现之间还存在着许多差异（部分原因是由于标准草案还不完善），故而如果运用了 ACL 技术，就会危及应用的可移植性。

本章将介绍 ACL，并就其用法提供简明教程。此外，还会讲解用来操纵和获取 ACL 的某些库函数。鉴于此类库函数数量众多，本章不会逐一对其做深入探讨。（详细信息可参考手册页。）

17.1 概述

一个 ACL 由一系列 ACL 记录（以下简称 ACE）组成，其中每条记录都针对单个用户或用户组定义了对文件的访问权限（如图 17-1 所示）。

标记类型	标记限定符	权限	
ACL_USER_OBJ	-	rwx	← 对应于传统的 属主（用户）权限
ACL_USER	1007	r--	
ACL_USER	1010	rwx	
ACL_GROUP_OBJ	-	rwx	← 对应于传统的 组权限
ACL_GROUP	102	r--	
ACL_GROUP	103	-w-	
ACL_GROUP	109	--x	
ACL_MASK	-	rwx	
ACL_OTHER	-	r--	← 对应于传统的 其他用户权限

按组分类记录

图 17-1: 访问控制列表示例

ACL 记录

每条 ACE 都由 3 部分组成。

- 标记类型：表示该记录作用于一个用户、组，还是其他类别的用户。
- 标记限定符（可选项）标识特定的用户或组（亦即，某个用户 ID 或组 ID）。
- 权限集合：本字段包含所授予的权限信息（读、写及执行）。

标记类型取值可为下列各值之一。

ACL_USER_OBJ

带有该标记的 ACE 记录了授予文件属主的权限。每个 ACL 只能包含一条该类型标记的记录。该记录与传统的文件属主（用户）权限相对应。

ACL_USER

携带该值的 ACE 记录了授予某用户（由标记限定符标识）的权限。一个 ACL 可包含零条或多条此标记类型的记录，但针对一个特定用户最多只能定义一条此类记录。

ACL_GROUP_OBJ

包含该值的 ACE 记录了授予文件组的权限。每个 ACL 只会包含一条此标记类型的记录。除非 ACL 还包含类型为“ACL_MASK”的记录，否则此类记录对应于传统的文件组权限。

ACL_GROUP

包含该值的 ACE 记录了授予某个组（由标记限定符标识）的权限。每个 ACL 可包含零条或多条此标记类型的记录，但针对一个特定组最多只能定义一条此类记录。

ACL_MASK

包含该值的 ACE 记录了可由 ACL_USER、ACL_GROUP_OBJ 以及 ACL_GROUP 型 ACE 所能授予的最高权限。一个 ACL 最多只能包含一条标记类型为 ACL_MASK 的 ACE。假如 ACL 含有标记类型为 ACL_USER 或 ACL_GROUP 的记录，那么就必须包含一条 ACL_MASK 型的 ACE。稍后会细述这一标记类型。

ACL_OTHER

对于不匹配任何其他 ACE 的用户，由包含该值的 ACE 授予权限。每个 ACL 只能包含一

条标记类型为“ACL_OTHER”的 ACE。该记录对应于传统的文件其他（other）用户权限。

只有标记类型为“ACL_USER”和“ACL_GROUP”的记录，才会采用标记限定符来指定用户 ID 和组 ID。

最小 ACL 和扩展 ACL

最小化（minimal）ACL 语义上等同于传统的文件权限集合，恰好由 3 条记录组成。每条标记的类型分别为 ACL_USER_OBJ、ACL_GROUP_OBJ 以及 ACL_OTHER。扩展 ACL 则是指除此之外，还包含标记类型为 ACL_USER、ACL_GROUP 和 ACL_MASK 的记录。

之所以要对最小化 ACL 和扩展 ACL 加以区分，原因之一是后者可对传统文件权限模型提供语义的扩展。而另一个原因则与 ACL 的 Linux 实现有关。Linux 系统是以系统扩展属性来实现 ACL 的（详见第 16 章）。用于维护文件访问型 ACL 的系统扩展属性名为 system.posix_acl_access。仅当文件具有扩展 ACL 时，才需使用这一扩展属性。可将针对最小化 ACL 的权限信息存储于传统的文件权限位中。

17.2 ACL 权限检查算法

与传统的文件权限模型（15.4.3 节）相比，对具有 ACL 的文件进行权限检查时，环境并没有什么不同。检查将按以下顺序执行，直至某一标准得到匹配。

1. 若进程具有特权，则拥有所有访问权限。与 15.4.3 节所述的传统文件权限模型相类似，这里也有一个例外。执行某文件时，仅当将可执行权限通过至少一条 ACL 记录授予该文件时，系统才会向特权级进程授予该权限。
2. 若某一进程的有效用户 ID 匹配文件的属主（用户 ID），则授予该进程标记类型为 ACL_USER_OBJ 的 ACE 所指定的权限。严谨的说法是在 Linux 系统中，本节所介绍的 ACL 权限检测使用的是进程的文件系统 ID（请参阅本书 9.5 节），而非其有效用户 ID。
3. 若进程的有效用户 ID 与某一 ACL_USER 类型记录的标记限定符相匹配，则授予该进程此记录所指定权限与 ACL_MASK 型记录值相与（&）的结果。
4. 若进程的组 ID（亦即，有效组 ID 或任一辅助组 ID）之一匹配于文件组（对应于标记类型为 ACL_GROUP_OBJ 的 ACE），或者任一 ACL_GROUP 型记录的标记限定符，则会依次进行如下检查，直至发现匹配项。
 - a) 若进程的组 ID 之一匹配于文件组，且标记类型为 ACL_GROUP_OBJ 的 ACE 授予了所请求的权限，则会依据此记录来判定对文件的访问权限。如果 ACL 中还包含了标记类型为 ACL_MASK 的 ACE，那么对该文件的访问权限将是两记录权限相与（&）后的结果。
 - b) 若进程的组 ID 之一匹配于该文件所辖 ACL_GROUP 型 ACE 的标记限定符，且该 ACE 授予了所请求的权限，那么会依据此记录来判定对文件的访问权限。如果 ACL 中包含了 ACL_MASK 型 ACE，那么对该文件的访问权限应为两记录权限相与（&）的结果。
 - c) 否则，拒绝对该文件的访问。
5. 否则，将以 ACL_OTHER 型 ACE 所记录的权限授予进程。

下面举例说明这些与组 ID 相关的文件授权规则。假定某文件的组 ID 为 100，并受图 17-1 所列 ACL 的保护。若组 ID 为 100 的某一进程发起系统调用 access(file, R_OK)，本次调用将会成功（亦即，返回 0）。（15.4.4 节介绍了 access()。）而另一方面，即便标记类型为

ACL_GROUP_OBJ 的 ACE 授予了所有权限，系统调用 `access(file, R_OK | W_OK | X_OK)` 仍将失败（亦即，返回-1，且将 `errno` 置为 `EACCES`），这是由于访问权限是该类型权限与 `ACL_MASK` 型记录权限相与（&）的结果，而这一结果禁用了对文件的执行权限。

再拿图 17-1 举个例子，假定某进程的组 ID 为 102，其附属组 ID 之一为 103。对该进程来说，调用 `access(file, R_OK)` 和 `access(file, W_OK)` 都会成功，因为这两次调用所请求的文件权限分别匹配标记类型为 `ACL_GROUP`，且标记限定符为 102 和 103 的 ACE 所记录的权限。另外，该进程调用 `access(file, R_OK | W_OK)` 将会失败，因为并无标记类型为 `ACL_GROUP` 的匹配记录同时包含读、写权限。

17.3 ACL 的长、短文本格式

执行 `setfacl` 和 `getfacl` 命令，或是使用某些 ACL 库函数操纵 ACL 时，需指明 ACE 的文本表现形式。ACE 的文本格式有两种。

- 长文本格式的 ACL：每行都包含一条 ACE，还可以包含注释，注释需以“#”开始，直至行尾结束。`getfacl` 命令的输出会以长文本格式显示 ACL。`getfacl` 命令的 `-M acl-file` 选项从指定文件中“提取”长文本格式的 ACL 定义。
- 短文本格式的 ACL：包含一系列以“,”分隔的 ACE。

无论是上述哪种格式，每条 ACE 都由以“:”分隔的 3 部分组成。

tag-type:[tag-qualifier]: permissions

标记类型字段的取值限于表 17-1 第一列所示范围之内。标记类型之后的标记限定符为可选项，采用名称或数字 ID 来标识用户或组。仅当标记类型为 `ACL_USER` 和 `ACL_GROUP` 时，才允许标记限定符的存在。

表 17-1: 对 ACE 文本格式的解释

标记文本格式	是否存在标记限定符	对应的标记类型	ACE 的用途
user	N	ACL_USER_OBJ	文件属主（用户）
u, user	Y	ACL_USER	特定用户
g, group	N	ACL_GROUP_OBJ	文件组
g, group	Y	ACL_GROUP	特定组
m, mask	N	ACL_MASK	组分类掩码
o, other	N	ACL_OTHER	其他用户

以下所示为短文本格式的 ACL，对应于传统权限掩码 0650:

```
u::rw-,g::r-x,o::---
u::rw,g::rx,o::-
user::rw,group::rx,other::-
```

下面这一短文本格式 ACL 则包含了两条命名用户 ACE、一条命名组 ACE 以及一条掩码 ACE。

```
u::rw,u:paulh:rw,u:annabel:rw,g::r,g:teach:rw,m::rwx,o::-
```

17.4 ACL_mask 型 ACE 和 ACL 组分类

如果一个 ACL 包含了标记类型为 ACL_USER 或 ACL_GROUP 的 ACE，那么也一定会包含标记类型为 ACL_MASK 的 ACE。若 ACL 未包含任何标记类型为 ACL_USER 或 ACL_GROUP 的 ACE，那么标记类型为 ACL_MASK 的 ACE 则为可选项。

对于 ACL_MASK 标记类型的 ACE，其作用在于是所谓“组分类”（group class）中 ACE 所能授予权限的上限。组分类是指在 ACL 中，由所有标记类型为 ACL_USER、ACL_GROUP 以及 ACL_GROUP_OBJ 的 ACE 所组成的集合。

提供标记类型为 ACL_MASK 的 ACE，其目的在于即使运行并无 ACL 概念的应用程序，也能保障其行为的一致性。作为这一论点的例证，假设与文件关联的 ACL 包含以下记录：

```
user::rwx                # ACL_USER_OBJ
user:paulh:r-x          # ACL_USER
group::r-x              # ACL_GROUP_OBJ
group:teach:--x        # ACL_GROUP
other:--x               # ACL_OTHER
```

若某程序针对该文件按以下方式调用 chmod()。

```
chmod(pathname, 0700);    /* Set permissions to rwx----- */
```

对于对 ACL 一无所知的应用程序而言，这意味着“除文件属主以外，不允许其他任何用户访问”。即便存在针对该文件的 ACL，这层意思也不会变。如果 ACL 中不含 ACL_MASK 型记录，那么可以有多种方法来实现这一行为，但每种方法都存在缺陷。

- 只是将 ACL_GROUP_OBJ 和 ACL_OTHER 型记录的掩码简单地修改为---是不足以解决问题的，因为用户 paulh 和组 teach 依旧对该文件拥有某些权限。
- 另一种可能是，将针对组和其他用户的权限新设置（即，全部屏蔽）应用于标记类型为 ACL_USER、ACL_GROUP、ACL_GROUP_OBJ 以及 ACL_OTHER 的记录。

```
user::rwx                # ACL_USER_OBJ
user:paulh:---          # ACL_USER
group:---                # ACL_GROUP_OBJ
group:teach:---        # ACL_GROUP
other:---                # ACL_OTHER
```

这一方法的问题在于，之前由具有 ACL 概念的应用所确立的文件权限语义会被对 ACL 一无所知的应用所“错杀”，因为如下调用（举例说明）不会将 ACL 中的 ACL_USER 和 ACL_GROUP 型记录恢复到其之前的状态：

```
chmod(pathname, 751);
```

- 要避免这些问题，可以考虑将标记类型为 ACL_GROUP_OBJ 的记录置为对所有 ACL_USER 和 ACL_GROUP 类记录的约束。然而，这也意味着总是需要将 ACL_GROUP_OBJ 型记录置为 ACL_USER 和 ACL_GROUP 型记录所允许权限的并集。而系统又会使用 ACL_GROUP_OBJ 型记录来判定赋予文件组的权限，这会引发冲突。

设计标记类型为 ACL_MASK 的记录，正是为了解决上述问题。这一机制在实现传统意义上的 chmod() 操作的同时，也无损于由具有 ACL 概念的应用所确立的文件权限语义。当 ACL 包含标记类型为 ACL_MASK 的 ACE 时：

- 调用 chmod() 对传统组权限所做的变更，会改变 ACL_MASK（而非 ACL_GROUP_OBJ）

标记类型 ACE 的设置。

- 调用 `stat()`，在 `st_mode` 字段（图 15-1）的组权限位中会返回 `ACL_MASK` 权限（而非 `ACL_GROUP_OBJ` 权限）。

尽管 `ACL_MASK` 型记录的出现保护了 ACL 信息，使其免遭并无 ACL 概念的应用的“误伤”，反之却并非如此。ACL 的优先级要高于对文件组权限的传统操作。例如，假设为某文件设置了如下 ACL：

```
user::rw-,group:---,mask:---,other::r--
```

若针对该文件执行 `chmod g+rw` 命令，则 ACL 将会变为：

```
user::rw-,group:---,mask::rw-,other::r--
```

这时，组用户仍无法访问该文件。一种迂回策略是修改针对组的 ACE，赋予其所有权限。结果，组用户总是能获得 `ACL_MASK` 型记录的所有权限。

17.5 getfacl 和 setfacl 命令

在 shell 中运行 `getfacl` 命令，可查看到应用于文件的 ACL。

```
$ umask 022                Set shell umask to known state
$ touch tfile              Create a new file
$ getfacl tfile
# file: tfile
# owner: mtk
# group: users
user::rw-
group::r--
other::r--
```

由 `getfacl` 命令的输出可知，新建文件具有最小的 ACL 权限。`getfacl` 命令会在输出 ACL 记录的文本格式之前，显示该文件的名称和属主、属组。执行 `getfacl` 命令时，如带有 `--omit-header` 选项，可省略上述内容。

接下来的例子则显示，执行传统的 `chmod` 命令来改变文件访问权限时，其效果贯穿到文件的 ACL 上。

```
$ chmod u=rwx,g=rx,o=x tfile
$ getfacl --omit-header tfile
user::rwx
group::r-x
other::-x
```

`setfacl` 命令可用来修改文件的 ACL。下例中执行 `setfacl -m` 命令，为文件的 ACL 追加标记类型为 `ACL_USER` 和 `ACL_GROUP` 的记录。

```
$ setfacl -m u:paulh:rx,g:teach:x tfile
$ getfacl --omit-header tfile
user::rwx
user:paulh:r-x                ACL_USER entry
group::r-x
group:teach:--x              ACL_GROUP entry
mask::r-x                    ACL_MASK entry
other::-x
```

带-m选项的 setfacl 命令可修改现有 ACE，或者，当给定标记类型和限定符的 ACE 不存在时，会追加新的 ACE。setfacl 命令还可使用-R 选项，将指定的 ACL “递归”应用于目录树中的所有文件。

由 getfacl 命令的输出可知，setfacl 自动为该 ACL 新建了一条标记类型为 ACL_MASK 的记录。

追加了 ACL_USER 和 ACL_GROUP 标记类型的记录会将该 ACL 转变为扩展 ACL。因此，在执行 ls -l 命令时，会在文件的传统权限掩码之后多一个加号（“+”）。

```
$ ls -l tfile
-rwxr-x--x+ 1 mtk      users          0 Dec 3 15:42 tfile
```

接下来继续执行 setfacl 命令，以禁用 ACL_MASK 标记类型记录中除执行权限以外的所有权限，然后再执行 getfacl 命令来查看文件的 ACL。

```
$ setfacl -m m::x tfile
$ getfacl --omit-header tfile
user::rwx
user:paulh:r-x          #effective:--x
group::r-x              #effective:--x
group:teach:--x
mask:--x
other:--x
```

在用户 paulh 和文件组输出后的 “#effective:” 注释是指在与 ACL_MASK 型记录相与（AND）后，由上述记录所赋予的权限实际上要小于记录中所描述的情况。

再次执行 ls -l 命令来观察文件的传统权限位，由输出可知，组分类权限位反映的是 ACL_MASK 型记录的权限 (--x)，而非 ACL_GROUP 型记录中的权限(r-x)。

```
$ ls -l tfile
-rwx--x--x+ 1 mtk      users          0 Dec 3 15:42 tfile
```

setfacl -x 则用来从 ACL 中删除记录。下例删除了用户 paulh 和组 teach 的记录（删除 ACE 时无需指定其权限）：

```
$ setfacl -x u:paulh,g:teach tfile
$ getfacl --omit-header tfile
user::rwx
group::r-x
mask::r-x
other:--x
```

请注意，在执行上述操作时，setfacl 命令会自动将掩码型 ACE 调整为所有组分类 ACE 权限的集合（只有一条此类 ACE：ACL_GROUP_OBJ）。若不想进行这种调整，执行 setfacl 命令时要带上-n 选项。

最后需要说明的是，执行带-b 选项的 setfacl 命令，可从 ACL 中删除所有扩展 ACE，而只保留最小化 ACE（亦即，用户、组及其他）。

17.6 默认 ACL 与文件创建

行文至此，对 ACL 的讨论所描述的均属访问型（access）ACL。顾名思义，当进程访问与该 ACL 相关的文件时，将使用访问型 ACL 来判定进程对文件的访问权限。针对目录，还可创建第二种 ACL：默认型（default）ACL。

访问目录时，默认型 ACL 并不参与判定所授予的权限。相反，默认型 ACL 的存在与否

决定了在目录下所创建文件或子目录的 ACL 和权限。(默认型 ACL 存储于名为 `system.posix_acl_default` 的扩展属性中。)

想查看和设置与目录相关的默认型 ACL, 需要执行带有 `-d` 选项的 `getfacl` 和 `setfacl` 命令。

```
$ mkdir sub
$ setfacl -d -m u::rwx,u:paulh:rx,g::rx,g:teach:rwx,o::- sub
$ getfacl -d --omit-header sub
user::rwx
user:paulh:r-x
group::r-x
group:teach:rwx
mask::rwx
other:----
                                setfacl generated ACL_MASK entry automatically
```

执行带有 `-k` 选项的 `setfacl` 命令, 可删除针对目录而设的默认型 ACL。

若针对目录设置了默认型 ACL, 则:

- 新建于目录下的子目录会将该目录的默认型 ACL 继承为其默认型 ACL。换言之, 默认型 ACL 会随子目录的创建而沿目录树传播开来。
- 新建于目录下的文件或子目录会将该目录的默认型 ACL 继承为其访问型 ACL。与传统文件权限位相对应的 ACL 记录将和创建文件或子目录时系统调用 (`open()`、`mkdir()`等等) 中的 `mode` 参数相与 (&)。所谓“对应的 ACL 记录”是指:
 - `L_USER_OBJ`;
 - `ACL_MASK`, 若不含 `ACL_MASK`, 则为 `ACL_GROUP_OBJ`;
 - `ACL_OTHER`。

一旦目录拥有默认型 ACL, 那么对于新建于该目录下的文件来说, 进程的 `umask` (15.4.6 节) 并不参与判定文件访问型 ACL 中所记录的权限。

试举一例, 演示一新建文件如何将其父目录的默认型 ACL 继承为自身的访问型 ACL。假设使用如下 `open()` 调用, 在前例所建目录下创建一新文件:

```
open("sub/tfile", O_RDWR | O_CREAT,
      S_IRWXU | S_IXGRP | S_IXOTH); /* rwx--x--x */
```

这一新文件的访问型 ACL 如下:

```
$ getfacl --omit-header sub/tfile
user::rwx
user:paulh:r-x          #effective:--x
group::r-x             #effective:--x
group:teach:rwx        #effective:--x
mask:--x
other:----
```

若该目录并无默认 ACL, 则:

- 新建于该目录下的子目录也不存在默认 ACL。
- 会沿用传统规则来设置目录下新建文件或目录的权限。除去按进程的 `umask` 而屏蔽权限位之外, 将文件权限置为 (`open()`、`mkdir()`等调用中) `mode` 参数的值。这时, 新文件将拥有最小化的 ACL。

17.7 ACL 在实现方面的限制

各类文件系统都对一 ACL 中所含记录的条数有所限制。

- 对于 ext2、ext3 以及 ext4 文件系统，某文件所含所有 ACL 记录的总和受制于如下要求：该文件扩展属性的所有名称与值所占字节必须位于同一逻辑磁盘块之内（见 16.2 节）。每条 ACL 记录需占 8 字节，因而一文件所含 ACE 的最大条数会略少于块大小的 1/8（因为 ACL 的扩展属性名称也有一定开销）。因此，大小为 4096 字节的块最多允许 500 条左右的 ACE。（2.6.11 版本之前，内核要求 ext2 和 ext3 文件系统中文件 ACL 的记录总数不得超过 32 条。）
- 对于 XFS 文件系统，每一 ACL 的记录数上限为 25 条。
- 对于 Reiserfs 和 JFS 文件系统，ACL 最多可含 8191 条记录。之所以如此，是由于 VFS 要求扩展属性的值大小不得超过 64KB（见 16.2 节）。

写作本书时，Btrfs 文件系统将 ACL 所含记录的条数限制在 500 条左右。但鉴于该文件系统的开发极其活跃，故而这一限制随时可能发生变化。

尽管上面提及的文件系统大多允许一 ACL 包含大量记录，但出于以下原因，还是应当避免：

- 冗长的 ACL 将增加维护工作的复杂程度，且容易出错；
- 扫描 ACL 寻找匹配记录（在执行组 ID 检查时，还将匹配多条记录）所需的时间，将随记录条数的增长而增长。

通常的做法是：在系统组文件（8.3 节）中定义适当的组，并在 ACL 中运用起来，从而将文件 ACL 的记录条数保持在一个较低的合理水平。

17.8 ACL API

POSIX.1e 标准草案围绕着操纵 ACL 定义了大量函数和数据结构。鉴于其规模庞大，要描述所有函数的细节是不现实的。本节会先对此类函数的用法进行概括，再以相关编程示例作为总结。

程序要使用 ACL API，就应包含 `<sys/acl.h>`。如果还用到了 POSIX.1e 标准草案中的各种 Linux 扩展（acl(5)手册页罗列了一系列 Linux 扩展），程序可能还需要包含 `<acl/libacl.h>`。为与 libacl 库链接，编译此类程序时需带有 `-lacl` 选项。

如前所述，在 Linux 上，ACL 是以扩展属性的方式来实现的，而将 ACL API 实现为一套操纵用户空间数据结构的库函数，并且会在必要时调用 `getxattr()` 和 `setxattr()`，来获取和修改持有 ACL 的持久层 system 扩展属性。此外，应用程序直接调用 `getxattr()` 和 `setxattr()` 去操纵 ACL 也是可行的，尽管并不推荐这一做法。

概述

组成 ACL API 的函数刊载于 acl(5)手册页中。乍看起来，此类函数及数据结构数量之巨，着实令人不得其门而入。图 17-2 概括了各种数据结构之间的关系，并标明了诸多 ACL 函数的用法。

由图 17-2 可知，ACL API 将 ACL 视为一层次化对象：

- 一个 ACL 包含一条或多条 ACL 记录；
- 每条记录均包含一标记类型、一标记限定符（可选），以及一权限集合。

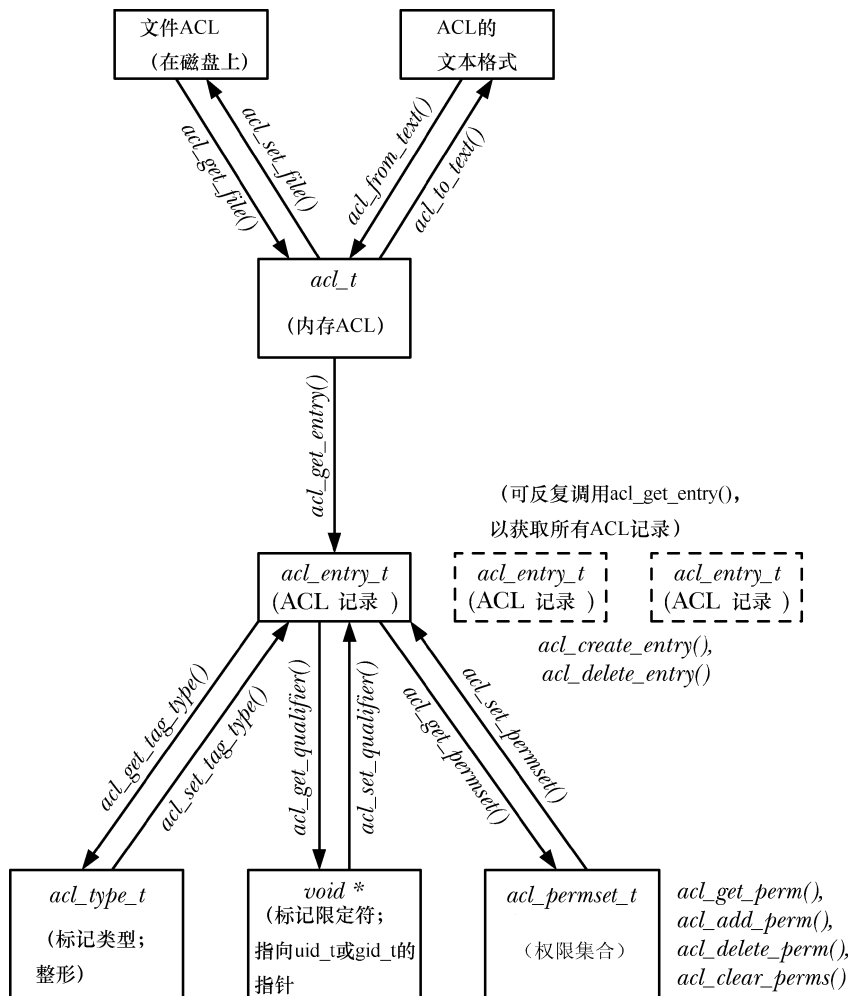


图 17-2: ACL 库函数及数据结构之间的关系

接下来，将简要介绍各种 ACL 函数。多数情况下，不会对每个函数的返回错误加以描述。函数返回整数（状态）时，通常以 0 表示成功，以 -1 表示错误。返回句柄（指针）的函数出错时将返回 NULL。诊断错误时，则可将检查 `errno` 作为常规手段。

句柄（handler）是一抽象术语，用以指代一对象或数据结构。句柄的表现方式由 API 实现决定，例如：可以是指针、数组索引，或者 hash 键。

将文件的 ACL 读入内存

`acl_get_file()`函数可用来获取（由 `pathname` 所标识）文件的 ACL 副本。

```
acl_t acl;
```

```
acl = acl_get_file(pathname, type);
```

取决于参数 `type` 的值（`ACL_TYPE_ACCESS` 或 `ACL_TYPE_DEFAULT`），可调用该函数来获取访问型 ACL 或默认型 ACL。`acl_get_file()`函数将返回一（类型为 `acl_t` 的）句柄，供其

他 ACL 函数使用。

从内存 ACL 中获取记录

`acl_get_entry()`函数会返回一个句柄，指向内存 ACL（由函数的 `acl` 参数指代）中的记录之一。句柄的返回位置由函数的最后一个参数指定。

```
acl_entry_t entry;

status = acl_get_entry(acl, entry_id, &entry);
```

`entry_id` 参数决定返回那条记录的句柄。若将其指定为 `ACL_FIRST_ENTRY`，则会返回的句柄指向 ACL 中的首条 ACE。若将该参数指定为 `ACL_NEXT_ENTRY`，则所返回的句柄将指向上次所获取记录之后的 ACE。因此，在首次调用 `acl_get_entry()`时，把 `type` 参数指定为 `ACL_FIRST_ENTRY`，在随后的调用中，再将其指定为 `ACL_NEXT_ENTRY`，如此这般，就可以遍历 ACL 的所有记录。

若成功获取到一条 ACE，`acl_get_entry()`函数将返回 1；如无记录可取，则返回 0；失败，则返回-1。

获取并修改 ACL 记录中的属性

函数 `acl_get_tag_type()`和 `acl_set_tag_type()`可分别用来获取和修改（由 `entry` 参数所指定）ACL 记录中的标记类型。

```
acl_tag_t tag_type;

status = acl_get_tag_type(entry, &tag_type);
status = acl_set_tag_type(entry, tag_type);
```

`tag_type` 参数类型为 `acl_type_t`（整型），取值可为 `ACL_USER_OBJ`、`ACL_USER`、`ACL_GROUP_OBJ`、`ACL_GROUP`、`ACL_OTHER` 或 `ACL_MASK` 之一。

函数 `acl_get_qualifier()`和 `acl_set_qualifier()`可分别用来获取和修改（由 `entry` 参数所指定）ACL 记录中的标记限定符。下面是两个函数的使用实例，这里假设通过对标记类型的检测，已然确定该记录属于 `ACL_USER`。

```
uid_t *qualp;          /* Pointer to UID */

qualp = acl_get_qualifier(entry);
status = acl_set_qualifier(entry, qualp);
```

仅当 ACE 的标记类型为 `ACL_USER` 或 `ACL_GROUP` 时，标记限定符才有意义。在上例中，`qualp` 是指向用户 ID (`uid_t *`)的一枚指针，在下例中，则是指向组 ID (`gid_t *`)的指针。

函数 `acl_get_permset()`和 `acl_set_permset()`则可分别用来获取和修改（由 `entry` 参数所指定）ACE 中的权限集合。

```
acl_permset_t permset;

status = acl_get_permset(entry, &permset);
status = acl_set_permset(entry, permset);
```

数据类型 `acl_permset_t` 是一个指代权限集合的句柄。

下列函数则用来操纵某一权限集合中的内容：

```
int is_set;

is_set = acl_get_perm(perms, perm);

status = acl_add_perm(perms, perm);
status = acl_delete_perm(perms, perm);
status = acl_clear_perms(perms);
```

在上述各个调用中，可将 `perm` 参数指定为 `ACL_READ`、`ACL_WRITE` 或 `ACL_EXECUTE`——顾名思义。上述函数的用法如下所述：

- 若在（由 `perms` 参数指代的）权限集合中成功激活由 `perm` 参数所指定的权限，`acl_get_perm()`函数将返回 1（真值），否则返回 0。该函数为 Linux 对 POSIX.1e 标准草案的扩展。
- `acl_add_perm()`函数用来向由 `perms` 参数所指代的权限集合中追加由 `perm` 参数所指定的权限。
- `acl_delete_perm()`函数用来从 `perms` 参数所指代的权限集合中删除由 `perm` 参数所指定的权限。（即便要删除的权限在权限集合中并不存在，函数也不会报错。）
- `acl_clear_perm()`函数用来从 `perms` 参数所指代的权限集合中删除所有权限。

创建和删除 ACE

`acl_create_entry()`函数用来在某一现有 ACL 中新建一条记录。该函数会将一个指代新建 ACE 的句柄返回到由其第二个参数所指定的内存位置。

```
acl_entry_t entry;

status = acl_create_entry(&acl, &entry);
```

然后，即可利用先前介绍过的函数来设置该记录。

`acl_delete_entry()`函数用来从 ACL 中删除一条 ACE。

```
status = acl_delete_entry(acl, entry);
```

更新文件的 ACL

`acl_set_file()`函数的作用与 `acl_get_file()`相反，将使用驻留于内存的 ACL 内容（由 `acl` 参数所指代）来更新磁盘上的 ACL。

```
int status;

status = acl_set_file(pathname, type, acl);
```

如欲更新访问型 ACL，需将该函数的 `type` 参数指定为 `ACL_TYPE_ACCESS`；如欲更新目录的默认型 ACL，则需将 `type` 指定为 `ACL_TYPE_DEFAULT`。

ACL 在内存和文本格式之间的转换

`acl_from_text()`函数可将包含文本格式 ACL（长短不拘）的字符串转换为内存 ACL，并返回一个句柄，用以在后续函数调用中指代该 ACL。

```
acl = acl_from_text(acl_string);
```

`acl_to_text()`则执行与上述函数相反的转换，并同时返回对应于 ACL（由 `acl` 参数指定）

的长文本格式字符串。

```
char *str;  
ssize_t len;
```

```
str = acl_to_text(acl, &len);
```

若参数 len 不为 NULL，那么会在该参数所指向的缓冲区中放置返回字符串的长度。

ACL API 中的其他函数

接下来将介绍几个未见诸于图 17-2 的常用 ACL 函数。

`acl_calc_mask(&acl)`函数用来计算并设置内存 ACL（其句柄由 `acl` 参数指定）中 `ACL_MASK` 型记录的权限。通常，只要是修改或创建 ACL，就会用到该函数。其会对所有 `ACL_USER`、`ACL_GROUP` 以及 `ACL_GROUP_OBJ` 型记录的权限并集进行计算，作为 `ACL_MASK` 型记录的权限。若 `ACL_MASK` 型记录不存在，则该函数会创建一个，这也算是该函数的妙用之一。也就是说，在将 `ACL_USER` 和 `ACL_GROUP` 型记录添加到前面提及的“最小化”ACL 时，调用该函数就能确保 `ACL_MASK` 型记录的创建。

若参数 `acl` 所指定的 ACL 有效，`acl_valid(acl)`函数将返回 0，否则，返回-1。若以下所有条件成立(为真)，则可判定该 ACL 有效。

- `ACL_USER_OBJ`、`ACL_GROUP_OBJ` 以及 `ACL_OTHER` 类型的记录均只能有一条。
- 若有任一 `ACL_USER` 或 `ACL_GROUP` 类型的记录存在，则也必然存在一条 `ACL_MASK` 型记录。
- 标记类型为 `ACL_MASK` 的 ACE 至多只有一条。
- 每条标记类型为 `ACL_USER` 的记录都有一唯一的用户 ID。
- 每条标记类型为 `ACL_GROUP` 的记录都有一唯一的组 ID。

`acl_check()`和 `acl_error()`函数(后者为 Linux 的扩展)与 `acl_valid()`函数有异曲同工之妙，尽管可移植性不强，但在处理畸形 ACL 时却能对错误提供更为精确的描述。欲知详情，请参考手册页。

`acl_delete_def_file(pathname)`函数用来删除目录（由参数 `pathname` 指定）的默认型 ACL。

`acl_init(count)`函数用来新建一个空的 ACL 结构，其空间足以容纳由参数 `count` 所指定的记录数。(参数 `count` 向系统传递的是编程者的柔性诉求，而非硬性要求。)函数将返回这一新建 ACL 的句柄。

`acl_dup(acl)`函数用来为由 `acl` 参数所指定的 ACL 创建副本，并以该 ACL 副本的句柄作为返回值。

`acl_free(handle)`函数用来释放由其他 ACL 函数所分配的内存。例如，必须使用该函数来释放由 `acl_from_text()`、`acl_to_text()`、`acl_get_file()`、`acl_init()`以及 `acl_dup()`调用所分配的内存。

程序示例

程序清单 17-1 对某些 ACL 库函数的使用做了演示。该程序可获取并展示与文件相关的 ACL（亦即，该程序提供了 `getfacl` 命令的部分功能）。若以 `-d` 命令行选项执行该程序，则将显示与目录相关的默认型 ACL，而非访问型 ACL。

以下为该程序的运行示例。

```

$ touch tfile
$ setfacl -m 'u:annie:r,u:paulh:rw,g:teach:r' tfile
$ ./acl_view tfile
user_obj          rw-
user              annie  r--
user              paulh  rw-
group_obj         r--
group             teach  r--
mask              rw-
other             r--

```

随本书发行的源码中还包含了另一程序：`acl/acl_update.c`，可用来更新 ACL（该程序提供了 `setfacl` 命令的部分功能）。

程序清单 17-1：显示与文件挂钩的访问或默认 ACL

```

acl/acl\_view.c
#include <acl/libacl.h>
#include <sys/acl.h>
#include "ugid_functions.h"
#include "tlpi_hdr.h"

static void
usageError(char *progName)
{
    fprintf(stderr, "Usage: %s [-d] filename\n", progName);
    exit(EXIT_FAILURE);
}

int
main(int argc, char *argv[])
{
    acl_t acl;
    acl_type_t type;
    acl_entry_t entry;
    acl_tag_t tag;
    uid_t *uidp;
    gid_t *gidp;
    acl_permset_t permset;
    char *name;
    int entryId, permVal, opt;

    type = ACL_TYPE_ACCESS;
    while ((opt = getopt(argc, argv, "d")) != -1) {
        switch (opt) {
            case 'd': type = ACL_TYPE_DEFAULT;        break;
            case '?': usageError(argv[0]);
        }
    }

    if (optind + 1 != argc)
        usageError(argv[0]);

    acl = acl_get_file(argv[optind], type);
    if (acl == NULL)
        errExit("acl_get_file");

    /* Walk through each entry in this ACL */

```

```

    for (entryId = ACL_FIRST_ENTRY; ; entryId = ACL_NEXT_ENTRY) {

        if (acl_get_entry(acl, entryId, &entry) != 1)
            break; /* Exit on error or no more entries */
/* Retrieve and display tag type */

if (acl_get_tag_type(entry, &tag) == -1)
    errExit("acl_get_tag_type");

printf("%-12s", (tag == ACL_USER_OBJ) ? "user_obj" :
        (tag == ACL_USER) ? "user" :
        (tag == ACL_GROUP_OBJ) ? "group_obj" :
        (tag == ACL_GROUP) ? "group" :
        (tag == ACL_MASK) ? "mask" :
        (tag == ACL_OTHER) ? "other" : "???");

/* Retrieve and display optional tag qualifier */

if (tag == ACL_USER) {
    uidp = acl_get_qualifier(entry);
    if (uidp == NULL)
        errExit("acl_get_qualifier");

    name = groupNameFromId(*uidp);
    if (name == NULL)
        printf("%-8d ", *uidp);
    else
        printf("%-8s ", name);

    if (acl_free(uidp) == -1)
        errExit("acl_free");
} else if (tag == ACL_GROUP) {
    gidp = acl_get_qualifier(entry);
    if (gidp == NULL)
        errExit("acl_get_qualifier");

    name = groupNameFromId(*gidp);
    if (name == NULL)
        printf("%-8d ", *gidp);
    else
        printf("%-8s ", name);

    if (acl_free(gidp) == -1)
        errExit("acl_free");
} else {
    printf(" ");
}

/* Retrieve and display permissions */

if (acl_get_permset(entry, &permset) == -1)
    errExit("acl_get_permset");

permVal = acl_get_perm(permset, ACL_READ);
if (permVal == -1)
    errExit("acl_get_perm - ACL_READ");

```



```

printf("%c", (permVal == 1) ? 'r' : '-');
permVal = acl_get_perm(permset, ACL_WRITE);
if (permVal == -1)
    errExit("acl_get_perm - ACL_WRITE");
printf("%c", (permVal == 1) ? 'w' : '-');
permVal = acl_get_perm(permset, ACL_EXECUTE);
if (permVal == -1)
    errExit("acl_get_perm - ACL_EXECUTE");
printf("%c", (permVal == 1) ? 'x' : '-');

printf("\n");
}

if (acl_free(acl) == -1)
    errExit("acl_free");

exit(EXIT_SUCCESS);
}

```

acl/acl_view.c

17.9 总结

自 2.6 版本起，Linux 开始支持 ACL。ACL 是对传统 UNIX 文件权限模型的扩展，籍此可在每用户或每组的基础上来控制对文件的访问。

进阶信息

访问 <http://wt.tuxomania.net/publications/posix.1e/>，可在线查看 POSIX.1e 和 POSIX.2c 标准草案的最后一稿（第 17 号草案）。

acl(5)手册页简要介绍了 ACL，并针对 Linux 平台所实现的各种 ACL 库函数，就其可移植性给出了指导。

ACL 及扩展属性的 Linux 实现细节刊载于[Grünbacher, 2003]。Andreas Grünbacher 所维护的 Web 站点也包含了与 ACL 有关的信息，链接为 <http://acl.bestbits.at/>。

17.10 练习

- 17-1.** 编写一个程序，根据与一特定用户或组相对应的 ACE 来显示权限。该程序应接受 2 个命令行参数。第一个参数可以为字母“u”或“g”，用以表明第二个参数是用户还是组。（利用定义于程序清单 8-1 中的函数，还可将第二个参数任意指定为数字或名称。）若与给定用户或组相对应的 ACE 隶属于组分类，则程序还需另外显示与 ACL 掩码型记录相与后的权限。

第 18 章

目录与链接

作为文件相关议题的结局篇，本章将讨论目录和链接。首先是对其系统级实现进行了回顾，之后则描述了用于创建和移除目录、链接的系统调用。接下来所探讨的库函数，可允许程序扫描单个目录下的内容并遍历一个目录树（即，检查目录树中的每个文件）。

每个进程都有两个目录相关属性根目录及当前工作目录，分别用于为解释绝对路径名和相对路径名提供参照点。本章将描述修改二者的系统调用。

最后，本章讨论了相关库函数，可用来解析路径名，并将其分解为目录和文件名两部分。

18.1 目录和（硬）链接

在文件系统中，目录的存储方式类似于普通文件。目录与普通文件的区别有二。

- 在其 `i-node` 条目中，会将目录标记为一种不同的文件类型（参见 14.4 节）。
- 目录是经特殊组织而成的文件。本质上说就是一个表格，包含文件名和 `i-node` 编号。

在大多数原生 Linux 文件系统中，文件名长度可达 255 个字符。图 18-1 所示为针对示例文件（`/etc/passwd`）所维护的文件系统 `i-node` 表以及相关目录文件的部分内容，展示了目录与 `i-node` 之间的关系。

虽然一个进程能够打开一个目录，但却不能使用 `read()` 去读取目录的内容。为了检索目录内容，进程必须使用本章后续讨论的系统调用和库函数。（在一些 UNIX 实现中，也可以对目录执行 `read()`，但这会给应用带来可移植性方面的问题。）进程同样也不能使用 `write()` 来改变一个目录的内容，仅能借助于诸如 `open()`（创建一个新文件）、`link()`、`mkdir()`、`symlink()`、`unlink()` 及 `rmdir()` 之类的系统调用（4.3 节描述了 `open()` 调用，本章稍后会介绍余下的系统调用）来间接（向内核请求）改变其内容。

`i-node` 表的编号始于 1，而非 0，因为若目录条目的 `i-node` 字段值为 0，则表明该条目尚未使用。`i-node 1` 用来记录文件系统的坏块。文件系统根目录（`/`）总是存储在 `i-node` 条目 2 中（如图 18-1 所示），所以内核在解析路径名时就知道该从哪里着手。

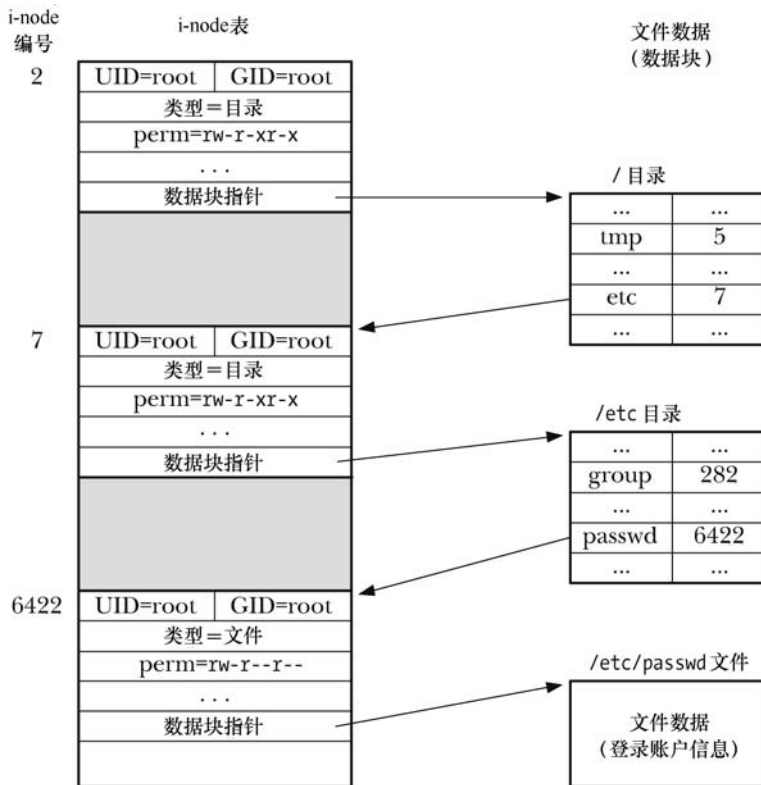


图 18-1: 以文件/etc/passwd 为例, 展示 i-node 和目录结构之间的关系

回顾文件 i-node (14.4 节) 中存储的信息列表, 会发现其中并未包含文件名, 而仅通过目录列表内的一个映射来定义文件名称。其妙用在于, 能够在相同或者不同目录中创建多个名称, 每个均指向相同的 i-node 节点。也将这些名称称为链接, 有时也称之为硬链接 (稍后介绍), 以示与符号链接有所区别。

所有的原生 Linux 和 UNIX 文件系统均支持硬链接, 然而, 许多非 UNIX 文件系统 (比如, 微软的 VFAT) 则不支持。(微软的 NTFS 文件系统支持硬链接。)

可在 shell 中利用 ln 命令为一个业已存在的文件创建新的硬链接, 正如下面 shell 会话日志所示。

```

$ echo -n 'It is good to collect things,' > abc
$ ls -li abc
122232 -rw-r--r--  1 mtk      users          29 Jun 15 17:07 abc
$ ln abc xyz
$ echo ' but it is better to go on walks.' >> xyz
$ cat abc
It is good to collect things, but it is better to go on walks.
$ ls -li abc xyz
122232 -rw-r--r--  2 mtk      users          63 Jun 15 17:07 abc
122232 -rw-r--r--  2 mtk      users          63 Jun 15 17:07 xyz

```

Cat 命令输出中一目了然的事, 经过 ls-li 命令所示 i-node 编码 (即第一列) 得到了进一步证实。名称 abc 和 xyz 指向相同的 i-node 条目, 因此均指向相同文件。ls-li 命令所示内容的第三列为对 i-node 链接的计数。执行 ln abc xyz 命令后, abc 所指向 i-node 的链接计数升至 2,

因为现在指向该文件的有两个名字。（由于指向相同的 i-node，针对文件 xyz 输出的链接计数也是 2。）

若移除其中一个文件名，另一文件名以及文件本身将继续存在。

```
$ rm abc
$ ls -li xyz
122232 -rw-r--r-- 1 mtk users 63 Jun 15 17:07 xyz
```

仅当 i-node 的链接计数降为 0 时，也就是移除了文件的所有名字时，才会删除（释放）文件的 i-node 记录和数据块。总结如下：`rm` 命令从目录列表中删除一文件名，将相应 i-node 的链接计数减一，若链接计数因此而降为 0，则还将释放该文件名所指代的 i-node 和数据块。

同一文件的所有名字（链接）地位平等——没有一个名字（比如，第一个）会优于其他名字。正如上例所示，在移除与文件相关的第一个名称后，物理文件继续存在，但只能通过另一文件名来访问其内容。

在线论坛上经常会有这样的问题出现：在程序中如何找到与文件描述符 X 相关联的文件名？简单的回答是不能，至少缺乏明确而又便于移植的手段，因为一个文件描述符指向一个 i-node，而指向这个 i-node 的文件名则可能有多个（或者甚至如 18.3 节所述，一个都没有）。

在 Linux 系统上，借助于 `readdir()` 对 Linux 特有 `/proc/PID/fd` 目录内容（内含符号链接指向进程当前打开的每个文件描述符）的扫描，可以获知一个进程当前打开了哪些文件。此外，已经移植到多个 UNIX 系统中的 `lsdf(1)` 和 `fuser(1)` 工具也精于此道。

对硬链接的限制有二，均可用符号链接来加以规避。

- 因为目录条目（硬链接）对文件的指代采用了 i-node 编号，而 i-node 编号的唯一性仅在一个文件系统之内才能得到保障，所以硬链接必须与其指代的文件驻留在同一文件系统中。
- 不能为目录创建硬链接，从而避免出现令诸多系统程序陷于混乱的链接环路。

早期的 UNIX 实现一度曾允许超级用户为目录创建硬链接。这在当时是必要的，因为这些实现并未提供 `mkdir()` 系统调用。相反，当时会使用 `mknod()` 调用创建一个目录，然后为 `.` 和 `..` 创建链接（[Vahalia, 1996]）。虽然这一特性已是昨日黄花，但一些现代 UNIX 实现出于向后兼容的目的仍对其加以保留。

使用绑定挂载（`bind mount`）可以获得与为目录创建硬链接相似的效果。

18.2 符号（软）链接

符号链接，有时也称为软链接，是一种特殊的文件类型，其数据是另一文件的名称。图 18-2 展示的情况是：两个硬链接——`/home/arena/this` 和 `/home/allyn/that`——指向同一个文件，而符号链接 `/home/kiran/other`，则指代文件名 `/home/arena/this`。

在 shell 中，符号链接是由 `ln -s` 命令创建的。`ls -F` 命令的输出结果中会在符号链接的尾部标记 `@`。

符号链接的内容既可以是绝对路径，也可以是相对路径。解释相对符号链接时以链接本身的位置作为参照点。

符号链接的地位不如硬链接。尤其是，文件的链接计数中并未将符号链接计算在内。（因此，图 18-2 中编号为 61 的 i-node，其链接计数为 2，而不是 3。）因此，如果移除了符号链接所指向

的文件名，符号链接本身还将继续存在，尽管无法再对其进行解引用（下溯）操作，也将此类链接称之为悬空链接。更有甚者，还可以为并不存在的文件名创建一个符号链接。

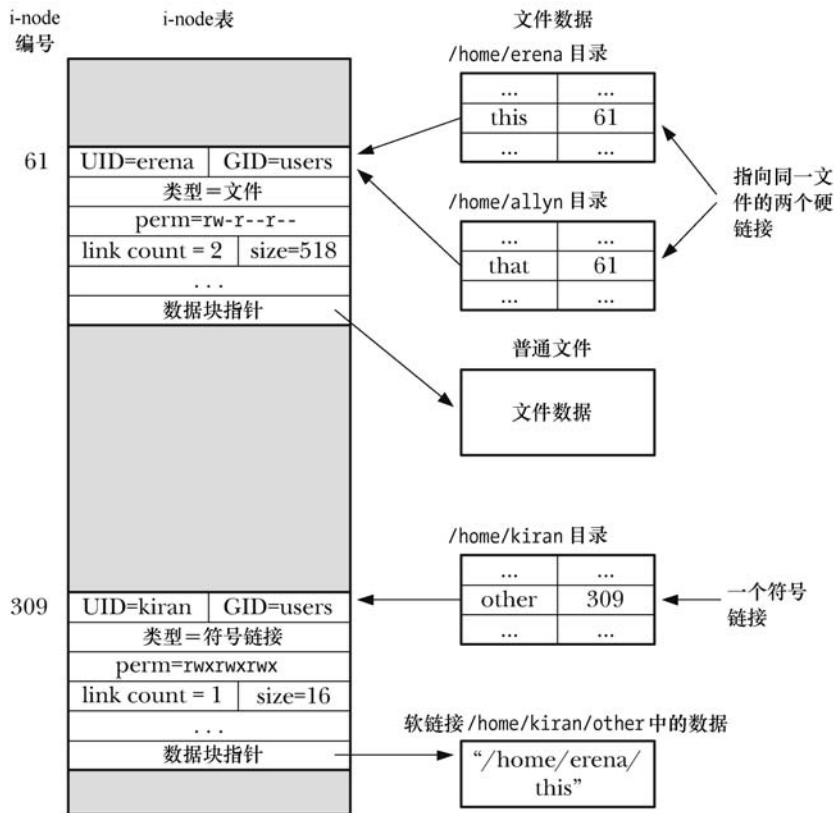


图 18-2: 对硬链接和符号链接的展现

引入符号链接的是 4.2BSD。虽然未获 POSIX.1-1990 接纳，但 SUSv1 和 SUSv3 随后还是将其纳入规范。

因为符号链接指代一个文件名，而非 i-node 编号，所以可以用其来链接不同文件系统中的文件。对硬链接的那些制约也就不会困扰到符号链接，可以为目录创建符号链接。诸如 `find` 和 `tar` 之类的工具命令有能力识别硬链接和符号链接之间的差异，要么会在默认情况下不对符号链接进行解引用，要么会避免因使用符号链接而陷入引用环路。

符号链接之间可能会形成链路（例如，`a` 是指向 `b` 的符号链接，而 `b` 是指向 `c` 的符号链接）。当在各个文件相关的系统调用中指定了符号链接时，内核会对一系列链接层层解去引用，直抵最终文件。

SUSv3 规定，针对路径名中的每个符号链接部件，系统实现应允许对其实施至少 `_POSIX_SYMLINK_MAX` 次解除引用操作。`_POSIX_SYMLINK_MAX` 的规定值为 8。然而，在内核 2.6.18 之前，Linux 将解析符号链接链路时的解引用操作次数限制为 5 次。始于版本 2.6.18，Linux 内核实现了 SUSv3 所规定的最小解引用次数：8 次。Linux 还将对一个完整路径名的解引用总数限制为 40 次。施加这些限制，意在应对超长符号链接链路以及符号链接环路，从而使内核代码在解析符号链接时免于引发堆栈溢出。

某些 UNIX 文件系统的优化举措，不但没有在正文中提及，而且也未见于图 18-2。如果构成符号链接内容的字符串总长度很小，足以放入 i-node 中通常用于存放数据指针的位置，那么就会将字符串直接存储在那里。这省去对一个磁盘块的分配，也加速了对符号链接信息的访问，因为获取信息时仅涉及到文件的 i-node。例如，ext2、ext3 及 ext4 采用这一技术，将 i-node 中通常用于存放数据块指针的 60 个字节转而用于存放长度合适的符号字符串。实践证明，这一优化措施卓有成效。在笔者检查过的某一系统中，符号链接共计 20 700 个，其中内容长度不超过 60 个字节的占 97%。

系统调用对符号链接的解释

许多系统调用都会对符号链接进行解引用处理（即下溯 follow），从而对链接所指向的文件展开操作。还有一些系统调用对符号链接则不作处理，直接操作于链接文件本身。书中会在论及每个系统调用的同时，描述其针对符号链接的行为。表 18-1 对此作了总结。

表 18-1: 各个函数对符号链接的解释

函 数	是否对链接解引用	备 注	
access()	•	UNIX 域套接字带有路径名	
acct()	•		
bind()	•		
chdir()	•		
chmod()	•		
chown()	•		
chroot()	•		
creat()	•		
exec()	•		
getxattr()	•		
lchown()			
lgetxattr()			
link()			参见 18.3 节
listxattr()	•		
lstat()			
lremovexattr()			
lsetxattr()			
lstat()			
lutimes()			
open()	•	除非指定了 O_NOFOLLOW 或者 O_EXCL O_CREAT	
opendir()	•		
pathconf()	•		

续表

函 数	是否对链接解引用	备 注
pivot_root()	•	无论哪个参数中的链接，都不会对其进行解引用 若参数为符号链接，则调用失败，并将 <code>errno</code> 置为 <code>ENOTDIR</code>
quotactl()	•	
readlink()		
removexattr()	•	
rename()		
rmdir()		
setxattr()	•	
stat()	•	
statfs(), statvfs()	•	
swapon(), swapoff()	•	
truncate()	•	
unlink()		
uselib()	•	
utime(), utimes()	•	

少数情况下，符号链接本身及其所指向的文件会需要类似的功能，系统这时就会提供两套系统调用：一套会对链接解除引用，另一套则反之，后者在命名时会冠以字母 `l`。 `stat()` 和 `lstat()` 就是一例。

有一点是约定俗成的：总是会对路径名中目录部分（即最后一个斜线字符前的所有组成部分）的符号链接进行解除引用操作。因此，在路径 `/somedir/somesubdir/file` 中，若 `somedir` 和 `somesubdir` 属于符号链接，则一定会解除对这两个目录的引用，而针对 `file` 是否进行解引用与否，则取决于路径名所传入的系统调用。

18.11 节描述了一组自版本 2.6.16 加入的系统调用，对表 18-1 所展示的部分接口功能有所扩展。对于其中的某些调用而言，可以利用 `flags` 参数来控制是否对符号链接进行解引用操作。

符号链接的文件权限和所有权

大部分操作会无视符号链接的所有权和权限（创建符号链接时会为其赋予所有权限）。是否允许操作反而是由符号链接所指代文件的所有权和权限来决定。仅当在带有粘性权限位（15.4.5 节）的目录中对符号链接进行移除或改名操作时，才会考虑符号链接自身的所有权。

18.3 创建和移除（硬）链接： `link()` 和 `unlink()`

`link()` 和 `unlink()` 系统调用分别创建和移除硬链接。

```
#include <unistd.h>

int link(const char *oldpath, const char *newpath);

Returns 0 on success, or -1 on error
```

若 `oldpath` 中提供的是一个已存在文件的路径名，则系统调用 `link()` 将以 `newpath` 参数所指定的路径名创建一个新链接。若 `newpath` 指定的路径名已然存在，则不会将其覆盖；相反，将产生一个错误（EEXIST）。

在 Linux 中，`link()` 系统调用不会对符号链接进行解引用操作。若 `oldpath` 属于符号链接，则会将 `newpath` 创建为指向相同符号链接文件的全新硬链接。（换言之，`newpath` 也是符号链接，指向 `oldpath` 所指向的同一文件。）这一行为有悖于 SUSv3 规范。SUSv3 要求，除非另行规定（`link()` 系统调用不在此列），否则所有执行路径名解析操作的函数都应对符号链接进行解引用。大多数其他 UNIX 实现的行事方式都与 SUSv3 相符。值得注意的是，Solaris 是个例外，默认情况下的行为与 Linux 相同。但若采用适当的编译器选项，又可提供符合 SUSv3 规范的行为。鉴于系统实现间的这种差异，应避免将 `oldpath` 参数指定为符号链接，以保障程序的可移植性。

SUSv4 承认现有实现间存在不一致性，同时规定 `link()` 调用对符号链接解引用与否由实现定义。SUSv4 还将 `linkat()` 纳入规范，在执行与 `link()` 相同任务的同时，可利用 `flag` 参数来控制调用是否解析符号链接。更多细节参见 18.11 节。

```
#include <unistd.h>

int unlink(const char *pathname);

Returns 0 on success, or -1 on error
```

`unlink()` 系统调用移除一个链接（删除一个文件名），且如果此链接是指向文件的最后一个链接，那么还将移除文件本身。若 `pathname` 中指定的链接不存在，则 `unlink()` 调用失败，并将 `errno` 置为 ENOENT。

`unlink()` 不能移除一个目录，完成这一任务需要使用 `rmdir()` 或 `remove()`，将于 18.6 节进行介绍。

SUSv3 规定，若 `pathname` 中指定的是一个目录，则 `unlink()` 调用失败，并将 `errno` 置为 EPERM。然而，在 Linux 中，`unlink()` 在这种情况下会将 `errno` 置为 EISDIR 值。（对于与 SUSv3 间的这一差别，LSB 倒也并不讳言。）为保障可移植性，应用程序在检查这种情况时应做两手准备。

`unlink()` 系统调用不会对符号链接进行解引用操作，若 `pathname` 为符号链接，则移除链接本身，而非链接指向的名称。

仅当关闭所有文件描述符时，方可删除一个已打开的文件

内核除了为每个 i-node 维护链接计数之外，还对文件的打开文件描述（参见图 5-2）计数。

当移除指向文件的最后一个链接时，如果仍有进程持有指代该文件的打开文件描述符，那么在关闭所有此类描述符之前，系统实际上将不会删除该文件。这一特性的妙用在于允许取消对文件的链接，而无需担心是否有其他进程已将其打开。（然而，对于链接数已降为 0 的打开文件，就无法将文件名与其重新关联起来。）此外，基于上述事实，还可以玩点小技巧：先创建并打开一个临时文件，随即取消对文件的链接（`unlink`），然后在程序中继续使用该文件。（这正是 5.12 节所述 `tmpfile()` 函数的所作所为。）

程序清单 18-1 对此现象做了展示。

程序清单 18-1：使用 `unlink()` 移除一个链接

```
----- dirs_links/t_unlink.c
#include <sys/stat.h>
#include <fcntl.h>
#include "t1pi_hdr.h"

#define CMD_SIZE 200
#define BUF_SIZE 1024

int
main(int argc, char *argv[])
{
    int fd, j, numBlocks;
    char shellCmd[CMD_SIZE];          /* Command to be passed to system() */
    char buf[BUF_SIZE];              /* Random bytes to write to file */

    if (argc < 2 || strcmp(argv[1], "--help") == 0)
        usageErr("%s temp-file [num-1kB-blocks] \n", argv[0]);

    numBlocks = (argc > 2) ? getInt(argv[2], GN_GT_0, "num-1kB-blocks")
        : 100000;

    fd = open(argv[1], O_WRONLY | O_CREAT | O_EXCL, S_IRUSR | S_IWUSR);
    if (fd == -1)
        errExit("open");

    if (unlink(argv[1]) == -1)        /* Remove filename */
        errExit("unlink");

    for (j = 0; j < numBlocks; j++)  /* Write lots of junk to file */
        if (write(fd, buf, BUF_SIZE) != BUF_SIZE)
            fatal("partial/failed write");

    snprintf(shellCmd, CMD_SIZE, "df -k `dirname %s`", argv[1]);
    system(shellCmd);                /* View space used in file system */

    if (close(fd) == -1)             /* File is now destroyed */
        errExit("close");
    printf("***** Closed file descriptor\n");

    system(shellCmd);                /* Review space used in file system */
    exit(EXIT_SUCCESS);
}
----- dirs_links/t_unlink.c
```

程序清单 18-1 中程序接受两个命令行参数。第一个参数标识程序应该创建的文件名称。

程序打开此文件后随即取消与文件名的链接。虽然文件名已消失，但是文件本身依然存在。然后程序向文件随机写入一些数据块，数据块数量由程序的第二个命令行参数（可选项）指定。这时，程序会利用 `df(1)` 命令显示文件系统的空间使用情况。程序接着会关闭文件描述符，系统因之而将文件移除，程序会再次使用 `df(1)` 命令来显示有所下降的磁盘使用情况。如下 shell 会话演示了运行程序清单 18-1 程序的情况：

```
$ ./t_unlink /tmp/tfile 1000000
Filesystem      1K-blocks      Used Available Use% Mounted on
/dev/sda10      5245020      3204044   2040976  62% /
***** Closed file descriptor
Filesystem      1K-blocks      Used Available Use% Mounted on
/dev/sda10      5245020      2201128   3043892  42% /
```

程序清单 18-1 使用 `system()` 函数来执行 shell 命令，27.6 节将对此函数做详细描述。

18.4 更改文件名：rename()

借助于 `rename()` 系统调用，既可以重命名文件，又可以将文件移至同一文件系统另一目录。

```
#include <stdio.h>

int rename(const char *oldpath, const char *newpath);

Returns 0 on success, or -1 on error
```

调用会将现有的一个路径名 `oldpath` 重命名为 `newpath` 参数所指定的路径名。

`rename()` 调用仅操作目录条目，而不移动文件数据。改名既不影响指向该文件的其他硬链接，也不影响持有该文件打开描述符的任何进程，因为这些文件描述符指向的是打开文件描述，（在调用 `open()` 之后）与文件名并无瓜葛。

以下规则适用与对 `rename()` 的调用。

- 若 `newpath` 已经存在，则将其覆盖。
- 若 `newpath` 与 `oldpath` 指向同一文件，则不发生变化（且调用成功）。这很不合常理。顺着上一条规则的思路，通常的推断是：如果两个文件名 `x` 和 `y` 都存在，那么调用 `rename("x","y")` 时应当把 `x` 移除才是。但如果 `x` 和 `y` 链接的是同一文件，事实却并非如此。

此规则源于早期的 BSD 实现，其动机可能是出于这样的考虑：要保障诸如 `rename("x", "x")`、`rename("x", "./x")` 以及 `rename("x", "somedir/./x")` 之类的调用不会移除文件，内核必须进行检查。而这种设计则有助于简化这一检查。

- `rename()` 系统调用对其两个参数中的符号链接均不进行解引用。如果 `oldpath` 是一符号链接，那么将重命名该符号链接。如果 `newpath` 是一符号链接，那么会将其视为由 `oldpath` 重命名而成的普通路径名（即移除已有的符号链接 `newpath`）。
- 如果 `oldpath` 指代文件，而非目录，那么就不能将 `newpath` 指定为一个目录的路径名（否则将 `errno` 置为 `EISDIR`）。要想重命名一个文件到某一目录中（亦即将文件移到另一目录），`newpath` 必须包含新的文件名。如下调用既将一个文件移动到另一目录中，同时又将其改名：

```
rename("sub1/x", "sub2/y");
```

- 若将 `oldpath` 指定为目录名，则意在重命名该目录。这种情况下，必须保证 `newpath` 要么不存在，要么是一个空目录的名称。无论 `newpath` 是一个已有文件还是一个非空目录，调用都将出错（分别将 `errno` 置为 `ENOTDIR` 和 `ENOTEMPTY`）。
- 若 `oldpath` 是一目录，则 `newpath` 不能包含 `oldpath` 作为其目录前缀。例如，不能将 `/home/mtk` 重命名为 `/home/mtk/bin`（错误为 `EINVAL`）。
- `oldpath` 和 `newpath` 所指代的文件必须位于同一文件系统。之所以如此，是因为目录内容由硬链接列表组成，且硬链接所指向的 `i-node` 与目录位于同一文件系统。如前所述，`rename()` 仅限于操作目录列表的内容。试图将一文件重命名至不同的文件系统将返回错误 `EXDEV`。（非要如此，必须从一个文件系统中将其文件内容复制到另一文件系统，然后再删除老文件。这正是 `mv` 命令的功能所在。）

18.5 使用符号链接：symlink()和 readlink()

现在来看看用于创建符号链接，以及检查其内容的系统调用。

`symlink()` 系统调用会针对由 `filepath` 所指定的路径名创建一个新的符号链接——`linkpath`。（想移除符号链接，需使用 `unlink()` 调用。）

```
#include <unistd.h>

int symlink(const char *filepath, const char *linkpath);

Returns 0 on success, or -1 on error
```

若 `linkpath` 中给定的路径名已然存在，则调用失败（且将 `errno` 置为 `EEXIST`）。由 `filepath` 指定的路径名可以是绝对路径，也可以是相对路径。

由 `filepath` 所命名的文件或目录在调用时无需存在。即便当时存在，也无法阻止后来将其删除。这时，`linkpath` 成为“悬空链接”，其他系统调用试图对其进行解引用操作都将出错（通常错误号为 `ENOENT`）。

如果指定一符号链接作为 `open()` 调用的 `pathname` 参数，那么将打开链接指向的文件。有时，倒宁愿获取链接本身的内容，即其所指向的路径名。这正是 `readlink()` 系统调用的本职工作，将符号链接字符串的一份副本置于 `buffer` 指向的字符数组中。

```
#include <unistd.h>

ssize_t readlink(const char *pathname, char *buffer, size_t bufsiz);

Returns number of bytes placed in buffer on success, or -1 on error
```

`bufsiz` 是一个整型参数，用以告知 `readlink()` 调用 `buffer` 中的可用字节数。

如果一切顺利，`readlink()` 将返回实际放入 `buffer` 中的字节数。若链接长度超过 `bufsiz`，则置于 `buffer` 中的是经截断处理的字符串（并返回字符串大小，亦即 `bufsiz`）。

由于 `buffer` 尾部并未放置终止空字符，故而也无法分辨 `readlink()` 所返回的字符串到底是经过截断处理，还是恰巧将 `buffer` 填满。验证后者的方法之一是重新分配一块更大的 `buffer`，并再次调用 `readlink()`。另外，还可以将 `pathname` 的长度定义为常量 `PATH_MAX`（参见 11.1

节)，该常量定义了程序可拥有的最长路径名长度。

程序清单 18-4 演示了 `readlink()` 的用法。

为限制符号链接中所能存储的最大字节数，SUSv3 定义了一个新限制 `SYMLINK_MAX`，要求系统实现对此加以定义，并规定其不得少于 255 个字节。写作本书时，Linux 尚未对此限制作出定义。而正文之所以建议使用 `PATH_MAX`，是因为该限制至少与 `SYMLINK_MAX` 大小相当。

SUSv2 将 `readlink()` 的返回类型定义为 `int`，而当前的许多实现（以及 Linux 中较老版本的 `glibc`）对此奉行不悖。SUSv3 则将 `readlink()` 的返回值类型改为 `ssize_t`。

18.6 创建和移除目录：`mkdir()`和`rmdir()`

`mkdir()` 系统调用创建一个新目录。

```
#include <sys/stat.h>

int mkdir(const char *pathname, mode_t mode);

Returns 0 on success, or -1 on error
```

`pathname` 参数指定了新目录的路径名。该路径名可以是相对路径，也可以是绝对路径。若具有该路径名的文件已经存在，则调用失败并将 `errno` 置为 `EEXIST`。

对新目录所有权的设置遵循 15.3.1 节所述规则。

`mode` 参数指定了新目录的权限。（15.3.1、15.3.2 和 15.4.5 各节描述了目录权限位的含义。）对该位掩码值的指定方式既可以与 `open()` 调用相同——对表 15-4 所列各常量进行或(`|`)操作，也可直接赋予八进制数值。既定的 `mode` 值还将于进程掩码相与 (`&`)（参见 15.4.6 节）。另外，`set-user-ID` 位始终处于关闭状态，因为该位对于目录而言毫无意义。

如果在 `mode` 中设置了粘滞位 (`S_ISVTX`)，那么将对新目录设置该权限。

调用还会忽略在 `mode` 中设置的 `set-group-ID` 位 (`S_ISGID`)。相反，若对其父目录设置了 `set-group-ID` 位，则同样会对新建目录设置该权限。15.3.1 节曾指出，对目录设置 `set-group-ID` 权限位将导致目录中新建文件的组 ID 取自目录组 ID，而非进程有效组 ID。`mkdir()` 系统调用按此处描述的方式来传播 `set-group-ID` 权限位，以保证目录下所有子目录的行为均保持一致。

SUSv3 明文规定，`mkdir()` 对 `set-user-ID`、`set-group-ID` 以及粘滞位的处理方式由实现定义。某些 UNIX 实现在新建目录时总是关闭这 3 个权限位。

新建目录包括两个条目：`.`（点），即指向目录自身的链接；`..`（点点），即指向父目录的链接。

SUSv3 并未要求目录中包括 `.` 和 `..` 条目，只是要求当路径中出现 `.` 和 `..` 时，实现应能正确解释。若要保证应用程序的可移植性，则不应假设目录中存在这些条目。

`mkdir()` 系统调用所创建的仅仅是路径名中的最后一部分。换言之，`mkdir("aaa/bbb/ccc", mode)` 仅当目录 `aaa` 和 `aaa/bbb` 已经存在的情况下才会成功。（这相当于 `mkdir(1)` 命令的默认行为，但 `mkdir(1)` 同时也提供 `-p` 选项，可将不存在的中间目录一一创建。）

GNU C 语言库提供有 `mkdtemp(template)` 函数，可谓 `mkstemp()` 函数的“目录”版。该函数会创建一个唯一命名的目录，在对所有者开启读、写和执行权限的同时，不对任何其他用户赋予任何权限。不过，`mkdtemp()` 所返回的并非一个文件描述符，而是一枚指针，指向经过修改处理的字符串，内含 `template` 中的实际目录名。该函数并未纳入 SUSv3 规范，也未获得所有 UNIX 实现的支持，但 SUSv4 对其作了定义。

`rmdir()` 系统调用移除由 `pathname` 指定的目录，该目录可以是绝对路径名，也可以是相对路径名。

```
#include <unistd.h>

int rmdir(const char *pathname);

Returns 0 on success, or -1 on error
```

要使 `rmdir()` 调用成功，则要删除的目录必须为空。如果 `pathname` 的最后一部分为符号链接，那么 `rmdir()` 调用将不对其进行解引用操作，并返回错误，同时将 `errno` 置为 `ENOTDIR`。

18.7 移除一个文件或目录：remove()

`remove()` 库函数移除一个文件或一个空目录。

```
#include <stdio.h>

int remove(const char *pathname);

Returns 0 on success, or -1 on error
```

如果 `pathname` 是一文件，那么 `remove()` 去调用 `unlink()`；如果 `pathname` 为一目录，那么 `remove()` 去调用 `rmdir()`。

与 `unlink()`、`rmdir()` 一样，`remove()` 不对符号链接进行解引用操作。若 `pathname` 是一符号链接，则 `remove()` 会移除链接本身，而非链接所指向的文件。

如果移除一个文件只是为创建同名新文件做准备，那么编码时使用 `remove()` 函数会更加简单，无需再去检查目录名所指是文件还是目录，然后再决定是调用 `unlink()` 还是 `rmdir()`。

`remove()` 属于 C 语言标准库函数，广为 UNIX 和非 UNIX 系统所支持。大多数非 UNIX 系统并不支持硬链接，所以将移除文件的函数命名为 `unlink()` 也不合情理。

18.8 读目录：opendir()和 readdir()

本节所述库函数可用于打开一个目录，并逐一获取其包含文件的名称。

读取目录的库函数均以 `getdents()` 系统调用（未纳入 SUSv3 规范）为基础，但其接口更易于使用。Linux 还提供了 `readdir(2)` 系统调用（相对于此处描述的 `readdir(3)` 库函数），所执行的任务类似于 `getdents()`，也因之而遭废止。

`opendir()`函数打开一个目录，并返回指向该目录的句柄，供后续调用使用。

```
#include <dirent.h>

DIR *opendir(const char *dirpath);

Returns directory stream handle, or NULL on error
```

`opendir()`函数打开由 `dirpath` 指定的目录，并返回指向 `DIR` 类型结构的指针。该结构即所谓目录流（directory stream），亦即调用者传递给下述其他函数的句柄。一旦从 `opendir()`返回，则将目录流指向目录列表的首条记录。

除了要创建的目录流所针对的目录由打开文件描述符指代之外，`fdopendir()`与 `opendir()`并无不同。

```
#include <dirent.h>

DIR *fdopendir(int fd);

Returns directory stream handle, or NULL on error
```

提供 `fdopendir()`函数，意在帮助应用程序免受 18.11 节所述各种竞态条件的困扰。

调用 `fdopendir()`成功后，文件描述符将处于系统的控制之下，且除了利用本节余下部分所描述的函数之外，程序不应采取任何其他方式对其进行访问。

SUSv4 定义了 `fdopendir()`函数（但 SUSv3 并未将其纳入规范）。

`readdir()`函数从一个目录流中读取连续的条目。

```
#include <dirent.h>

struct dirent *readdir(DIR *dirp);

Returns pointer to a statically allocated structure describing
next directory entry, or NULL on end-of-directory or error
```

每调用 `readdir()`一次，就会从 `dirp` 所指代的目录流中读取下一目录条目，并返回一枚指针，指向经静态分配而得的 `dirent` 类型结构，内含与该条目相关的如下信息：

```
struct dirent {
    ino_t d_ino;          /* File i-node number */
    char d_name[];       /* Null-terminated name of file */
};
```

每次调用 `readdir()`都会覆盖该结构。

出于对程序可移植性的考虑，上述定义略去了 Linux `dirent` 结构中的各种非标准字段。这其中最令人感兴趣的当属 `d_type`，它同时获得了 BSD 流派的支持，但并未在其他 UNIX 系统中实现。该属性值用于标识命名于 `d_name` 之中文件的类型，诸如 `DT_REG`（普通文件）、`DT_DIR`（目录）、`DT_LNK`（符号链接）或 `DT_FIFO`（FIFO）。（这些名称类似于表 15-1 所列诸宏。）利用该属性值可省去为确定文件类型而对 `lstat()`的调用。注意，写作本书时，该属性仅获得 Btrfs、ext2、ext3 以及 ext4 的全面支持。

调用 `lstat()`（或者 `stat()`，如果应对符号链接解引用时）可获得 `d_name` 所指向文件的更多信息，其中，路径名由之前调用 `opendir()`时指定的 `dirpath` 参数与 “/” 字符以及 `d_name` 字段的

返回值拼接组成。

`readdir()`返回时并未对文件名进行排序，而是按照文件在目录中出现的天然次序（这取决于文件系统向目录添加文件时所遵循的次序，及其在删除文件后对目录列表中空隙的填补方式）。（命令 `ls-f` 对文件列表的排列与调用 `readdir()` 时一样，均未做排序处理。）

使用 `scandir(3)` 函数可以获得经过排序处理的文件列表，且排列规则可由程序员定义，具体细节请参考手册页。尽管该函数未获 SUSv3 接纳，但得到了大多数 UNIX 实现的支持。SUSv4 也对 `scandir()` 作了定义。

一旦遇到目录结尾或是出错，`readdir()` 将返回 `NULL`，针对后一种情况，还会设置 `errno` 以示具体错误。为了区别这两种情况，可编码如下：

```
errno = 0;
direntp = readdir(dirp);
if (direntp == NULL) {
    if (errno != 0) {
        /* Handle error */
    } else {
        /* We reached end-of-directory */
    }
}
```

如果目录内容恰逢应用调用 `readdir()` 扫描该目录时发生变化，那么应用程序可能无法观察到这些变动。SUSv3 明确指出，对于 `readdir()` 是否会返回自上次调用 `opendir()` 或 `rewinddir()` 后在目录中增减的文件，规范不做要求。至于最后一次执行上述调用前就存在的文件，应确保其全部返回。

`rewinddir()` 函数可将目录流回移到起点，以便对 `readdir()` 的下一调用将从目录的第一个文件开始。

```
#include <dirent.h>

void rewinddir(DIR *dirp);
```

`closedir()` 函数将由 `dirp` 指代、处于打开状态的目录流关闭，同时释放流所使用的资源。

```
#include <dirent.h>

int closedir(DIR *dirp);
```

Returns 0 on success, or -1 on error

SUSv3 还定义了两个高级函数：`telldir()` 和 `seekdir()`，允许随机访问目录流。有关这些函数的深入信息请参考手册页。

目录流与文件描述符

有一个目录流，就有一个文件描述符与之关联。`dirfd()` 函数返回与 `dirp` 目录流相关联的文件描述符。

```
#include <dirent.h>

int dirfd(DIR *dirp);
```

Returns file descriptor on success, or -1 on error

例如，将 `dirfd()` 返回的文件描述符传递给 `fchdir()`（参见 18.10 节），就可以把进程的当前工作目录改成相应目录。此外，还可以将其传递给 18.11 节所述各函数的 `dirfd` 参数。

`dirfd()` 函数还见于 BSD 系统，但在其他实现中则鲜有踪迹。该函数未获 SUSv3 接纳，但 SUSv4 则对其做了规范。

这里值得一提的是，`opendir()` 会为与目录流相关联的文件描述符自动设置 `close-on-exec` 标志 (`FD_CLOEXEC`)，以确保当执行 `exec()` 时自动关闭该文件描述符。（SUSv3 要求这一行为。）`close-on-exec` 标志将在 27.4 节加以描述。

示例程序

程序清单 18-2 使用 `opendir()`、`readdir()` 和 `closedir()` 函数来列出由命令行参数所指定各目录的内容（若未提供参数则为当前工作目录）。以下是运行该程序的一个例子：

```
$ mkdir sub                                Create a test directory
$ touch sub/a sub/b                        Make some files in the test directory
$ ./list_files sub                          List contents of directory
sub/a
sub/b
```

程序清单 18-2：扫描一个目录

```
----- dirs_links/list_files.c
#include <dirent.h>
#include "tspi_hdr.h"

static void          /* List all files in directory 'dirPath' */
listFiles(const char *dirpath)
{
    DIR *dirp;
    struct dirent *dp;
    Boolean isCurrent;          /* True if 'dirpath' is "." */

    isCurrent = strcmp(dirpath, ".") == 0;

    dirp = opendir(dirpath);
    if (dirp == NULL) {
        errMsg("opendir failed on '%s'", dirpath);
        return;
    }

    /* For each entry in this directory, print directory + filename */

    for (;;) {
        errno = 0;              /* To distinguish error from end-of-directory */
        dp = readdir(dirp);
        if (dp == NULL)
            break;

        if (strcmp(dp->d_name, ".") == 0 || strcmp(dp->d_name, "..") == 0)
            continue;          /* Skip . and .. */

        if (!isCurrent)
            printf("%s/", dirpath);
            printf("%s\n", dp->d_name);
    }
}
```



```

    if (errno != 0)
        errExit("readdir");

    if (closedir(dirp) == -1)
        errMsg("closedir");
}

int
main(int argc, char *argv[])
{
    if (argc > 1 && strcmp(argv[1], "--help") == 0)
        usageErr("%s [dir...]\n", argv[0]);

    if (argc == 1)                /* No arguments - use current directory */
        listFiles(".");
    else
        for (argv++; *argv; argv++)
            listFiles(*argv);

    exit(EXIT_SUCCESS);
}

```

dirs_links/list_files.c

readdir_r()函数

readdir_r()函数是 readdir()的变体。二者之间语义上的关键差异在于前者是可重入的，而后者不是。这是因为 readdir_r()对文件条目的返回利用的是由调用者分配的 entry 参数，而 readdir()则是将信息置于静态分配的结构并返回其指针。21.1.2 节和 31.1 节讨论了可重入性 (reentrancy)。

```

#include <dirent.h>

int readdir_r(DIR *dirp, struct dirent *entry, struct dirent **result);

```

Returns 0 on success, or a positive error number on error

针对既定 dirp，亦即之前调用 opendir()所打开的目录流，readdir_r()将下一项目录条目录置于由 entry 指向的 dirent 结构中。另外，还会在 result¹中放置指向该结构的一枚指针。如果抵达目录流尾部，那么会在 result²中返回 NULL(且 readdir_r()返回 0)。当出现错误时，readdir_r()不会返回-1，而是返回一个对应于 errno 的正整型值。

在 Linux 中，dirent 结构的 d_name 字段是大小为 256 字节的一个数组，足以容纳可能出现的最长文件名。虽然有几个其他的 UNIX 实现也为 d_name 定义了相同的大小，但 SUSv3 对此却并做规定，而另一些 UNIX 实现则将该字段定义为 1 字节的数组，并将正确分配结构大小的工作交给调用程序。这时，应将 d_name 字段大小设定为常量 NAME_MAX+1 (考虑终止空字节)。为确保可移植性，应用程序应以如下方式分配 dirent 结构：

```

struct dirent *entryp;
size_t len;

len = offsetof(struct dirent, d_name) + NAME_MAX + 1;
entryp = malloc(len);
if (entryp == NULL)
    errExit("malloc");

```

1 译者注：kernel.org 的在线手册页则为*result，译者倾向于后者，但未与作者沟通。请读者自行验证。

2 译者注：同上，疑为*result。

鉴于 `dirent` 结构中 `d_name` 字段（该属性总是位于结构的最后）之前各类属性的数量和大小在不同系统中实现不一，采用 `offsetof()` 宏（定义于 `<stddef.h>` 中）可避免程序对此产生依赖。

`offsetof()` 宏接受两个参数——结构类型和该结构中某一字段的名称——并返回一 `size_t` 类型值，亦即该字段距结构起点的字节偏移量。这个宏之所以必要，是由于编译器为满足诸如 `int` 之类的类型的对齐要求，可能在结构中插入填充字节。这会导致结构中某一字段的偏移量可能要大于该属性之前所有字段的长度总和。

18.9 文件树遍历：`nftw()`

`nftw()` 函数允许程序对整个目录子树进行递归遍历，并为子树中的每个文件执行某些操作（即，调用由程序员定义的函数）。

`nftw()` 函数是对执行类似功能的老函数 `ftw()` 的加强。由于提供了更多功能，对符号链接的处理也更易于把握（SUSv3 规定，`ftw()` 的函数实现无论是否对符号链接进行解引用，均符合规范），故而新近开发的应用程序应考虑采用 `nftw()`（new `ftw`）。SUSv3 将 `nftw()` 和 `ftw()` 均纳入规范，但 SUSv4 将后者标记为“已废止”。

GNU C 语言函数库也提供了派生自 BSD 分支的 `fts` API (`fts_open()`、`fts_read()`、`fts_children()`、`fts_set()` 和 `fts_close()`)。这些函数执行的任务类似于 `ftw()` 和 `nftw()`，但在遍历树方面为应用程序提供了更大的灵活性。然而，因为这些 API 目前尚未获得业界标准的接纳，也鲜有见诸于 BSD 后裔之外的其他 UNIX 实现，所以在此略而不论。

`nftw()` 函数遍历由 `dirpath` 指定的目录树，并为目录树中的每个文件调用一次由程序员定义的 `func` 函数。

```
#define _XOPEN_SOURCE 500
#include <ftw.h>

int nftw(const char *dirpath,
         int (*func) (const char *pathname, const struct stat *statbuf,
                     int typeflag, struct FTW *ftwbuf),
         int nopenfd, int flags);

Returns 0 after successful walk of entire tree, or -1 on error,
or the first nonzero value returned by a call to func
```

默认情况下，`nftw()` 会针对给定的树执行未排序的前序遍历，即对各目录的处理要先于各目录下的文件和子目录。

当 `nftw()` 遍历目录树时，最多会为树的每一层级打开一个文件描述符。参数 `nopenfd` 指定了 `nftw()` 可使用文件描述符数量的最大值。如果目录树深度超过这一最大值，那么 `nftw()` 会在做好记录的前提下，关闭并重新打开描述符，从而避免同时持有的描述符数目突破上限 `nopenfd`（从而导致运行越来越慢）。在较老的 UNIX 实现中，有的系统要求每个进程可打开的文件描述符数量不得超过 20 个，这更突显出这一参数的必要性。现代 UNIX 实现允许进程打开大量的文件描述符，因此，在指定该数目时出手可以大方一些（比如，10 或者更多）。

`nftw()` 的 `flags` 参数由 0 个或多个下列常量相或(`|`)组成，这些常量可对函数的操作做出修正。

FTW_CHDIR

在处理目录内容之前先调用 `chdir()` 进入每个目录。如果打算让 `func` 在 `pathname` 参数所指定文件的驻留目录下展开某些工作，那么就应当使用这一标志。

FTW_DEPTH

对目录树执行后序遍历。这意味着，`nftw()` 会在对目录本身执行 `func` 之前先对目录中的所有文件（及子目录）执行 `func` 调用。（这一标志名称容易引起误会——`nftw()` 遍历目录树遵循的是深度优先原则，而非广度优先。而这一标志的作用其实就是将先序遍历改为后序遍历。）

FTW_MOUNT

不会越界进入另一文件系统。因此，如果树中某一子目录是挂载点，那么不会对其进行遍历。

FTW_PHYS

默认情况下，`nftw()` 对符号链接进行解引用操作。而使用该标志则告知 `nftw()` 函数不要这么做。相反，函数会将符号链接传递给 `func` 函数，并将 `typeflag` 值置为 `FTW_SL`，如下所述。

`nftw()` 为每个文件调用 `func` 时传递 4 个参数。第一个参数 `pathname` 是文件的路径名。这个路径名可以是绝对路径，也可以是相对路径。如果指定 `dirpath` 时使用的是绝对路径，那么 `pathname` 就可能是绝对路径。反之，如果指定 `dirpath` 时使用的是相对路径名，则 `pathname` 中的路径可能是相对于进程调用 `nftw()` 时的当前工作目录而言。第二个参数 `statbuf` 是一枚指针，指向 `stat` 结构（参见 15.1 节），内含该文件的相关信息。第三个参数 `typeflag` 提供了有关该文件的深入信息，并具有如下特征值之一。

FTW_D

这是一个目录。

FTW_DNR

这是一个不能读取的目录（所以 `nftw()` 不能遍历其后代）。

FTW_DP

正在对一个目录进行后序遍历，当前项是一个目录，其所包含的文件和子目录已经处理完毕。

FTW_F

该文件的类型是除目录和符号链接以外的任何类型。

FTW_NS

对该文件调用 `stat()` 失败，可能是因为权限限制。`Statbuf` 中的值未定义。

FTW_SL

这是一个符号链接。仅当使用 `FTW_PHYS` 标志调用 `nftw()` 函数时才返回该值。

FTW_SLN

这是一个悬空的符号链接。仅当未在 `flags` 参数中指定 `FTW_PHYS` 标志时才会出现该值。

`Func` 的第四个参数 `ftwbuf` 是一枚指针，所指向结构定义如下：

```
struct FTW {
    int base;          /* Offset to basename part of pathname */
    int level;        /* Depth of file within tree traversal */
};
```

该结构的 `base` 字段是指 `func` 函数中 `pathname` 参数内文件名部分（最后一个 “/” 字符之

后的部分)的整型偏移量。level 字段是指该条目相对于遍历起点(其 level 为 0)的深度。

每次调用 func 都必须返回一个整型值,由 nftw()加以解释。如果返回 0, nftw()会继续对树进行遍历,如果所有对 func 的调用均返回 0,那么 nftw()本身也将返回 0 给调用者。若返回非 0 值,则通知 nftw()立即停止对树的遍历,这时 nftw()也会返回相同的非 0 值。

由于 nftw()使用的数据结构是动态分配的,故而应用程序提前终止目录树遍历的唯一方法就是让 func 调用返回一个非 0 值。调用 longjmp() (6.8 节)从 func 退出会导致不可预期的结果——至少会引起内存泄漏。

示例程序

程序清单 18-3 展示了 nftw()的使用。

程序清单 18-3: 使用 nftw()遍历目录树

```
----- dirs_links/nftw_dir_tree.c
#define _XOPEN_SOURCE 600      /* Get nftw() and S_IFSOCK declarations */
#include <ftw.h>
#include "tspi_hdr.h"

static void
usageError(const char *progName, const char *msg)
{
    if (msg != NULL)
        fprintf(stderr, "%s\n", msg);
    fprintf(stderr, "Usage: %s [-d] [-m] [-p] [directory-path]\n", progName);
    fprintf(stderr, "\t-d Use FTW_DEPTH flag\n");
    fprintf(stderr, "\t-m Use FTW_MOUNT flag\n");
    fprintf(stderr, "\t-p Use FTW_PHYS flag\n");
    exit(EXIT_FAILURE);
}

static int
dirTree(const char *pathname, const struct stat *sbuf, int type,
        struct FTW *ftwb)
{
    switch (sbuf->st_mode & S_IFMT) {      /* Print file type */
    case S_IFREG:  printf("-"); break;
    case S_IFDIR:  printf("d"); break;
    case S_IFCHR:  printf("c"); break;
    case S_IFBLK:  printf("b"); break;
    case S_IFLNK:  printf("l"); break;
    case S_IFIFO:  printf("p"); break;
    case S_IFSOCK: printf("s"); break;
    default:       printf("?"); break;    /* Should never happen (on Linux) */
    }

    printf(" %s ",
           (type == FTW_D) ? "D " : (type == FTW_DNR) ? "DNR " :
           (type == FTW_DP) ? "DP " : (type == FTW_F) ? "F " :
           (type == FTW_SL) ? "SL " : (type == FTW_SLN) ? "SLN " :
           (type == FTW_NS) ? "NS " : " ");

    if (type != FTW_NS)
        printf("%7ld ", (long) sbuf->st_ino);
    else

```

```

        printf("        ");

    printf("%*s", 4 * ftwb->level, "");          /* Indent suitably */
    printf("%s\n", &pathname[ftwb->base]);      /* Print basename */
    return 0;                                   /* Tell nftw() to continue */
}

int
main(int argc, char *argv[])
{
    int flags, opt;

    flags = 0;
    while ((opt = getopt(argc, argv, "dmp")) != -1) {
        switch (opt) {
            case 'd': flags |= FTW_DEPTH;    break;
            case 'm': flags |= FTW_MOUNT;    break;
            case 'p': flags |= FTW_PHYS;     break;
            default: usageError(argv[0], NULL);
        }
    }

    if (argc > optind + 1)
        usageError(argv[0], NULL);

    if (nftw((argc > optind) ? argv[optind] : ".", dirTree, 10, flags) == -1) {
        perror("nftw");
        exit(EXIT_FAILURE);
    }
    exit(EXIT_SUCCESS);
}
----- dirs_links/nftw_dir_tree.c

```

程序清单 18-3 中程序以层级缩进方式显示了一个目录树中的文件。每行显示一个文件，内容包括文件名、文件类型及 i-node 编号。可通过命令行选项来指定 nftw()调用中的 flags 参数值。下面的 shell 会话展示了运行程序的示例结果。首先创建一个新的空目录，并在其中填充各种类型的文件。

```

$ mkdir dir
$ touch dir/a dir/b           Create some plain files
$ ln -s a dir/s1             and a symbolic link
$ ln -s x dir/ds1            and a dangling symbolic link
$ mkdir dir/sub               and a subdirectory
$ touch dir/sub/x             with a file of its own
$ mkdir dir/sub2              and another subdirectory
$ chmod 0 dir/sub2            that is not readable

然后使用该程序调用 nftw()函数，其 flags 参数为 0:
$ ./nftw_dir_tree dir
d D  2327983  dir
- F  2327984   a
- F  2327985   b
- F  2327984  s1           The symbolic link s1 was resolved to a
l SLN 2327987  ds1
d D  2327988  sub
- F  2327989   x
d DNR 2327994  sub2

从以上输出可见，对符号链接 s1 进行了解析。

```

然后再使用该程序来调用 `nftw()` 函数，令 `flags` 参数包含 `FTW_PHYS` 和 `FTW_DEPTH` 标志：

```
$ ./nftw_dir_tree -p -d dir
- F 2327984 a
- F 2327985 b
l SL 2327986 s1 The symbolic link s1 was not resolved
l SL 2327987 ds1
- F 2327989 x
d DP 2327988 sub
d DNR 2327994 sub2
d DP 2327983 dir
```

从以上输出可见，未对符号链接 `s1` 进行解析。

nftw() 的 FTW_ACTIONRETVAL 标识

始于 2.3.3 版本，`glibc` 允许在 `nftw()` 的 `flags` 参数中指定一个额外的非标准标志 `FTW_ACTIONRETVAL`。此标志改变了 `nftw()` 函数对 `func()` 返回值的解释方式。当指定该标识时，`func()` 应返回下列值之一。

FTW_CONTINUE

与传统 `func()` 函数返回 0 时一样，继续处理目录树中的条目。

FTW_SKIP_SIBLINGS.

不再进一步处理当前目录中的条目，恢复对父目录的处理。

FTW_SKIP_SUBTREE

如果 `pathname` 是目录（即 `typeflag` 为 `FTW_D`），那么就不对该目录下的条目调用 `func()`。恢复进行对该目录的下一个同级目录的处理。

FTW_STOP

与传统 `func()` 函数返回非 0 值时一样，不再进一步处理目录树下的任何条目。`nftw()` 将返回 `FTW_STOP` 给调用者。

想从 `<ftw.h>` 文件中获得对 `FTW_ACTIONRETVAL` 的定义，必须定义 `_GNU_SOURCE` 特性测试宏。

18.10 进程的当前工作目录

一个进程的当前工作目录（`current working directory`）定义了该进程解析相对路径名的起点。新进程的当前工作目录继承自其父进程。

获取当前工作目录

进程可使用 `getcwd()` 来获取当前工作目录。

```
#include <unistd.h>

char *getcwd(char *cwbdbuf, size_t size);

Returns cwbdbuf on success, or NULL on error
```

`getcwd()` 函数将内含当前工作目录绝对路径的字符串（包括结尾空字符）置于 `cwbdbuf` 指

向的已分配缓冲区中。调用者必须为 `cwdbuf` 缓冲区分配至少 `size` 个字节的空间。（通常，`cwdbuf` 的大小与 `PATH_MAX` 常量相当。）

一旦调用成功，`getcwd()`将返回一枚指向 `cwdbuf` 的指针。如果当前工作目录的路径名长度超过 `size` 个字节，那么 `getcwd()`会返回 `NULL`，并将 `errno` 置为 `ERANGE`。

在 Linux/x86-32 系统中，`getcwd()`返回指针所指向的字符串最大长度可达 4096 个字节。如果当前工作目录（以及 `cwdbuf` 和 `size`）突破了这一限制，那么就会直接对路径名做截断处理，移去始于起点的整个目录前缀（字符串仍以空字符结尾）。换言之，当当前工作目录的绝对路径超出这一限制时，`getcwd()`的行为也不再可靠。

实际上，Linux 的 `getcwd()`系统调用为要返回的路径名在内部分配了一个虚拟内存页。x86-32 架构的页大小为 4096 字节，而在页尺寸更大的架构中（比如，Alpha 的页大小为 8192 字节），`getcwd()`能返回更长的路径名。

若 `cwdbuf` 为 `NULL`，且 `size` 为 0，则 `glibc` 封装函数会为 `getcwd()`按需分配一个缓冲区，并将指向该缓冲区的指针作为函数的返回值。为避免内存泄漏，调用者之后必须调用 `free()`来释放这一缓冲区。对可移植性有所要求的应用程序应当避免依赖该特性。大多数其他实现则针对 `SUSv3` 规范提供了一个更为简单的扩展。如果 `cwdbuf` 是 `NULL`，那么 `getcwd()`将分配一个大小为 `size` 字节的缓冲区，用于向调用者返回结果。`glibc` 的 `getcwd()`也实现了这一特性。

GNU C 函数库还为获取当前工作目录提供了另外两个函数。派生自 BSD 的 `getwd(path)`函数容易引起缓冲区溢出，因为该函数无法为返回的路径名长度设定上限。`get_current_dir_name()`函数也会返回包含当前工作目录名的一个字符串。虽然该函数易于使用，但却不具有可移植性。考虑到安全性和可移植性，`getcwd()`无疑是不二之选（前提是避免使用 GNU 扩展功能）。

只要具有合适的权限（大体要求是，身为进程属主或者具有 `CAP_SYS_PTRACE` 能力），就可通过读取（`readlink()`）Linux 专有符号链接 `/proc/PID/cwd` 的内容来确定任何进程的当前工作目录。

改变当前工作目录

`chdir()`系统调用将调用进程的当前工作目录改变为由 `pathname` 指定的相对或绝对路径名（如属于符号链接，还会对其解除引用）。

```
#include <unistd.h>

int chdir(const char *pathname);

Returns 0 on success, or -1 on error
```

`fchdir()`系统调用与 `chdir()`作用相同，只是在指定目录时使用了文件描述符，而该描述符是之前调用 `open()`打开相应目录时获得的。

```
#define _XOPEN_SOURCE 500 /* Or: #define _BSD_SOURCE */
#include <unistd.h>

int fchdir(int fd);

Returns 0 on success, or -1 on error
```

以下代码片段所示为使用 `fchdir()` 将进程的当前工作目录变为另一位置，然后再改回原始位置：

```
int fd;

fd = open(".", O_RDONLY);      /* Remember where we are */
chdir(somepath);             /* Go somewhere else */
fchdir(fd);                  /* Return to original directory */
close(fd);
```

使用 `chdir()` 达到同等效果的代码如下所示：

```
char buf[PATH_MAX];

getcwd(buf, PATH_MAX);      /* Remember where we are */
chdir(somepath);           /* Go somewhere else */
chdir(buf);                 /* Return to original directory */
```

18.11 针对目录文件描述符的相关操作

始于版本 2.6.16，Linux 内核提供了一系列新的系统调用，在执行与传统系统调用相似任务的同时，还提供了一些附加功能，对某些应用程序非常有用。表 18-2 对这些调用进行了归纳。之所以在本章介绍这些系统调用，是因为它们对进程当前工作目录的传统语义做了改动。

表 18-2：系统调用使用目录文件描述来解释相对路径

新 接 口	类似的传统接口	备 注
<code>faccessat()</code>	<code>access()</code>	支持 <code>AT_EACCESS</code> 和 <code>AT_SYMLINK_NOFOLLOW</code> 标志
<code>fchmodat()</code>	<code>chmod()</code>	
<code>fchownat()</code>	<code>chown()</code>	支持 <code>AT_SYMLINK_NOFOLLOW</code> 标志
<code>fstatat()</code>	<code>stat()</code>	支持 <code>AT_SYMLINK_NOFOLLOW</code> 标志
<code>linkat()</code>	<code>link()</code>	支持（始于 Linux 2.6.18） <code>AT_SYMLINK_FOLLOW</code> 标志
<code>mkdirat()</code>	<code>mkdir()</code>	
<code>mkfifoat()</code>	<code>mkfifo()</code>	基于 <code>mknodat()</code> 的库函数
<code>mknodat()</code>	<code>mknod()</code>	
<code>openat()</code>	<code>open()</code>	
<code>readlinkat()</code>	<code>readlink()</code>	
<code>renameat()</code>	<code>rename()</code>	
<code>symlinkat()</code>	<code>symlink()</code>	
<code>unlinkat()</code>	<code>unlink()</code>	支持 <code>AT_REMOVEDIR</code> 标志
<code>utimensat()</code>	<code>utimes()</code>	支持 <code>AT_SYMLINK_NOFOLLOW</code> 标志

为便于描述这些系统调用，这里就以 `openat()` 为例。

```
#define _XOPEN_SOURCE 700      /* Or define _POSIX_C_SOURCE >= 200809 */
#include <fcntl.h>

int openat(int dirfd, const char *pathname, int flags, ... /* mode_t mode */);
           Returns file descriptor on success, or -1 on error
```


`openat()`系统调用类似于传统的 `open()`系统调用，只是添加了一个 `dirfd` 参数，其作用如下。

- 如果 `pathname` 中为一相对路径名，那么对其解释则以打开文件描述符 `dirfd` 所指向的目录为参照点，而非进程的当前工作目录。
- 如果 `pathname` 中为一相对路径，且 `dirfd` 中所含为特殊值 `AT_FDCWD`，那么对 `pathname` 的解释则相对与进程当前工作目录（即与 `open(2)`行为一致）而言。
- 如果 `pathname` 中为绝对路径，那么将忽略 `dirfd` 参数。

`openat()`的 `flag` 参数目的与 `open()`相同。然而，部分表 18-2 中所列系统调用还支持 `flags` 参数，这是相应的传统系统调用所不具备的，其目的在于修改调用语义。出现频率最高的标志为 `AT_SYMLINK_NOFOLLOW`，其含义是如果 `pathname` 为符号链接，那么系统调用将操作于符号链接本身，而非符号链接所指向的文件。（`linkat()`系统调用提供了 `AT_SYMLINK_FOLLOW` 标志，其作用正好相反，即改变 `linkat()`的默认行为，当 `oldpath` 属于符号链接时对其进行解引用操作。）有关其他标志的详情，请参考相应手册页。

之所以要支持表 18-2 中所列的系统调用，其原因有二（此处再以 `openat()`为例）。

- 当调用 `open()`打开位于当前工作目录之外的文件时，可能会发生某些竞态条件。而使用 `openat()`就能够避免这一问题。在调用 `open()`的同时，如果 `pathname` 目录前缀的某些部分发生了改变，就可能导致竞争。要想避免这类竞态，可以针对目标目录打开一个文件描述符，然后将该描述符传递给 `openat()`。
- 如第 29 章所述，工作目录是进程的属性之一，为进程中所有线程所共享。而对某些应用程序而言，需要针对不同线程拥有不同的“虚拟”工作目录。将 `openat()`与应用所维护的目录文件描述符相结合，就可以模拟出这一功能。

SUSv3 并未对这些系统调用加以规范，但 SUSv4 将其包括在内。为了获得对这些系统调用的声明，必须在包含相应头文件之前（比如定义 `open()`的 `<fcntl.h>`）将 `_XOPEN_SOURCE` 特性测试宏定义为大于或等于 700 的值。另外，将 `_POSIX_C_SOURCE` 宏的值定义为大于或等于 200809 也能收到同样效果。（在 2.10 版本之前的 `glibc` 中，要获得对这些系统调用的声明还需要定义 `_ATFILE_SOURCE` 宏。）

Solaris 9 及其更高版本也提供了一些表 18-2 所列接口的版本，只是语义略有不同。

18.12 改变进程的根目录：chroot()

每个进程都有一个根目录，该目录是解释绝对路径（即那些以/开始的目录）时的起点。默认情况下，这是文件系统的真实根目录。（新进程从其父进程处继承根目录。）有些场合需要改变一个进程的根目录，而特权级（`CAP_SYS_CHROOT`）进程通过 `chroot()`系统调用能够做到这一点。

```
#define _BSD_SOURCE
#include <unistd.h>

int chroot(const char *pathname);
```

Returns 0 on success, or -1 on error

`chroot()`系统调用将进程的根目录改为由 `pathname` 指定的目录（如果 `pathname` 是符号链接，还将对其解引用）。自此，对所有绝对路径名的解释都将以该文件系统的这一位置作为起

点。鉴于这会将应用程序限定于文件系统的特定区域，有时也将此称为设立了一个 `chroot` 监禁区。

SUSv2 包含了对 `chroot()` 的定义（标记为 `LEGACY`），但 SUSv3 又将其删去。无论如何，`chroot()` 还是获得了大多数 UNIX 实现的支持。

借助于 `chroot()` 系统调用，`chroot` 命令可在 `chroot` 监禁区中执行 `shell` 命令。

而通过读取（`readlink()`）Linux 专有 `/proc/PID/root` 符号链接的内容，可以获取任何进程的根目录。

`ftp` 程序就是应用 `chroot()` 的典型实例之一。作为一种安全措施，当用户匿名登录 `ftp` 时，`ftp` 程序将使用 `chroot()` 为新进程设置根目录——一个专门预留给匿名登录用户的目录。调用 `chroot()` 后，用户将受困于文件系统中新根目录下的子树中，无法在整个文件系统中信马由缰。（这里所依赖的事实是根目录是其自身的父目录。也就是说 `./` 是 `/` 的一个链接，所以改变目录到后再执行 `cd ..` 命令时，用户依然会待在同一目录下。）

一些 UNIX 实现（不包括 Linux）允许多个硬链接指向同一目录，这时就有可能在一个子目录中创建指向其父目录（或者更高层级远祖目录）的硬链接。在这种系统实现中，如果存在指向监禁区目录树之外的硬链接，那么监禁区的安全将受到威胁。而指向监禁区之外目录的符号链接则不是问题，因为对这些符号链接的解释将在进程新根目录的框架之内进行，所以是无法染指到 `chroot` 监禁区之外的。

通常情况下，不是随便什么程序都可以在 `chroot` 监禁区中运行的，因为大多数程序与共享库之间采取的是动态链接的方式。因此，要么只能局限于运行静态链接程序，要么就在监禁区中复制一套标准的共享库系统目录（比如，包括 `/lib` 和 `/usr/lib`）（针对这一点，14.9.4 节描述的绑定挂载特性就派上了用场）。

`chroot()` 系统调用从未被视为一个完全安全的监禁机制。首先，特权级程序可以在随后对 `chroot()` 的进一步调用中利用种种手段而越狱成功。例如，特权级（`CAP_MKNOD`）程序能够使用 `mknod()` 来创建一个内存设备文件（类似于 `/dev/mem`），并通过该设备来访问 RAM 的内容，到那时，就一切皆有可能了。通常，最好不要在 `chroot` 监禁区文件系统内放置 `set-user-ID-root` 程序。

即便是对于无特权程序，也必须小心防范如下几条可能的越狱路线。

- 调用 `chroot()` 并未改变进程的当前工作目录。因此，通常应在调用 `chroot()` 之前或者之后调用一次 `chdir()` 函数（例如，`chroot()` 调用之后执行 `chdir("/")`）。如果没有这么做，那么进程就能够使用相对路径去访问监狱之外的文件和目录。（一些 BSD 的衍生系统杜绝了这一可能性——如果当前工作目录位于新的根目录树之外，那么 `chroot()` 调用会将其修改为与根目录一致。）
- 如果进程针对监禁区之外的某一目录持有一打开文件描述符，那么结合 `fchdir()` 和 `chroot()` 即可越狱成功，如下面代码所示：

```
int fd;

fd = open("/", O_RDONLY);
chroot("/home/mtk");          /* Jailed */
fchdir(fd);
chroot(".");                  /* Out of jail */
```

为了防止这种可能性，必须关闭所有指向监禁区外目录的文件描述符。（其他一些 UNIX 实现提供了 `fchroot()` 系统调用，可用于获得与上述代码片段类似的结果。）

- 即使针对上述可能性采取了防范措施，仍不足以阻止任意非特权程序（即无法控制其操作的程序）越狱成功。遭到囚禁的进程仍然能够利用 UNIX 域套接字来接受（自另一进程处）指向监禁区外目录的文件描述符。（61.13.3 节简要描述了进程间利用套接字来传递文件描述符的概念。）将这一文件描述指定为 `fchdir()` 调用的入参，程序即可将其当前工作目录置于监禁区外，之后再通过相对路径来随意访问文件和目录。

一些 BSD 衍生系统提供的 `jail()` 系统调用解决了包括上述问题在内的不少问题，其所创建的监禁区即使针对特权级进程也是安全的。

18.13 解析路径名: `realpath()`

`realpath()` 库函数对 `pathname`（以空字符结尾的字符串）中的所有符号链接一一解除引用，并解析其中所有对 `/` 和 `./` 的引用，从而生成一个以空字符结尾的字符串，内含相应的绝对路径名。

```
#include <stdlib.h>

char *realpath(const char *pathname, char *resolved_path);

Returns pointer to resolved pathname on success, or NULL on error
```

生成的字符串将置于 `resolved_path` 指向的缓冲区中，该字符串应当是一个字符数组，长度至少为 `PATH_MAX` 个字节。一旦调用成功，`realpath()` 将返回指向该字符串的一枚指针。

`glibc` 的 `realpath()` 实现允许调用者将 `resolved_path` 参数指定为空。这时，`realpath()` 会为经解析生成的路径名分配一个多达 `PATH_MAX` 个字节的缓冲区，并将指向该缓冲区的指针作为结果返回。（调用者必须自行调用 `free()` 来释放该缓冲区。）`SUSv3` 并未将该扩展功能纳入规范，但 `SUSv4` 对其进行了定义。

程序清单 18-4 中的程序采用 `readlink()` 和 `realpath()` 来读取符号链接的内容，并将该链接解析为一个绝对路径名。下面是运行该程序的一个示例：

```
$ pwd                               Where are we?
/home/mtk
$ touch x                             Make a file
$ ln -s x y                           and a symbolic link to it
$ ./view_symlink y
readlink: y --> x
realpath: y --> /home/mtk/x
```

程序清单 18-4：读取并解析一个符号链接

```
----- dirs_links/view_symlink.c
#include <sys/stat.h>
#include <limits.h>                /* For definition of PATH_MAX */
#include "tspi_hdr.h"

#define BUF_SIZE PATH_MAX

int
main(int argc, char *argv[])
```

```

{
    struct stat statbuf;
    char buf[BUF_SIZE];
    ssize_t numBytes;

    if (argc != 2 || strcmp(argv[1], "--help") == 0)
        usageErr("%s pathname\n", argv[0]);
    if (lstat(argv[1], &statbuf) == -1)
        errExit("lstat");

    if (!S_ISLNK(statbuf.st_mode))
        fatal("%s is not a symbolic link", argv[1]);

    numBytes = readlink(argv[1], buf, BUF_SIZE - 1);
    if (numBytes == -1)
        errExit("readlink");
    buf[numBytes] = '\0'; /* Add terminating null byte */
    printf("readlink: %s --> %s\n", argv[1], buf);

    if (realpath(argv[1], buf) == NULL)
        errExit("realpath");
    printf("realpath: %s --> %s\n", argv[1], buf);

    exit(EXIT_SUCCESS);
}

```

----- dirs_links/view_symlink.c

18.14 解析路径名字符串：dirname()和basename()

dirname()和basename()函数将一个路径名字符串分解成目录和文件名两部分。（这些函数执行的任务与dirname(1)和basename(1)命令相类似。）

```
#include <libgen.h>
```

```
char *dirname(char *pathname);
char *basename(char *pathname);
```

Both return a pointer to a null-terminated (and possibly
statically allocated) string

比如，给定路径名为/home/britta/prog.c，dirname()将返回/home/britta，而basename()将返回prog.c。将dirname()返回的字符串与一斜线字符(/)以及basename()返回的字符串拼接起来，将生成一条完整的路径名。

关于dirname()和basename()的操作请注意以下几点。

- 将忽略pathname中尾部的斜线字符。
- 如果pathname中未包含斜线字符，那么dirname()将返回字符串.(点)，而basename()将返回pathname。
- 如果pathname仅由一个斜线字符组成，那么dirname()和basename()均将返回字符串/。将其应用于上述的拼接规则，所创建的路径名字符串为///。该路径名属于有效路径名。因为多个连续斜线字符相当于单个斜线字符，所以路径名///就相当于路径名/。
- 如果pathname为空指针或者空字符串，那么dirname()和basename()均将返回字符串.(点)。（拼接这些字符串将生成路径名./，对等于.，即当前目录。）

表 18-3 所示为 `dirname()` 和 `basename()` 针对各种示例路径名所返回的字符串。

表 18-3: `dirname()` 和 `basename()` 返回的字符串示例

路径名字符串	<code>dirname()</code>	<code>basename()</code>
/	/	/
/usr/bin/zip	/usr/bin	zip
/etc/passwd///	/etc	passwd
/etc///passwd	/etc	passwd
etc/passwd	etc	passwd
passwd	.	passwd
passwd/	.	passwd
..	.	..
NULL	.	.

程序清单 18-5: `dirname()` 和 `basename()` 的应用

```
----- dirs_links/t_dirbasename.c
#include <libgen.h>
#include "tlpi_hdr.h"

int
main(int argc, char *argv[])
{
    char *t1, *t2;
    int j;

    for (j = 1; j < argc; j++) {
        t1 = strdup(argv[j]);
        if (t1 == NULL)
            errExit("strdup");
        t2 = strdup(argv[j]);
        if (t2 == NULL)
            errExit("strdup");

        printf("%s ==> %s + %s\n", argv[j], dirname(t1), basename(t2));

        free(t1);
        free(t2);
    }

    exit(EXIT_SUCCESS);
}
----- dirs_links/t_dirbasename.c
```

`dirname()` 和 `basename()` 均可修改 `pathname` 所指向的字符串。因此，如果希望保留原有的路径名字符串，那么就必须向 `dirname()` 和 `basename()` 传递该字符串的副本，如程序清单 18-5 所示。该程序使用 `strdup()`（该函数调用了 `malloc()`）来制作传递给 `dirname()` 和 `basename()` 的字符串副本，然后再使用 `free()` 将其释放。

最后需要指出的是，`dirname()` 和 `basename()` 所返回的指针均可指向经由静态分配的字符

串，对相同函数的后续调用可能会修改这些字符串的内容。

18.15 总结

i-node 中并不包含文件的名称。相反，对文件的命名利用的是目录条目，而目录则是列出文件名和 i-node 编号之间对应关系的一个表格。也将这些目录条目称作（硬）链接。一个文件可以有多个链接，这些链接之间的地位是平等的。可使用 `link()`和 `unlink()`来创建和移除链接，对文件的重命名则使用系统调用 `rename()`。

调用 `symlink()`可创建符号（或者软）链接。符号链接在某些方面与硬链接相类似，其差异则在于符号链接可以跨越文件系统边界，还可指代目录。符号链接只是一个内容包含了另一文件名称的文件，可通过 `readlink()`来获取该文件的名称。（目标）i-node 的链接计数中并未包含符号链接，如果将该链接所指向的文件移除，那么此链接将处于悬空状态。一些系统调用会自动对符号链接进行解引用（下溯），其余的则不会。有时系统会提供两种版本的系统调用，一种会解引用符号链接，另一种则不会，例如 `stat()`和 `lstat()`。

创建目录使用的是 `mkdir()`，移除目录则使用 `rmdir()`。而扫描一个目录的内容则可使用 `opendir()`、`readdir()`以及相关函数。`nftw()`函数允许程序遍历一棵完整的目录树，并为树中每个文件调用由程序员定义的函数。

`remove()`函数可以用来移除一个文件（即一个链接）或者一个空目录。

每个进程都拥有一个根目录和一个当前工作目录，分别作为解释绝对路径和相对路径的参照点。可通过 `chroot()`和 `chdir()`系统调用来修改这些属性。而 `getcwd()`函数则返回进程的当前工作目录。

Linux 还提供了一套新的系统调用（如：`openat()`），其行为与其传统同行（如：`open()`）相类似，不同之处则在于可利用新的系统调用来提供一个指向目录的文件描述符（而非进程的当前工作目录），用于作为解释相对路径名的参照点。这将有助于避免特定类型的竞态条件，以及为每个线程实现虚拟工作目录。

`realpath()`函数解析一个路径名——解引用所有的符号链接，并将所有的 `.`和 `..`解析为相应目录——从而生成相应的绝对路径名。`dirname()`和 `basename()`函数可用来将路径名分解为目录和文件名两部分。

18.16 练习

- 18-1.** 4.3.2 节曾指出，如果一个文件正处于执行状态，那么要将其打开以执行写操作是不可能的（`open()`调用返回-1，且将 `errno` 置为 `ETXTBSY`）。然而，在 shell 中执行如下操作却是可能的：

```
$ cc -o longrunner longrunner.c
$ ./longrunner &                               Leave running in background
$ vi longrunner.c                               Make some changes to the source code
$ cc -o longrunner longrunner.c
```

最后一条命令覆盖了现有的同名可执行文件。原因何在？（提示：在每次编译后调用 `ls -li` 命令来查看可执行文件的 i-node 编号。）

- 18-2.** 以下代码中对 `chmod()`的调用为什么会失败？

```
mkdir("test", S_IRUSR | S_IWUSR | S_IXUSR);
chdir("test");
fd = open("myfile", O_RDWR | O_CREAT, S_IRUSR | S_IWUSR);
symlink("myfile", "../mylink");
chmod("../mylink", S_IRUSR);
```

- 18-3.** 实现 `realpath()`。
- 18-4.** 修改程序清单 18-2 中的程序，用 `readdir_r()`来取代 `readdir()`。
- 18-5.** 实现一个功能与 `getcwd()`相当的函数。提示：要获得当前工作目录的名称，可调用 `opendir()`和 `readdir()`来遍历其父目录（..）中的各个条目，查找其中与当前工作目录具有相同 i-node 编号及设备编号（即，分别为 `stat()`和 `lstat()`调用所返回 `stat` 结构中的 `st_ino` 和 `st_dev` 属性）的一项。如此这般，沿着目录树层层拾级而上（`chdir("../")`）并进行扫描，就能构建出完整的目录路径。当父目录与当前工作目录相同时（回忆 `../`与 `/`相同的情况），就结束遍历。无论调用该函数成功与否，都应将调用者遣回其起始目录（使用 `open()`和 `fchdir()`能很方便地实现这一功能）。
- 18-6.** 使用 `FTW_DEPTH` 标志来修改程序清单 18-3(`nftw_dir_tree.c`)中的程序。注意目录树遍历顺序的差异。
- 18-7.** 编写一程序，使用 `nftw()`来遍历目录树，并打印出树中各类文件（普通文件、目录、符号链接等）的总和及百分比。
- 18-8.** 实现 `nftw()`。（需要使用 `opendir()`、`readdir()`、`closedir()`和 `stat()`等系统调用。）
- 18-9.** 18.10 节展示了两种技术（分别为 `fchdir()`和 `chdir()`），用于在将当前工作目录转到另一位置后，再返回之前的当前工作目录。假设需要反复执行这一操作，哪种方法更为高效？原因何在？请写一段程序加以验证。

第 19 章

监控文件事件

某些应用程序需要对文件或目录进行监控，已侦测其是否发生了特定事件。例如，当把文件加入或移出一目录时，图形化文件管理器应能判定此目录是否在其当前显示之列，而守护进程可能也想要监控自己的配置文件，以了解其是否被修改。

自内核 2.6.13 起，Linux 开始提供 `inotify` 机制，以允许应用程序监控文件事件。本章将介绍 `inotify` 的用法。

`inotify` 机制所取代的是 `dnotify`，后者一来较为陈旧，二来仅具备前者的部分功能。本章会在结尾处对 `dnotify` 做简要介绍，着重突出 `inotify` 的优势所在。

`inotify` 和 `dnotify` 都是 Linux 专有机制。（仅有少数其他系统提供类似机制。比如，BSD 所提供的 `kqueue` API。）

有几个函数库提供的（类似）API，比之 `inotify` 和 `dnotify`，既更为抽象，在可移植性方面也略胜一筹。某些应用可能更倾向于采用这样的函数库。而这其中的一些库也会调用 `inotify` 和 `dnotify`，前提是要获得操作系统的支持。FAM（文件变更监控 <http://oss.sgi.com/projects/fam/>）和 Gamin (<http://www.gnome.org/~veillard/gamin/>)便是此类库的两个例子。

19.1 概述

使用 `inotify` API 有以下几个关键步骤。

1. 应用程序使用 `inotify_init()` 来创建一个 `inotify` 实例，该系统调用所返回的文件描述符用于在后续操作中指代该实例。
2. 应用程序使用 `inotify_add_watch()` 向 `inotify` 实例（由步骤 1 创建）的监控列表添加条目，藉此告知内核哪些文件是自己的兴趣所在。每个监控项都包含一个路径名以及相关的位掩码。位掩码针对路径名指明了所要监控的事件集合。作为函数结果，`inotify_add_watch()` 将返回一监控描述符，用于在后续操作中指代该监控项。（系统调用 `inotify_rm_watch()` 执行其逆向操作，将之前添加入 `inotify` 实例的监控项移除。）
3. 为获得事件通知，应用程序需针对 `inotify` 文件描述符执行 `read()` 操作。每次对 `read()` 的成

功调用，都会返回一个或多个 `inotify_event` 结构，其中各自记录了处于 `inotify` 实例监控之下的某一路径名所发生的事件。

4. 应用程序在结束监控时会关闭 `inotify` 文件描述符。这会自动清除与 `inotify` 实例相关的所有监控项。

inotify 机制可用于监控文件或目录。当监控目录时，与路径自身及其所含文件相关的事件都会通知给应用程序。

`inotify` 监控机制为非递归。若应用程序有意监控整个目录子树内的事件，则需对该树中的每个目录发起 `inotify_add_watch()` 调用。

可使用 `select()`、`poll()`、`epoll` 以及由信号驱动的 I/O（自 Linux 2.6.25 起）来监控 `inotify` 文件描述符。只要有事件可供读取，上述 API 便会将 `inotify` 文件描述符标记为可读。关乎这些编程接口的详细信息请见第 63 章。

`inotify` 机制属可选的 Linux 内核组件，可通过 `CONFIG_INOTIFY` 和 `CONFIG_INOTIFY_USER` 选项进行配置。

19.2 inotify API

`inotify_init()` 系统调用可创建一新的 `inotify` 实例。

```
#include <sys/inotify.h>

int inotify_init(void);

Returns file descriptor on success, or -1 on error
```

作为函数结果，`inotify_init()` 会返回一个文件描述符（句柄），用于在后续操作中指代此 `inotify` 实例。

Linux 自内核 2.6.27 开始支持一个新的、非标准的系统调用 `inotify_init1()`。该系统调用所执行的任务与 `inotify_init()` 相同，但提供了一个额外的参数 `flag`，用于修改系统调用的行为。该参数支持的标志有二：`IN_CLOEXEC` 标志会使内核针对新文件描述符激活 `close-on-exec` 标志 (`FD_CLOEXEC`)。引入该标志的原因正如 4.3.1 节所述 `open()` 的 `O_CLOEXEC` 标志一样。`IN_NONBLOCK` 标志会导致内核激活底层打开文件描述的 `O_NONBLOCK` 标志，如此一来，未来的读操作将是非阻塞式的，省得还要额外调用 `fcntl()` 来获得相同效果。

针对文件描述符 `fd` 所指代 `inotify` 实例的监控列表，系统调用 `inotify_add_watch()` 既可以追加新的监控项，也可以修改现有监控项。（请参考图 19-1。）

```
#include <sys/inotify.h>

int inotify_add_watch(int fd, const char *pathname, uint32_t mask);

Returns watch descriptor on success, or -1 on error
```

参数 `pathname` 标识欲创建或修改的监控项所对应的文件。调用程序必须对该文件具有读权限（调用 `inotify_add_watch()` 时，会对文件权限做一次性检查。只要监控项继续存在，即便有

人更改了文件权限，使调用程序不再对文件具有读权限，调用程序依然会继续收到文件的通
知消息)。

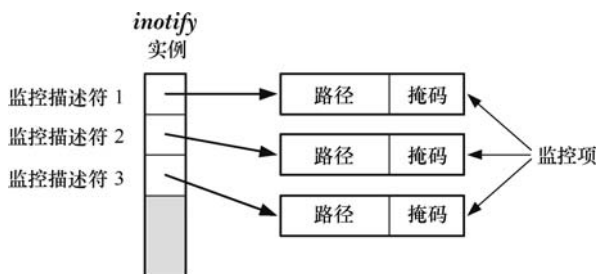


图 19-1: 一个 inotify 实例及与之相关的内核数据结构

参数 `mask` 为一位掩码，针对 `pathname` 定义了意欲监控的事件。稍后会论及可在掩码中指定的各种位值。

如果先前未将 `pathname` 加入 `fd` 的监控列表，那么 `inotify_add_watch()` 会在列表中创建一
新的监控项，并返回一新的、非负监控描述符，用来在后续操作中指代此监控项。对 `inotify` 实例
来说，该监控描述符是唯一的。

若先前已将 `pathname` 加入 `fd` 的监控列表，则 `inotify_add_watch()` 会修改现有 `pathname` 监控
项的掩码，并返回其监控描述符。(此描述符就是最初将 `pathname` 加入该监控列表的系统调用
`inotify_add_watch()` 所返回的监控描述符。) 下节在讨论 `IN_MASK_ADD` 标志时会就掩码的修
改过程做进一步描述。

系统调用 `inotify_rm_watch()` 会从文件描述符 `fd` 所指代的 `inotify` 实例中，删除由 `wd` 所定义
的监控项。

```
#include <sys/inotify.h>

int inotify_rm_watch(int fd, uint32_t wd);

Returns 0 on success, or -1 on error
```

参数 `wd` 是一监控描述符，由之前对 `inotify_add_watch()` 的调用返回。(`uint32_t` 数据类型
为一符号 32 位整数。)

删除监控项会为该监控描述符生成 `IN_IGNORED` 事件。稍后将讨论该事件。

19.3 inotify 事件

使用 `inotify_add_watch()` 删除或修改监控项时，位掩码参数 `mask` 标识了针对给定路径名
(`pathname`) 而要监控的事件。表 19-1 的“in”列列出了可在 `mask` 中定义的事件位。

表 19-1: inotify 事件

位 值	In	Out	描 述
<code>IN_ACCESS</code>	●	●	文件被访问 (<code>read()</code>)
<code>IN_ATTRIB</code>	●	●	文件元数据改变
<code>IN_CLOSE_WRITE</code>	●	●	关闭为了写入而打开的文件

位 值	In	Out	描 述
IN_CLOSE_NOWRITE	●	●	关闭以只读方式打开的文件
IN_CREATE	●	●	在受监控目录内创建了文件/目录
IN_DELETE	●	●	在受监控目录内删除文件/目录
IN_DELETE_SELF	●	●	删除受监控目录/文件本身
IN_MODIFY	●	●	文件被修改
IN_MOVE_SELF	●	●	移动受监控目录/文件本身
IN_MOVED_FROM	●	●	文件移出到受监控目录之外
IN_MOVED_TO	●	●	将文件移入受监控目录
IN_OPEN	●	●	文件被打开
IN_ALL_EVENTS	●		以上所有输出事件的统称
IN_MOVE	●		IN_MOVED_FROM IN_MOVED_TO 事件的统称
IN_CLOSE	●		IN_CLOSE_WRITE IN_CLOSE_NOWRITE 事件的统称
IN_DONT_FOLLOW	●		不对符号链接解引用（始于 Linux 2.6.15）
IN_MASK_ADD	●		将事件追加到 <code>pathname</code> 的当前监控掩码
IN_ONESHOT	●		只监控 <code>pathname</code> 的一个事件
IN_ONLYDIR	●		<code>pathname</code> 不为目录时会失败（始于 Linux 2.6.15）
IN_IGNORED		●	监控项为内核或应用程序所移除
IN_ISDIR		●	<code>name</code> 中所返回的文件名为路径
IN_Q_OVERFLOW		●	事件队列溢出
IN_UNMOUNT		●	包含对象的文件系统遭卸载

对于表 19-1 所列出的绝大多数位而言，顾名思义可知义。以下是对一些细节的澄清。

- 当文件的元数据（比如，权限、所有权、链接计数、扩展属性、用户 ID 或组 ID 等）改变时，会发生 `IN_ATTRIB` 事件。
- 删除受监控对象（即，一个文件或目录）时，发生 `IN_DELETE_SELF` 事件。当受监控对象是一个目录，并且该目录所含文件之一遭删除时，发生 `IN_DELETE` 事件。
- 重命名受监控对象时，发生 `IN_MOVE_SELF` 事件。重命名受监控目录内的对象时，发生 `IN_MOVED_FROM` 和 `IN_MOVED_TO` 事件。其中，前一事件针对包含旧对象名的目录，后一事件则针对包含新对象名的目录。
- `IN_DONT_FOLLOW`、`IN_MASK_ADD`、`IN_ONESHOT` 和 `IN_ONLYDIR` 位并非对监控事件的定义，而是意在控制 `inotify_add_watch()` 系统调用的行为。
- `IN_DONT_FOLLOW` 则规定，若 `pathname` 为符号链接，则不对其解引用。其作用在于令应用程序可以监控符号链接，而非符号连接所指代的文件。
- 倘若对已为同一 `inotify` 描述符所监控的同一路径名再次执行 `inotify_add_watch()` 调用，那么默认情况下会用给定的 `mask` 掩码来替换该监控项的当前掩码。如果指定了 `IN_MASK_ADD`，那么则会将 `mask` 值与当前掩码相或。
- `IN_ONESHOT` 允许应用只监控 `pathname` 的一个事件。事件发生后，监控项会自动从

监控列表中消失。

- 只有当 `pathname` 为目录时，`IN_ONLYDIR` 才允许应用程序对其进行监控。如果 `pathname` 并非目录，那么调用 `inotify_add_watch()` 失败，报错为 `ENOTDIR`。如要确保监控对象为一目录，则使用该标志可以规避竞争条件的发生。

19.4 读取 inotify 事件

将监控项在监控列表中登记后，应用程序可用 `read()` 从 `inotify` 文件描述符中读取事件，以判定发生了哪些事件。若时至读取时尚未发生任何事件，`read()` 会阻塞下去，直至有事件产生（除非对该文件描述符设置了 `O_NONBLOCK` 状态标志，这时若无任何事件可读，`read()` 将立即失败，并报错 `EAGAIN`）。

事件发生后，每次调用 `read()` 会返回一个缓冲区，内含一个或多个如下类型的结构（请见图 19-2）：

```
struct inotify_event {
    int      wd;          /* Watch descriptor on which event occurred */
    uint32_t mask;       /* Bits describing event that occurred */
    uint32_t cookie;     /* Cookie for related events (for rename()) */
    uint32_t len;        /* Size of 'name' field */
    char     name[];     /* Optional null-terminated filename */
};
```

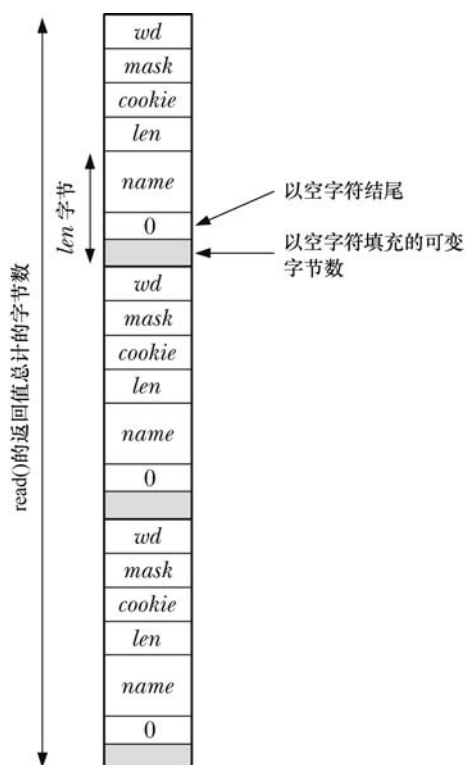


图 19-2: 包含 3 个 `inotify_event` 结构的输入缓冲区

字段 `wd` 指明发生事件的是那个监控描述符。该字段值由之前对 `inotify_add_watch()` 的调用返回。当应用程序要监控同一 `inotify` 文件描述符下的多个文件和目录时，字段 `wd` 就派上用

场。应用利用其所提供的线索来判定发生事件的特定文件或目录。（要做到这一点，应用程序必须维护专有数据结构，记录监控描述符与路径名之间的关系。）

`mask` 字段会返回描述该事件的位掩码。由表 19-1 所示的 `Out` 列展示了可出现于 `mask` 中的位范围。还要注意下列与特殊位相关的更多细节。

- 移除监控项时，会产生 `IN_IGNORED` 事件。起因可能有两个：其一，应用程序使用了 `inotify_rm_watch()` 系统调用显式移除监控项；其二，因受监控对象被删除或其所驻留的文件系统遭卸载，致使内核隐式删除监控项。以 `IN_ONESHOT` 而建立的监控项因事件触发而遭自动移除时，不会产生 `IN_IGNORED` 事件。
- 如果事件的主体为路径，那么除去其他位以外，在 `mask` 中还会设置 `IN_ISDIR` 位。
- `IN_UNMOUNT` 事件会通知应用程序包含受监控对象的文件系统已遭卸载。该事件发生之后，还会产生包含 `IN_IGNORED` 置位的附加事件。
- 19.5 节将介绍 `IN_Q_OVERFLOW`，并讨论对排队 `inotify` 事件的限制。

使用 `cookie` 字段可将相关事件联系在一起。目前，只有在对文件重命名时才会用到该字段。当这种情况发生时，系统会针对待重命名文件所在目录产生 `IN_MOVED_FROM` 事件，然后，还会针对重命名后文件的所在目录生成 `IN_MOVED_TO` 事件。（若仅是在同一目录内为文件改名，系统则会针对同一目录产生上述两个事件。）两个事件的 `cookie` 字段值相等，故而应用程序得以将它们关联起来。

当受监控目录中有文件发生事件时，`name` 字段返回一个以空字符结尾的字符串，以标识该文件。若受监控对象自身有事件发生，则不使用 `name` 字段，将 `len` 字段置 0。

`len` 字段用于表示实际分配给 `name` 字段的字节数。在 `read()` 所返回的缓冲区中，存储于 `name` 内的字符串结尾与下一个 `inotify_event` 结构的开始（请参见 19.2 节）之间，可能会有额外填充字节，故而 `len` 字段不可或缺。单个 `inotify` 事件的长度是 `sizeof(struct inotify_event)+ len`。

如果传递给 `read()` 的缓冲区过小，无法容纳下一个 `inotify_event` 结构，那么 `read()` 调用将以失败告终，并以 `EINVAL` 错误向应用程序报告这一情况。（在 2.6.21 之前版本的内核中，这种情况下 `read()` 将返回 0。在改为报告 `EINVAL` 错误之后，则对编程错误的提示更为清晰。）应用程序可再次以更大的缓冲区执行 `read()` 操作。然而，只要确保缓冲区足以容纳至少一个事件，这一问题将得以完全规避：传给 `read()` 的缓冲区应至少为 `sizeof(struct inotify_event)+ NAME_MAX + 1` 字节，其中 `NAME_MAX` 是文件名的最大长度，此外在加上终止空字符使用的 1 个字节。

采用的缓冲区大小如大于最小值，则可自单个 `read()` 中读取多个事件，效率极高。对 `inotify` 文件描述符所执行的 `read()`，将在已发生事件数量与缓冲区可容纳事件数量间取最小值并返回之。

针对文件描述符 `fd` 调用 `ioctl(fd, FIONREAD, &numbytes)`，会返回其所指代的 `inotify` 实例中的当前可读字节数。

从 `inotify` 文件描述符中读取的事件形成了一个有序队列。打个比方，这样一来，对文件重命名时，便可保证在 `IN_MOVED_TO` 事件之前能读取到 `IN_MOVED_FROM` 事件。

在事件队列的末尾追加一个新事件时，如果此新事件与队列当前的尾部事件拥有相同的 `wd`、`mask`、`cookie` 和 `mask` 值，那么内核会将两者合并（以避免对新事件排队）。之所以这么做，是因为很多应用程序都并不关注同一事件的反复出现，而丢弃多余的事件能降低内核维护事件队列所需的内存总量。然而，这也意味着使用 `inotify` 将无法可靠判定出周期性事件的发生次数或频率。

程序示例

虽然在前文中描述了 inotify API 的诸多细节，但实际上，该 API 使用起来却颇为简单。程序清单 19-1 展示了对 inotify 的运用。

程序清单 19-1：运用 inotify API

```
----- inotify/demo_inotify.c
#include <sys/inotify.h>
#include <limits.h>
#include "tspi_hdr.h"

static void          /* Display information from inotify_event structure */
displayInotifyEvent(struct inotify_event *i)
{
    printf("    wd =%2d; ", i->wd);
    if (i->cookie > 0)
        printf("cookie =%4d; ", i->cookie);

    printf("mask = ");
    if (i->mask & IN_ACCESS)        printf("IN_ACCESS ");
    if (i->mask & IN_ATTRIB)       printf("IN_ATTRIB ");
    if (i->mask & IN_CLOSE_NOWRITE) printf("IN_CLOSE_NOWRITE ");
    if (i->mask & IN_CLOSE_WRITE)  printf("IN_CLOSE_WRITE ");
    if (i->mask & IN_CREATE)       printf("IN_CREATE ");
    if (i->mask & IN_DELETE)       printf("IN_DELETE ");
    if (i->mask & IN_DELETE_SELF)  printf("IN_DELETE_SELF ");
    if (i->mask & IN_IGNORED)      printf("IN_IGNORED ");
    if (i->mask & IN_ISDIR)        printf("IN_ISDIR ");
    if (i->mask & IN_MODIFY)       printf("IN_MODIFY ");
    if (i->mask & IN_MOVE_SELF)    printf("IN_MOVE_SELF ");
    if (i->mask & IN_MOVED_FROM)   printf("IN_MOVED_FROM ");
    if (i->mask & IN_MOVED_TO)    printf("IN_MOVED_TO ");
    if (i->mask & IN_OPEN)         printf("IN_OPEN ");
    if (i->mask & IN_Q_OVERFLOW)   printf("IN_Q_OVERFLOW ");
    if (i->mask & IN_UNMOUNT)     printf("IN_UNMOUNT ");
    printf("\n");

    if (i->len > 0)
        printf("        name = %s\n", i->name);
}
#define BUF_LEN (10 * (sizeof(struct inotify_event) + NAME_MAX + 1))

int
main(int argc, char *argv[])
{
    int inotifyFd, wd, j;
    char buf[BUF_LEN];
    ssize_t numRead;
    char *p;
    struct inotify_event *event;

    if (argc < 2 || strcmp(argv[1], "--help") == 0)
        usageErr("%s pathname... \n", argv[0]);

    ① inotifyFd = inotify_init();          /* Create inotify instance */
    if (inotifyFd == -1)
```

```

    errExit("inotify_init");

    for (j = 1; j < argc; j++) {
②      wd = inotify_add_watch(inotifyFd, argv[j], IN_ALL_EVENTS);
        if (wd == -1)
            errExit("inotify_add_watch");

        printf("Watching %s using wd %d\n", argv[j], wd);
    }

    for (;;) {
③      numRead = read(inotifyFd, buf, BUF_LEN);
        if (numRead == 0)
            fatal("read() from inotify fd returned 0!");

        if (numRead == -1)
            errExit("read");

        printf("Read %ld bytes from inotify fd\n", (long) numRead);

        /* Process all of the events in buffer returned by read() */

        for (p = buf; p < buf + numRead; ) {
④      event = (struct inotify_event *) p;
            displayInotifyEvent(event);

            p += sizeof(struct inotify_event) + event->len;
        }
    }

    exit(EXIT_SUCCESS);
}

```

inotify/demo_inotify.c

程序清单 19-1 中程序将执行以下步骤。

- 使用 `inotify_init()`，创建 `inotify` 文件描述符①。
- 使用 `inotify_add_watch()`，将程序命令行参数中指定的每个文件加入监控项②。每个监控项都将监控所有可能发生的事件。
- 执行无限循环。
 - 从 `inotify` 描述符读取事件缓冲区③。
 - 调用 `displayInotifyEvent()` 函数，以显示上述缓冲区中各 `inotify_event` 结构的内容④。

以下 shell 会话演示了对程序清单 19-1 所列程序的使用。首先，在后台运行该程序的实例，对两个目录进行监控。

```

$ ./demo_inotify dir1 dir2 &
[1] 5386
Watching dir1 using wd 1
Watching dir2 using wd 2

```

然后，执行某些命令，从而在两个目录中产生事件。先使用 `cat(1)` 创建一个文件：

```

$ cat > dir1/aaa
Read 64 bytes from inotify fd
  wd = 1; mask = IN_CREATE
  name = aaa
  wd = 1; mask = IN_OPEN
  name = aaa

```

由后台程序所生成的上述输出表明，`read()`读取了包含两个事件的缓冲区。继续在该文件中执行某些输入操作，然后输入 `end-of-file` 字符串：

```
Hello world
Read 32 bytes from inotify fd
  wd = 1; mask = IN_MODIFY
  name = aaa
Type Control-D
Read 32 bytes from inotify fd
  wd = 1; mask = IN_CLOSE_WRITE
  name = aaa
```

接下来，将该文件转移至另一个受监控的目录，同时对其重新命名。这会产生两个事件，一个对应于文件的源目录（监控描述符 1），另一个对应于文件的目标目录（监控描述符 2）。

```
$ mv dir1/aaa dir2/bbb
Read 64 bytes from inotify fd
  wd = 1; cookie = 548; mask = IN_MOVED_FROM
  name = aaa
  wd = 2; cookie = 548; mask = IN_MOVED_TO
  name = bbb
```

以上两个事件共享相同的 `cookie` 值，允许应用程序将它们联系起来。

当在其中一个受监控目录下创建子目录时，由此产生的事件掩码会置 `IN_ISDIR` 位，以示该事件的对象是一目录。

```
$ mkdir dir2/ddd
Read 32 bytes from inotify fd
  wd = 1; mask = IN_CREATE IN_ISDIR
  name = ddd
```

此处，再次提醒大家，`inotify` 监控是非递归的。如果应用程序有意对新创建的子目录进行监控，则需进一步执行 `inotify_add_watch()` 系统调用，并指明子目录的路径名。

最后，将其中一个受监控目录删除：

```
$ rmdir dir1
Read 32 bytes from inotify fd
  wd = 1; mask = IN_DELETE_SELF
  wd = 1; mask = IN_IGNORED
```

系统会生成最后一个事件，以通知应用程序，内核已从监控列表中删除该监控项。

19.5 队列限制和/proc 文件

对 `inotify` 事件做排队处理，需要消耗内核内存。正因如此，内核会对 `inotify` 机制的操作施以各种限制。超级用户可配置 `/proc/sys/fs/inotify` 路径中的 3 个文件来调整这些限制：

`max_queued_events`

调用 `inotify_init()` 时，使用该值来为新 `inotify` 实例队列中的事件数量设置上限。一旦超出这一上限，系统将生成 `IN_Q_OVERFLOW` 事件，并丢弃多余的事件。溢出事件的 `wd` 字段值为 -1。

`max_user_instances`

对由每个真实用户 ID 创建的 `inotify` 实例数的限制值。

`max_user_watches`

对由每个真实用户 ID 创建的监控项数量的限制值。

这 3 个文件的典型默认值分别为 16384、128 和 8192。

19.6 监控文件的旧有系统：dnotify

Linux 还为监控文件事件提供了另一种机制。该机制名为 `dnotify`，问世于内核 2.4 版本，但 `inotify` 已令其“落伍”。相较于 `inotify`，`dnotify` 存在如下局限性。

- `dnotify` 机制通过向应用程序发送信号来通告事件。使用信号作为通告机制，会使应用程序的设计复杂化（请参见 22.12 节）。这也使得在函数库中使用 `dnotify` 变得困难，因为调用该函数的程序可能会改变对通告信号的处置（disposition）。而 `inotify` 机制则不使用信号。
- `dnotify` 的监控单元为目录。只要对该目录下的任一文件执行了任何操作，系统都会通知应用程序。相形之下，`inotify` 的监控对象则既可以是单个文件，也能是目录。
- 为监控目录，`dnotify` 需要应用程序为该目录打开文件描述符。使用文件描述符会导致两个问题。其一，由于程序处于运行中，将无法卸载包含此目录的文件系统。其二，因为每个目录都需要一个文件描述符，所以应用程序最终可能会消耗大量文件描述符。而 `inotify` 不使用文件描述符，故而可以避免上述问题。
- 与 `inotify` 相比，由 `dnotify` 提供的与文件事件相关的信息不够精确。当位于受监控目录下的文件发生改变时，`dnotify` 只会通知有事件发生，但不会说明事件具体涉及了哪个文件。因此，应用程序必须通过缓存目录内容来进行判断。此外，针对已发生事件的类型，`inotify` 所提供的信息也比 `dnotify` 更详细。
- 在某些情况下，`dnotify` 不支持可靠的文件事件通告机制。

在 `fcntl(2)` 手册页中对 `F_NOTIFY` 操作的描述部分，以及内核源码 `Documentation/dnotify.txt` 中，都可找到有关 `dnotify` 的更多信息。

19.7 总结

当一组受监控的文件或目录有事件发生（对文件的打开、关闭、创建、删除、修改以及重命名等操作）时，Linux 专有的 `inotify` 机制可让应用程序获得通知。`inotify` 机制取代了较老的 `dnotify` 机制。

19.8 练习

- 19-1. 编写一个程序，针对其命令行参数所指定的目录，记录所有的文件创建、删除和改名操作。该程序应能够监控指定目录下所有子目录中的事件。获得所有子目录的列表需使用 `nftw()`（参见 18.9 节）。当在目录树下添加或删除了子目录时，受监控的子目录集合应能保持同步更新。

第 20 章

信号：基本概念

本章和接下来的两章将讨论信号。虽然基本概念较为简单，但因为要涵盖大量细节，所以篇幅较长。

本章包括以下主题。

- 各种不同信号及其用途。
- 内核可能为进程产生信号的环境，以及某一进程向另一进程发送信号所使用的系统调用。
- 进程在默认情况下对信号的响应方式，以及进程改变对信号响应方式的手段，特别是借助于信号处理器程序的手段，即程序收到信号时去自动调用的函数，由程序员定义。
- 使用进程信号掩码来阻塞信号，以及等待信号的相关概念。
- 如何暂停进程的执行，并等待信号的到达。

20.1 概念和概述

信号是事件发生时对进程的通知机制。有时也称之为软件中断。信号与硬件中断的相似之处在于打断了程序执行的正常流程，大多数情况下，无法预测信号到达的精确时间。

一个（具有合适权限的）进程能够向另一进程发送信号。信号的这一用法可作为一种同步技术，甚至是进程间通信（IPC）的原始形式。进程也可以向自身发送信号。然而，发往进程的诸多信号，通常都是源于内核。引发内核为进程产生信号的各类事件如下。

- 硬件发生异常，即硬件检测到一个错误条件并通知内核，随即再由内核发送相应信号给相关进程。硬件异常的例子包括执行一条异常的机器语言指令，诸如，被 0 除，或者引用了无法访问的内存区域。
- 用户键入了能够产生信号的终端特殊字符。其中包括中断字符（通常是 Control-C）、暂停字符（通常是 Control-Z）。
- 发生了软件事件。例如，针对文件描述符的输出变为有效，调整了终端窗口大小，定时器到期，进程执行的 CPU 时间超限，或者该进程的某个子进程退出。

针对每个信号，都定义了一个唯一的（小）整数，从 1 开始顺序展开。<signal.h>以 SIGxxxx 形式的符号名对这些整数做了定义。由于每个信号的实际编号随系统不同而不同，所以在程序中

总是使用这些符号名。例如，当用户键入中断字符时，将传递给进程 SIGINT 信号（信号编号为 2）。

信号分为两大类。第一组用于内核向进程通知事件，构成所谓传统或者标准信号。Linux 中标准信号的编号范围为 1~31。本章将描述这些标准信号。另一组信号由实时信号构成，其与标准信号的差异将在 22.8 节中描述。

信号因某些事件而产生。信号产生后，会于稍后被传递给某一进程，而进程也会采取某些措施来响应信号。在产生和到达期间，信号处于等待（pending）状态。

通常，一旦（内核）接下来要调度该进程运行，等待信号会马上送达，或者如果进程正在运行，则会立即传递信号（例如，进程向自身发送信号）。然而，有时需要确保一段代码不为传递来的信号所中断。为了做到这一点，可以将信号添加到进程的信号掩码中——目前会阻塞该组信号的到达。如果所产生的信号属于阻塞之列，那么信号将保持等待状态，直至稍后对其解除阻塞（从信号掩码中移除）。进程可使用各种系统调用对其信号掩码添加和移除信号。

信号到达后，进程视具体信号执行如下默认操作之一。

- 忽略信号：也就是说，内核将信号丢弃，信号对进程没有产生任何影响（进程永远都不知道曾经出现过该信号）。
- 终止（杀死）进程：这有时是指进程异常终止，而不是进程因调用 `exit()` 而发生的正常终止。
- 产生核心转储文件，同时进程终止：核心转储文件包含对进程虚拟内存的镜像，可将其加载到调试器中以检查进程终止时的状态。
- 停止进程：暂停进程的执行。
- 于之前暂停后再度恢复进程的执行。

除了根据特定信号而采取默认行为之外，程序也能改变信号到达时的响应行为。也将此称之为对信号的处置（disposition）设置。程序可以将对信号的处置设置为如下之一。

- 采取默认行为。这适用于撤销之前对信号处置的修改、恢复其默认处置的场景。
- 忽略信号。这适用于默认行为为终止进程的信号。
- 执行信号处理器程序。

信号处理器程序是由程序员编写的函数，用于为响应传递来的信号而执行适当任务。例如，shell 为 SIGINT 信号（由中断字符串 Control-C 产生）提供了一个处理器程序，令其停止当前正在执行的工作，并将控制返回到（shell 的）主输入循环，并再次向用户呈现 shell 提示符。通知内核应当去调用某一处理器程序的行为，通常称之为安装或者建立信号处理器程序。调用信号处理器程序以响应传递来的信号，则称之为信号已处理（handled），或者已捕获（caught）。

请注意，无法将信号处置设置为终止进程或者转储核心（除非这是对信号的默认处置）。效果最为近似的是为信号安装一个处理器程序，并于其中调用 `exit()` 或者 `abort()`。`abort()` 函数（21.2.2 节）为进程产生一个 SIGABRT 信号，该信号将引发进程转储核心文件并终止。

Linux 特有的 `/proc/PID/status` 文件包含有各种位掩码字段，通过检查这些掩码可以确定进程对信号的处理。位掩码以十六进制数形式显示，最低有效位代表信号 1，相临的左边一位代表信号 2，以此类推。这些字段分别为 `SigPnd`（基于线程的等待信号）、`ShdPnd`（进程级等待信号，始于 Linux 2.6）、`SigBlk`（阻塞信号）、`SigIgn`（忽略信号）和 `SigCgt`（捕获信号）。（33.2 节阐述了多线程进程对信号的处理，这将有助于澄清 `SigPnd` 与 `ShdPnd` 之间的差异。）使用 `ps(1)` 命令的各种选项也能获得相同信息。

信号在 UNIX 实现中出现很早，诞生之后又历经变革。在早期实现中，信号在特定场景下有可能会丢失（即，没有传递到目标进程）。此外，尽管系统提供了执行关键代码时阻塞信号传递的机制，但阻塞有时也不大可靠。4.2BSD 利用所谓可靠信号解决了这些问题。（BSD 在创新上还更进一步，增加了额外信号来支持 shell 作业控制，请参考 34.7 节。）

System V 后来也为信号增加了可靠语义，但采用的模型与 BSD 无法兼容。这一不兼容性直到 POSIX.1-1990 标准出台后才得以解决。该标准针对可靠信号所采取规范主要基于 BSD 模型。

22.7 节将就可靠和不可靠信号的细节展开讨论，22.13 节则简要说明了老版 BSD 和 System V 的信号 API。

20.2 信号类型和默认行为

此前曾提及，Linux 对标准信号的编号为 1~31。然而，Linux 于 `signal(7)` 手册页中列出的信号名称却超出了 31 个。名称超出的原因有多种。有些名称只是其他名称的同义词，之所以定义是为了与其他 UNIX 实现保持源码兼容。其他名称虽然有定义，但却并未使用。以下列表介绍了各种信号。

SIGABRT

当进程调用 `abort()` 函数（21.2.2 节）时，系统向进程发送该信号。默认情况下，该信号会终止进程，并产生核心转储文件。这实现了调用 `abort()` 的预期目标，产生核心转储文件用于调试。

SIGALRM

经调用 `alarm()` 或 `setitimer()` 而设置的实时定时器一旦到期，内核将产生该信号。实时定时器是根据挂钟时间进行计时的（即人类对逝去时间的概念）。更多细节参见 23.1 节。

SIGBUS

产生该信号（总线错误，bus error）即表示发生了某种内存访问错误。如 49.4.3 节所述，当使用由 `mmap()` 所创建的内存映射时，如果试图访问的地址超出了底层内存映射文件的结尾，那么将产生该错误。

SIGCHLD

当父进程的某一子进程终止（或者因为调用了 `exit()`，或者因为被信号杀死）时，（内核）将向父进程发送该信号。当父进程的某一子进程因收到信号而停止或恢复时，也可能向父进程发送该信号。详情请参考 26.3 节。

SIGCLD

与 SIGCHLD 信号同义。

SIGCONT

将该信号发送给已停止的进程，进程将会恢复运行（即在之后某个时间点重新获得调度）。当接收信号的进程当前不处于停止状态时，默认情况下将忽略该信号。进程可以捕获该信号，以便在恢复运行时可以执行某些操作。关于该信号的更多细节请参考 22.2 节和 34.7 节。

SIGEMT

UNIX 系统通常用该信号来标识一个依赖于实现的硬件错误。Linux 系统仅在 Sun SPARC

实现中使用了该信号。后缀 EMT 源自仿真器陷阱 (emulator trap)，Digital PDP-11 的汇编程序助记符之一。

SIGFPE

该信号因特定类型的算术错误而产生，比如除以 0。后缀 FPE 是浮点异常的缩写，不过整型算术错误也能产生该信号。该信号于何时产生的精确细节取决于硬件架构和对 CPU 控制寄存器的设置。例如，在 x86-32 架构中，整数除以 0 总是产生 SIGFPE 信号，但是对浮点数除以 0 的处理则取决于是否启用了 FE_DIVBYZERO 异常。如果启用了该异常（使用 `feenableexcept()`），那么浮点数除以 0 也将产生 SIGFPE 信号，否则，将为操作数产生符合 IEEE 标准的结果（无穷大的浮点表示形式）。更多信息请参考 `fenv(3)` 手册页和 `<fenv.h>` 文件。

SIGHUP

当终端断开（挂机）时，将发送该信号给终端控制进程。34.6 节将描述控制进程的概念以及产生 SIGHUP 信号的各种环境。SIGHUP 信号还可用于守护进程（比如，`init`、`httpd` 和 `inetd`）。许多守护进程会在收到 SIGHUP 信号时重新进行初始化并重读配置文件。借助于显式执行 `kill` 命令或者运行同等功效的程序或脚本，系统管理员可向守护进程手工发送 SIGHUP 信号来触发这些行为。

SIGILL

如果进程试图执行非法（即格式不正确）的机器语言指令，系统将向进程发送该信号。

SIGINFO

在 Linux 中，该信号名与 SIGPWR 信号名同义。在 BSD 系统中，键入 Control-T 可产生 SIGINFO 信号，用于获取前台进程组的状态信息。

SIGINT

当用户键入终端中断字符（通常为 Control-C）时，终端驱动程序将发送该信号给前台进程组。该信号的默认行为是终止进程。

SIGIO

利用 `fcntl()` 系统调用，即可于特定类型（诸如终端和套接字）的打开文件描述符发生 I/O 事件时产生该信号。63.3 节将就此特性做进一步说明。

SIGIOT

在 Linux 中，该信号名与 SIGABRT 信号同义。在其他一些 UNIX 实现中，该信号表示发生了由实现定义的硬件错误。

SIGKILL

此信号为“必杀 (sure kill)”信号，处理器程序无法将其阻塞、忽略或者捕获，故而“一击必杀”，总能终止进程。

SIGLOST

Linux 中存在该信号名，但并未加以使用。在其他一些 UNIX 实现中，如果远端 NFS 服务器在崩溃之后重新恢复，而 NFS 客户端却未能重新获得由本地进程所持有的锁，那么 NFS 客户端将向这些进程发送此信号。（NFS 规范并未对该特性进行标准化。）

SIGPIPE

当某一进程试图向管道、FIFO 或套接字写入信息时，如果这些设备并无相应的阅读进程，

那么系统将产生该信号。之所以如此，通常是因为阅读进程已经关闭了其作为 IPC 通道的文件描述符。更多细节请参考 44.2 节。

SIGPOLL

该信号从 System V 派生而来，与 Linux 中的 SIGIO 信号同义。

SIGPROF

由 `setitimer()` 调用（参见 23.1 节）所设置的性能分析定时器刚一过期，内核就将产生该信号。性能分析定时器用于记录进程所使用的 CPU 时间。与虚拟定时器不同（参见下面的 SIGVTALRM 信号），性能分析定时器在对 CPU 时间计数时会将用户态与内核态都包含在内。

SIGPWR

这是电源故障信号。当系统配备有不间断电源（UPS）时，可以设置守护进程来监控电源发生故障时备用电池的剩余电量。如果电池电量行将耗尽（长时间停电之后），那么监控进程会将该信号发往 `init` 进程，而后者则将其解读为快速、有序关闭系统的一个请求。

SIGQUIT

当用户在键盘上键入退出字符（通常为 `Control-\`）时，该信号将发往前台进程组。默认情况下，该信号终止进程，并生成可用于调试的核心转储文件。进程如果陷入无限循环，或者不再响应时，使用 SIGQUIT 信号就很合适。键入 `Control-\`，再调用 `gdb` 调试器加载刚才生成的核心转储文件，接着用 `backtrace` 命令来获取堆栈跟踪信息，就能发现正在执行的是程序的哪部分代码。（[Matloff, 2008]描述了 `gdb` 的用法。）

SIGSEGV

这一信号非常常见，当应用程序对内存的引用无效时，就会产生该信号。引起对内存无效引用的原因很多，可能是因为要引用的页不存在（例如，该页位于堆和栈之间的未映射区域），或者进程试图更新只读内存（比如，程序文本段或者标记为只读的一块映射内存区域）中某一位置的内容，又或者进程企图在用户态（参见 2.1 节）去访问内核的部分内存。C 语言中引发这些事件的往往是解引用的指针里包含了错误地址（例如，未初始化的指针），或者传递了一个无效参数供函数调用。该信号的命名源于术语“段违例”。

SIGSTKFLT

`signal(7)` 手册页中将其记载为“协处理器栈错误”，Linux 对该信号作了定义，但并未加以使用。

SIGSTOP

这是一个必停（`sure stop`）信号，处理器程序无法将其阻塞、忽略或者捕获，故而总是能停止进程。

SIGSYS

如果进程发起的系统调用有误，那么将产生该信号。这意味着系统将进程执行的指令视为一个系统调用陷阱（`trap`），但相关的系统调用编号却是无效的（参见 3.1 节）。

SIGTERM

这是用来终止进程的标准信号，也是 `kill` 和 `killall` 命令所发送的默认信号。用户有时会使用 `kill-KILL` 或者 `kill-9` 显式向进程发送 SIGKILL 信号。然而，这一做法通常是错误的。精心设计的应用程序应当为 SIGTERM 信号设置处理器程序，以便于其能够预先清除临时文件和释

放其他资源，从而全身而退。发送 SIGKILL 信号可以杀掉某个进程，从而绕开了 SIGTERM 信号的处理程序。因此，总是应该首先尝试使用 SIGTERM 信号来终止进程，而把 SIGKILL 信号作为最后手段，去对付那些不响应 SIGTERM 信号的失控进程。

SIGTRAP

该信号用来实现断点调试功能以及 `strace(1)` 命令（附录 A）所执行的跟踪系统调用功能。更多信息参见 `ptrace(2)` 手册页。

SIGTSTP

这是作业控制的停止信号，当用户在键盘上输入挂起字符（通常是 Control-Z）时，将发送该信号给前台进程组，使其停止运行。第 34 章详细描述了进程组（作业）和作业控制，以及程序应在何时以及如何去处理该信号。该信号名源自“终端停止（terminal stop）”的术语。

SIGTTIN

在作业控制 shell 下运行时，若后台进程组试图对终端进行 `read()` 操作，终端驱动程序则将向该进程组发送此信号。该信号默认将停止进程。

SIGTTOU

该信号的目的与 SIGTTIN 信号类似，但所针对的是后台作业的终端输出。在作业控制 shell 下运行时，如果对终端启用了 TOSTOP（终端输出停止）选项（可能是通过 `stty tostop` 命令），而某一后台进程组试图对终端进行 `write()` 操作（参见 34.7.1 节），那么终端驱动程序将向该进程组发送 SIGTTOU 信号。该信号默认将停止进程。

SIGUNUSED

顾名思义，该信号没有使用。在 Linux 2.4 及其后续版本中，该信号名在很多架构中与 SIGSYS 信号同义。换言之，尽管信号名还保持向后兼容，但信号编号在这些架构中不再处于未使用状态。

SIGURG

系统发送该信号给一个进程，表示套接字上存在带外（也称作紧急）数据（参见 61.13.1 节）。

SIGUSR1

该信号和 SIGUSR2 信号供程序员自定义使用。内核绝不会为进程产生这些信号。进程可以使用这些信号来相互通知事件的发生，或是彼此同步。在早期的 UNIX 实现中，这是可供应用随意使用的仅有的两个信号。（实际上，进程间可以相互发送任何信号，但如果内核也为进程产生了同类信号，这两种情况就有可能产生混淆。）现代 UNIX 实现则提供了很多实时信号，也可用于程序员自定义的目的（参见 22.8 节）。

SIGUSR2

参见对 SIGUSR1 信号的描述。

SIGVTALRM

调用 `setitimer()`（参见 23.1 节）设置的虚拟定时器刚一到期，内核就会产生该信号。虚拟定时器记录的是进程在用户态所使用的 CPU 时间。

SIGWINCH

在窗口环境中，当终端窗口尺寸发生变化时（如 62.9 节所述，要么是由于用户手动调

整了大小，要么是因为程序调用 `ioctl()` 对大小做了调整)，会向前台进程组发送该信号。借助于为该信号安装的处理程序，诸如 `vi` 和 `less` 之类的程序会在窗口尺寸调整后重新绘制输出。

SIGXCPU

当进程的 CPU 时间超出对应的资源限制时（参见 36.3 节对 `RLIMIT_CPU` 的描述），将发送此信号给进程。

SIGXFSZ

如果进程因试图增大文件（调用 `write()` 或 `truncate()`）而突破对进程文件大小的资源限制（参见 36.3 节对 `RLIMIT_FSIZE` 的描述）时，那么将发送此信号给进程。

表 20-1 总结了 Linux 下与信号相关的一系列信息。关于此表，请注意以下几点。

- 信号编号列所示为在不同硬件架构下对信号的编号。除非另有说明，信号在所有架构中编号相同。信号编号在架构上的差异会在括号中予以说明，所涉及的架构包括 Sun SPARC、SPARC64 (S)、HP/Compaq/Digital Alpha (A)、MIPS (M) 和 HP PA-RISC (P)。此列中的 `undef` 表示此符号在所示架构中未定义。
- `SUSv3` 列则表示 SUSv3 是否定义了该信号。
- 默认列显示了信号的默认行为。`term` 表示信号终止进程，`core` 表示进程产生核心转储文件并退出，`ignore` 表示忽略该信号，`stop` 表示信号停止了进程，`cont` 表示信号恢复了一个已停止的进程。

某些前面列出的信号并未见诸于表 20-1，如 `SIGCLD`（`SIGCHLD` 信号的同义词）、`SIGINFO`（未使用）、`SIGIOT`（`SIGABRT` 信号的同义词）、`SIGLOST`（未使用）和 `SIGUNUSED`（在许多架构中是 `SIGSYS` 信号的同义词）。

表 20-1: Linux 信号

名称	信号值	描述	SUSv3	默认
<code>SIGABRT</code>	6	中止进程	•	<code>core</code>
<code>SIGALRM</code>	14	实时定时器过期	•	<code>term</code>
<code>SIGBUS</code>	7 (SAMP=10)	内存访问错误	•	<code>core</code>
<code>SIGCHLD</code>	17 (SA=20, MP=18)	终止或者停止子进程	•	<code>ignore</code>
<code>SIGCONT</code>	18 (SA=19, M=25, P=26)	若停止则继续	•	<code>cont</code>
<code>SIGEMT</code>	<code>undef</code> (SAMP=7)	硬件错误		<code>term</code>
<code>SIGFPE</code>	8	算术异常	•	<code>core</code>
<code>SIGHUP</code>	1	挂起	•	<code>term</code>
<code>SIGILL</code>	4	非法指令	•	<code>core</code>
<code>SIGINT</code>	2	终端中断	•	<code>term</code>
<code>SIGIO /</code>	29 (SA=23, MP=22)	I/O 时可能产生	•	<code>term</code>
<code>SIGPOLL</code>				

续表

名 称	信 号 值	描 述	SUSv3	默认
SIGKILL	9	必杀（确保杀死）	●	term
SIGPIPE	13	管道断开	●	term
SIGPROF	27 (M=29, P=21)	性能分析定时器过期	●	term
SIGPWR	30 (SA=29, MP=19)	电量行将耗尽		term
SIGQUIT	3	终端退出	●	core
SIGSEGV	11	无效的内存引用	●	core
SIGSTKFLT	16 (SAM=undef, P=36)	协处理器栈错误		term
SIGSTOP	19 (SA=17, M=23, P=24)	确保停止	●	stop
SIGSYS	31 (SAMP=12)	无效的系统调用	●	core
SIGTERM	15	终止进程	●	term
SIGTRAP	5	跟踪/断点陷阱	●	core
SIGTSTP	20 (SA=18, M=24, P=25)	终端停止	●	stop
SIGTTIN	21 (M=26, P=27)	BG ¹ 从终端读取	●	stop
SIGTTOU	22 (M=27, P=28)	BG 向终端写	●	stop
SIGURG	23 (SA=16, M=21, P=29)	套接字上的紧急数据	●	ignore
SIGUSR1	10 (SA=30, MP=16)	用户自定义信号 1	●	term
SIGUSR2	12 (SA=31, MP=17)	用户自定义信号 2	●	term
SIGVTALRM	26 (M=28, P=20)	虚拟定时器过期	●	term
SIGWINCH	28 (M=20, P=23)	终端窗口尺寸发生变化		ignore
SIGXCPU	24 (M=30, P=33)	突破对 CPU 时间的限制	●	core
SIGXFSZ	25 (M=31, P=34)	突破对文件大小的限制	●	core

针对表 20-1 中某些信号的默认行为，要注意以下几点。

- 在 Linux 2.2 中，信号 SIGXCPU、SIGXFSZ、SIGSYS 和 SIGBUS 的默认行为是终止进程，但不会产生核心转储文件。自内核 2.4 以后，Linux 实现满足了 SUSv3 的要求，这些信号不但会引发进程终止，也将生成核心转储文件。在其他几个 UNIX 实现中，对信号 SIGXCPU 和 SIGXFSZ 的处理方式与 Linux 2.2 相同。
- 在其他的 UNIX 实现中，对 SIGPWR 信号的默认行为通常是将其忽略。
- 几个 UNIX 实现（特别是 BSD 衍生系统）默认情况下将忽略 SIGIO 信号。
- 虽然 SIGEMT 信号尚未获得任何标准的接纳，但却得到大多数 UNIX 实现的支持。然而，在其他实现中，该信号通常会导致进程终止并产生核心转储文件。
- SUSv1 将 SIGURG 信号的默认行为定义为终止进程，这也是一些较老 UNIX 实现的默认做法。而 SUSv2 则采用了现行规范（将其忽略）。

¹ 译者注：即后台进程组。

20.3 改变信号处置: signal()

UNIX 系统提供了两种方法来改变信号处置: `signal()`和 `sigaction()`。本节描述的 `signal()`系统调用, 是设置信号处置的原始 API, 所提供的接口比 `sigaction()`简单。另一方面, `sigaction()`提供了 `signal()`所不具备的功能。进一步而言, `signal()`的行为在不同 UNIX 实现间存在差异 (22.7 节), 这也意味着对可移植性有所追求的程序绝不能使用此调用来建立信号处理器函数。故此, `sigaction()`是建立信号处理器的首选 API (强力推荐)。自 20.13 节介绍了 `sigaction()`调用的用法之后, 本书示例将一律采用该调用来建立信号处理器程序。

`signal()`函数虽然记录在 Linux 手册页的第 2 部分, 但实际却被实现为一个基于 `sigaction()`系统调用的 `glibc` 库函数。

```
#include <signal.h>

void ( *signal(int sig, void (*handler)(int)) ) (int);

Returns previous signal disposition on success, or SIG_ERR on error
```

这里需要对 `signal()`函数的原型做一些解释。第一个参数 `sig`, 标识希望修改处置的信号编号, 第二个参数 `handler`, 则标识信号抵达时所调用函数的地址。该函数无返回值 (`void`), 并接收一个整型参数。因此, 信号处理器函数一般具有以下形式:

```
void
handler(int sig)
{
    /* Code for the handler */
}
```

20.4 节将描述处理器函数中 `sig` 参数的目的。

`signal()`的返回值是之前的信号处置。像 `handler` 参数一样, 这是一枚指针, 所指向的是带有一个整型参数且无返回值的函数。换言之, 编写如下代码, 可以暂时为信号建立一个处理器函数, 然后再将信号处置重置为其本来面目:

```
void (*oldHandler)(int);

oldHandler = signal(SIGINT, newHandler);
if (oldHandler == SIG_ERR)
    errExit("signal");

/* Do something else here. During this time, if SIGINT is
   delivered, newHandler will be used to handle the signal. */

if (signal(SIGINT, oldHandler) == SIG_ERR)
    errExit("signal");
```

使用 `signal()`, 将无法在不改变信号处置的同时, 还能获取到当前的信号处置。要想做到这一点, 必须使用 `sigaction()`。

针对信号处理器函数指针做如下类型定义, 将有助于理解 `signal()`的原型:

```
typedef void (*sig_handler_t)(int);
```

signal()原型可以改写成如下形式:

```
sig_handler_t signal(int sig, sig_handler_t handler);
```

如果定义了_GNU_SOURCE 特性测试宏, 那么 glibc 将在<signal.h>头文件中暴露非标准的 sig_handler_t 数据类型。

在为 signal()指定 handler 参数时, 可以以如下值来代替函数地址:

SIG_DFL

将信号处置重置为默认值 (表 20-1)。这适用于将之前 signal()调用所改变的信号处置还原。

SIG_IGN

忽略该信号。如果信号专为此进程而生, 那么内核会默默将其丢弃。进程甚至从未知道曾经产生了该信号。

调用 signal()成功将返回先前的信号处置, 有可能是先前安装的处理器函数地址, 也可能是常量 SIG_DFL 和 SIG_IGN 之一。如果调用失败, signal()将返回 SIG_ERR。

20.4 信号处理器简介

信号处理器程序 (也称为信号捕捉器) 是当指定信号传递给进程时将会调用的一个函数。本节描述了信号处理器的基本原理, 而第 21 章将继续做详细介绍。

调用信号处理器程序, 可能会随时打断主程序流程; 内核代表进程来调用处理器程序, 当处理器返回时, 主程序会在处理器打断的位置恢复执行。这一工作序列可用图 20-1 来加以说明。

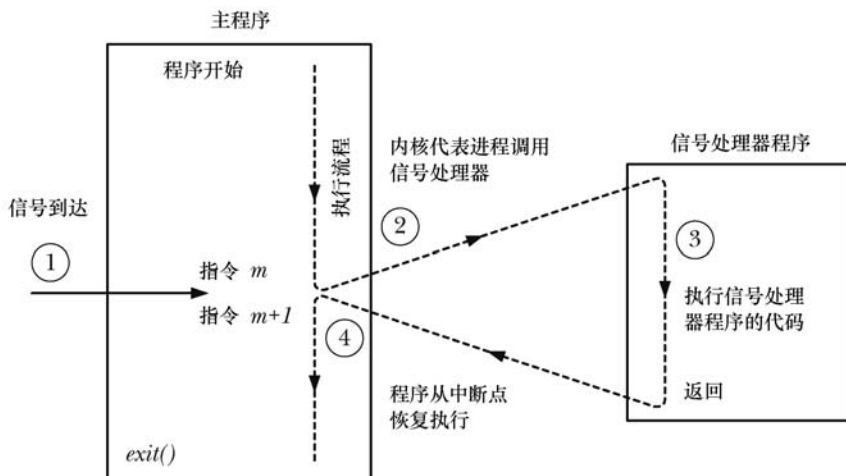


图 20-1: 信号到达并执行处理器程序

虽然信号处理器程序几乎可以为所欲为, 但一般而言, 设计应力求简单。21.1 节将对这一点展开论述。

程序清单 20-1: 为 SIGINT 信号安装一个处理器程序

```
----- signals/ouch.c
#include <signal.h>
#include "tspi_hdr.h"

static void
sigHandler(int sig)
{
    printf("Ouch!\n");                /* UNSAFE (see Section 21.1.2) */
}

int
main(int argc, char *argv[])
{
    int j;

    if (signal(SIGINT, sigHandler) == SIG_ERR)
        errExit("signal");

    for (j = 0; ; j++) {
        printf("%d\n", j);
        sleep(3);                    /* Loop slowly... */
    }
}
----- signals/ouch.c
```

程序清单 20-1 所示为一个简单的信号处理器函数，由主程序为 SIGINT 信号而建立。当键入中断字符（通常为 Control-C）时，终端驱动程序将产生该信号。处理器只是简单打印一条消息，随即返回。

主程序会持续循环。每次迭代，程序都将递增计数器值并将其打印出来，然后休眠几秒钟。（为了按这种方式休眠，程序使用了 sleep() 函数，该函数会令调用者处于暂停状态，持续时间则由指定的秒数决定。该函数将在 23.4.1 节中进行描述。）

运行程序清单 20-1 中程序的结果如下：

```
$ ./ouch
0                               Main program loops, displaying successive integers
Type Control-C
Ouch!                            Signal handler is executed, and returns
1                               Control has returned to main program
2
Type Control-C again
Ouch!
3
Type Control-\ (the terminal quit character)
Quit (core dumped)
```

内核在调用信号处理器程序时，会将引发调用的信号编号作为一个整型参数传递给处理器函数。（就是程序清单 20-1 中处理器函数的 sig 参数）。如果信号处理器程序只捕获一种类型的信号，那么这个参数几乎无用。然而，如果安装相同的处理器来捕获不同类型的信号，那么就可以利用此参数来判定引发对处理器调用的是何种信号。

程序清单 20-2 中程序展示了这一思路，为 SIGINT 和 SIGQUIT 信号建立了同一处理器程序。（当键入终端退出字符时，通常为 Control-\，终端驱动程序将产生 SIGQUIT 信号。）处理器程序代码通过检查 sig 参数来区分这两种信号，并为每种信号采取不同措施。main() 函数则

使用 `pause()` 函数（参见 20.14 节的描述）来阻塞进程，直至捕获到信号。

如下 shell 会话日志演示了对该程序的使用：

```
$ ./intquit
Type Control-C
Caught SIGINT (1)
Type Control-C again
Caught SIGINT (2)
and again
Caught SIGINT (3)
Type Control-\
Caught SIGQUIT - that's all folks!
```

程序清单 20-1 和程序清单 20-2 都在信号处理器程序中使用了 `printf()` 函数来显示消息。现实世界的应用程序一般绝不会在信号处理器程序中使用 `stdio` 函数，21.1.2 节将就其原因进行讨论。然而，本书各种示例仍然会在信号处理器程序中调用 `printf()` 函数，作为观察处理器程序调用的一种简单手段。

程序清单 20-2：为两个不同信号建立同一处理器函数

```
-----signals/intquit.c
#include <signal.h>
#include "tlpi_hdr.h"

static void
sigHandler(int sig)
{
    static int count = 0;

    /* UNSAFE: This handler uses non-async-signal-safe functions
       (printf(), exit()); see Section 21.1.2) */

    if (sig == SIGINT) {
        count++;
        printf("Caught SIGINT (%d)\n", count);
        return;          /* Resume execution at point of interruption */
    }

    /* Must be SIGQUIT - print a message and terminate the process */

    printf("Caught SIGQUIT - that's all folks!\n");
    exit(EXIT_SUCCESS);
}

int
main(int argc, char *argv[])
{
    /* Establish same handler for SIGINT and SIGQUIT */

    if (signal(SIGINT, sigHandler) == SIG_ERR)
        errExit("signal");
    if (signal(SIGQUIT, sigHandler) == SIG_ERR)
        errExit("signal");

    for (;;)                /* Loop forever, waiting for signals */
        pause();           /* Block until a signal is caught */
}
-----signals/intquit.c
```

20.5 发送信号：kill()

与 shell 的 kill 命令相类似，一个进程能够使用 kill() 系统调用向另一进程发送信号。（之所以选择 kill 作为术语，是因为早期 UNIX 实现中大多数信号的默认行为是终止进程。）

```
#include <signal.h>

int kill(pid_t pid, int sig);

Returns 0 on success, or -1 on error
```

pid 参数标识一个或多个目标进程，而 sig 则指定了要发送的信号。如何解释 pid，要视以下 4 种情况而定。

- 如果 pid 大于 0，那么会发送信号给由 pid 指定的进程。
- 如果 pid 等于 0，那么会发送信号给与调用进程同组的每个进程，包括调用进程自身。（SUSv3 声明，除去“一组未予明确的系统进程”¹之外，应将信号发送给同一进程组中的所有进程，且这一排除条件同样适用于余下的两种情况。）
- 如果 pid 小于 -1，那么会向组 ID 等于该 pid 绝对值的进程组内所有下属进程发送信号。向一个进程组的所有进程发送信号在 shell 作业控制中有特殊用途（参见 34.7 节）。
- 如果 pid 等于 -1，那么信号的发送范围是：调用进程有权将信号发往的每个目标进程，除去 init（进程 ID 为 1）和调用进程自身。如果特权级进程发起这一调用，那么会发送信号给系统中的所有进程，上述两个进程除外。显而易见，有时也将这种信号发送方式称之为广播信号。（SUSv3 并未要求将调用进程排除在信号接收范围之外，Linux 此处所遵循的是 BSD 系统的语义。）

如果并无进程与指定的 pid 相匹配，那么 kill() 调用失败，同时将 errno 置为 ESRCH（“查无此进程”）。

进程要发送信号给另一进程，还需要适当的权限，其权限规则如下。

- 特权级（CAP_KILL）进程可以向任何进程发送信号。
- 以 root 用户和组运行的 init 进程（进程号为 1），是一种特例，仅能接收已安装了处理器函数的信号。这可以防止系统管理员意外杀死 init 进程——这一系统运作的基石。
- 如图 20-2 所示，如果发送者的实际或有效用户 ID 匹配于接受者的实际用户 ID 或者保存设置用户 ID(saved set-user-id)，那么非特权进程也可以向另一进程发送信号。利用这一规则，用户可以向由他们启动的 set-user-ID 程序发送信号，而无需考虑目标进程有效用户 ID 的当前设置。将目标进程有效用户 ID 排除在检查范围之外，这一举措的辅助作用在于防止用户某甲向用户某乙的进程发送信号，而该进程正在执行的 set-user-ID 程序又属于用户某甲。（SUSv3 要求强制执行图 20-2 所示的规则，但如 kill(2) 手册页所述，Linux 内核在 2.0 版本之前所遵循的规则略有不同。）
- SIGCONT 信号需要特殊处理。无论对用户 ID 的检查结果如何，非特权进程可以向同一会话中的任何其他进程发送这一信号。利用这一规则，运行作业控制的 shell 可以重启已停止的作业（进程组），即使作业进程已经修改了它们的用户 ID。（亦即，使用

¹ 译者注：参考 APUEv2，其实就是由系统实现决定的意思，再次鄙视 SUS 的学究笔法。

9.7 节所述系统调用来改变其凭据，进而成为特权级进程。)

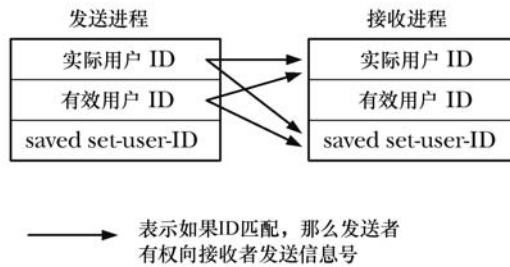


图 20-2: 非特权进程发送信号所需的权限

如果进程无权发送信号给所请求的 pid，那么 kill()调用将失败，且将 errno 置为 EPERM。若 pid 所指为一系列进程（即 pid 是负值）时，只要可以向其中之一发送信号，则 kill()调用成功。

程序清单 20-3 中展示了 kill()的用法。

20.6 检查进程的存在

kill()系统调用还有另一重功用。若将参数 sig 指定为 0（即所谓空信号），则无信号发送。相反，kill()仅会去执行错误检查，查看是否可以向目标进程发送信号。从另一角度来看，这意味着，可以使用空信号来检测具有特定进程 ID 的进程是否存在。若发送空信号失败，且 errno 为 ESRCH，则表明目标进程不存在。如果调用失败，且 errno 为 EPERM（表示进程存在，但无权向目标进程发送信号）或者调用成功（有权向进程发送信号），那么就表示进程存在。

验证一个特定进程 ID 的存在并不能保证特定程序仍在运行。因为内核会随着进程的生死而循环使用进程 ID。而一段时间之后，同一进程 ID 所指恐怕是另一进程了。此外，特定进程 ID 可能存在，但是一个僵尸（亦即，进程已死，但其父进程尚未执行 wait()来获取其终止状态，如 26.2 节所述）。

还可使用各种其他技术来检查某一特定进程是否正在运行，其中包括如下技术。

- wait()系统调用：第 26 章将描述这些调用。这些调用仅用于监控调用者的子进程。
- 信号量和排他文件锁：如果进程持续持有某一信号量或文件锁，并且一直处于被监控状态，那么如能获取到信号量或锁时，即表明该进程已经终止。第 47 章和第 53 章将描述信号量，第 55 章将描述文件锁。
- 诸如管道和 FIFO 之类的 IPC 通道：可对监控目标进程进行设置，令其在自身生命周期内持有对通道进行写操作的打开文件描述符。同时，令监控进程持有针对通道进行读操作的打开文件描述符，且当通道写入端关闭时（遭遇文件结束符），即可获知监控目标进程已经终止。监控进程对此情况的判定，既可借助于对自身文件描述符的读取，也可采用第 63 章所述的描述符监控技术之一。
- /proc/PID 接口：例如，如果进程 ID 为 12345 的进程存在，那么目录/proc/12345 将存在，可以发起诸如 stat()之类的调用来进行检查。

除去最后一项之外，循环使用进程 ID 不会影响上述所有技术。

程序清单 20-3 展示了 kill()的用法。该程序接受两个命令行参数，分别为信号编号和进程 ID，并使用 kill()将该信号发送给指定进程。如果指定了信号 0（空信号），那么程序将报告目

标进程是否存在。

20.7 发送信号的其他方式：raise()和 killpg()

有时，进程需要向自身发送信号（34.7.3 节就有此一例）。raise()函数就执行了这一任务。

```
#include <signal.h>

int raise(int sig);

Returns 0 on success, or nonzero on error
```

在单线程程序中，调用 raise()相当于对 kill()的如下调用：

```
kill(getpid(), sig);
```

支持线程的系统会将 raise(sig)实现为：

```
pthread_kill(pthread_self(), sig)
```

33.2.3 节描述了 pthread_kill()函数，但目前仅需了解一点就已足够，该实现意味着将信号传递给调用 raise()的特定线程。相比之下，kill(getpid(), sig)调用会发送一个信号给调用进程，并可将该信号传递给该进程的任一线程。

raise()函数起源于 C89。C 语言标准不包含诸如进程 ID 之类的操作系统细节，raise()函数之所以得以定义，是因为该函数不需要引用进程 ID。

当进程使用 raise()（或者 kill()）向自身发送信号时，信号将立即传递（即，在 raise()返回调用者之前）。

注意，raise()出错将返回非 0 值（不一定为-1）。调用 raise()唯一可能发生的错误为 EINVAL，即 sig 无效。因此，在任何指定了某一 SIGxxxx 常量的位置，都未检查该函数的返回状态。

程序清单 20-3：使用 kill()系统调用

```
signals/t_kill.c

#include <signal.h>
#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    int s, sig;

    if (argc != 3 || strcmp(argv[1], "--help") == 0)
        usageErr("%s sig-num pid\n", argv[0]);

    sig = getInt(argv[2], 0, "sig-num");

    s = kill(getLong(argv[1], 0, "pid"), sig);

    if (sig != 0) {
        if (s == -1)
            errExit("kill");
    } else {
        /* Null signal: process existence check */
        if (s == 0) {
```



```

        printf("Process exists and we can send it a signal\n");
    } else {
        if (errno == EPERM)
            printf("Process exists, but we don't have "
                "permission to send it a signal\n");
        else if (errno == ESRCH)
            printf("Process does not exist\n");
        else
            errExit("kill");
    }
}

exit(EXIT_SUCCESS);
}

```

signals/t_kill.c

killpg()函数向某一进程组的所有成员发送一个信号。

```

#include <signal.h>

int killpg(pid_t pgrp, int sig);

```

Returns 0 on success, or -1 on error

killpg()调用相当于对 kill()的如下调用：

```
kill(-pgrp, sig);
```

如果指定 pgrp 的值为 0，那么会向调用者所属进程组的所有进程发送此信号。SUSv3 对此未作规范，但大多数 UNIX 实现对该情况的处理方式与 Linux 相同。

20.8 显示信号描述

每个信号都有一串与之相关的可打印说明。这些描述位于数组 sys_siglist 中。例如，可以用 sys_siglist[SIGPIPE]来获取对 SIGPIPE 信号（管道断开）的描述。然而，较之于直接引用 sys_siglist 数组，还是推荐使用 strsignal()函数。

```

#define _BSD_SOURCE
#include <signal.h>

extern const char *const sys_siglist[];

#define _GNU_SOURCE
#include <string.h>

char *strsignal(int sig);

```

Returns pointer to signal description string

strsignal()函数对 sig 参数进行边界检查，然后返回一枚指针，指向针对该信号的可打印描述字符串，或者是当信号编号无效时指向错误字符串。（在其他一些 UNIX 实现中，strsignal()函数会在 sig 无效时返回空值。）

除去边界检查之外，strsignal()函数较之于直接引用 sys_siglist 数组的另一优势是对本地

(locale) 设置敏感 (10.4 节), 所以显示信号描述时会使用本地语言。

程序清单 20-4 中所示为使用 `strsignal()` 的例子之一。

`psignal()` 函数 (在标准错误设备上) 所示为 `msg` 参数所给定的字符串, 后面跟有一个冒号, 随后是对应于 `sig` 的信号描述。和 `strsignal()` 一样, `psignal()` 函数也对本地设置敏感。

```
#include <signal.h>

void psignal(int sig, const char *msg);
```

尽管 SUSv3 并未将 `psignal()`、`strsignal()` 和 `sys_siglist` 纳入标准, 但还是有许多 UNIX 实现支持它们。(SUSv4 中加入了对于 `psignal()` 和 `strsignal()` 的规范。)

20.9 信号集

许多信号相关的系统调用都需要能表示一组不同的信号。例如, `sigaction()` 和 `sigprocmask()` 允许程序指定一组将由进程阻塞的信号, 而 `sigpending()` 则返回一组目前正在等待送达给一进程的信号。(稍后将描述这些系统调用。)

多个信号可使用一个称之为信号集的数据结构来表示, 其系统数据类型为 `sigset_t`。SUSv3 规定了一系列函数来操纵信号集, 现在将描述这些函数。

像在大多数 UNIX 实现中一样, `sigset_t` 数据类型在 Linux 中是一个位掩码。然而, SUSv3 对此并无要求。使用其他一些类型的结构来表示信号集也是有可能的。SUSv3 仅要求可对 `sigset_t` 类型赋值即可。因此, 必须使用某些标量类型 (比如一个整数) 或者一个 C 语言结构 (也许包含了一个整型数组) 来实现该类型。

`sigemptyset()` 函数初始化一个未包含任何成员的信号集。`sigfillset()` 函数则初始化一个信号集, 使其包含所有信号 (包括所有实时信号)。

```
#include <signal.h>

int sigemptyset(sigset_t *set);
int sigfillset(sigset_t *set);
```

Both return 0 on success, or -1 on error

必须使用 `sigemptyset()` 或者 `sigfillset()` 来初始化信号集。这是因为 C 语言不会对自动变量进行初始化, 并且, 借助于将静态变量初始化为 0 的机制来表示空信号集的作法在可移植性上存在问题, 因为有可能使用位掩码之外的结构来实现信号集。(出于同一原因, 为将信号集标记为空而使用 `memset(3)` 函数将其内容清零的做法也不正确。)

信号集初始化后, 可以分别使用 `sigaddset()` 和 `sigdelset()` 函数向一个集合中添加或者移除单个信号。

```
#include <signal.h>

int sigaddset(sigset_t *set, int sig);
int sigdelset(sigset_t *set, int sig);
```

Both return 0 on success, or -1 on error

在 `sigaddset()` 和 `sigdelset()` 中, `sig` 参数均表示信号编号。

sigismember()函数用来测试信号 sig 是否是信号集 set 的成员。

```
#include <signal.h>

int sigismember(const sigset_t *set, int sig);
                                Returns 1 if sig is a member of set, otherwise 0
```

如果 sig 是 set 的一个成员，那么 sigismember()函数将返回 1 (true)，否则返回 0 (false)。GNU C 库还实现了 3 个非标准函数，是对上述信号集标准函数的补充。

```
#define _GNU_SOURCE
#include <signal.h>

int sigandset(sigset_t *set, sigset_t *left, sigset_t *right);
int sigorset(sigset_t *dest, sigset_t *left, sigset_t *right);
                                Both return 0 on success, or -1 on error

int sigisemptyset(const sigset_t *set);
                                Returns 1 if sig is empty, otherwise 0
```

这些函数执行了如下任务。

- sigandset()将 left 集和 right 集的交集置于 dest 集。
- sigorset()将 left 集和 right 集的并集置于 dest 集。
- 若 set 集内未包含信号，则 sigisemptyset()返回 true。

示例程序

程序清单 20-4 所示为使用本节介绍的函数来编写的函数，供本书后续各程序调用。第一个函数 printSigset()显示了指定信号集的成员信号。该函数使用了定义于<signal.h>文件中的 NSIG 常量，其值等于信号最大编号加 1。当获取信号集成员时，会在测试所有信号编号的循环中将该值作为循环上限。

虽然 SUSv3 并未定义 NSIG，但是大多数 UNIX 实现都支持这一常量。只不过，要想使其可见，可能需要使用特定于实现的编译器选项。例如，在 Linux 中，就必须定义如下功能测试宏之一：BSD_SOURCE、_SVID_SOURCE 或者 _GNU_SOURCE。

利用 printSigset()函数，printSigMask()和 printPendingSigs()函数分别用于显示进程的信号掩码和当前处于等待状态的信号集。这两个函数还分别使用了 sigprocmask()和 sigpending()系统调用。sigprocmask()和 sigpending()系统调用将分别在 20.10 节和 20.11 节中予以描述。

程序清单 20-4：显示信号集的函数

```
----- signals/signal_functions.c

#define _GNU_SOURCE
#include <string.h>
#include <signal.h>
#include "signal_functions.h"          /* Declares functions defined here */
#include "tspi_hdr.h"

/* NOTE: All of the following functions employ fprintf(), which
   is not async-signal-safe (see Section 21.1.2). As such, these
```

```

functions are also not async-signal-safe (i.e., beware of
indiscriminately calling them from signal handlers). */

void          /* Print list of signals within a signal set */
printSigset(FILE *of, const char *prefix, const sigset_t *sigset)
{
    int sig, cnt;

    cnt = 0;
    for (sig = 1; sig < NSIG; sig++) {
        if (sigismember(sigset, sig)) {
            cnt++;
            fprintf(of, "%s%d (%s)\n", prefix, sig, strsignal(sig));
        }
    }

    if (cnt == 0)
        fprintf(of, "%s<empty signal set>\n", prefix);
}

int          /* Print mask of blocked signals for this process */
printSigMask(FILE *of, const char *msg)
{
    sigset_t currMask;

    if (msg != NULL)
        fprintf(of, "%s", msg);

    if (sigprocmask(SIG_BLOCK, NULL, &currMask) == -1)
        return -1;

    printSigset(of, "\t\t", &currMask);

    return 0;
}

int          /* Print signals currently pending for this process */
printPendingSigs(FILE *of, const char *msg)
{
    sigset_t pendingSigs;

    if (msg != NULL)
        fprintf(of, "%s", msg);

    if (sigpending(&pendingSigs) == -1)
        return -1;

    printSigset(of, "\t\t", &pendingSigs);

    return 0;
}

```

signals/signal_functions.c

20.10 信号掩码（阻塞信号传递）

内核会为每个进程维护一个信号掩码，即一组信号，并将阻塞其针对该进程的传递。如

果将遭阻塞的信号发送给某进程，那么对该信号的传递将延后，直至从进程信号掩码中移除该信号，从而解除阻塞为止。（由 33.2.1 节可知，信号掩码实际属于线程属性，在多线程进程中，每个线程都可使用 `pthread_sigmask()` 函数来独立检查和修改其信号掩码。）

向信号掩码中添加一个信号，有如下几种方式。

- 当调用信号处理器程序时，可将引发调用的信号自动添加到信号掩码中。是否发生这一情况，要视 `sigaction()` 函数在安装信号处理器程序时所使用的标志而定。
- 使用 `sigaction()` 函数建立信号处理器程序时，可以指定一组额外信号，当调用该处理器程序时会将其阻塞。
- 使用 `sigprocmask()` 系统调用，随时可以显式向信号掩码中添加或移除信号。

对前两种情况的讨论将推迟到 20.13 节对 `sigaction()` 函数的介绍之后，现在先来讨论 `sigprocmask()` 函数。

```
#include <signal.h>

int sigprocmask(int how, const sigset_t *set, sigset_t *oldset);
                Returns 0 on success, or -1 on error
```

使用 `sigprocmask()` 函数既可修改进程的信号掩码，又可获取现有掩码，或者两重功效兼具。**how 参数指定了 `sigprocmask()` 函数想给信号掩码带来的变化。**

SIG_BLOCK

将 `set` 指向信号集内的指定信号添加到信号掩码中。换言之，将信号掩码设置为其当前值和 `set` 的并集。

SIG_UNBLOCK

将 `set` 指向信号集中的信号从信号掩码中移除。即使要解除阻塞的信号当前并未处于阻塞状态，也不会返回错误。

SIG_SETMASK

将 `set` 指向的信号集赋给信号掩码。

上述各种情况下，若 `oldset` 参数不为空，则其指向一个 `sigset_t` 结构缓冲区，用于返回之前的信号掩码。

如果想获取信号掩码而又对其不作改动，那么可将 `set` 参数指定为空，这时将忽略 `how` 参数。

要想暂时阻止信号的传递，可以使用程序清单 20-5 中所示的一系列调用来阻塞信号，然后再将信号掩码重置为先前的状态以解除对信号的锁定。

程序清单 20-5：暂时阻塞信号传递

```
sigset_t blockSet, prevMask;

/* Initialize a signal set to contain SIGINT */

sigemptyset(&blockSet);
sigaddset(&blockSet, SIGINT);

/* Block SIGINT, save previous signal mask */

if (sigprocmask(SIG_BLOCK, &blockSet, &prevMask) == -1)
```

```

    errExit("sigprocmask1");

    /* ... Code that should not be interrupted by SIGINT ... */

    /* Restore previous signal mask, unblocking SIGINT */

    if (sigprocmask(SIG_SETMASK, &prevMask, NULL) == -1)
        errExit("sigprocmask2");

```

SUSv3 规定，如果有任何等待信号因对 `sigprocmask()` 的调用而解除了锁定，那么在此调用返回前至少会传递一个信号。换言之，如果解除了对某个等待信号的锁定，那么会立刻将该信号传递给进程。

系统将忽略试图阻塞 `SIGKILL` 和 `SIGSTOP` 信号的请求。如果试图阻塞这些信号，`sigprocmask()` 函数既不会予以关注，也不会产生错误。这意味着，可以使用如下代码来阻塞除 `SIGKILL` 和 `SIGSTOP` 之外的所有信号：

```

sigfillset(&blockSet);
if (sigprocmask(SIG_BLOCK, &blockSet, NULL) == -1)
    errExit("sigprocmask");

```

20.11 处于等待状态的信号

如果某进程接受了一个该进程正在阻塞的信号，那么会将该信号添加到进程的等待信号集中。当（且如果）之后解除了对该信号的锁定时，会随之将信号传递给此进程。**为了确定进程中处于等待状态的是哪些信号**，可以使用 `sigpending()`。

```

#include <signal.h>

int sigpending(sigset_t *set);

```

Returns 0 on success, or -1 on error

`sigpending()` 系统调用为调用进程返回处于等待状态的信号集，并将其置于 `set` 指向的 `sigset_t` 结构中。随后可以使用 20.9 节描述的 `sigismember()` 函数来检查 `set`。

如果修改了对等待信号的处置，那么当后来解除对信号的锁定时，将根据新的处置来处理信号。这项技术虽然不经常使用，但还是存在一个应用场景，即将对信号的处处置为 `SIG_IGN`，或者 `SIG_DFL`（如果信号的默认行为是忽略），从而阻止传递处于等待状态的信号。因此，会将信号从进程的等待信号集中移除，从而不传递该信号。

20.12 不对信号进行排队处理

等待信号集只是一个掩码，仅表明一个信号是否发生，而未表明其发生的次数。换言之，如果同一信号在阻塞状态下产生多次，那么会将该信号记录在等待信号集中，并在稍后仅传递一次。（标准信号和实时信号之间的差异之一在于，如 22.8 节所述，对实时信号进行了排队处理。）

程序清单 20-6 和程序清单 20-7 显示了两个程序，可用于观察未作排队处理的信号。清单 20-6 的程序可接受多达四个命令行参数，如下所示：

```
$ ./sig_sender PID num-sigs sig-num [sig-num-2]
```

第一个参数是程序发送信号的目标进程 ID。第二个参数则指定发送给目标进程的信号数量。第三个参数指定发往目标进程的信号编号。如果还提供了信号编号作为第四个参数，那么当程序发送完之前参数所指定的信号之后，将发送该信号的一个实例。在如下 shell 会话示例中，就使用了最后一个参数向目标进程发送一个 SIGINT 信号，发送该信号的目的将在稍后揭晓。

程序清单 20-6：发送多个信号

```
----- signals/sig_sender.c
#include <signal.h>
#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    int numSigs, sig, j;
    pid_t pid;

    if (argc < 4 || strcmp(argv[1], "--help") == 0)
        usageErr("%s pid num-sigs sig-num [sig-num-2]\n", argv[0]);
    pid = getLong(argv[1], 0, "PID");
    numSigs = getInt(argv[2], GN_GT_0, "num-sigs");
    sig = getInt(argv[3], 0, "sig-num");

    /* Send signals to receiver */

    printf("%s: sending signal %d to process %ld %d times\n",
           argv[0], sig, (long) pid, numSigs);

    for (j = 0; j < numSigs; j++)
        if (kill(pid, sig) == -1)
            errExit("kill");

    /* If a fourth command-line argument was specified, send that signal */

    if (argc > 4)
        if (kill(pid, getInt(argv[4], 0, "sig-num-2")) == -1)
            errExit("kill");

    printf("%s: exiting\n", argv[0]);
    exit(EXIT_SUCCESS);
}
----- signals/sig_sender.c
```

程序清单 20-7 中程序则被设计为去捕获程序清单 20-6 程序所发送的信号并汇总其统计数据。该程序执行了以下步骤。

- 该程序建立了单个处理器程序来捕获所有信号。（捕获 SIGKILL 和 SIGSTOP 信号是不可能的，不过将忽略在尝试为这些信号建立处理器时所发生的错误。）对于大多数类型的信号，处理器程序只是简单地使用一个数组来对信号计数。如果收到的信号为 SIGINT，那么处理器程序将对标志 (gotSigint) 置位，从而使程序退出主循环（下面所描述的 while 循环）。（至于 volatile 修饰符以及声明 gotSigint 变量的 sig_atomic_t 数据类型，将在 21.1.3 节中解释其用途。）
- 如果提供有一个命令行参数给程序，那么程序对所有信号的阻塞秒数将由该参数指

定，并且在解除阻塞之前会显示待处理的信号集，从而使用户在进程执行下面的步骤前向其发送信号。

- 程序执行 while 循环以消耗 CPU 时间，直至将 gotSigint 标志置位。（20.14 节和 22.9 节描述了 pause() 和 sigsuspend() 的用法，二者在等待信号到来期间对 CPU 的使用方式都颇为高效。）
- 退出 while 循环后，程序显示对所有接收信号的计数。

首先使用这两个程序来展示的是遭阻塞的信号无论产生了多少次，仅会传递一次。这里为接收者指定了一个睡眠间隔，并在醒来之前发送所有信号。

```
$ ./sig_receiver 15 &                               Receiver blocks signals for 15 secs
[1] 5368
./sig_receiver: PID is 5368
./sig_receiver: sleeping for 15 seconds
$ ./sig_sender 5368 1000000 10 2                   Send SIGUSR1 signals, plus a SIGINT
./sig_sender: sending signal 10 to process 5368 1000000 times
./sig_sender: exiting
./sig_receiver: pending signals are:
                2 (Interrupt)
                10 (User defined signal 1)
./sig_receiver: signal 10 caught 1 time
[1]+ Done ./sig_receiver 15
```

发送程序的命令行参数指定了 SIGUSR1 和 SIGINT 信号，其在 Linux/x86 中的编号分别为 10 和 2。

从以上输出可知，即使一个信号发送了一百万次，但仅会传递一次给接收者。

即使进程没有阻塞信号，其所收到的信号也可能比发送给它的要少得多。如果信号发送速度如此之快，以至于在内核考虑将执行权调度给接收进程之前，这些信号就已经到达，这时就会发生上述情况，从而导致多次发送的信号在进程等待信号集中只记录了一次。如果不带任何命令行参数来执行程序清单 20-7 中程序（因此，进程没有阻塞信号，也没有睡眠），那么将看到如下情况：

```
$ ./sig_receiver &
[1] 5393
./sig_receiver: PID is 5393
$ ./sig_sender 5393 1000000 10 2
./sig_sender: sending signal 10 to process 5393 1000000 times
./sig_sender: exiting
./sig_receiver: signal 10 caught 52 times
[1]+ Done ./sig_receiver
```

在所发送的一百万次信号之中，接收进程仅捕获到 52 次。（捕获信号的精确数目每每不同，这取决于内核调度算法变幻莫测的决策结果。）之所以如此，原因在于，发送程序会在每次获得调度而运行时发送多个信号给接收者。然而，当接收进程得以运行时，传递来的信号只有一个，因为只会将这些信号中的一个标记为等待状态。

程序清单 20-7：捕获信号并对其计数

```
----- signals/sig_receiver.c
#define _GNU_SOURCE
#include <signal.h>
#include "signal_functions.h" /* Declaration of printSigset() */
#include "tlpi_hdr.h"

static int sigCnt[NSIG]; /* Counts deliveries of each signal */
```



```

static volatile sig_atomic_t gotSigint = 0;
                                /* Set nonzero if SIGINT is delivered */

static void
① handler(int sig)
{
    if (sig == SIGINT)
        gotSigint = 1;
    else
        sigCnt[sig]++;
}

int
main(int argc, char *argv[])
{
    int n, numSecs;
    sigset_t pendingMask, blockingMask, emptyMask;

    printf("%s: PID is %ld\n", argv[0], (long) getpid());

    ② for (n = 1; n < NSIG; n++)          /* Same handler for all signals */
        (void) signal(n, handler);      /* Ignore errors */

    /* If a sleep time was specified, temporarily block all signals,
       sleep (while another process sends us signals), and then
       display the mask of pending signals and unblock all signals */

    ③ if (argc > 1) {
        numSecs = getInt(argv[1], GN_GT_0, NULL);

        sigfillset(&blockingMask);
        if (sigprocmask(SIG_SETMASK, &blockingMask, NULL) == -1)
            errExit("sigprocmask");

        printf("%s: sleeping for %d seconds\n", argv[0], numSecs);
        sleep(numSecs);

        if (sigpending(&pendingMask) == -1)
            errExit("sigpending");

        printf("%s: pending signals are: \n", argv[0]);
        printSigset(stdout, "\t\t", &pendingMask);

        sigemptyset(&emptyMask);        /* Unblock all signals */
        if (sigprocmask(SIG_SETMASK, &emptyMask, NULL) == -1)
            errExit("sigprocmask");
    }

    ④ while (!gotSigint)                /* Loop until SIGINT caught */
        continue;

    ⑤ for (n = 1; n < NSIG; n++)          /* Display number of signals received */
        if (sigCnt[n] != 0)
            printf("%s: signal %d caught %d time%s\n", argv[0], n,
                sigCnt[n], (sigCnt[n] == 1) ? "" : "s");

    exit(EXIT_SUCCESS);
}

```

signals/sig_receiver.c

20.13 改变信号处置: sigaction ()

除去 `signal()` 之外, `sigaction()` 系统调用是设置信号处置的另一选择。虽然 `sigaction()` 的用法比之 `signal()` 更为复杂, 但作为回报, 也更具灵活性。尤其是, `sigaction()` 允许在获取信号处置的同时无需将其改变, 并且, 还可设置各种属性对调用信号处理器程序时的行为施以更加精准的控制。此外, 如 22.7 节所述, 在建立信号处理器程序时, `sigaction()` 较之 `signal()` 函数可移植性更佳。

```
#include <signal.h>

int sigaction(int sig, const struct sigaction *act, struct sigaction *oldact);

Returns 0 on success, or -1 on error
```

`sig` 参数标识想要获取或改变的信号编号。**该参数可以是除去 SIGKILL 和 SIGSTOP 之外的任何信号。**

`act` 参数是一枚指针, 指向描述信号新处置的数据结构。如果仅对信号的现有处置感兴趣, 那么可将该参数指定为 `NULL`。`oldact` 参数是指向同一结构类型的指针, 用来返回之前信号处置的相关信息。如果无意获取此类信息, 那么可将该参数指定为 `NULL`。`act` 和 `oldact` 所指向的结构类型如下所示:

```
struct sigaction {
    void (*sa_handler)(int); /* Address of handler */
    sigset_t sa_mask; /* Signals blocked during handler
                       invocation */
    int sa_flags; /* Flags controlling handler invocation */
    void (*sa_restorer)(void); /* Not for application use */
};
```

`sigaction` 结构实际要比此处所展示的更为复杂, 更多细节请参见 21.4 节。

`sa_handler` 字段对应于 `signal()` 的 `handler` 参数。其所指定的值为信号处理器函数的地址, 亦或是常量 `SIG_IGN`、`SIG_DFL` 之一。仅当 `sa_handler` 是信号处理程序的地址时, 亦即 `sa_handler` 的取值在 `SIG_IGN` 和 `SIG_DFL` 之外, 才会对 `sa_mask` 和 `sa_flags` 字段 (稍后讨论) 加以处理。余下的字段 `sa_restorer`, 则不适用于应用程序 (SUSv3 未予规定)。

`sa_restorer` 字段仅供内部使用, 用以确保当信号处理器程序完成后, 会去调用专用的 `sigreturn()` 系统调用, 借此来恢复进程的执行上下文, 以便于进程从信号处理器中断的位置继续执行。这一用法的实例可见诸于 `glibc` 源文件 `sysdeps/unix/sysv/linux/i386/sigaction.c` 中。

`sa_mask` 字段定义了一组信号, 在调用由 `sa_handler` 所定义的处理程序时将阻塞该组信号。当调用信号处理器程序时, 会在调用信号处理器之前, 将该组信号中当前未处于进程掩码之列的任何信号自动添加到进程掩码中。这些信号将保留在进程掩码中, 直至信号处理器函数返回, 届时将自动删除这些信号。利用 `sa_mask` 字段可指定一组信号, 不允许它们中断此处理器程序的执行。此外, 引发对处理器程序调用的信号将自动添加到进程信号掩码中。这意味着, 当正在执行处理器程序时, 如果同一个信号实例第二次抵达, 信号处理器程序将不会递归中断自己。由于不会对遭阻塞的信号进行排队处理, 如果在处理器程序执行过程中

重复产生这些信号中的任何信号，（稍后）对信号的传递将是一次性的。

`sa_flags` 字段是一个位掩码，指定用于控制信号处理过程的各种选项。该字段包含的位如下（可以相或（|））。

SA_NOCLDSTOP

若 `sig` 为 `SIGCHLD` 信号，则当因接受一信号而停止或恢复某一子进程时，将不会产生此信号。参见 26.3.2 节。

SA_NOCLDWAIT

（始于 Linux 2.6）若 `sig` 为 `SIGCHLD` 信号，则当子进程终止时不会将其转化为僵尸。更多细节参见 26.3.3 节。

SA_NODEFER

捕获该信号时，不会在执行处理器程序时将该信号自动添加到进程掩码中。`SA_NOMASK` 历史上曾是 `SA_NODEFER` 的代名词。之所以建议使用后者，是因为 SUSv3 将其纳入规范。

SA_ONSTACK

针对此信号调用处理器函数时，使用了由 `sigaltstack()` 安装的备选栈。参见 21.3 节。

SA_RESETHAND

当捕获该信号时，会在调用处理器函数之前将信号处置重置为默认值（即 `SIG_DFL`）（默认情况下，信号处理器函数保持建立状态，直至进一步调用 `sigaction()` 将其显式解除。）`SA_ONESHOT` 历史上曾是 `SA_RESETHAND` 的代名词，之所以建议使用后者，是因为 SUSv3 将其纳入规范。

SA_RESTART

自动重启由信号处理器程序中断的系统调用。参见 21.5 节。

SA_SIGINFO

调用信号处理器程序时携带了额外参数，其中提供了关于信号的深入信息。对该标志的描述参见 21.4 节。

SUSv3 定义了上述所有选项。

程序清单 21-1 展现了对 `sigaction()` 的使用。

20.14 等待信号：pause()

调用 `pause()` 将暂停进程的执行，直至信号处理器函数中断该调用为止（或者直至一个未处理信号终止进程为止）。

```
#include <unistd.h>
```

```
int pause(void);
```

Always returns -1 with *errno* set to `EINTR`

处理信号时，`pause()` 遭到中断，并总是返回 -1，并将 `errno` 置为 `EINTR`。（21.5 节描述了关于 `EINTR` 错误的更多信息。）

程序清单 20-2 提供了应用 `pause()` 的一例子。

在 22.9 节、22.10 节及 22.11 节中，可以看到程序等待信号时暂停执行的各种其他方式。

20.15 总结

信号是发生某种事件的通知机制，可以由内核、另一进程或进程自身发送给进程。存在一系列的标准信号类型，每种都有唯一的编号和目的。

信号传递通常是异步行为，这意味着信号中断进程执行的位置是不可预测的。有时（比如，硬件产生的信号），信号也可同步传递，这意味着在程序执行的某一点可以预期并重现信号的传递。

默认情况下，要么忽略信号，要么终止进程（生成或者不生成核心转储文件），要么停止一个正在运行的进程，要么重启一个已停止的进程。特定的默认行为取决于信号类型。此外，程序可以使用 `signal()` 或者 `sigaction()` 来显式忽略一个信号，或者建立一个由程序员自定义的信号处理器程序，以供信号到达时调用。出于可移植性考虑，最好使用 `sigaction()` 来建立信号处理器函数。

一个（具有适当权限的）进程可以使用 `kill()` 向另一进程发送信号。发送空信号（0）是判定特定进程 ID 是否在用的方式之一。

每个进程都具有一个信号掩码，代表当前传递遭到阻塞的一组信号。使用 `sigprocmask()` 可从信号掩码中添加或者移除信号。

如果接收的信号当前遭到阻塞，那么该信号将保持等待状态，直至解除对其阻塞。系统不会对标准信号进行排队处理，也就是说，将信号标记为等待状态（以及后续的传递）只会发生一次。进程能够使用 `sigpending()` 系统调用来获取等待信号集（用以描述多个不同信号的数据结构）。

与 `signal()` 相比，`sigaction()` 系统调用在设置信号处置方面提供了更多控制，且更具灵活性。首先，可以指定一组调用处理器函数时将阻塞的额外信号。此外，可以使用各种标志来控制调用信号处理器时所发生的行为。例如，启用某些标志即可选择旧有的不可靠信号语义（不阻塞引发处理器调用的信号，在调用信号处理器之前就将信号处置重置为默认值）。

借助于 `pause()`，进程可暂停执行，直至信号到达为止。

更多信息

[Bovet & Cesati, 2005] 和 [Maxwell, 1999] 提供了 Linux 信号实现的背景资料。[Goodheart & Cox, 1994] 详细说明了 System V 版本 4 对信号的实现。GNU C 函数库手册（在线网访问：<http://www.gnu.org/>）包含了对信号的全面描述。

20.16 练习

- 20-1. 如 20.3 节所指，比之 `signal()`，`sigaction()` 函数在建立信号处理器时可移植性更佳。请用 `sigaction()` 替换程序清单 20-7 程序（`sig_receiver.c`）中对 `signal()` 的调用。
- 20-2. 编写一程序来展示当将对等待信号的处置改为 `SIG_IGN` 时，程序绝不会看到（捕获）信号。
- 20-3. 编写一程序，以 `sigaction()` 来建立信号处理器函数，请验证 `SA_RESETHAND` 和 `SA_NODEFER` 标志的效果。
- 20-4. 请用 `sigaction()` 调用来实现 `siginterrupt()`。

第 21 章

信号：信号处理器函数

承接上一章，本章将继续介绍信号。本章的描述重点是信号处理器函数（handler），同时会拓展 20.4 节所发起的讨论。本章涵盖主题如下。

- 如何设计信号处理器函数，其中对可重入性以及异步信号安全函数的探讨必不可少。
- 从信号处理器函数中正常返回的各种途径，特别是非本地跳转技术的运用。
- 利用备选栈处理信号。
- 借助于带有 SA_SIGINFO 标志的 sigaction() 函数，处理器函数能够获取引发其调用信号的更多详细信息。
- 信号处理器函数如何中断处于阻塞状态的系统调用，以及如何重启该系统调用。

21.1 设计信号处理器函数

一般而言，将信号处理器函数设计得越简单越好。其中的一个重要原因就在于，这将降低引发竞争条件的风险。下面是针对信号处理器函数的两种常见设计。

- 信号处理器函数设置全局性标志变量并退出。主程序对此标志进行周期性检查，一旦置位随即采取相应动作。（主程序若因监控一个或多个文件描述符的 I/O 状态而无法进行这种周期性检查时，则可令信号处理器函数向一专用管道写入一个字节的数据，同时将该管道的读取端置于主程序所监控的文件描述符范围之内。63.5.2 节展示了这一技术的运用。）
- 信号处理器函数执行某种类型的清理动作，接着终止进程或者使用非本地跳转（21.2.1 节）将栈解开并将控制返回到主程序中的预定位置。

后续各节将会探讨这些设计理念，以及信号处理器函数设计中的其他一些重要概念。

21.1.1 再论信号的非队列化处理

20.10 节已然提及，在执行某信号的处理函数时会阻塞同类信号的传递（除非在调用 sigaction() 时指定了 SA_NODEFER 标志）。如果在执行处理器函数时（再次）产生同类信号，那

么会将该信号标记为等待状态并在处理器函数返回之后再行传递。前一章还曾指出，不会对信号进行排队处理。在处理器函数执行期间，如果多次产生同类信号，那么仍然会将其标记为等待状态，但稍后只会传递一次。

信号的这种“失踪”方式无疑将影响对信号处理器函数的设计。首先，无法对信号的产生次数进行可靠计数。其次，在为信号处理器函数编码时可能需要考虑处理同类信号多次产生的情况。26.3.1 节在讨论 SIGCHLD 信号时会有相关示例。

21.1.2 可重入函数和异步信号安全函数

在信号处理器函数中，并非所有系统调用以及库函数均可予以安全调用。要了解来龙去脉，就需要解释一下以下两种概念：可重入（reentrant）函数和异步信号安全（async-signal-safe）函数。

可重入和非可重入函数

要解释可重入函数为何物，首先需要区分单线程程序和多线程程序。典型 UNIX 程序都具有一条执行线程，贯穿程序始终，CPU 围绕单条执行逻辑来处理指令。而对于多线程程序而言，同一进程却存在多条独立、并发的执行逻辑流。

第 29 章将会展示如何显式创建一个包含多条执行线程的程序。不过，多执行线程的概念与使用了信号处理器函数的程序也有关联。因为信号处理器函数可能会在任一时点异步中断程序的执行，从而在同一个进程中实际形成了两条（即主程序和信号处理器函数）独立（虽然不是并发）的执行线程。

如果同一个进程的多条线程可以同时安全地调用某一函数，那么该函数就是可重入的。此处，“安全”意味着，无论其他线程调用该函数的执行状态如何，函数均可产生预期结果。

SUSv3 对可重入函数的定义是：函数由两条或多条线程调用时，即便是交叉执行，其效果也与各线程以未定义¹顺序依次调用时一致。

更新全局变量或静态数据结构的函数可能是不可重入的。（只用到本地变量的函数肯定是可重入的。）如果对函数的两个调用（例如：分别由两条执行线程发起）同时试图更新同一全局变量或数据类型，那么二者很可能会相互干扰并产生不正确的结果。例如，假设某线程正在为一链表数据结构添加一个新的链表项，而另一线程也正试图更新同一链表。由于为链表添加新项涉及对多枚指针的更新，一旦另一线程中断这些步骤并修改了相同的指针，结果就会产生混乱。

在 C 语言标准函数库中，这种可能性非常普遍。例如，7.1.3 节所提及的 malloc() 和 free() 就维护有一个针对已释放内存块的链表，用于从堆中重新分配内存。如果主程序在调用 malloc() 期间为一个同样调用 malloc() 的信号处理器函数所中断，那么该链表可能会遭到破坏。因此，malloc() 函数族以及使用它们的其他库函数都是不可重入的。

还有一些函数库之所以不可重入，是因为它们使用了经静态分配的内存来返回信息。此类函数（本书其他地方也有论及）的例子包括 crypt()、getpwnam()、gethostbyname() 以及 getservbyname()。如果信号处理器用到了这类函数，那么将会覆盖主程序中上次调用同一函数所返回的信息（反之亦然）。

将静态数据结构用于内部记账的函数也是不可重入的。其中最明显的例子就是 stdio 函数库成员（printf()、scanf() 等），它们会为缓冲区 I/O 更新内部数据结构。所以，如果在信号处

¹ 译者注：任意。

处理器函数中调用了 `printf()`，而主程序又在调用 `printf()` 或其他 `stdio` 函数期间遭到了处理器函数的中断，那么有时就会看到奇怪的输出，甚至导致程序崩溃或者数据的损坏。

即使并未使用不可重入的库函数，可重入问题依然不容忽视。如果信号处理器函数和主程序都要更新由程序员自定义的全局性数据结构，那么对于主程序而言，这种信号处理器函数就是不可重入的。

如果函数是不可重入的，那么其手册页通常会或明或暗地给出提示。对于其中那些使用或返回静态分配变量的函数，需要特别留意。

示例程序

程序清单 21-1 展示了函数 `crypt()` (8.5 节) 不可重入的本来面目。该程序接受两个字符串作为命令行参数，执行步骤如下。

1. 调用 `crypt()` 加密第 1 个命令行参数中的字符串，并使用 `strdup()` 将结果复制到独立缓冲区中。
2. 为 `SIGINT` 信号（按下 `Ctrl-C` 产生）创建处理器函数。处理器函数调用 `crypt()` 加密第 2 个命令行参数所提供的字符串。
3. 进入无限 `for` 循环，使用 `crypt()` 加密第 1 个命令行参数中的字符串，并检查其返回字符串与第 1 步保存的结果是否一致。

在不产生信号的情况下，第 3 步中的检查结果将总是匹配。然而，一旦收到 `SIGINT` 信号，而主程序又恰在 `for` 循环内的 `crypt()` 调用之后，字符串的匹配检查之前遭到信号处理器函数的中断，这时就会发生字符串不匹配的情况。程序运行结果如下：

```
$ ./non_reentrant abc def
Repeatedly type Control-C to generate SIGINT
Mismatch on call 109871 (mismatch=1 handled=1)
Mismatch on call 128061 (mismatch=2 handled=2)
Many lines of output removed
Mismatch on call 727935 (mismatch=149 handled=156)
Mismatch on call 729547 (mismatch=150 handled=157)
Type Control-\ to generate SIGQUIT
Quit (core dumped)
```

由对上述输出 `mismatch` 和 `handled` 值的比较可知，在大多数情况下，处理器函数会在 `main()` 中的 `crypt()` 调用与字符串比较之间去覆盖静态分配的缓冲区。

程序清单 21-1：在 `main()` 以及信号处理函数中调用不可重入的函数

```
----- signals/nonreentrant.c

#define _XOPEN_SOURCE 600
#include <unistd.h>
#include <signal.h>
#include <string.h>
#include "tlpi_hdr.h"

static char *str2;          /* Set from argv[2] */
static int handled = 0;    /* Counts number of calls to handler */

static void
handler(int sig)
{
    crypt(str2, "xx");
    handled++;
}
```

```

int
main(int argc, char *argv[])
{
    char *cr1;
    int callNum, mismatch;
    struct sigaction sa;

    if (argc != 3)
        usageErr("%s str1 str2\n", argv[0]);

    str2 = argv[2];
    cr1 = strdup(crypt(argv[1], "xx")); /* Make argv[2] available to handler */
                                        /* Copy statically allocated string
                                        to another buffer */

    if (cr1 == NULL)
        errExit("strdup");

    sigemptyset(&sa.sa_mask);
    sa.sa_flags = 0;
    sa.sa_handler = handler;
    if (sigaction(SIGINT, &sa, NULL) == -1)
        errExit("sigaction");

    /* Repeatedly call crypt() using argv[1]. If interrupted by a
       signal handler, then the static storage returned by crypt()
       will be overwritten by the results of encrypting argv[2], and
       strcmp() will detect a mismatch with the value in 'cr1'. */

    for (callNum = 1, mismatch = 0; ; callNum++) {
        if (strcmp(crypt(argv[1], "xx"), cr1) != 0) {
            mismatch++;
            printf("Mismatch on call %d (mismatch=%d handled=%d)\n",
                callNum, mismatch, handled);
        }
    }
}

```

signals/nonreentrant.c

标准的异步信号安全函数

异步信号安全的函数是指当从信号处理器函数调用时，可以保证其实现是安全的。如果某一函数是可重入的，又或者信号处理器函数无法将其中断时，就称该函数是异步信号安全的。

表 21-1 所列为各种标准要求实现为异步信号安全的函数。其中，名称后未跟 v2 或 v3 字符串的函数是由 POSIX.1-1990 规定为异步信号安全的。带有 v2 标记的函数由 susv2 加入，带有 v3 标记的则由 susv3 加入。个别 UNIX 实现可能会将其他某些函数实现为异步信号安全的，但所有符合标准的 UNIX 实现都必须保证至少表中这些函数是异步信号安全的（假设由实现来提供这些函数，Linux 并未实现所有这些函数）。

SUSv4 对表 21-1 做了如下修改。

- 移除如下函数：fpathconf()、pathconf()和 sysconf()。
- 添加如下函数：execl()、execv()、faccessat()、fchmodat()、fchownat()、fexecve()、fstatat()、futimens()、linkat()、mkdirat()、mkfifoat()、mknod()、mknodat()、openat()、readlinkat()、renameat()、symlinkat()、unlinkat()、utimensat()和 utimes()。

表 21-1: POSIX.1-1990、SUSv2 和 SUSv3 规定为异步信号安全的函数

<code>_Exit()</code> (v3)	<code>getpid()</code>	<code>sigdelset()</code>
<code>_exit()</code>	<code>getppid()</code>	<code>sigemptyset()</code>
<code>abort()</code> (v3)	<code>getsockname()</code> (v3)	<code>sigfillset()</code>
<code>accept()</code> (v3)	<code>getsockopt()</code> (v3)	<code>sigismember()</code>
<code>access()</code>	<code>getuid()</code>	<code>signal()</code> (v2)
<code>aio_error()</code> (v2)	<code>kill()</code>	<code>sigpause()</code> (v2)
<code>aio_return()</code> (v2)	<code>link()</code>	<code>sigpending()</code>
<code>aio_suspend()</code> (v2)	<code>listen()</code> (v3)	<code>sigprocmask()</code>
<code>alarm()</code>	<code>lseek()</code>	<code>sigqueue()</code> (v2)
<code>bind()</code> (v3)	<code>lstat()</code> (v3)	<code>sigset()</code> (v2)
<code>cfgetispeed()</code>	<code>mkdir()</code>	<code>sigsuspend()</code>
<code>cfgetospeed()</code>	<code>mkfifo()</code>	<code>sleep()</code>
<code>cfsetispeed()</code>	<code>open()</code>	<code>socket()</code> (v3)
<code>cfsetospeed()</code>	<code>pathconf()</code>	<code>socketatmark()</code> (v3)
<code>chdir()</code>	<code>pause()</code>	<code>socketpair()</code> (v3)
<code>chmod()</code>	<code>pipe()</code>	<code>stat()</code>
<code>chown()</code>	<code>poll()</code> (v3)	<code>symlink()</code> (v3)
<code>clock_gettime()</code> (v2)	<code>posix_trace_event()</code> (v3)	<code>sysconf()</code>
<code>close()</code>	<code>pselect()</code> (v3)	<code>tcdrain()</code>
<code>connect()</code> (v3)	<code>raise()</code> (v2)	<code>tcflow()</code>
<code>creat()</code>	<code>read()</code>	<code>tcflush()</code>
<code>dup()</code>	<code>readlink()</code> (v3)	<code>tcgetattr()</code>
<code>dup2()</code>	<code>recv()</code> (v3)	<code>tcgetpgrp()</code>
<code>execle()</code>	<code>recvfrom()</code> (v3)	<code>tcsendbreak()</code>
<code>execve()</code>	<code>recvmsg()</code> (v3)	<code>tcsetattr()</code>
<code>fchmod()</code> (v3)	<code>rename()</code>	<code>tcsetpgrp()</code>
<code>fchown()</code> (v3)	<code>rmdir()</code>	<code>time()</code>
<code>fcntl()</code>	<code>select()</code> (v3)	<code>timer_getoverrun()</code> (v2)
<code>fdatasync()</code> (v2)	<code>sem_post()</code> (v2)	<code>timer_gettime()</code> (v2)
<code>fork()</code>	<code>send()</code> (v3)	<code>timer_settime()</code> (v2)
<code>fpathconf()</code> (v2)	<code>sendmsg()</code> (v3)	<code>times()</code>
<code>fstat()</code>	<code>sendto()</code> (v3)	<code>umask()</code>
<code>fsync()</code> (v2)	<code>setgid()</code>	<code>uname()</code>
<code>ftruncate()</code> (v3)	<code>setpgid()</code>	<code>unlink()</code>
<code>getegid()</code>	<code>setsid()</code>	<code>utime()</code>
<code>geteuid()</code>	<code>setsockopt()</code> (v3)	<code>wait()</code>
<code>getgid()</code>	<code>setuid()</code>	<code>waitpid()</code>
<code>getgroups()</code>	<code>shutdown()</code> (v3)	<code>write()</code>
<code>getpeername()</code> (v3)	<code>sigaction()</code>	
<code>getpgrp()</code>	<code>sigaddset()</code>	

SUSv3 强调，表 21-1 之外的所有函数对于信号而言都是不安全的，但同时指出，仅当信号处理器函数中断了不安全函数的执行，且处理器函数自身也调用了这个不安全函数时，该函数才是不安全的。换言之，编写信号处理器函数有如下两种选择。

- 确保信号处理器函数代码本身是可重入的，且只调用异步信号安全的函数。
- 当主程序执行不安全函数或是去操作信号处理器函数也可能更新的全局数据结构时，阻塞信号的传递。

第 2 种方法的问题是，在一个复杂程序中，要想确保主程序对不安全函数的调用不为信号处理器函数所中断，这有些困难。出于这一原因，通常就将上述规则简化为在信号处理器函数中绝不调用不安全的函数。

如果使用同一处理器函数来处理多个不同信号，或者在调用 `sigaction()` 时设置了 `SA_NODEFER` 标志，那么处理器函数就有可能中断自己。因此，处理器函数如果更新了全局（或静态）变量，即便主程序不使用这些变量，那么它们依然可能是不可重入的。

信号处理器函数内部对 `errno` 的使用

由于可能会更新 `errno`，调用表 21-1 中函数依然会导致信号处理器函数不可重入，因为它们可能会覆盖之前由主程序调用函数时所设置的 `errno` 值。有一种变通方法，即当信号处理器函数使用了表 21-1 所列函数时，可在其入口处保存 `errno` 值，并在其出口处恢复 `errno` 的旧有值，请看下面的例子：

```
void
handler(int sig)
{
    int savedErrno;

    savedErrno = errno;

    /* Now we can execute a function that might modify errno */

    errno = savedErrno;
}
```

在本书示例程序中使用不安全函数

虽然 `printf()` 不是异步信号安全的函数，但却频频现身于本书各种示例的信号处理器函数中。之所以如此，是因为在展示对信号处理器的调用，以及显示处理器相关变量的内容时，`printf()` 都不失为一种简明而又便捷的方式。出于类似原因，在信号处理器函数中偶尔也会用到其他一些不安全函数，包括其他的 `stdio` 函数以及 `strsignal()`。

真正的应用程序应当避免在信号处理器函数中调用非异步信号安全的函数。为了明确这一点，每当示例的信号处理器调用这些函数时，代码注释中都会注明这一用法是不安全的。

```
printf("Some message\n");          /* UNSAFE */
```

21.1.3 全局变量和 `sig_atomic_t` 数据类型

尽管存在可重入问题，有时仍需要在主程序和信号处理器函数之间共享全局变量。信号处理器函数可能会随时修改全局变量——只要主程序能够正确处理这种可能性，共享全局变量就是安全的。例如，一种常见的设计是，信号处理器函数只做一件事情，设置全局标志。主程序则会周期性地检查这一标志，并采取相应动作来响应信号传递（同时清除标志）。当信号处理器函数以此方式来访问全局变量时，应该总是在声明变量时使用 `volatile` 关键字，从而防止编译器将其优化到寄存器中。

对全局变量的读写可能不止一条机器指令，而信号处理器函数就可能会在这些指令序列之间将主程序中中断（也将此类变量访问称为非原子操作）。因此，C 语言标准以及 SUSv3 定义了一种整型数据类型 `sig_atomic_t`，意在保证读写操作的原子性。因此，所有在主程序与信号处理器函数之间共享的全局变量都应声明如下：

```
volatile sig_atomic_t flag;
```

程序清单 22-5 提供了使用 `sig_atomic_t` 数据类型的一个例子。

注意，C 语言的递增（++）和递减（--）操作符并不在 `sig_atomic_t` 所提供的保障范围之内。这些操作在某些硬件架构上可能不是原子操作（更多细节请参考 30.1 节）。在使用 `sig_atomic_t` 变量时唯一所能做的就是信号处理器中进行设置，在主程序中进行检查（反之亦可）。

C99 和 SUSv3 规定，实现应当（在 `<stdint.h>` 中）定义两个常量 `SIG_ATOMIC_MIN` 和 `SIG_ATOMIC_MAX`，用于规定可赋给 `sig_atomic_t` 类型的值范围。标准要求，如果将 `sig_atomic_t` 表示为有符号值，其范围至少应该在 -127~127 之间，如果作为无符号值，则应该在 0~255 之间。在 Linux 中，这两个常量分别等于有符号 32 位整型数的负、正极限值。

21.2 终止信号处理器函数的其他方法

目前为止所看到的信号处理器函数都是以返回主程序而终结。不过，只是简单地从信号处理器函数中返回并不能满足需要，有时候甚至没什么用处。（22.4 节在讨论硬件产生的信号时会举出这方面的例子。）

以下是从信号处理器函数中终止的其他一些方法。

- 使用 `_exit()` 终止进程。处理器函数事先可以做一些清理工作。注意，不要使用 `exit()` 来终止信号处理器函数，因为它不在表 21-1 所列的安全函数中。之所以不安全，是因为如 25.1 节所述，该函数会在调用 `_exit()` 之前刷新 `stdio` 的缓冲区。
- 使用 `kill()` 发送信号来杀掉进程（即，信号的默认动作是终止进程）。
- 从信号处理器函数中执行非本地跳转。
- 使用 `abort()` 函数终止进程，并产生核心转储。

以下各节将会对最后两点做深入讨论。

21.2.1 在信号处理器函数中执行非本地跳转

6.8 节曾论及使用 `setjmp()` 和 `longjmp()` 来执行非本地跳转，以便从一个函数跳转至该函数的某个调用者。在信号处理器函数中也可以使用这种技术。这也是因硬件异常（例如内存访问错误）而导致信号传递之后的一条恢复途径，允许将信号捕获并把控制返回到程序中某个特定位置。例如，一旦收到 `SIGINT` 信号（通常由键入 `Ctrl-C` 产生），`shell` 执行一个非本地跳转，将控制返回到主输入循环中（以便读取下一条命令）。

然而，使用标准 `longjmp()` 函数从处理器函数中退出存在一个问题。之前曾经提及，在进入信号处理器函数时，内核会自动将引发调用的信号以及由 `act.sa_mask` 所指定的任意信号添加到进程的信号掩码中，并在处理器函数正常返回时再将它们从掩码中清除。

如果使用 `longjmp()` 来退出信号处理器函数，那么信号掩码会发生什么情况呢？这取决于特定 UNIX 实现的血统。在 System V 一脉中，`longjmp()` 不会将信号掩码恢复，亦即在离开处理器函数时不会对遭阻塞的信号解除阻塞。Linux 遵循 System V 的这一特性。（这通常并非所

希望的行为，因为引发对信号处理器调用的信号仍将保持阻塞状态。) 在源于 BSD 一脉的实现中，`setjmp()`将信号掩码保存在其 `env` 参数中，而信号掩码的保存值由 `longjmp()`恢复。(继承自 BSD 的实现还提供另外两个拥有 System V 语义的函数：`_setjmp()`和`_longjmp()`。) 换言之，使用 `longjmp()`来退出信号处理器函数将有损于程序的可移植性。

如果编译程序时定义了 `_BSD_SOURCE` 特性检测宏，那么 (glibc 的) `setjmp()`将遵循 BSD 语义。

鉴于两大 UNIX 流派之间的差异，POSIX.1-1990 选择不对 `setjmp()`和 `longjmp()`的信号掩码处理进行规范，而是定义了一对新函数：`sigsetjmp()`和 `siglongjmp()`，针对执行非本地跳转时的信号掩码进行显式控制。

```
#include <setjmp.h>

int sigsetjmp(sigjmp_buf env, int savesigs);
    Returns 0 on initial call, nonzero on return via siglongjmp()
void siglongjmp(sigjmp_buf env, int val);
```

函数 `sigsetjmp()`和 `siglongjmp()`的操作与 `setjmp()`和 `longjmp()`类似。唯一的区别是参数 `env` 的类型不同 (是 `sigjmp_buf` 而不是 `jmp_buf`)，并且 `sigsetjmp()`多出一个参数 `savesigs`。如果指定 `savesigs` 为非 0，那么会将调用 `sigsetjmp()`时进程的当前信号掩码保存于 `env` 中，之后通过指定相同 `env` 参数的 `siglongjmp()`调用进行恢复。如果 `savesigs` 为 0，则不会保存和恢复进程的信号掩码。

函数 `longjmp()`和 `siglongjmp()`都不在表 21-1 所列异步信号安全函数的范围之内。因为与在信号处理器中调用这些函数一样，在执行非本地跳转之后去调用任何非异步信号安全的函数也需要冒同样的风险。此外，如果信号处理器函数中断了正在更新数据结构的主程序，那么执行非本地跳转退出处理器函数后，这种不完整的更新动作很可能会将数据结构置于不一致状态。规避这一问题的一种技术是在程序对敏感数据进行更新时，借助于 `sigprocmask()`临时将信号阻塞起来。

示例程序

程序清单 21-2 展示了两种类型的非本地跳转在处理信号掩码上的差异。该程序为 `SIGINT` 创建处理器函数，并允许选择 `setjmp()+longjmp()`组合或者 `sigsetjmp()+siglongjmp()`组合的方式来退出信号处理器函数，具体采用何种函数组合则取决于程序编译时是否对宏 `USE_SIGSETJMP` 进行了定义。程序会分别在进入信号处理器函数时，以及非本地跳转将控制从信号处理器交还给主程序后，显示信号掩码的当前设置。

如果利用 `longjmp()`来退出信号处理器函数，其结果如下：

```
$ make -s sigmask_longjmp          Default compilation causes setjmp() to be used
$ ./sigmask_longjmp
Signal mask at startup:
    <empty signal set>
Calling setjmp()
Type Control-C to generate SIGINT
Received signal 2 (Interrupt), signal mask is:
    2 (Interrupt)
After jump from handler, signal mask is:
    2 (Interrupt)
(At this point, typing Control-C again has no effect, since SIGINT is blocked)
Type Control-\ to kill the program
Quit
```

由程序输出结果可知，信号处理器函数调用 `longjmp()` 之后的信号掩码设置与进入处理器函数时保持一致。

上述 shell 会话中构建程序所使用的 `makefile` 由随本书发布的源码提供。选项 `-s` 告知 `make` 程序不要显示正在执行的命令。使用该选项意在避免对会话日志的显示产生干扰。（[Mecklenbug, 2005]对 GNU `make` 程序做了说明。）

如果编译同一源文件来创建利用 `siglongjmp()` 退出信号处理器函数的程序，则结果如下：

```
$ make -s sigmask_siglongjmp      Compiles using cc -DUSE_SIGSETJMP
$ ./sigmask_siglongjmp x
Signal mask at startup:
    <empty signal set>
Calling sigsetjmp()
Type Control-C
Received signal 2 (Interrupt), signal mask is:
    2 (Interrupt)
After jump from handler, signal mask is:
    <empty signal set>
```

在这里，没有将 `SIGINT` 信号阻塞，因为 `siglongjmp()` 恢复了原来的信号掩码。接着，再次按下 `Ctrl-C`，会再次调用该信号处理器函数。

```
Type Control-C
Received signal 2 (Interrupt), signal mask is:
    2 (Interrupt)
After jump from handler, signal mask is:
    <empty signal set>
Type Control-\ to kill the program
Quit
```

由上述输出可知，`siglongjmp()` 将信号掩码恢复到调用 `sigsetjmp()` 时的值（即一个空信号集）。

程序清单 21-2 还展示了信号处理器函数执行非本地跳转时的一种实用技术。信号随时可能产生，所以有可能发生于 `sigsetjmp()`（或 `setjmp()`）设置跳转目标之前。为杜绝这种可能（这将导致处理器函数使用未初始化的 `env` 缓冲区来执行非本地跳转），程序启用了守卫变量 `canJump`，来表征 `env` 缓冲区的初始化与否。如果 `canJump` 不为真（`false`），处理器函数将不执行跳转而直接返回。另一种方法是调整程序代码，在创建信号处理器函数之前去调用 `sigsetjmp()`（或 `setjmp()`）。不过对于复杂的程序而言，苛求这样的步骤执行顺序可能会有困难，而使用守卫变量也许会更简单一些。

注意，在编写程序清单 21-2 程序时使用 `#ifndef` 是使其编码风格符合标准的最简单的手段。特别是当无法用下面的运行时检查代码来取代 `#ifdef` 时。

```
if (useSiglongjmp)
    s = sigsetjmp(senv, 1);
else
    s = setjmp(env);
if (s == 0)
    ...
```

这一做法有违规范，因为 `SUSv3` 不允许在赋值语句（6.8 节）中调用 `setjmp()` 和 `sigsetjmp()`。

程序清单 21-2：在信号处理器函数中执行非本地跳转

```
----- signals/sigmask_longjmp.c
#define _GNU_SOURCE      /* Get strsignal() declaration from <string.h> */
#include <string.h>
#include <setjmp.h>
```

```

#include <signal.h>
#include "signal_functions.h"          /* Declaration of printSigMask() */
#include "tlpi_hdr.h"

static volatile sig_atomic_t canJump = 0;
                                /* Set to 1 once "env" buffer has been
                                initialized by [sig]setjmp() */

#ifdef USE_SIGSETJMP
static sigjmp_buf senv;
#else
static jmp_buf env;
#endif

static void
handler(int sig)
{
    /* UNSAFE: This handler uses non-async-signal-safe functions
       (printf(), strsignal(), printSigMask()); see Section 21.1.2) */

    printf("Received signal %d (%s), signal mask is:\n", sig,
           strsignal(sig));
    printSigMask(stdout, NULL);

    if (!canJump) {
        printf("'env' buffer not yet set, doing a simple return\n");
        return;
    }

#ifdef USE_SIGSETJMP
    siglongjmp(senv, 1);
#else
    longjmp(env, 1);
#endif
}

int
main(int argc, char *argv[])
{
    struct sigaction sa;

    printSigMask(stdout, "Signal mask at startup:\n");

    sigemptyset(&sa.sa_mask);
    sa.sa_flags = 0;
    sa.sa_handler = handler;
    if (sigaction(SIGINT, &sa, NULL) == -1)
        errExit("sigaction");

#ifdef USE_SIGSETJMP
    printf("Calling sigsetjmp()\n");
    if (sigsetjmp(senv, 1) == 0)
#else
    printf("Calling setjmp()\n");
    if (setjmp(env) == 0)
#endif
        canJump = 1;                                /* Executed after [sig]setjmp() */

    else                                            /* Executed after [sig]longjmp() */
        printSigMask(stdout, "After jump from handler, signal mask is:\n");
}

```

```
for (;;)                                /* Wait for signals until killed */
    pause();
}
```

signals/sigmask_longjmp.c

21.2.2 异常终止进程：abort()

函数 `abort()` 终止其调用进程，并生成核心转储。

```
#include <stdlib.h>

void abort(void);
```

函数 `abort()` 通过产生 `SIGABRT` 信号来终止调用进程。对 `SIGABRT` 的默认动作是产生核心转储文件并终止进程。调试器可以利用核心转储文件来检测调用 `abort()` 时的程序状态。

`SUSv3` 要求，无论阻塞或者忽略 `SIGABRT` 信号，`abort()` 调用均不受影响。同时规定，除非进程捕获 `SIGABRT` 信号后信号处理器函数尚未返回，否则 `abort()` 必须终止进程。后一句话值得三思。21.2 节所描述的信号处理器函数终止方法中，与此相关的就是使用非本地跳转退出处理器函数。这一做法将抵消 `abort()` 的效果。否则，`abort()` 将总是终止进程。在大多数实现中，终止时可确保发生如下事件：若进程在发出一次 `SIGABRT` 信号后仍未终止（即，处理器捕获信号并返回，以便恢复执行 `abort()`），则 `abort()` 会将对 `SIGABRT` 信号的处理重置为 `SIG_DFL`，并再度发出 `SIGABRT` 信号，从而确保将进程杀死。

如果 `abort()` 成功终止了进程，那么还将刷新 `stdio` 流并将其关闭。

程序清单 3-3 在错误处理函数中提供了使用 `abort()` 的一个例子。

21.3 在备选栈中处理信号：sigaltstack()

在调用信号处理器函数时，内核通常会在进程栈中为其创建一帧。不过，如果进程对栈的扩展突破了对栈大小的限制时，这种做法就不大可行。例如，栈的增长过大，以至于会触及到一片映射内存（48.5 节）或者向上增长的堆，又或者栈的大小已经直逼 `RLIMIT_STACK`（36.3 节）资源限制，这些都会造成这种情况的发生。

当进程对栈的扩展试图突破其上限时，内核将为该进程产生 `SIGSEGV` 信号。不过，因为栈空间已然耗尽，内核也就无法为进程已经安装的 `SIGSEGV` 处理器函数创建栈帧。结果是，处理器函数得不到调用，而进程也就终止了（`SIGSEGV` 的默认动作）。

如果希望在这种情况下确保对 `SIGSEGV` 信号处理器函数的调用，就需要做如下工作。

1. 分配一块被称为“备选信号栈”的内存区域，作为信号处理器函数的栈帧。
2. 调用 `sigaltstack()`，告之内核该备选信号栈的存在。
3. 在创建信号处理器函数时指定 `SA_ONSTACK` 标志，亦即通知内核在备选栈上为处理器函数创建栈帧。

利用系统调用 `sigaltstack()`，既可以创建一个备选信号栈，也可以将已创建备选信号栈的相关信息返回。

```
#include <signal.h>
```

```
int sigaltstack(const stack_t *sigstack, stack_t *old_sigstack);
```

Returns 0 on success, or -1 on error

参数 `sigstack` 所指向的数据结构描述了新备选信号栈的位置及属性。参数 `old_sigstack` 指向的结构则用于返回上一备选信号栈的相关信息（如果存在）。两个参数之一均可为 `NULL`。例如，将参数 `sigstack` 设为 `NULL` 可以发现现有备选信号栈，并且不用将其改变。不为 `NULL` 时，这些参数所指向的数据结构类型如下：

```
typedef struct {  
    void *ss_sp;          /* Starting address of alternate stack */  
    int   ss_flags;      /* Flags: SS_ONSTACK, SS_DISABLE */  
    size_t ss_size;     /* Size of alternate stack */  
} stack_t;
```

字段 `ss_sp` 和 `ss_size` 分别指定了备选信号栈的位置和大小。在实际使用信号栈时，内核会将 `ss_sp` 值自动对齐为与硬件架构相适宜的地址边界。

备选信号栈通常既可以静态分配，也可以在堆上动态分配。`SUSv3` 规定将常量 `SIGSTKSZ` 作为划分备选栈大小的典型值，而将 `MINSSIGSTKSZ` 作为调用信号处理器函数所需的最小值。在 `Linux/x86-32` 系统上，分别将这两个值定义为 8192 和 2048。

内核不会重新划分备选栈的大小。如果栈溢出了分配给它的空间，就会产生混乱（例如，写变量超出了对栈的限制）。这通常不是一个问题，因为一般情况下会利用备选栈来处理标准栈溢出的特殊情况，常常只在这个栈上分配为数不多的几帧。`SIGSEGV` 处理器函数的工作不是在执行清理动作后终止进程，就是使用非本地跳转解开标准栈。

`ss_flags` 可以包含如下值之一：

SS_ONSTACK

如果在获取已创建备选信号栈的当前信息时该标志已然置位，就表明进程正在备选信号栈上执行。当进程已经在备选信号栈上运行时，试图调用 `sigaltstack()` 来创建一个新的备选信号栈将会产生一个错误（`EPERM`）。

SS_DISABLE

在 `old_sigstack` 中返回，表示当前不存在已创建的备选信号栈。如果在 `sigstack` 中指定，则会禁用当前已创建的备选信号栈。

程序清单 21-3 演示了备选信号栈的创建和使用。在创建一个新的备选信号栈以及 `SIGSEGV` 的信号处理器函数之后，程序将调用一个无限递归函数，这会导致栈溢出，同时系统会向进程发送 `SIGSEGV` 信号。运行该程序的结果如下。

```
$ ulimit -s unlimited  
$ ./t_sigaltstack  
Top of standard stack is near 0xbffff6b8  
  
Alternate stack is at          0x804a948-0x804cfff  
Call   1 - top of stack near 0xbff0b3ac  
Call   2 - top of stack near 0xbfe1714c  
Many intervening lines of output removed  
Call 2144 - top of stack near 0x4034120c  
Call 2145 - top of stack near 0x4024cfac  
Caught signal 11 (Segmentation fault)  
Top of handler stack near      0x804c860
```


在这一 shell 会话中，命令 ulimit 负责移除 shell 之前可能设置的任何 RLIMIT_STACK 资源限制。36.3 节会解释这种资源限制。

程序清单 21-3: 使用 sigaltstack()

```
-----signals/t_sigaltstack.c
#define _GNU_SOURCE          /* Get strsignal() declaration from <string.h> */
#include <string.h>
#include <signal.h>
#include "tspi_hdr.h"

static void
sigsegvHandler(int sig)
{
    int x;

    /* UNSAFE: This handler uses non-async-signal-safe functions
       (printf(), strsignal(), fflush()); see Section 21.1.2) */

    printf("Caught signal %d (%s)\n", sig, strsignal(sig));
    printf("Top of handler stack near    %10p\n", (void *) &x);
    fflush(NULL);

    _exit(EXIT_FAILURE);          /* Can't return after SIGSEGV */
}

static void          /* A recursive function that overflows the stack */
overflowStack(int callNum)
{
    char a[100000];          /* Make this stack frame large */

    printf("Call %4d - top of stack near %10p\n", callNum, &a[0]);
    overflowStack(callNum+1);
}

int
main(int argc, char *argv[])
{
    stack_t sigstack;
    struct sigaction sa;
    int j;

    printf("Top of standard stack is near %10p\n", (void *) &j);

    /* Allocate alternate stack and inform kernel of its existence */

    sigstack.ss_sp = malloc(SIGSTKSZ);
    if (sigstack.ss_sp == NULL)
        errExit("malloc");
    sigstack.ss_size = SIGSTKSZ;
    sigstack.ss_flags = 0;
    if (sigaltstack(&sigstack, NULL) == -1)
        errExit("sigaltstack");
    printf("Alternate stack is at    %10p-%p\n",
           sigstack.ss_sp, (char *) sbrk(0) - 1);

    sa.sa_handler = sigsegvHandler;          /* Establish handler for SIGSEGV */
    sigemptyset(&sa.sa_mask);
}
```

```

    sa.sa_flags = SA_ONSTACK;          /* Handler uses alternate stack */
    if (sigaction(SIGSEGV, &sa, NULL) == -1)
        errExit("sigaction");

    overflowStack(1);
}

```

signals/t_sigaltstack.c

21.4 SA_SIGINFO 标志

如果在使用 `sigaction()` 创建处理器函数时设置了 `SA_SIGINFO` 标志，那么在收到信号时处理器函数可以获取该信号的一些附加信息。为获取这一信息，需要将处理器函数声明如下：

```
void handler(int sig, siginfo_t *siginfo, void *ucontext);
```

如同标准信号处理器函数一样，第 1 个参数 `sig` 表示信号编号。第 2 个参数 `siginfo` 是用于提供信号附加信息的一个结构。该结构会与最后一个参数 `ucontext` 一起，在下面做详细说明。

因为上述信号处理器函数的原型不同于标准处理器函数，依照 C 语言的类型规则，将无法利用 `sigaction` 结构的 `sa_handler` 字段来指定处理器函数地址。此时需要使用另一个字段：`sa_sigaction`。换言之，`sigaction` 结构比 20.13 节所展示的要稍微复杂一些。其完整定义如下：

```

struct sigaction {
    union {
        void (*sa_handler)(int);
        void (*sa_sigaction)(int, siginfo_t *, void *);
    } __sigaction_handler;
    sigset_t  sa_mask;
    int       sa_flags;
    void      (*sa_restorer)(void);
};

/* Following defines make the union fields look like simple fields
   in the parent structure */

#define sa_handler __sigaction_handler.sa_handler
#define sa_sigaction __sigaction_handler.sa_sigaction

```

结构 `sigaction` 使用联合体来合并 `sa_sigaction` 和 `sa_handler`。（大部分其他 UNIX 实现也采用相同的方式。）之所以使用联合体，是因为对 `sigaction()` 的特定调用只会用到其中的一个字段。（不过，如果天真地认为可以彼此独立地设置 `sa_handler` 和 `sa_sigaction`，就有可能导致一些奇怪的 bug。可能的原因是在为不同的信号创建处理器函数时，多次对 `sigaction()` 的调用复用了同一个 `sigaction` 结构。）

这里是使用 `SA_SIGINFO` 创建信号处理器函数的一个例子：

```

struct sigaction act;

sigemptyset(&act.sa_mask);
act.sa_sigaction = handler;
act.sa_flags = SA_SIGINFO;

if (sigaction(SIGINT, &act, NULL) == -1)
    errExit("sigaction");

```

至于使用 `SA_SIGINFO` 标志的完整例子，请参考程序清单 22-3 和程序清单 23-5。

结构 `siginfo_t`

在以 `SA_SIGINFO` 标志创建的信号处理器函数中，结构 `siginfo_t` 是其第 2 个参数，格式如下：

```
typedef struct {
    int    si_signo;        /* Signal number */
    int    si_code;        /* Signal code */
    int    si_trapno;      /* Trap number for hardware-generated signal
                           (unused on most architectures) */
    union sigval si_value; /* Accompanying data from sigqueue() */
    pid_t  si_pid;        /* Process ID of sending process */
    uid_t  si_uid;        /* Real user ID of sender */
    int    si_errno;      /* Error number (generally unused) */
    void   *si_addr;      /* Address that generated signal
                           (hardware-generated signals only) */

    int    si_overrun;     /* Overrun count (Linux 2.6, POSIX timers) */
    int    si_timerid;     /* (Kernel-internal) Timer ID
                           (Linux 2.6, POSIX timers) */

    long   si_band;       /* Band event (SIGPOLL/SIGIO) */
    int    si_fd;         /* File descriptor (SIGPOLL/SIGIO) */
    int    si_status;     /* Exit status or signal (SIGCHLD) */
    clock_t si_utime;     /* User CPU time (SIGCHLD) */
    clock_t si_stime;     /* System CPU time (SIGCHLD) */
} siginfo_t;
```

要获取 `<signal.h>` 对 `siginfo_t` 的声明，必须将特性测试宏 `_POSIX_C_SOURCE` 的值定义为大于或等于 199309。

如同大部分 UNIX 实现一样，在 Linux 系统中，`siginfo_t` 结构的很多字段都是联合体，因为对每个信号而言，并非所有字段都有必要。（参考 `<bits/siginfo.h>` 中的细节。）

一旦进入信号处理器函数，对结构 `siginfo_t` 中字段的设置如下。

`si_signo`

需要为所有信号设置。内含引发处理器函数调用的信号编号——与处理器函数 `sig` 参数的值相同。

`si_code`

需要为所有信号设置。如表 21-2 所示，所含代码提供了关于信号来源的深入信息。

`si_value`

该字段包含调用 `sigqueue()` 发送信号时的伴随数据。22.8.1 节将讨论 `sigqueue()`。

`si_pid`

对于经由 `kill()` 或 `sigqueue()` 发送的信号，该字段保存了发送进程的进程 ID。

`si_uid`

对于经由 `kill()` 或 `sigqueue()` 发送的信号，该字段保存了发送进程的真实用户 ID。系统之所以提供真实用户 ID，是因为其信息量比之有效用户 ID 更为丰富。回忆 20.5 节所述关于信号发送的权限规则，如果有效用户 ID 授予发送者发送信号的权力，那么发送方的用户 ID 必须要么为 0（特权级用户），要么与接收进程的真实用户 ID 或者保存设置用户 ID（`saved set-user-ID`）相同。这时，接收者了解发送者的真实用户 ID 就很有用，因为它有可能不同于有效用户 ID（例如，如果发送者是一个 `set-user-ID` 程序）。

`si_errno`

如果将该字段置为非 0 值，则其所包含为一错误号（类似 `errno`），标志信号的产生原因。

Linux 通常不使用该字段。

si_addr

仅针对由硬件产生的 SIGBUG、SIGSEGV、SIGILL 和 SIGFPE 信号设置该字段。对于 SIGBUS 和 SIGSEGV 而言，该字段内含引发无效内存引用的地址。对于 SIGILL 和 SIGFPE 信号，则包含导致信号产生的程序指令地址。

以下各字段均属非标准的 Linux 扩展，仅当 POSIX 定时器（23.6 节）到期而产生信号传递时设置：

si_timerid

内含供内核内部使用的 ID，用以标识定时器。

si_overrun

设置该字段为定时器的溢出次数。

仅当收到 SIGIO 信号（63.3 节）时，才会设置下面两个字段。

si_band

该字段包含与 I/O 事件相关的“带事件”值。（直到 glibc 2.3.2，si_band 的类型都是 int 型。）

si_fd

该字段包含与 I/O 事件相关的文件描述符编号。SUSv3 并未定义这一字段，不过许多其他实现都予以了支持。

仅当收到 SIGCHLD 信号（26.3 节）时，才会对以下各字段进行设置。

si_status

该字段包含子进程的退出状态（当 si_code=CLD_EXITED 时）或者发给子进程的信号编号（即 26.1.3 节所述终止或停止子进程的信号编号）。

si_utime

该字段包含子进程使用的用户 CPU 时间。在版本 2.6 以前，以及 2.6.27 以后的内核版本中，对该字段的度量以系统时钟滴答除以 sysconf (_SC_CLK_TCK) 的返回值作为基本单位。而在版本 2.6.27 之前的 2.6 内核中则存在 bug，该字段在报告时间时采用的度量单位为（可由用户配置的）jiffy（10.6 节）。SUSv3 没有定义该字段，但许多其他实现都予以支持。

si_stime

该字段包含了子进程使用的系统 CPU 时间。可参考对 si_utime 的描述。同样，SUSv3 并未定义该字段，不过许多其他实现都予以支持。

si_code 字段提供了关于信号来源的更多信息，其值如表 21-2 所示。表中第 2 列列出的信号特有值（特别是由硬件产生的 4 种信号：SIGBUS、SIGSEGV、SIGILL 和 SIGFPE）不会悉数现身于所有的 UNIX 实现以及硬件架构之上——尽管 Linux 定义了所有常量，而且 SUSv3 也定义了其中的大部分。

关于表 21-2 中所示各值，还需注意以下几点附加说明。

- 值 SI_KERNEL 和 SI_SIGIO 为 Linux 所特有，既未获 SUSv3 定义，也未获其他 UNIX 实现支持。
- SI_SIGIO 仅在 Linux2.2 中用到。自内核 2.4 起，Linux 转而采用表中的 POLL_* 常量。

表 21-2: 结构 `siginfo_t` 中 `si_code` 字段返回值一览表

信 号	si_code 的值	信 号 来 源
任意 (所有)	SI_ASYNCIO	异步 I/O (AIO) 操作已经完成
	SI_KERNEL	从内核发送 (例如, 来自于终端驱动程序信号)
	SI_MESGQ	消息到达 POSIX 消息队列 (自 Linux 2.6.6)
	SI_QUEUE	利用 <code>sigqueue()</code> 从用户进程发出的实时信号
	SI_SIGIO	SIGIO 信号 (仅 Linux 2.2 支持)
	SI_TIMER	POSIX (实时) 定时器到期
	SI_TKILL	调用 <code>tkill()</code> 或 <code>tgkill()</code> 的用户进程 (自 Linux 2.4.19)
	SI_USER	调用 <code>kill()</code> 或 <code>raise()</code> 的用户进程
SIGBUS	BUS_ADRALN	无效的地址对齐
	BUS_ADRERR	不存在的物理地址
	BUS_MCEERR_AO	硬件内存错误, 动作为可选 (自 Linux 2.6.32)
	BUS_MCEERR_AR	硬件内存错误, 动作为必需 (自 Linux 2.6.32)
	BUS_OBJERR	对象特有的硬件错误
SIGCHLD	CLD_CONTINUED	因 SIGCONT 信号, 子进程得以继续执行 (自 Linux 2.6.9)
	CLD_DUMPED	子进程异常终止, 并产生核心转储
	CLD_EXITED	子进程退出
	CLD_KILLED	子进程异常终止, 且不产生核心转储
	CLD_STOPPED	子进程停止
	CLD_TRAPPED	受到跟踪的子进程停止
SIGFPE	FPE_FLTDIV	浮点除 0
	FPE_FLTINV	无效的浮点操作
	FPE_FLTOVF	浮点溢出
	FPE_FLTRES	浮点结果不精确
	FPE_FLTUND	浮点下溢
	FPE_INTDIV	整型除 0
	FPE_INTOVF	整型溢出
	FPE_SUB	下标超出范围
SIGILL	ILL_BADSTK	内部栈错误
	ILL_COPROC	协处理器错误
	ILL_ILLADR	非法地址模式
	ILL_ILLOPC	非法操作码
	ILL_ILLOPN	非法操作数
	ILL_ILLTRP	非法陷入
	ILL_PRVOPC	特权级操作码
	ILL_PRIVREG	特权级寄存器

信 号	si_code 的值	信 号 来 源
SIGPOLL/ SIGIO	POLL_ERR	I/O 错误
	POLL_HUP	设备断开
	POLL_IN	输入数据有效
	POLL_MSG	输入消息有效
	POLL_OUT	输出缓冲区有效
	POLL_PRI	高优先级输入有效
SIGSEGV	SEGV_ACCERR	映射对象的无效权限
	SEGV_MAPERR	未映射为对象的地址
SIGTRAP	TRAP_BRANCH	进程分支陷入
	TRAP_BRKPT	进程断点
	TRAP_HWBKPT	硬件断点/监测点
	TRAP_TRACE	进程跟踪陷入

SUSv4 定义了功用与 `psignal()` (20.8 节) 相仿的 `psiginfo()` 函数。函数 `psiginfo()` 带有两个参数，分别是指向 `siginfo_t` 结构的指针和一个消息字符串。该函数在标准错误设备上输出字符串消息，接着显示描述于 `siginfo_t` 结构中的信号信息。`glibc` 自 2.10 版开始提供 `psiginfo()` 函数。`glibc` 实现会显示信号的描述信息及来源（根据 `si_code` 字段所示），对于某些信号，还会列出 `siginfo_t` 结构中的其他字段。函数 `psiginfo()` 是 SUSv4 中的新丁，并非所有系统都予以支持。

参数 `ucontext`

以 `SA_SIGINFO` 标志所创建的信号处理器函数，其最后一个参数是 `ucontext`，一个指向 `ucontext_t` 类型结构（定义于 `<ucontext.h>`）的指针。（因为 SUSv3 并未规定该参数的任何细节，所以将其定义为 `void` 类型指针。）该结构提供了所谓的用户上下文信息，用于描述调用信号处理器函数前的进程状态，其中包括上一个进程信号掩码以及寄存器的保存值，例如程序计数器（`cp`）和栈指针寄存器（`sp`）。信号处理器函数很少用到此类信息，所以此处也略而不论。

使用结构 `ucontext_t` 的其他函数有 `getcontext()`、`makecontext()`、`setcontext()` 和 `swapcontext()`，分别对应的功能是允许进程去接收、创建、改变以及交换执行上下文。（这些操作有点类似于 `setjmp()` 和 `longjmp()`，但更为通用。）可以使用这些函数来实现协程（`coroutines`），令进程的执行线程在两个（或多个）函数之间交替。SUSv3 规定了这些函数，但将它们标记为已废止。SUSv4 则将其删去，并建议使用 POSIX 线程来重写旧有的应用程序。`glibc` 手册页提供了关于这些函数的深入信息。

21.5 系统调用的中断和重启

考虑如下场景。

1. 为某信号创建处理器函数。
2. 发起一个阻塞的系统调用 (blocking system call)，例如，从终端设备调用的 `read()` 就会阻塞到有数据输入为止。
3. 当系统调用遭到阻塞时，之前创建了处理器函数的信号传递了过来，随即引发对处理器函数的调用。

信号处理器返回后又会发生什么？默认情况下，系统调用失败，并将 `errno` 置为 `EINTR`。这是一种有用的特性。23.3 节将会描述如何使用定时器（会产生 `SIGALRM` 信号）来设置像 `read()` 之类阻塞系统调用的超时。

不过，更为常见的情况是希望遭到中断的系统调用得以继续运行。为此，可在系统调用遭信号处理器中断的事件中，利用如下代码来手动重启系统调用。

```
while ((cnt = read(fd, buf, BUF_SIZE)) == -1 && errno == EINTR)
    continue;          /* Do nothing loop body */

if (cnt == -1)         /* read() failed with other than EINTR */
    errExit("read");
```

如果需要频繁使用上述代码，那么定义成如下宏会很方便：

```
#define NO_EINTR(stmt) while ((stmt) == -1 && errno == EINTR);
```

使用该宏，可以将原先对 `read()` 的调用改写如下：

```
NO_EINTR(cnt = read(fd, buf, BUF_SIZE));
```

```
if (cnt == -1)         /* read() failed with other than EINTR */
    errExit("read");
```

GNU C 库提供了一个（非标准）宏，其作用与定义于 `<unistd.h>` 中的 `NO_EINTR()` 相同。该宏名为 `TEMP_FAILURE_RETRY()`，定义特性测试宏 `_GNU_SOURCE` 后即可使用。

即使采用了类似 `NO_EINTR()` 这样的宏，让信号处理器来中断系统调用还是颇为不便，因为只要有意重启阻塞的调用，就需要为每个阻塞的系统调用添加代码。反之，可以调用指定了 `SA_RESTART` 标志的 `sigaction()` 来创建信号处理器函数，从而令内核代表进程自动重启系统调用，还无需处理系统调用可能返回的 `EINTR` 错误。

标志 `SA_RESTART` 是针对每个信号的设置。换言之，允许某些信号的处理器函数中断阻塞的系统调用，而其他系统调用则可以自动重启。

SA_RESTART 标志对哪些系统调用（和库函数）有效

不幸的是，并非所有的系统调用都可以通过指定 `SA_RESTART` 来达到自动重启的目的。究其原因，有部分历史因素。

- 4.2BSD 引入了重启系统调用的概念，包括中断对 `wait()` 和 `waitpid()` 的调用，以及如下 I/O 系统调用：`read()`、`readv()`、`write()` 和阻塞的 `ioctl()` 操作。I/O 系统调用都是可中断的，所以只有在操作“慢速 (slow)”设备时，才可以利用 `SA_RESTART` 标志来自动重启调用。慢速设备包括终端 (terminal)、管道 (pipe)、FIFO 以及套接字 (socket)。对于这

些文件类型，各种 I/O 操作都有可能堵塞。（相比之下，磁盘文件并未沦入慢速设备之列，因为借助于缓冲区高速缓存，磁盘 I/O 请求一般都可以立即得到满足。当出现磁盘 I/O 请求时，内核会令该进程休眠，直至完成 I/O 动作为止。）

- 其他大量阻塞的系统调用则继承自 System V，在其初始设计中并未提供重启系统调用的功能。

在 Linux 中，如果采用 SA_RESTART 标志来创建系统处理器函数，则如下阻塞的系统调用（以及构建于其上的库函数）在遭到中断时是可以自动重启的。

- 用来等待子进程（26.1 节）的系统调用：wait()、waitpid()、wait3()、wait4()和 waitid()。
- 访问慢速设备时的 I/O 系统调用：read()、readv()、write()、writev()和 ioctl()。如果在收到信号时已经传递了部分数据，那么还是会中断输入输出系统调用，但会返回成功状态：一个整型值，表示已成功传递数据的字节数。
- 系统调用 open()，在可能阻塞的情况下（例如，如 44.7 节所述，在打开 FIFO 时）。
- 用于套接字的各种系统调用：accept()、accept4()、connect()、send()、sendmsg()、sendto()、recv()、recvfrom()和 recvmsg()。（在 Linux 中，如果使用 setsockopt()来设置超时，这些系统调用就不会自动重启。更多细节请参考 signal(7)手册页。）
- 对 POSIX 消息队列进行 I/O 操作的系统调用：mq_receive()、mq_timedreceive()、mq_send()和 mq_timedsend()。
- 用于设置文件锁的系统调用和库函数：flock()、fcntl()和 lockf()。
- Linux 特有系统调用 futex()的 FUTEX_WAIT 操作。
- 用于递减 POSIX 信号量的 sem_wait()和 sem_timedwait()函数。（在一些 UNIX 实现上，如果设置了 SA_RESTART 标志，sem_wait()就会重启。）
- 用于同步 POSIX 线程的函数：pthread_mutex_lock()、pthread_mutex_trylock()、pthread_mutex_timedlock()、pthread_cond_wait()和 pthread_cond_timedwait()。

内核 2.6.22 之前，不管是否设置了 SA_RESTART 标志，futex()、sem_wait()和 sem_timedwait()遭到中断时总是产生 EINTR 错误。

以下阻塞的系统调用（以及构建于其上的库函数）则绝不会自动重启（即便指定了 SA_RESTART）。

- poll()、ppoll()、select()和 pselect()这些 I/O 多路复用调用。（SUSv3 明文规定，无论设置 SA_RESTART 标志与否，都不对 select()和 pselect()遭处理器函数中断时的行为进行定义。）
- Linux 特有的 epoll_wait()和 epoll_pwait()系统调用。
- Linux 特有的 io_getevents()系统调用。
- 操作 System V 消息队列和信号量的阻塞系统调用：semop()、semopop()、msgrecv()和 msgsnd()。（虽然 System V 原本并未提供自动重启系统调用的功能，但在某些 UNIX 实现上，如果设置了 SA_RESTART 标志，这些系统调用还是会自动重启。）
- 对 inotify 文件描述符发起的 read()调用。
- 用于将进程挂起指定时间的系统调用和库函数：sleep()、nanosleep()和 clock_nanosleep()。
- 特意设计用来等待某一信号到达的系统调用：pause()、sigsuspend()、sigtimedwait()和 sigwaitinfo()。

为信号修改 SA_RESTART 标志

函数 siginterrupt()用于改变信号的 SA_RESTART 设置。


```
#include <signal.h>

int siginterrupt(int sig, int flag);
```

Returns 0 on success, or -1 on error

若参数 `flag` 为真 (1)，则针对信号 `sig` 的处理器函数将会中断阻塞的系统调用的执行。如果 `flag` 为假 (0)，那么在执行了 `sig` 的处理器函数之后，会自动重启阻塞的系统调用。

函数 `siginterrupt()` 的工作原理是：调用 `sigaction()` 获取信号当前处置的副本，调整自结构 `oldact` 中返回的 `SA_RESTART` 标志，接着再次调用 `sigaction()` 来更新信号处置。

SUSv4 标记 `sigterrupt()` 为已废止，并推荐使用 `sigaction()` 加以替代。

对于某些 Linux 系统调用，未处理的停止信号会产生 EINTR 错误

在 Linux 上，即使没有信号处理器函数，某些阻塞的系统调用也会产生 `EINTR` 错误。如果系统调用遭到阻塞，并且进程因信号 (`SIGSTOP`、`SIGTSTP`、`SIGTTIN` 或 `SIGTTOU`) 而停止，之后又因收到 `SIGCONT` 信号而恢复执行时，就会发生这种情况。

以下系统调用和函数具有这一行为：`epoll_pwait()`、`epoll_wait()`、对 `inotify` 文件描述符执行的 `read()` 调用、`semop()`、`semtimedop()`、`sigtimedwait()` 和 `sigwaitinfo()`。

内核 2.6.24 之前，`poll()` 也曾存在这种行为，2.6.22 之前的 `sem_wait()`、`sem_timedwait()`、`futex(FUTEX_WAIT)`，2.6.9 之前的 `msgrcv()` 和 `msgsnd()`，以及 Linux 2.4 及其之前的 `nanosleep()` 也同样如此。

在 Linux 2.4 及其之前的版本中，也可以以这种方式来中断 `sleep()`，但是不会返回错误值，而是返回休眠所剩余的秒数。

这种行为的结果是，如果程序可能因信号而停止和重启，那么就需要添加代码来重新启动这些系统调用，即便该程序并未为停止信号设置处理器函数。

21.6 总结

本章讨论了影响信号处理器函数操作与设计的一系列因素。

由于没有对信号排队，故而在为处理器编码时，有时必须要考虑特定类型信号多次发生的可能性，即使之前信号只产生过一次。可重入问题会影响到对全局变量的修改方式，还限制了可从信号处理器函数中安全调用的函数范围。

除了返回之外，信号处理器函数的终止还存在多种其他方法，其中包括：调用 `_exit()`，发送信号来终止进程 (`kill()`、`raise()` 或 `abort()`)，或者执行非本地跳转。借助于 `sigsetjmp()` 和 `siglongjmp()`，可以在执行非本地跳转时为程序提供处理信号掩码的显式控制手段。

可以使用 `sigaltstack()` 来为进程定义备选信号栈。这是调用信号处理器函数时，用来替代标准进程栈的一块内存。当标准栈因增长过大（内核会在此时向进程发送 `SIGSEGV` 信号）而消耗殆尽时，备选栈就特别有用。

如果在调用 `sigaction()` 时设置了 `SA_SIGINFO` 标志，那么所创建的信号处理器函数就能接收信号的附加信息。`siginfo_t` 结构提供了这些信息，其地址则传递给信号处理器作为参数。

如果信号处理器函数中断了阻塞的系统调用，系统调用会产生 `EINTR` 错误。利用这种特性，就可以为阻塞的系统调用设置一个定时器。如果有意，可以手动重启遭到中断的系统调

用。另外，在调用 `sigaction()` 创建信号处理器函数时，如果设置了 `SA_RESTART` 标志，那么大部分（但非全部）系统调用都将会自动重启。

更多信息

参考 20.15 节所列信息来源。

21.7 练习

21.1. 实现 `abort()`。

第 22 章

信号：高级特性

本章是“信号”主题系列讨论（始于第 20 章）的完结篇，涵盖了一些更为高级的议题，如下所示。

- 核心转储文件。
- 与信号的传递、处置及处理相关的特殊情况。
- 信号的同步产生和异步产生。
- 信号的传递时机及传递顺序。
- 信号处理器函数对系统调用的中断，以及如何自动重启遭到中断的系统调用。
- 实时信号。
- 用 `sigsuspend()` 来设置进程信号掩码并等待信号到达。
- 用 `sigwaitinfo()`（和 `sigtimedwait()`）同步等待信号到达。
- 用 `signalfd()` 从一个文件描述符中接收信号。
- 较老的 BSD 版信号 API 和 System V 版信号 API。

22.1 核心转储文件

特定信号会引发进程创建一个核心转储文件并终止运行（参考表 20-1）。所谓核心转储是内含进程终止时内存映像的一个文件。（术语 `core` 源于一种老迈的内存技术。）将该内存映像加载到调试器中，即可查明信号到达时程序代码和数据的状态。

引发程序生成核心转储文件的方式之一是键入退出字符（通常为 `Control-\`），从而生成 `SIGQUIT` 信号。

```
$ ulimit -c unlimited           Explained in main text
$ sleep 30
Type Control-\
Quit (core dumped)
$ ls -l core                     Shows core dump file for sleep(1)
-rw----- 1 mtk users 57344 Nov 30 13:39 core
```

本例中，当检测出子进程（运行 `sleep` 命令的进程）为 `SIGQUIT` 信号所杀，并生成核心

转储文件时，shell 会显示 “Quit (core dump)” 消息。

核心转储文件创建于进程的工作目录中，名为 `core`。这是核心转储文件的默认位置和名称。稍后，将解释如何改变这些默认值。

借助于许多实现所提供的工具（例如 FreeBSD 和 Solaris 中的 `gcore`），可获取某一正在运行进程的核心转储文件。Linux 系统也有类似功能，使用 `gdb` 去连接（`attach`）一个正在运行的进程，然后运行 `gcore` 命令。

不产生核心转储文件的情况

以下情况不会产生核心转储文件。

- 进程对于核心转储文件没有写权限。造成这种情况的原因有进程对将要创建核心转储文件的所在目录可能没有写权限，或者是因为存在同名（且不可写，亦或非正规类型，例如，目录或符号链接）的文件。
- 存在一个同名、可写的普通文件，但指向该文件的（硬）链接数超过一个。
- 将要创建核心转储文件的所在目录并不存在。
- 把进程“核心转储文件大小”这一资源限制置为 0。36.3 节将就这一限制（`RLIMIT_CORE`）进行详细讨论。上例就使用了 `ulimit` 命令（C shell 中为 `limit` 命令）来取消对核心转储文件大小的任何限制。
- 将进程“可创建文件的大小”这一资源限制设置为 0。36.3 节将描述这一限制（`RLIMIT_FSIZE`）。
- 对进程正在执行的二进制可执行文件没有读权限。这样¹就防止了用户借助于核心转储文件来获取本无法读取的程序代码。
- 以只读方式挂载当前工作目录所在的文件系统，或者文件系统空间已满，又或者 `i-node` 资源耗尽。还有一种情况，即用户已经达到其在该文件系统上的配额限制。
- `Set-user-ID`（`set-group-ID`）程序在由非文件属主（或属组）执行时，不会产生核心转储文件。这可以防止恶意用户将一个安全程序的内存转储出来，再针对诸如密码之类的敏感信息进行刺探。

借助于 Linux 专有系统调用 `prctl()` 的 `PR_SET_DUMPABLE` 操作，可以为进程设置 `dumpable` 标志。当非文件属主（或属组）运行 `set-user-ID`（`set-group-ID`）程序时，如设置该标志即可生成核心转储文件。`PR_SET_DUMPABLE` 操作始见于 Linux 2.4，更多详细信息参见 `prctl(2)` 手册页。另外，始于内核版本 2.6.13，针对 `set-user-ID` 和 `set-group-ID` 进程是否产生核心转储文件，`/proc/sys/fs/suid_dumpable` 文件开始提供系统级控制。详情参见 `proc(5)` 手册页。

始于内核版本 2.6.23，利用 Linux 特有的 `/proc/PID/coredump_filter`，可以对写入核心转储文件的内存映射类型（第 49 章将解释内存映射）施以进程级控制。该文件中的值是一个 4 位掩码，分别对应于 4 种类型的内存映射：私有匿名映射、私有文件映射、共享匿名映射以及共享文件映射。文件默认值提供了传统的 Linux 行为：仅对私有匿名映射和共享匿名映射进行转储。详情参见 `core(5)` 手册页。

¹ 译者注：指不生成核心转储文件。

为核心转储文件命名：/proc/sys/kernel/core_pattern

从 Linux 版本 2.6 开始，可以根据 Linux 特有的 /proc/sys/kernel/core_pattern 文件所包含的格式字符串来控制对系统上生成的所有核心转储文件的命名。默认情况下，该文件所含字符串为 `core`。特权级用户可以将该文件内容定义为包含表 22-1 所列的任一格式说明符，待实际命名时再以表中右列所示相应值加以替换。此外，允许字符串中包含斜线 (/)。换言之，处在控制范围之内的，不仅包括核心文件的名称，还包括核心文件的所在（绝对或相对）目录。替换所有格式说明符后，由此生成的路径名字符串长度至多可达 128 个字符（Linux 2.6.19 之前为 64 个字符），超出部分将予以截断。

Linux 从内核版本 2.6.19 开始支持 `core_pattern` 文件的另一种语法。如果该文件包含一个以管道符 (|) 为首的字符串，那么会将该文件的剩余字符串视为一个程序，其可选参数可包含表 22-1 所示的 % 说明符——当进程转储核心文件时，将执行该程序。并且会将核心转储至该程序的标准输入，而非一个文件。详情请参考 `core(5)` 手册页。

其他一些 UNIX 实现也提供了类似于 `core_pattern` 的机制。例如，在 BSD 一派中，会将程序名追加到文件名尾部，形如 `core.progname`。Solaris 提供了一个工具 (`coreadm`)，允许由用户来选择核心转储文件的名称和存放目录。

表 22-1: 服务于 /proc/sys/kernel/core_pattern 的文件说明符

说 明 符	替 代 为
%c	对核心文件大小的资源软限制（字节数；始于 Linux 2.6.24）
%e	可执行文件名（不含路径前缀）
%g	遭转储进程的实际组 ID
%h	主机系统的名称
%p	遭转储进程的进程 ID
%s	导致进程终止的信号编号
%t	转储时间，始于 Epoch，以秒为单位
%u	遭转储进程的实际用户 ID
%%	单个 % 字符

22.2 传递、处置及处理的特殊情况

本节讨论了针对特定信号，适用于其传递、处置以及处理方面的特殊规则。

SIGKILL 和 SIGSTOP

SIGKILL 信号的默认行为是终止一个进程，SIGSTOP 信号的默认行为是停止一个进程，二者的默认行为均无法改变。当试图用 `signal()` 和 `sigaction()` 来改变对这些信号的处置时，将总是返回错误。同样，也不能将这两个信号阻塞。这是一个深思熟虑的设计决定。不允许修改这些信号的默认行为，这也意味着总是可以利用这些信号来杀死或者停止一个

失控进程。

SIGCONT 和停止信号

如前所述，可使用 SIGCONT 信号来使某些（因接收 SIGSTOP、SIGTSTP、SIGTTIN 和 SIGTTOU 信号而）处于停止状态的进程得以继续运行。由于这些停止信号具有独特目的，所以在某些情况下内核对它们的处理方式将有别于其他信号。

如果一个进程处于停止状态，那么一个 SIGCONT 信号的到来总是会促使其恢复运行，即使该进程正在阻塞或者忽略 SIGCONT 信号。该特性之所以必要，是因为如果要恢复这些处于停止状态的进程，舍此之外别无他法。（如果处于停止状态的进程正在阻塞 SIGCONT 信号，并且已经为 SIGCONT 信号建立了处理器函数，那么在进程恢复运行后，只有当取消了对 SIGCONT 的阻塞时，进程才会去调用相应的处理器函数。）

如果有任一其他信号发送给了一个已经停止的进程，那么在进程收到 SIGCONT 信号而恢复运行之前，信号实际上并未传递。SIGKILL 信号则属于例外，因为该信号总是会杀死进程，即使进程目前处于停止状态。

每当进程收到 SIGCONT 信号时，会将处于等待状态的停止信号丢弃（即进程根本不知道这些信号）。相反，如果任何停止信号传递给了进程，那么进程将自动丢弃任何处于等待状态的 SIGCONT 信号。之所以采取这些步骤，意在防止之前发送的一个停止信号会在随后撤销 SIGCONT 信号的行为，反之亦然。

由终端产生的信号若已被忽略，则不应改变其信号处置

如果程序在执行时发现，已将对由终端产生信号的处处置置为了 SIG_IGN（忽略），那么程序通常不应试图去改变信号处置。这并非系统的硬性规定，而是编写应用程序时所应遵循的惯例，34.7.3 节将解释其理由。与之相关的信号有：SIGHUP、SIGINT、SIGQUIT、SIGTTIN、SIGTTOU 和 SIGTSTP。

22.3 可中断和不可中断的进程睡眠状态

前文曾指出，SIGKILL 和 SIGSTOP 信号对进程的作用是立竿见影的。对于这一论断，此处要加入一条限制。内核经常需要令进程进入休眠，而休眠状态又分为两种。

- **TASK_INTERRUPTIBLE**：进程正在等待某一事件。例如，正在等待终端输入，等待数据写入当前的空管道，或者等待 System V 信号量值的增加。进程在该状态下所耗费的时间可长可短。如果为这种状态下的进程产生一个信号，那么操作将中断，而传递来的信号将唤醒进程。ps(1)命令在显示处于 TASK_INTERRUPTIBLE 状态的进程时，会将其 STAT（进程状态）字段标记为字母 S。
- **TASK_UNINTERRUPTIBLE**：进程正在等待某些特定类型的事件，比如磁盘 I/O 的完成。如果为这种状态下的进程产生一个信号，那么在进程摆脱这种状态之前，系统不会把信号传递给进程。ps(1)命令在显示处于 TASK_UNINTERRUPTIBLE 状态的进程时，会将其 STAT 字段标记为字母 D。

因为进程处于 TASK_UNINTERRUPTIBLE 状态的时间通常转瞬即逝，所以系统在进程脱离该状态时传递信号的现象也不易于被发现。然而，在极少数情况下，进程可能会因硬件故

障、NFS 问题或者内核缺陷而在该状态下保持挂起。这时，SIGKILL 将不会终止挂起进程。如果问题诱因无法得到解决，那么就只能通过重启系统来消灭该进程。

大多数 UNIX 系统实现都支持 TASK_INTERRUPTIBLE 和 TASK_UNINTERRUPTIBLE 状态。从内核 2.6.25 开始，Linux 加入第三种状态来解决上述挂起进程的问题。

- **TASK_KILLABLE**: 该状态类似于 TASK_UNINTERRUPTIBLE，但是会在进程收到一个致命信号（即一个杀死进程的信号）时将其唤醒。在对内核代码的相关部分进行改造后，就可使用该状态来避免各种因进程挂起而重启系统的情况。这时，向进程发送一个致命信号就能杀死进程。为使用 TASK_KILLABLE 而进行代码改造的首个内核模块是 NFS。

22.4 硬件产生的信号

硬件异常可以产生 SIGBUS、SIGFPE、SIGILL，和 SIGSEGV 信号，调用 kill() 函数来发送此类信号是另一种途径，但较为少见。SUSv3 规定，在硬件异常的情况下，如果进程从此类信号的处理函数中返回，亦或进程忽略或阻塞了此类信号，那么进程的行为未定义。原因如下。

- **从信号处理器中返回**: 假设机器语言指令产生了上述信号之一，并因此而调用了信号处理器函数。当从处理器函数正常返回后，程序会尝试从其中断处恢复执行。可当初引发信号产生的恰恰正是这条指令，所以信号会再次“光临”。故事的结局通常是，程序进入无限循环，重复调用信号处理器函数。
- **忽略信号**: 忽略因硬件而产生的信号于情理不合，试想算术异常之后，程序应当如何继续执行呢？无法明确。当由于硬件异常而产生上述信号之一时，Linux 会强制传递信号，即使程序已经请求忽略此类信号。
- **阻塞信号**。与上一种情况一样，阻塞因硬件而产生的信号也不合情理：不清楚程序随后应当如何继续执行。在 2.4 以及更早的版本中，Linux 内核仅会将阻塞硬件产生信号的种种企图一一忽略，信号无论如何都会传递给进程，随后要么进程终止，要么信号处理器会捕获信号——在程序安装有信号处理器的情况下。始于 Linux 2.6，如果信号遭到阻塞，那么该信号总是会立刻杀死进程，即使进程已经为此信号安装了处理器函数。（对于因硬件而产生的信号，Linux 2.6 之所以会改变对其处于阻塞状态下的处理方式，是由于 Linux 2.4 的行为中隐藏有缺陷，并可能在多线程程序中引起死锁。）

随本书发布源码中的 signals/demo_SIGFPE.c 程序就展示了忽略或者阻塞 SIGFPE 信号的后果，或者可正常返回的处理器将其捕获的结果。

正确处理硬件产生信号的方法有二：要么接受信号的默认行为（进程终止）；要么为其编写不会正常返回的处理器函数。除了正常返回之外，终结处理器执行的手段还包括调用 _exit() 以终止进程，或者调用 siglongjmp()（21.2.1 节），确保将控制传递回程序中（产生信号的指令位置之外）的某一位置。

22.5 信号的同步生成和异步生成

前文已然论及，进程一般无法预测其接收信号的时间。要证实这一点，需要对信号的同

步生成和异步生成加以区分。

截止目前所探讨的均属于信号的异步生成，即引发信号产生（无论信号发送者是内核还是另一进程）的事件，其发生与进程的执行无关。（例如，用户输入中断字符，或者子进程终止。）对于异步产生的信号，本节起始处的论断并非虚言。

然而，有时候信号的产生是由进程本身的执行造成的，前面就曾提及两个这样的例子。

- 执行特定的机器语言指令，可导致硬件异常，并因此而产生 22.4 节所述的硬件产生信号（SIGBUS、SIGFPE、SIGILL、SIGSEGV 和 SIGEMT）。
- 进程可以使用 `raise()`、`kill()` 或者 `killpg()` 向自身发送信号。

在这些情况下，信号的产生就是同步的——会立即传递信号（除非该信号遭到阻塞，但还要参考 22.4 节就阻塞硬件产生信号而展开的讨论）。换言之，本节开始处的论断则并不成立。对于同步产生的信号而言，其传递不但可以预测，而且可以重现。

注意，同步是对信号产生方式的描述，并不针对信号本身。所有的信号既可以同步产生（例如，进程使用 `kill()` 向自身发送信号），也可以异步产生（例如，由另一进程使用 `kill()` 来发送信号）。

22.6 信号传递的时机与顺序

本节的主题有二。其一，具体于何时去传递一个处于等待状态的信号；其二，对于多个遭到阻塞，且处于等待状态的信号一旦同时解除阻塞，将会发生什么情况？

何时传递一个信号？

如 22.5 节所述，同步产生的信号会立即传递。例如，硬件异常会触发一个即时信号，而当进程使用 `raise()` 向自身发送信号时，信号会在 `raise()` 调用返回前就已经发出。

当异步产生一个信号时，即使并未将其阻塞，在信号产生和实际传递之间仍可能会存在一个（瞬时）延迟。在此期间，信号处于等待状态。这是因为内核将等待信号传递给进程的时机是，该进程正在执行，且发生由内核态到用户态的下一切换时。实际上，这意味着在以下时刻才会传递信号。

- 进程在前度超时后，再度获得调度时（即，在一个时间片的开始处）。
- 系统调用完成时（信号的传递可能引起正在阻塞的系统调用过早完成）。

解除对多个信号的阻塞时，信号的传递顺序

如果进程使用 `sigprocmask()` 解除了对多个等待信号的阻塞，那么所有这些信号会立即传递给该进程。

就目前的 Linux 实现而言，Linux 内核按照信号编号的升序来传递信号。例如，如果对处于等待状态的信号 `SIGINT`（信号编号为 2）和 `SIGQUIT`（信号编号为 3）同时解除阻塞，那么无论这两个信号的产生次序如何，`SIGINT` 都将先于 `SIGQUIT` 而传递。

然而，也不能对传递（标准）信号的特定顺序产生任何依赖，因为 SUSv3 规定，多个信号的传递顺序由系统实现决定。（该条款仅适用于标准信号。如 22.8 节所述，实时信号的相关标准规定，对于解除阻塞的实时信号而言，其传递顺序必须得到保障。）

当多个解除了阻塞的信号正在等待传递时，如果在信号处理器函数执行期间发生了内核态和用户态之间的切换，那么将中断此处理器函数的执行，转而去调用第二个信号处理器函数（如此递进），如图 22-1 所示。

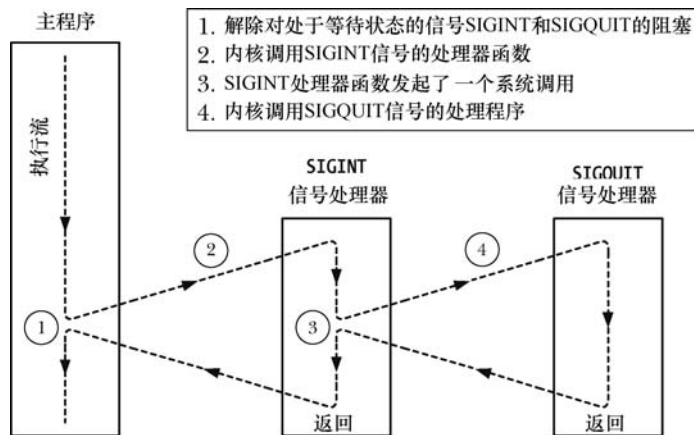


图 22-1: 对多个解除阻塞信号的传递

22.7 signal()的实现及可移植性

本节展示了如何使用 `sigaction()` 来实现 `signal()`。实现虽然简单明了，但还需要顾及这一事实，由于历史沿革和 UNIX 实现之间的差异，`signal()` 曾具有各种不同的语义。尤其是，信号的早期实现并不可靠，这意味着：

- 刚一进入信号处理器，会将信号处置重置为其默认行为。（这对应于 20.13 节描述的 `SA_RESETHAND` 标志。）要想在同一信号“再度光临”时再次调用该信号处理器函数，程序员必须在信号处理器内部调用 `signal()`，以显式重建处理器函数。这种情况存在一个问题：在进入信号处理器和重建处理器之间存在一个短暂的窗口期，而如果同一信号在此期间再度来袭，那么将只能按照其默认处置来进行处理。
- 在信号处理器执行期间，不会对新产生的信号进行阻塞。（这对应于 20.12 节描述的 `SA_NODEFER` 标志。）这意味着，如果在某一信号处理器函数执行期间，同类信号再度光顾，那么将对对该处理器函数进行递归调用。假定一串信号中彼此的时间间隔足够短，那么对处理器函数的递归调用将可能导致堆栈溢出。

除了不可靠之外，早期的 UNIX 实现并未提供系统调用的自动重启功能（即，21.5 节所述 `SA_RESTART` 标志的相关行为）。

4.2BSD 针对可靠信号的实现纠正了这些限制，其他一些 UNIX 实现也纷纷效仿。然而，时至今日，这些早期语义依然存在于 System V 的 `signal()` 实现之中。更有甚者，诸如 SUSv3 和 C99 之类的当代标准对 `signal()` 的这些方面也有意不予规范。

整合上述信息，对 `signal()` 的实现如程序清单 22-1 所示。该实现默认将提供信号的现代语义。如果编译时带有 `-DOLD_SIGNAL` 选项，那么将提供早期的不可靠信号语义，且不能启用系统调用的自动重启功能。

程序清单 22-1: `signal()` 的实现之一

```

-----signals/signal.c
#include <signal.h>

typedef void (*sighandler_t)(int);

```

```

sighandler_t
signal(int sig, sighandler_t handler)
{
    struct sigaction newDisp, prevDisp;

    newDisp.sa_handler = handler;
    sigemptyset(&newDisp.sa_mask);
#ifdef OLD_SIGNAL
    newDisp.sa_flags = SA_RESETHAND | SA_NODEFER;
#else
    newDisp.sa_flags = SA_RESTART;
#endif

    if (sigaction(sig, &newDisp, &prevDisp) == -1)
        return SIG_ERR;
    else
        return prevDisp.sa_handler;
}

```

—signals/signal.c

glibc 的一些细节

随着时间推移，glibc 对 signal() 库函数的实现也历经变化。较新版本（glibc 2 及更高版本）的函数库默认提供现代语义。而老版本则提供早期的不可靠（System V-兼容）语义。

Linux 内核将 signal() 实现为系统调用，并提供较老的、不可靠语义。然而，glibc 库则利用 sigaction() 实现了 signal() 库函数，从而将 signal() 系统调用旁路。

如果执意在现代 glibc 版本中使用不可靠信号语义，那么可以显式以（非标准的）sysv_signal() 函数来替代对 signal() 的调用。

```

#define _GNU_SOURCE
#include <signal.h>

void ( *sysv_signal(int sig, void (*handler)(int)) ) (int);

Returns previous signal disposition on success, or SIG_ERR on error

```

sysv_signal() 函数的参数与 signal() 函数相同。

若编译程序时并未定义 _BSD_SOURCE 特性测试宏，则 glibc 会隐式将所有 signal() 调用重新定义为 sysv_signal() 调用，亦即启用 signal() 的不可靠语义。默认情况下会定义 _BSD_SOURCE，但是（除非显式定义了 _BSD_SOURCE）如果编译程序时定义了诸如 _SVID_SOURCE 或 _XOPEN_SOURCE 之类的其他特性测试宏，那么对 _BSD_SOURCE 的默认定义将会失效。

sigaction() 是建立信号处理器的首选 API

鉴于上述 System V 与 BSD 之间（以及 glibc 新老版本之间）的可移植性问题，应当坚持使用 sigaction() 而非 signal() 来建立信号处理器，这不失为一种稳妥之举。本书剩下部分都将遵循这一做法。（另一种选择是，编写类似于程序清单 22-1 的 signal() 版本，精确设定所需要的标志，供应用程序内部使用。）不过，还应注意，使用 signal() 将信号处置设置为 SIG_IGN 或者 SIG_DFL 的手法具有良好的可移植性（程序也更为简短），所以也很常用。

22.8 实时信号

定义于 POSIX.1b 中的实时信号，意在弥补对标准信号的诸多限制。较之于标准信号，其优势如下所示。

- 实时信号的信号范围有所扩大，可应用于应用程序自定义的目的。而标准信号中可供应用随意使用的信号仅有两个：SIGUSR1 和 SIGUSR2。
- 对实时信号所采取的是队列化管理。如果将某一实时信号的多个实例发送给一进程，那么将会多次传递信号。相反，如果某一标准信号已经在等待某一进程，而此时即使再次向该进程发送此信号的实例，信号也只会传递一次。
- 当发送一个实时信号时，可为信号指定伴随数据（一整型数或者指针值），供接收进程的信号处理器获取。
- 不同实时信号的传递顺序得到保障。如果有多个不同的实时信号处于等待状态，那么将率先传递具有最小编号的信号。换言之，信号的编号越小，其优先级越高。如果是同一类型的多个信号在排队，那么信号（以及伴随数据）的传递顺序与信号发送来时的顺序保持一致。

SUSv3 要求，实现所提供的各种实时信号不得少于 `_POSIX_RTSIG_MAX`（定义为 8）个。Linux 内核则定义了 32 个不同的实时信号，编号范围为 32~63。<signal.h>头文件所定义的 `RTSIG_MAX` 常量则表征实时信号的可用数量，而此外所定义的常量 `SIGRTMIN` 和 `SIGRTMAX` 则分别表示可用实时信号编号的最小值和最大值。

采用 LinuxThreads 线程实现的系统将 `SIGRTMIN` 定义为 35（而非 32），这是因为 LinuxThreads 内部使用了前三个实时信号。而采用 NPTL 线程实现的系统则将 `SIGRTMIN` 定义为 34，因为 NPTL 内部使用了前两个实时信号。

对实时信号的区分方式有别于标准信号，不再依赖于所定义常量的不同。然而，程序员不应将实时信号编号的整型值在应用程序代码中写死，因为实时信号的范围因 UNIX 实现的不同而各异。与之相反，指代实时信号编号则可以采用 `SIGRTMIN+x` 的形式。例如，表达式 `(SIGRTMIN + 1)` 就表示第二个实时信号。

注意，SUSv3 并未要求 `SIGRTMAX` 和 `SIGRTMIN` 是简单的整数值。可以将其定义为函数（就像 Linux 中那样）。这也意味着，不能编写如下代码以供预处理器处理：

```
#if SIGRTMIN+100 > SIGRTMAX          /* WRONG! */
#error "Not enough realtime signals"
#endif
```

相反，必须在运行时执行等效检查。

对排队实时信号的数量限制

排队的实时信号（及其相关数据）需要内核维护相应的数据结构，用于罗列每个进程的排队信号。由于这些数据结构会消耗内核内存，故而内核对排队实时信号的数量设置了限制。

SUSv3 允许实现为每个进程中可排队的（各类）实时信号数量设置上限，并要求其不得少于 `_POSIX_SIGQUEUE_MAX`（定义为 32）。实现可借助于对 `SIGQUEUE_MAX` 常量的定义来表示其所允许的排队实时信号数量。发起如下调用也能获得这一信息：

```
lim = sysconf(_SC_SIGQUEUE_MAX);
```

若系统使用的 glibc 库版本在 2.4 之前，则该调用返回-1。从 glibc 2.4 开始，其返回值由内核版本决定。在 Linux 2.6.8 之前，调用将返回 Linux 专有文件/proc/sys/kernel/rtsig-max 中的值。该文件所定义为针对所有进程中可能排队的实时信号总数的系统级限制。默认值为 1024，不过特权级进程可以对其进行修改。至于当前的排队实时信号总数，可以从 Linux 专有的 /proc/sys/kernel/rtsig-nr 文件中读取。

从版本 2.6.8 开始，Linux 取消了这些/proc 文件。取而代之的是资源限制 RLIMIT_SIGPENDING（36.3 节）。针对某个特定实际用户 ID 下辖的所有进程，该限制限制了其可排队的信号总数。sysconf()调用从 glibc2.10 版本开始返回 RLIMIT_SIGPENDING 限制。（至于正在等待某一进程的实时信号数量，可以从 Linux 专有文件/proc/PID/status 中的 SigQ 字段读取。）

使用实时信号

为了能让一对进程收发实时信号，SUSv3 提出以下几点要求。

- 发送进程使用 sigqueue()系统调用来发送信号及其伴随数据。

使用 kill()、killpg()和 raise()调用也能发送实时信号。然而，至于系统是否会对利用此类接口所发送的信号进行排队处理，SUSv3 规定，由具体实现决定。这些接口在 Linux 中会对实时信号进行排队，但在其他许多 UNIX 实现中，情况则不然。

- 要为该信号建立了一个处理器函数，接收进程应以 SA_SIGINFO 标志发起对 sigaction()的调用。因此，调用信号处理器时就会附带额外参数，其中之一是实时信号的伴随数据。

在 Linux 中，即使接收进程在建立信号处理器时并未指定 SA_SIGINFO 标志，也能对实时信号进行队列化管理（但在这种情况下，将不可能获得信号的伴随数据）。然而，SUSv3 也不要实现确保这一行为，所以依赖这一点将有损于应用的可移植性。

22.8.1 发送实时信号

系统调用 sigqueue()将由 sig 指定的实时信号发送给由 pid 指定的进程。

```
#define _POSIX_C_SOURCE 199309
#include <signal.h>

int sigqueue(pid_t pid, int sig, const union signal value);

Returns 0 on success, or -1 on error
```

使用 sigqueue()发送信号所需要的权限与 kill()（参见 20.5 节）的要求一致。也可以发送空信号（即信号 0），其语义与 kill()中的含义相同。（不同于 kill()，sigqueue()不能通过将 pid 指定为负值而向整个进程组发送信号。）

程序清单 22-2：使用 sigqueue()发送实时信号

```
signals/t_sigqueue.c
#define _POSIX_C_SOURCE 199309
#include <signal.h>
#include "tspi_hdr.h"
```

```

int
main(int argc, char *argv[])
{
    int sig, numSigs, j, sigData;
    union sigval sv;

    if (argc < 4 || strcmp(argv[1], "--help") == 0)
        usageErr("%s pid sig-num data [num-sigs]\n", argv[0]);

    /* Display our PID and UID, so that they can be compared with the
       corresponding fields of the siginfo_t argument supplied to the
       handler in the receiving process */

    printf("%s: PID is %ld, UID is %ld\n", argv[0],
           (long) getpid(), (long) getuid());

    sig = getInt(argv[2], 0, "sig-num");
    sigData = getInt(argv[3], GN_ANY_BASE, "data");
    numSigs = (argc > 4) ? getInt(argv[4], GN_GT_0, "num-sigs") : 1;

    for (j = 0; j < numSigs; j++) {
        sv.sival_int = sigData + j;
        if (sigqueue(getLong(argv[1], 0, "pid"), sig, sv) == -1)
            errExit("sigqueue %d", j);
    }

    exit(EXIT_SUCCESS);
}

```

signals/t_sigqueue.c

参数 `value` 指定了信号的伴随数据，具有以下形式：

```

union sigval {
    int    sival_int;        /* Integer value for accompanying data */
    void *sival_ptr;        /* Pointer value for accompanying data */
};

```

对该参数的解释则取决于应用程序，由其选择对联合体（union）中的 `sival_int` 属性还是 `sival_ptr` 属性进行设置。`sigqueue()` 中很少使用 `sival_ptr`，因为指针的作用范围在进程内部，对于另一进程几乎没有意义。该字段得以一展身手之处，应该是在使用 `sigval` 联合体的其他函数中，诸如 23.6 节的 POSIX 计时器和 52.6 节的 POSIX 消息队列通知。

包括 Linux 在内的几个 UNIX 实现定义了与 `union sigval` 同义的数据类型 `sigval_t`。然而，该类型既未获得 SUSv3 接纳，也没有得到其他实现的支持。对可移植性有所要求的应用程序应当避免使用。

一旦触及对排队信号的数量限制，`sigqueue()` 调用将会失败。同时将 `errno` 置为 `EAGAIN`，以示需要再次发送该信号（在当前队列中某些信号传递之后的某一时间点）。

程序清单 22-2 提供了 `sigqueue()` 的应用示例。该程序最多接受 4 个参数，其中前 3 项为必填项：目标进程 ID、信号编号以及伴随实时信号的整型值。如果需要为指定信号发送多个实例，那么可以用可选的第 4 个参数来指定实例数量。在这种情况下，会为每个信号的伴随整型值依次加 1。22.8.2 节将展示该程序的用法。

22.8.2 处理实时信号

可以像标准信号一样，使用常规（单参数）信号处理器来处理实时信号。此外，也可

以用带有 3 个参数的信号处理器函数来处理实时信号，其建立则会用到 SA_SIGINFO 标志（参见 21.4 节）。以下为使用 SA_SIGINFO 标志为第六个实时信号建立处理器函数的代码示例：

```
struct sigaction act;

sigemptyset(&act.sa_mask);
act.sa_sigaction = handler;
act.sa_flags = SA_RESTART | SA_SIGINFO;

if (sigaction(SIGRTMIN + 5, &act, NULL) == -1)
    errExit("sigaction");
```

一旦采用了 SA_SIGINFO 标志，传递给信号处理器函数的第二个参数将是一个 siginfo_t 结构，内含实时信号的附加信息。21.4 节详细描述了这一数据结构。对于一个实时信号而言，会在 siginfo_t 结构中设置如下字段。

- si_signo 字段，其值与传递给信号处理器函数的第一个参数相同。
- si_code 字段表示信号来源，内容为表 21-2 中所示各值之一。对于通过 sigqueue() 发送的实时信号来说，该字段值总是为 SI_QUEUE。
- si_value 字段所含数据，由进程于使用 sigqueue() 发送信号时在 value 参数 (sigval union) 中指定。正如前文指出，对该数据的解释由应用程序决定。（若信号由 kill() 发送，则 si_value 字段所含信息无效。）
- si_pid 和 si_uid 字段分别包含信号发送进程的进程 ID 和实际用户 ID。

程序清单 22-3 提供了处理实时信号的一个例子。该程序捕获信号，并针对传递给信号处理器函数的 siginfo_t 结构，一一显示其中的各个字段值。该程序可接收两个整型命令行参数，均为可选项。如果提供了第一个参数，那么主程序将阻塞所有信号并进入休眠，休眠秒数由该参数指定。在此期间，将对进程的实时信号进行排队处理，并可观察解除对信号阻塞时所发生的情况。第二个参数指定了信号处理器函数在返回前所应休眠的秒数。指定一个非 0 值（默认为 1 秒）将有助于放缓程序的执行，便于看清处理多个信号时所发生的情况。

可以将程序清单 22-3 中程序与程序清单 22-2 中程序 (t_sigqueue.c) 结合起来探索实时信号的行为，正如以下 shell 会话日志所示：

```
$ ./catch_rtsigs 60 &
[1] 12842
$ ./catch_rtsigs: PID is 12842          Shell prompt mixed with program output
./catch_rtsigs: signals blocked - sleeping 60 seconds
Press Enter to see next shell prompt
$ ./t_sigqueue 12842 54 100 3          Send signal three times
./t_sigqueue: PID is 12843, UID is 1000
$ ./t_sigqueue 12842 43 200
./t_sigqueue: PID is 12844, UID is 1000
$ ./t_sigqueue 12842 40 300
./t_sigqueue: PID is 12845, UID is 1000
```

最终，catch_rtsigs 程序结束休眠，随着信号处理器捕获到各种信号而一一显示消息。（之所以看到 shell 提示符和程序的下一行输出混杂在一起，是因为 catch_rtsigs 程序正在后台输出信息。）可以看出，实时信号在传递时遵循低编号优先的原则，并且在传递给处理器函数的 siginfo_t 结构中包含了发送进程的进程 ID 和用户 ID。

```

$ ./catch_rtsigs: sleep complete
caught signal 40
  si_signo=40, si_code=-1 (SI_QUEUE), si_value=300
  si_pid=12845, si_uid=1000
caught signal 43
  si_signo=43, si_code=-1 (SI_QUEUE), si_value=200
  si_pid=12844, si_uid=1000

```

接下来的输出由同一实时信号的 3 个实例产生。由 `si_value` 值可知，这些信号的传递顺序与发送顺序相同。

```

caught signal 54
  si_signo=54, si_code=-1 (SI_QUEUE), si_value=100
  si_pid=12843, si_uid=1000
caught signal 54
  si_signo=54, si_code=-1 (SI_QUEUE), si_value=101
  si_pid=12843, si_uid=1000
caught signal 54
  si_signo=54, si_code=-1 (SI_QUEUE), si_value=102
  si_pid=12843, si_uid=1000

```

继续使用 `shell` 的 `kill` 命令向程序 `catch_rtsigs` 发送信号。一如既往，处理器函数接收到的 `siginfo_t` 结构中包含了发送进程的进程 ID 和用户 ID，但此时的 `si_code` 值为 `SI_USER`。

```

Press Enter to see next shell prompt
$ echo $$                                Display PID of shell
12780
$ kill -40 12842                          Uses kill(2) to send a signal
$ caught signal 40
  si_signo=40, si_code=0 (SI_USER), si_value=0
  si_pid=12780, si_uid=1000              PID is that of the shell
Press Enter to see next shell prompt
$ kill 12842                              Kill catch_rtsigs by sending SIGTERM
Caught 6 signals
Press Enter to see notification from shell about terminated background job
[1]+ Done                                ./catch_rtsigs 60

```

程序清单 22-3: 处理实时信号

```

signals/catch_rtsigs.c
#define _GNU_SOURCE
#include <string.h>
#include <signal.h>
#include "tspi_hdr.h"

static volatile int handlerSleepTime;
static volatile int sigCnt = 0;          /* Number of signals received */
static volatile int allDone = 0;

static void          /* Handler for signals established using SA_SIGINFO */
siginfoHandler(int sig, siginfo_t *si, void *ucontext)
{
    /* UNSAFE: This handler uses non-async-signal-safe functions
       (printf()); see Section 21.1.2) */

    /* SIGINT or SIGTERM can be used to terminate program */

    if (sig == SIGINT || sig == SIGTERM) {
        allDone = 1;
    }
}

```

```

        return;
    }

    sigCnt++;
    printf("caught signal %d\n", sig);

    printf("    si_signo=%d, si_code=%d (%s), ", si->si_signo, si->si_code,
           (si->si_code == SI_USER) ? "SI_USER" :
           (si->si_code == SI_QUEUE) ? "SI_QUEUE" : "other");
    printf("si_value=%d\n", si->si_value.sival_int);
    printf("    si_pid=%ld, si_uid=%ld\n", (long) si->si_pid, (long) si->si_uid);

    sleep(handlerSleepTime);
}
int
main(int argc, char *argv[])
{
    struct sigaction sa;
    int sig;
    sigset_t prevMask, blockMask;

    if (argc > 1 && strcmp(argv[1], "--help") == 0)
        usageErr("%s [block-time [handler-sleep-time]]\n", argv[0]);

    printf("%s: PID is %ld\n", argv[0], (long) getpid());

    handlerSleepTime = (argc > 2) ?
        getInt(argv[2], GN_NONNEG, "handler-sleep-time") : 1;

    /* Establish handler for most signals. During execution of the handler,
       mask all other signals to prevent handlers recursively interrupting
       each other (which would make the output hard to read). */

    sa.sa_sigaction = siginfoHandler;
    sa.sa_flags = SA_SIGINFO;
    sigfillset(&sa.sa_mask);

    for (sig = 1; sig < NSIG; sig++)
        if (sig != SIGSTP && sig != SIGQUIT)
            sigaction(sig, &sa, NULL);

    /* Optionally block signals and sleep, allowing signals to be
       sent to us before they are unblocked and handled */

    if (argc > 1) {
        sigfillset(&blockMask);
        sigdelset(&blockMask, SIGINT);
        sigdelset(&blockMask, SIGTERM);

        if (sigprocmask(SIG_SETMASK, &blockMask, &prevMask) == -1)
            errExit("sigprocmask");

        printf("%s: signals blocked - sleeping %s seconds\n", argv[0], argv[1]);
        sleep(getInt(argv[1], GN_GT_0, "block-time"));
        printf("%s: sleep complete\n", argv[0]);

        if (sigprocmask(SIG_SETMASK, &prevMask, NULL) == -1)
            errExit("sigprocmask");
    }
}

```



```
while (!allDone)                                /* Wait for incoming signals */
    pause();
}
```

signals/catch_rtsigs.c

22.9 使用掩码来等待信号：sigsuspend()

在解释 sigsuspend()的功用之前，先介绍一下它的一种使用场景。在对信号编程时偶尔会遇到如下情况。

1. 临时阻塞一个信号，以防止其信号处理器不会将某些关键代码片段的执行中断。
2. 解除对信号的阻塞，然后暂停执行，直至有信号到达。

为达到这一目的，可能会尝试使用程序清单 22-4 中代码所示方法。

程序清单 22-4：解除阻塞并等待信号的错误做法

```
sigset_t prevMask, intMask;
struct sigaction sa;

sigemptyset(&intMask);
sigaddset(&intMask, SIGINT);

sigemptyset(&sa.sa_mask);
sa.sa_flags = 0;
sa.sa_handler = handler;

if (sigaction(SIGINT, &sa, NULL) == -1)
    errExit("sigaction");

/* Block SIGINT prior to executing critical section. (At this
   point we assume that SIGINT is not already blocked.) */

if (sigprocmask(SIG_BLOCK, &intMask, &prevMask) == -1)
    errExit("sigprocmask - SIG_BLOCK");

/* Critical section: do some work here that must not be
   interrupted by the SIGINT handler */

/* End of critical section - restore old mask to unblock SIGINT */

if (sigprocmask(SIG_SETMASK, &prevMask, NULL) == -1)
    errExit("sigprocmask - SIG_SETMASK");

/* BUG: what if SIGINT arrives now... */

pause();                                /* Wait for SIGINT */
```

程序清单 22-4 中代码存在一个问题。假设 SIGINT 信号的传递发生在第二次调用 sigprocmask() 之后，调用 pause() 之前。（实际上，该信号可能产生于执行关键片段期间的任一时刻，仅当解除对信号的阻塞后才会随之而传递。）SIGINT 信号的传递将导致对处理器函数的调用，而当处理器返回后，主程序恢复执行，pause() 调用将陷入阻塞，直到 SIGINT 信号的第二个实例到达为止。这有违代码的本意：解除对 SIGINT 阻塞并等待其第一次出现。

即使在关键片段的起始点（即首次调用 `sigprocmask()`）和 `pause()`调用之间产生 `SIGINT` 信号的可能性不大，但这确实是上述代码的一处缺陷。这种取决于时间的缺陷是竞态条件（5.1 节）的例子之一。通常，竞态条件发生于两个进程或线程共享资源时。然而，此处的竞态条件却发生在主程序和其自身的信号处理器之间。

要避免这一问题，需要将解除信号阻塞和挂起进程这两个动作封装成一个原子操作。这正是 `sigsuspend()`系统调用的目的所在。

```
#include <signal.h>

int sigsuspend(const sigset_t *mask);

(Normally) returns -1 with errno set to EINTR
```

`sigsuspend()`系统调用将以 `mask` 所指向的信号集来替换进程的信号掩码，然后挂起进程的执行，直到其捕获到信号，并从信号处理器中返回。一旦处理器返回，`sigsuspend()`会将进程信号掩码恢复为调用前的值。

调用 `sigsuspend()`，相当于以不可中断方式执行如下操作：

```
sigprocmask(SIG_SETMASK, &mask, &prevMask);    /* Assign new mask */
pause();
sigprocmask(SIG_SETMASK, &prevMask, NULL);      /* Restore old mask */
```

虽然恢复老的信号掩码乍看起来似乎麻烦，但为了在需要反复等待信号的情况下避免竞态条件，这一做法就至关重要。在这种情况下，除非是在 `sigsuspend()`调用期间，否则信号必须保持阻塞状态。如果稍后需要对在调用 `sigsuspend()`之前遭到阻塞的信号解除阻塞，可以进一步调用 `sigprocmask()`。

若 `sigsuspend()`因信号的传递而中断，则将返回-1，并将 `errno` 置为 `EINTR`。如果 `mask` 指向的地址无效，则 `sigsuspend()`调用失败，并将 `errno` 置为 `EFAULT`。

示例程序

程序清单 22-5 展示了对 `sigsuspend()`的使用。该程序执行如下步骤。

- 调用 `printSigMask()`函数（程序清单 20-4）来显示进程信号掩码的初始值。
- 阻塞 `SIGINT` 和 `SIGQUIT` 信号，并保存原始的进程信号掩码。
- 为 `SIGINT` 和 `SIGQUIT` 信号建立相同的处理器函数。该处理器显示一条消息，且若对其调用因 `SIGQUIT` 信号的传递而引起，则设置全局变量 `gotSigquit`。
- 循环执行，直至对 `gotSigquit` 进行了设置。每次循环都执行如下步骤。
 - 使用 `printSigMask()`函数显示信号掩码的当前值。
 - 令 CPU 忙于循环并持续数秒钟，以此来模拟对一个关键片段的执行。
 - 使用 `printPendingSigs()`函数来显示等待信号的掩码（程序清单 20-4）。
 - 使用 `sigsuspend()`来解除对 `SIGINT` 和 `SIGQUIT` 信号的阻塞，并等待信号（如果尚未有信号处于等待状态）。
- 使用 `sigprocmask()`将进程信号掩码恢复为原始状态，然后再使用 `printSigMask()`来显示信号掩码。

程序清单 22-5：使用 `sigsuspend()`

```
signals/t_sigsuspend.c
#define _GNU_SOURCE    /* Get strsignal() declaration from <string.h> */
```

```

#include <string.h>
#include <signal.h>
#include <time.h>
#include "signal_functions.h"          /* Declarations of printSigMask()
                                       and printPendingSigs() */

#include "tspi_hdr.h"

static volatile sig_atomic_t gotSigquit = 0;

static void
handler(int sig)
{
    printf("Caught signal %d (%s)\n", sig, strsignal(sig));
    /* UNSAFE (see Section 21.1.2) */
    if (sig == SIGQUIT)
        gotSigquit = 1;
}

int
main(int argc, char *argv[])
{
    int loopNum;
    time_t startTime;
    sigset_t origMask, blockMask;
    struct sigaction sa;

①    printSigMask(stdout, "Initial signal mask is:\n");

    sigemptyset(&blockMask);
    sigaddset(&blockMask, SIGINT);
    sigaddset(&blockMask, SIGQUIT);
②    if (sigprocmask(SIG_BLOCK, &blockMask, &origMask) == -1)
        errExit("sigprocmask - SIG_BLOCK");

    sigemptyset(&sa.sa_mask);
    sa.sa_flags = 0;
    sa.sa_handler = handler;
③    if (sigaction(SIGINT, &sa, NULL) == -1)
        errExit("sigaction");
    if (sigaction(SIGQUIT, &sa, NULL) == -1)
        errExit("sigaction");

④    for (loopNum = 1; !gotSigquit; loopNum++) {
        printf("=== LOOP %d\n", loopNum);

        /* Simulate a critical section by delaying a few seconds */

        printSigMask(stdout, "Starting critical section, signal mask is:\n");
        for (startTime = time(NULL); time(NULL) < startTime + 4; )
            continue;          /* Run for a few seconds elapsed time */

        printPendingSigs(stdout,
            "Before sigsuspend() - pending signals:\n");
        if (sigsuspend(&origMask) == -1 && errno != EINTR)
            errExit("sigsuspend");
    }

⑤    if (sigprocmask(SIG_SETMASK, &origMask, NULL) == -1)
        errExit("sigprocmask - SIG_SETMASK");
}

```

```

⑥ printSigMask(stdout, "=== Exited loop\nRestored signal mask to:\n");

/* Do other processing... */

exit(EXIT_SUCCESS);
}

```

signals/t_sigsuspend.c

以下 shell 会话日志所示为程序清单 22-5 中程序的运行结果示例:

```

$ ./t_sigsuspend
Initial signal mask is:
    <empty signal set>
=== LOOP 1
Starting critical section, signal mask is:
    2 (Interrupt)
    3 (Quit)
Type Control-C; SIGINT is generated, but remains pending because it is blocked
Before sigsuspend() - pending signals:
    2 (Interrupt)
Caught signal 2 (Interrupt)      sigsuspend() is called, signals are unblocked

```

程序调用 sigsuspend()解除了对 SIGINT 信号的阻塞，还显示了最后一行输出。正是在那一点，调用了信号处理器，并显示了那一行输出。

主程序会继续循环。

```

=== LOOP 2
Starting critical section, signal mask is:
    2 (Interrupt)
    3 (Quit)
Type Control-\ to generate SIGQUIT
Before sigsuspend() - pending signals:
    3 (Quit)
Caught signal 3 (Quit)          sigsuspend() is called, signals are unblocked
=== Exited loop                Signal handler set gotSigquit
Restored signal mask to:
    <empty signal set>

```

此时按下 Control-\，将导致信号处理器去设置 gotSigquit 标志，并转而引发主程序终止循环。

22.10 以同步方式等待信号

22.9 节描述了如何结合信号处理器和 sigsuspend()来挂起一个进程的执行，直至传来一个信号。然而，这需要编写信号处理器函数，还需要应对信号异步传递所带来的复杂性。对于某些应用而言，这种方法过于繁复。作为替代方案，**可以利用 sigwaitinfo()系统调用来同步接收信号。**

```

#define _POSIX_C_SOURCE 199309
#include <signal.h>

int sigwaitinfo(const sigset_t *set, siginfo_t *info);

Returns number of delivered signal on success, or -1 on error

```

sigwaitinfo()系统调用挂起进程的执行，直至 set 指向信号集中的某一信号抵达。如果调用 sigwaitinfo()时，set 中的某一信号已经处于等待状态，那么 sigwaitinfo()将立即返回。传递来

的信号就此从进程的等待信号队列中移除，并且将返回信号编号作为函数结果。`info` 参数如果不为空，则会指向经过初始化处理的 `siginfo_t` 结构，其中所含信息与提供给信号处理器函数的 `siginfo_t` 参数（21.4 节）相同。

`sigwaitinfo()`所接受信号的传递顺序和排队特性与信号处理器所捕获的信号相同，就是说，不对标准信号进行排队处理，对实时信号进行排队处理，并且对实时信号的传递遵循低编号优先的原则。

除了卸去编写信号处理器的负担之外，使用 `sigwaitinfo()`来等待信号也要比信号处理器外加 `sigsuspend()`的组合要稍快一些（见练习 22-3）。

将对 `set` 中信号集的阻塞与调用 `sigwaitinfo()`结合起来，这当属明智之举。（即便某一信号遭到阻塞，仍然可以使用 `sigwaitinfo()`来获取等待信号。）如果没有这么做，而信号在首次调用 `sigwaitinfo()`之前，或者两次连续调用 `sigwaitinfo()`之间到达，那么对信号的处理将只能依照其当前处置。

SUSv3 规定，调用 `sigwaitinfo()`而不阻塞 `set` 中的信号将导致不可预知的行为（其行为未定义）。

程序清单 22-6 所示为使用 `sigwaitinfo()`的例子之一。程序首先阻塞所有信号，然后延迟数秒时间，具体秒数由可选命令行参数来指定，从而允许在调用 `sigwaitinfo()`之前向程序发送信号。程序随即持续循环调用 `sigwaitinfo()`来接收输入信号，直至收到 `SIGINT` 或 `SIGTERM` 信号。

如下 shell 会话日志展示了程序清单 22-6 中程序的运行情况。程序在后台运行，并指定在执行 `sigwaitinfo()`前需延迟 60 秒，随后再向进程发送两个信号：

```
$ ./t_sigwaitinfo 60 &
./t_sigwaitinfo: PID is 3837
./t_sigwaitinfo: signals blocked
./t_sigwaitinfo: about to delay 60 seconds
[1] 3837
$ ./t_sigqueue 3837 43 100          Send signal 43
./t_sigqueue: PID is 3839, UID is 1000
$ ./t_sigqueue 3837 42 200        Send signal 42
./t_sigqueue: PID is 3840, UID is 1000
```

最终，程序完成睡眠，`sigwaitinfo()`调用循环接收排队信号。（由于 `t_sigwaitinfo` 程序正在后台输出信息，故而可以观察到 shell 提示符和程序的下一行输出混杂在一起。）至于处理器所捕获到的实时信号，可以看出，编号低的信号率先传递，而且，借助于传递给信号处理器函数的 `siginfo_t` 结构，还可以获得发送进程的进程 ID 和用户 ID。

```
$ ./t_sigwaitinfo: finished delay
got signal: 42
  si_signo=42, si_code=-1 (SI_QUEUE), si_value=200
  si_pid=3840, si_uid=1000
got signal: 43
  si_signo=43, si_code=-1 (SI_QUEUE), si_value=100
  si_pid=3839, si_uid=1000
```

继续使用 shell 的 `kill` 命令向进程发送信号。可以观察到，这次将 `si_code` 字段置为 `SI_USER`（而非 `SI_QUEUE`）。

```
Press Enter to see next shell prompt
$ echo $$          Display PID of shell
3744
$ kill -USR1 3837  Shell sends SIGUSR1 using kill()
$ got signal: 10   Delivery of SIGUSR1
  si_signo=10, si_code=0 (SI_USER), si_value=100
```

```

    si_pid=3744, si_uid=1000                3744 is PID of shell
Press Enter to see next shell prompt
$ kill %1                                Terminate program with SIGTERM
$
Press Enter to see notification of background job termination
[1]+ Done                               ./t_sigwaitinfo 60

```

收到 SIGUSR1 信号，由其输出可知，si_value 字段值为 100。该值是由 sigqueue() 发送的前一信号初始化而成。前文曾指出，仅对由 sigqueue() 所发送的信号，si_value 字段所包含的信息才是可靠的。

程序清单 22-6：使用 sigwaitinfo() 来同步等待信号

```

                                                                    signals/t_sigwaitinfo.c
#define _GNU_SOURCE
#include <string.h>
#include <signal.h>
#include <time.h>
#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    int sig;
    siginfo_t si;
    sigset_t allSigs;

    if (argc > 1 && strcmp(argv[1], "--help") == 0)
        usageErr("%s [delay-secs]\n", argv[0]);

    printf("%s: PID is %ld\n", argv[0], (long) getpid());

    /* Block all signals (except SIGKILL and SIGSTOP) */

    sigfillset(&allSigs);
    if (sigprocmask(SIG_SETMASK, &allSigs, NULL) == -1)
        errExit("sigprocmask");
    printf("%s: signals blocked\n", argv[0]);

    if (argc > 1) {
        /* Delay so that signals can be sent to us */
        printf("%s: about to delay %s seconds\n", argv[0], argv[1]);
        sleep(getInt(argv[1], GN_GT_0, "delay-secs"));
        printf("%s: finished delay\n", argv[0]);
    }

    for (;;) {
        /* Fetch signals until SIGINT (^C) or SIGTERM */
        sig = sigwaitinfo(&allSigs, &si);
        if (sig == -1)
            errExit("sigwaitinfo");

        if (sig == SIGINT || sig == SIGTERM)
            exit(EXIT_SUCCESS);

        printf("got signal: %d (%s)\n", sig, strsignal(sig));
        printf("    si_signo=%d, si_code=%d (%s), si_value=%d\n",
            si.si_signo, si.si_code,
            (si.si_code == SI_USER) ? "SI_USER" :

```

```

        (si.si_code == SI_QUEUE) ? "SI_QUEUE" : "other",
        si.si_value.sival_int);
printf("    si_pid=%ld, si_uid=%ld\n",
       (long) si.si_pid, (long) si.si_uid);
    }
}

```

signals/t_sigwaitinfo.c

sigtimedwait()系统调用是 sigwaitinfo()调用的变体。唯一的区别是 sigtimedwait()允许指定等待时限。

```

#define _POSIX_C_SOURCE 199309
#include <signal.h>

int sigtimedwait(const sigset_t *set, siginfo_t *info,
                const struct timespec *timeout);

Returns number of delivered signal on success,
or -1 on error or timeout (EAGAIN)

```

timeout 参数指定了允许 sigtimedwait()等待一个信号的最大时长，是指向如下类型结构的一枚指针：

```

struct timespec {
    time_t tv_sec;        /* Seconds ('time_t' is an integer type) */
    long tv_nsec;        /* Nanoseconds */
};

```

填写 timespec 结构的所属字段，也就指定了允许 sigtimedwait()等待的最大秒数和纳秒数。如果将这两个字段均指定为 0，那么函数将立刻超时，就是说，会去轮询检查是否有指定信号集中的任一信号处于等待状态。¹如果调用超时而又没有收到信号，sigtimedwait()将调用失败，并将 errno 置为 EAGAIN。

如果将 timeout 参数指定为 NULL，那么 sigtimedwait()将完全等同于 sigwaitinfo()。SUSv3 对于 timeout 的 NULL 值含义也语焉不详，而某些 UNIX 实现则将该值视为轮询请求并立即将其返回。

22.11 通过文件描述符来获取信号

始于内核 2.6.22，Linux 提供了（非标准的）signalfd()系统调用；利用该调用可以创建一个特殊文件描述符，发往调用者的信号都可从该描述符中读取。signalfd 机制为同步接受信号提供了 sigwaitinfo()之外的另一种选择。

```

#include <sys/signalfd.h>

int signalfd(int fd, const sigset_t *mask, int flags);

Returns file descriptor on success, or -1 on error

```

mask 参数是一个信号集，指定了有意通过 signalfd 文件描述符来读取的信号。如同 sigwaitinfo()一样，通常也应该使用 sigprocmask()阻塞 mask 中的所有信号，以确保在有机会读

¹ 译者注：有，则将该信号的信息返回。作者的书似乎没有 man 手册页写得明了，请读者自行比较。

取这些信号之前，不会按照默认处置对它们进行处理。

如果指定 `fd` 为 `-1`，那么 `signalfd()` 会创建一个新的文件描述符，用于读取 `mask` 中的信号；否则，将修改与 `fd` 相关的 `mask` 值，且该 `fd` 一定是由之前对 `signalfd()` 的一次调用创建而成。

早期实现将 `flag` 参数保留下来供将来使用，且必须将其指定为 `0`。然而，Linux 从版本 2.6.27 开始支持下面两个标志。

SFD_CLOEXEC

为新的文件描述符设置 `close-on-exec` (`FD_CLOEXEC`) 标志。该标志之所以必要，与 4.3.1 节中描述的 `open(O_CLOEXEC)` 标志的设置理由相同。

SFD_NONBLOCK

为底层的打开文件描述设置 `O_NONBLOCK` 标志，以确保不会阻塞未来的读操作。既省去了一个额外的 `fcntl()` 调用，又获得了相同的结果。

创建文件描述符之后，可以使用 `read()` 调用从中读取信号。提供给 `read()` 的缓冲区必须足够大，至少应能够容纳一个 `signalfd_siginfo` 结构。<sys/signalfd.h> 文件定义了该结构，如下所示：

```
struct signalfd_siginfo {
    uint32_t  ssi_signo;      /* Signal number */
    int32_t   ssi_errno;     /* Error number (generally unused) */
    int32_t   ssi_code;      /* Signal code */
    uint32_t  ssi_pid;       /* Process ID of sending process */
    uint32_t  ssi_uid;       /* Real user ID of sender */
    int32_t   ssi_fd;        /* File descriptor (SIGPOLL/SIGIO) */
    uint32_t  ssi_tid;       /* Kernel timer ID (POSIX timers) */
    uint32_t  ssi_band;      /* Band event (SIGPOLL/SIGIO) */
    uint32_t  ssi_tid;       /* (Kernel-internal) timer ID (POSIX timers) */
    uint32_t  ssi_overrun;   /* Overrun count (POSIX timers) */
    uint32_t  ssi_trapno;    /* Trap number */
    int32_t   ssi_status;    /* Exit status or signal (SIGCHLD) */
    int32_t   ssi_int;       /* Integer sent by sigqueue() */
    uint64_t  ssi_ptr;       /* Pointer sent by sigqueue() */
    uint64_t  ssi_utime;     /* User CPU time (SIGCHLD) */
    uint64_t  ssi_stime;     /* System CPU time (SIGCHLD) */
    uint64_t  ssi_addr;      /* Address that generated signal
                               (hardware-generated signals only) */
};
```

该结构中字段所返回的信息与传统 `siginfo_t` 结构 (21.4 节) 中类似命名的字段信息相同。

`read()` 每次调用都将返回与等待信号数目相等的 `signalfd_siginfo` 结构，并填充到已提供的缓冲区中。如果调用时并无信号正在等待，那么 `read()` 将阻塞，直到有信号到达。也可以使用 `fcntl()` 的 `F_SETFL` 操作 (5.3 节) 来为文件描述符设置 `O_NONBLOCK` 标志，使得读操作不再阻塞，且若无信号等待，则调用失败，`errno` 为 `EAGAIN`。

当从 `signalfd` 文件描述符中读取到一信号时，该信号获得接纳，且不再为该进程而等待。

程序清单 22-7：使用 `signalfd()` 来读取信号

```
signals/signalfd_signal.c
#include <sys/signalfd.h>
#include <signal.h>
#include "tspi_hdr.h"
```



```

int
main(int argc, char *argv[])
{
    sigset_t mask;
    int sfd, j;
    struct signalfd_siginfo fdsi;
    ssize_t s;

    if (argc < 2 || strcmp(argv[1], "--help") == 0)
        usageErr("%s sig-num...\n", argv[0]);

    printf("%s: PID = %ld\n", argv[0], (long) getpid());

    sigemptyset(&mask);
    for (j = 1; j < argc; j++)
        sigaddset(&mask, atoi(argv[j]));

    if (sigprocmask(SIG_BLOCK, &mask, NULL) == -1)
        errExit("sigprocmask");

    sfd = signalfd(-1, &mask, 0);
    if (sfd == -1)
        errExit("signalfd");

    for (;;) {
        s = read(sfd, &fdsi, sizeof(struct signalfd_siginfo));
        if (s != sizeof(struct signalfd_siginfo))
            errExit("read");

        printf("%s: got signal %d", argv[0], fdsi.ssi_signal);
        if (fdsi.ssi_code == SI_QUEUE) {
            printf("; ssi_pid = %d; ", fdsi.ssi_pid);
            printf("ssi_int = %d", fdsi.ssi_int);
        }
        printf("\n");
    }
}

```

signals/signalfd_sigval.c

select()、poll()和 epoll（参见第 63 章）可以将 signalfd 描述符和其他描述符混合起来进行监控。撇开其他用途不提，该特性可成为 63.5.2 节所述 self-pipe 技巧之外的另一选择。如果有信号正在等待，那么这些技术将文件描述符指示为可读取。

当不再需要 signalfd 文件描述符时，应当关闭 signalfd 以释放相关内核资源。

程序清单 22-7 展示了 signalfd()的用法。程序为在命令行参数中指定的信号创建掩码，阻塞这些信号，然后创建用来读取这些信号的 signalfd 文件描述符，之后循环从该文件描述符中读取信号，并显示返回的 signalfd_siginfo 结构中的部分信息。如下 shell 会话在后台运行了程序清单 22-7 中程序，并使用程序清单 22-2 中程序 (t_sigqueue.c) 向该进程发送实时信号及伴随数据：

```

$ ./signalfd_sigval 44 &
./signalfd_sigval: PID = 6267
[1] 6267
$ ./t_sigqueue 6267 44 123          Send signal 44 with data 123 to PID 6267
./t_sigqueue: PID is 6269, UID is 1000
./signalfd_sigval: got signal 44; ssi_pid=6269; ssi_int=123
$ kill %1                          Kill program running in background

```

22.12 利用信号进行进程间通信

从某种角度，可将信号视为进程间通信（IPC）的方式之一。然而，信号作为一种 IPC 机制却也饱受限制。首先，与后续各章描述的其他 IPC 方法相比，对信号编程既繁且难，具体原因如下。

- 信号的异步本质就意味着需要面对各种问题，包括可重入性需求、竞态条件及在信号处理器中正确处理全局变量。（如果用 `sigwaitinfo()` 或者 `siglfd()` 来同步获取信号，这些问题中的大部分都不会遇到。）
- 没有对标准信号进行排队处理。即使是对于实时信号，也存在对信号排队数量的限制。这意味着，为了避免丢失信息，接收信号的进程必须想方设法通知发送者，自己为接受另一个信号做好了准备。要做到这一点，最显而易见的方法是由接收者向发送者发送信号。

还有一个更深层次的问题，信号所携带的信息量有限：信号编号以及实时信号情况下一字之长的附加数据（一个整数或者一枚指针值）。与诸如管道之类的其他 IPC 方法相比，过低的带宽使得信号传输极为缓慢。

由于上述种种限制，很少将信号用于 IPC。

22.13 早期的信号 API（System V 和 BSD）

之前对信号的讨论一直着眼于 POSIX 信号 API。本节将简要回顾一下 System V 和 BSD 提供的历史 API。虽然所有的新应用程序都应当使用 POSIX API，但是在从其他 UNIX 实现移植（通常较为老旧的）应用时，可能还是会碰到这些过时的 API。当移植这些使用老旧 API 的程序时，因为 Linux（像许多其他 UNIX 实现一样）提供了与 System V 和 BSD 兼容的 API，所以通常所要做的全部工作不过是在 Linux 平台上重新进行编译而已。

System V 信号 API

如前所述，System V 中的信号 API 存在一个重要差异：当使用 `signal()` 建立处理器函数时，得到的是老版、不可靠的信号语义。这意味着不会将信号添加到进程的信号掩码中，调用信号处理器时会将信号处置重置为默认行为，以及不会自动重启系统调用。

下面，简单介绍一些 System V 信号 API 中的函数。手册页提供有全部的细节。SUSv3 定义了所有这些函数，但指出应优先使用现代版的 POSIX 等价函数。SUSv4 将这些函数标记为已废止。

```
#define _XOPEN_SOURCE 500
#include <signal.h>
```

```
void (*sigset(int sig, void (*handler)(int)))(int);
```

On success: returns the previous disposition of *sig*, or `SIG_HOLD` if *sig* was previously blocked; on error `-1` is returned

为了建立一个具有可靠语义的信号处理器，System V 提供了 `sigset()` 调用（原型类似于

signal())。与 signal() 一样，可以将 sigset() 的 handler 参数指定为 SIG_IGN、SIG_DFL 或者信号处理器函数的地址。此外，还可以将其指定为 SIG_HOLD，在将信号添加到进程信号掩码的同时保持信号处置不变。

如果指定 handler 参数为 SIG_HOLD 之外的其他值，那么会将 sig 从进程信号掩码中移除（即，如果 sig 遭到阻塞，那么将解除对其阻塞）。

```
#define _XOPEN_SOURCE 500
#include <signal.h>

int sighold(int sig);
int sigrelse(int sig);
int sigignore(int sig);

int sigpause(int sig);
```

All return 0 on success, or -1 on error

Always returns -1 with *errno* set to EINTR

sighold() 函数将一个信号添加到进程信号掩码中。sigrelse() 函数则是从信号掩码中移除一个信号。sigignore() 函数设定对某信号的处置为“忽略 (ignore)”。sigpause() 函数类似于 sigsuspend() 函数，但仅从进程信号掩码中移除一个信号，随后将暂停进程，直到有信号到达。

BSD 信号 API

POSIX 信号 API 从 4.2BSD API 中汲取了很多灵感，所以 BSD 函数与 POSIX 函数大体相仿。

如同前文对 System V 信号 API 中函数的描述一样，首先给出 BSD 信号 API 中各函数的原型，随后简单解释一下每个函数的操作。再啰嗦一句，手册页提供有全部细节。

```
#define _BSD_SOURCE
#include <signal.h>

int sigvec(int sig, struct sigvec *vec, struct sigvec *ovec);
```

Returns 0 on success, or -1 on error

sigvec() 类似于 sigaction()。vec 和 ovec 参数是指向如下类型结构的指针：

```
struct sigvec {
    void (*sv_handler)();
    int sv_mask;
    int sv_flags;
};
```

sigvec 结构中的字段与 sigaction 结构中的那些字段紧密对应。第一个显著差异是 sv_mask（类似与 sa_mask）字段是一个整型，而非 sigset_t 类型。这意味着，在 32 位架构中，最多支持 31 个不同信号。另一个不同之处则在于在 sv_flags（类似与 a_flags）字段中使用了 SV_INTERRUPT 标志。因为重启系统调用是 4.2BSD 的默认行为，该标志是用来指定应使用信号处理器来中断慢速系统调用。（这与 POSIX API 截然相反，在使用 sigaction() 建立信号处理器时，如果希望启用系统调用重启功能，就必须显式指定 SA_RESTART 标志。）

```

#define _BSD_SOURCE
#include <signal.h>

int sigblock(int mask);
int sigsetmask(int mask);

int sigpause(int sigmask);

int sigmask(sig);

```

Both return previous signal mask

Always returns -1 with *errno* set to EINTR

Returns signal mask value with bit *sig* set

`sigblock()`函数向进程信号掩码中添加一组信号。这类似于 `sigprocmask()` 的 `SIG_BLOCK` 操作。`sigsetmask()`调用则为信号掩码指定了一个绝对值。这类似于 `sigprocmask()` 的 `SIG_SETMASK` 操作。

`sigpause()`类似于 `sigsuspend()`。注意，对该函数的定义在 System V 和 BSD API 中具有不同的调用签名。GNU C 函数库默认提供 System V 版本，除非在编译程序时指定了特性测试宏 `_BSD_SOURCE`。

`sigmask()`宏将信号编号转换成相应的 32 位掩码值。此类位掩码可以彼此相或，一起创建一组信号，如下所示：

```
sigblock(sigmask(SIGINT) | sigmask(SIGQUIT));
```

22.14 总结

某些信号会引发进程创建一个核心转储文件，并终止进程。核心转储所包含的信息可供调试器检查进程终止时的状态。默认情况下，对核心转储文件的命名为 `core`，但 Linux 提供了 `/proc/sys/kernel/core_pattern` 文件来控制对核心转储文件的命名。

信号的产生方式既可以是异步的，也可以是同步的。当由内核或者另一进程发送信号给进程时，信号可能是异步产生的。进程无法精确预测异步产生信号的传递时间。（文中曾指出，异步信号通常会在接收进程第二次从内核态切换到用户态时进行传递。）因进程自身执行代码而直接产生的信号则属于是同步产生的，例如，执行了一个引发硬件异常的指令，或者去调用 `raise()`。同步生成的信号，其传递可以精确预测（立即传递）。

实时信号是 POSIX 对原始信号模型的扩展，不同之处包括对实时信号进行队列化管理，具有特定的传递顺序，并且还可以伴随少量数据一同发送。设计实时信号，意在供应用程序自定义使用。实时信号的发送使用 `sigqueue()`系统调用，并且还向信号处理器函数提供了一个附加参数（`siginfo_t` 结构），以便其获得信号的伴随数据，以及发送进程的进程 ID 和实际用户 ID。

`sigsuspend()`系统调用在自动修改进程信号掩码的同时，还将挂起进程的执行直到信号到达，且二者属于同一原子操作。为了避免执行上述功能时出现竞态条件，确保 `sigsuspend()` 的原子性至关重要。

可以使用 `sigwaitinfo()`和 `sigtimedwait()`来同步等待一个信号。这省去了对信号处理器的设计和编码工作。对于以等待信号的传递为唯一目的的程序而言，使用信号处理器纯属多此一举。

像 `sigwaitinfo()` 和 `sigtimedwait()` 一样，可以使用 Linux 特有的 `signalfd()` 系统调用来同步等待一个信号。这一接口的独特之处在于可以通过文件描述符来读取信号。还可以使用 `select()`、`poll()` 和 `epoll` 来对其进行监控。

尽管可以将信号视为 IPC 的方式之一，但诸多制约因素令其常常无法胜任这一目的，其中包括信号的异步本质、不对信号进行排队处理的事实，以及较低的传递带宽。信号更为常见的应用场景是用于进程同步，或是各种其他目的（比如，事件通知、作业控制以及定时器到期）。

信号的基本概念虽然简单，但因为涉及的细节很多，所以对其讨论用去了 3 章的篇幅。信号在系统调用 API 的各部分中都扮演着重要角色，后面几章还将重温对信号的使用。此外，还有各种信号相关的函数是针对线程的（比如，`pthread_kill()` 和 `pthread_sigmask()`），将延后至 33.2 节进行讨论。

更多信息

参见 20.15 节所列的信息来源。

22.15 练习

- 22-1. 22.2 节曾指出，假设进程为 SIGCONT 信号建立了处理器函数并将其阻塞，如果该进程已停止 (stopped) 后因收到一个 SIGCONT 信号而恢复执行，那么仅当解除了对 SIGCONT 信号的阻塞时才会去调用信号处理器函数。编写一个程序来验证这一点。回忆一下，按下终端暂停字符（通常为 Control-Z）可以停止进程，使用 `kill-CONT` 命令（或者隐蔽一点，使用 shell 的 `fg` 命令）可以发送 SIGCONT 信号。
- 22-2. 如果实时信号和标准信号在同时等待一个进程，那么 SUSv3 对信号的传递顺序未予定义。编写一程序来展示 Linux 是如何处理这一情况的。（令程序为所有信号设置处理器函数，阻塞这些信号并持续一段时间，以便于向其发送各种信号，最后解除对所有信号的阻塞。）
- 22-3. 22.10 节指出，接收信号时，利用 `sigwaitinfo()` 调用要比信号处理器外加 `sigsuspend()` 调用的方法来得快。随本书发布的源码中提供的 `signals/sig_speed_ sigsuspend.c` 程序使用 `sigsuspend()` 在父、子进程之间交替发送信号。请对两进程间交换一百万次信号所花费的时间进行计时。（信号交换次数可通过程序命令行参数来提供。）使用 `sigwaitinfo()` 作为替代技术来对程序进行修改，并度量该版本的耗时。两个程序间的速度差异在哪里？
- 22-4. 使用 POSIX 信号 API 来实现 System V 函数 `sigset()`、`sighold()`、`sigrelse()`、`sigignore()` 和 `sigpause()`。

第 23 章

定时器与休眠

定时器是进程规划自己在未来某一时刻接获通知的一种机制。休眠则能使进程（或线程）暂停执行一段时间。本章讨论了定时器设置以及休眠的接口，涵盖主题如下。

- 针对间隔式定时器设置的传统 UNIX API（`setitimer()`和 `alarm()`），一经设定，会在特定的一段时间后通知进程。
- 允许进程休眠特定时间的 API 接口。
- POSIX.1b 时钟和定时器 API 接口。
- Linux 特有的 `timerfd` 功能，允许所创建定时器的到期信息可从文件描述符中读取。

23.1 间隔定时器

系统调用 `setitimer()`创建一个间隔式定时器（interval timer），这种定时器会在未来某个时间点到期，并于此后（可选择地）每隔一段时间到期一次。

```
#include <sys/time.h>

int setitimer(int which, const struct itimerval *new_value,
              struct itimerval *old_value);

Returns 0 on success, or -1 on error
```

通过在调用 `setitimer()`时为 `which` 指定以下值，进程可以创建 3 种不同类型的定时器。

ITIMER_REAL

创建以真实时间倒计时的定时器。到期时会产生 `SIGALARM` 信号并发送给进程。

ITIMER_VIRTUAL

创建以进程虚拟时间（用户模式下的 CPU 时间）倒计时的定时器。到期时会产生信号 `SIGVTALRM`。

ITIMER_PROF

创建一个 `profiling` 定时器，以进程时间（用户态与内核态 CPU 时间的总和）倒计时。到

期时，则会产生 SIGPROF 信号。

对所有这些信号的默认处置 (disposition) 均会终止进程。除非真地期望如此，否则就需要针对这些定时器信号创建处理器函数。

参数 `new_value` 和 `old_value` 均为指向结构 `itimerval` 的指针，结构的定义如下：

```
struct itimerval {
    struct timeval it_interval;    /* Interval for periodic timer */
    struct timeval it_value;      /* Current value (time until
                                next expiration) */
};
```

结构 `itimerval` 中的字段类型均为 `timeval` 结构，`timeval` 又由秒和微秒两部分组成：

```
struct timeval {
    time_t tv_sec;                /* Seconds */
    suseconds_t tv_usec;        /* Microseconds (long int) */
};
```

参数 `new_value` 的下属结构 `it_value` 指定了距离定时器到期的延迟时间。另一下属结构 `it_interval` 则说明该定时器是否为周期性定时器。如果 `it_interval` 的两个字段值均为 0，那么该定时器就属于在 `it_value` 所指定的时间间隔后到期的一次性定时器。只要 `it_interval` 中的任一字段非 0，那么在每次定时器到期之后，都会将定时器重置为在指定间隔后再次到期。

进程只能拥有上述 3 种定时器中的一种。当第 2 次调用 `setitimer()` 时，修改已有定时器的属性要符合参数 `which` 中的类型。如果调用 `setitimer()` 时将 `new_value.it_value` 的两个字段均置为 0，那么会屏蔽任何已有的定时器。

若参数 `old_value` 不为 NULL，则以其所指向的 `itimerval` 结构来返回定时器的前一设置。如果 `old_value.it_value` 的两个字段值均为 0，那么该定时器之前处于屏蔽状态。如果 `old_value.it_interval` 的两个字段值均为 0，那么该定时器之前被设置为历经 `old_value.it_value` 指定时间而到期的一次性定时器。对于需要在新定时器到期后将其还原的情况而言，获取定时器的前一设置就很重要。如果不关心定时器的前一设置，可以将 `old_value` 置为 NULL。

定时器会从初始值 (`it_value`) 倒计时一直到 0 为止。递减为 0 时，会将相应信号发送给进程，随后，如果时间间隔值 (`it_interval`) 非 0，那么会再次将 `it_value` 加载至定时器，重新开始向 0 倒计时。

可以在任何时刻调用 `getitimer()`，以了解定时器的当前状态、距离下次到期的剩余时间。

```
#include <sys/time.h>

int getitimer(int which, struct itimerval *curr_value);

Returns 0 on success, or -1 on error
```

系统调用 `getitimer()` 返回由 `which` 指定定时器的当前状态，并置于由 `curr_value` 所指向的缓冲区中。这与 `setitimer()` 借参数 `old_value` 所返回的信息完全相同，区别则在于 `getitimer()` 无需为了获取这些信息而改变定时器的设置。子结构 `curr_value.it_value` 返回距离下一次到期所剩余的总时间。该值会随定时器倒计时而变化，如果设置定时器时将 `it_interval` 置为非 0 值，那么会在定时器到期时将其重置。子结构 `curr_value.it_interval` 返回定时器的间隔时间，除非再次调用 `setitimer()`，否则该值一直保持不变。

使用 `setitimer()` (和 `alarm()`，稍后讨论) 创建的定时器可以跨越 `exec()` 调用而得以保存，但由 `fork()` 创建的子进程并不继承该定时器。

示例程序

程序清单 23-1 演示了 `setitimer()` 和 `getitimer()` 的使用，所执行的步骤如下。

- 为 `SIGALRM` 信号创建处理器函数③。
- 利用命令行参数为实时（`ITIMER_REAL`）定时器设置到期值及间隔时间④。若未提供命令行参数，程序默认创建一个两秒到期的一次性定时器。
- 进入一个循环⑤，消耗 CPU 时间并周期性地调用函数 `displayTimes()`①。该函数会显示自程序启动以来逝去的真实时间，以及 `ITIMER_REAL` 定时器的当前状态。

每当定时器到期时，都会调用 `SIGALRM` 处理器函数，其中会去设置全局标志 `gotAlarm`②。一旦设置了这一标志，主程序循环就会调用 `displayTimes()` 来显示处理器函数的调用时点以及定时器状态⑥。（之所以采用这一方式来设计信号处理器函数，意在避免从处理器函数内部去调用非异步信号安全的函数，究其原因可参考 21.1.2 节。）如果定时器的时间间隔为 0，那么程序会在第 1 次收到信号时即行退出。否则，程序会在捕获到 3 个信号后终止。⑦

运行程序清单 23-1 中的程序，可以看到如下结果：

```
$ ./real_timer 1 800000 1 0          Initial value 1.8 seconds, interval 1 second
      Elapsed  Value  Interval
START:  0.00
Main:    0.50   1.30   1.00          Timer counts down until expiration
Main:    1.00   0.80   1.00
Main:    1.50   0.30   1.00
ALARM:   1.80   1.00   1.00          On expiration, timer is reloaded from interval
Main:    2.00   0.80   1.00
Main:    2.50   0.30   1.00
ALARM:   2.80   1.00   1.00
Main:    3.00   0.80   1.00
Main:    3.50   0.30   1.00
ALARM:   3.80   1.00   1.00
That's all folks
```

程序清单 23-1：实时定时器的使用

```
timers/real_timer.c

#include <signal.h>
#include <sys/time.h>
#include <time.h>
#include "tlpi_hdr.h"

static volatile sig_atomic_t gotAlarm = 0;
/* Set nonzero on receipt of SIGALRM */

/* Retrieve and display the real time, and (if 'includeTimer' is
TRUE) the current value and interval for the ITIMER_REAL timer */

static void
① displayTimes(const char *msg, Boolean includeTimer)
{
    struct itimerval itv;
    static struct timeval start;
    struct timeval curr;
    static int callNum = 0;          /* Number of calls to this function */
```



```

if (callNum == 0) /* Initialize elapsed time meter */
    if (gettimeofday(&start, NULL) == -1)
        errExit("gettimeofday");

if (callNum % 20 == 0) /* Print header every 20 lines */
    printf("      Elapsed Value Interval\n");
if (gettimeofday(&curr, NULL) == -1)
    errExit("gettimeofday");
printf("%-7s %6.2f", msg, curr.tv_sec - start.tv_sec +
        (curr.tv_usec - start.tv_usec) / 1000000.0);

if (includeTimer) {
    if (getitimer(ITIMER_REAL, &itv) == -1)
        errExit("getitimer");
    printf(" %6.2f %6.2f",
        itv.it_value.tv_sec + itv.it_value.tv_usec / 1000000.0,
        itv.it_interval.tv_sec + itv.it_interval.tv_usec / 1000000.0);
}

printf("\n");
callNum++;
}

static void
sigalrmHandler(int sig)
{
    ② gotAlarm = 1;
}

int
main(int argc, char *argv[])
{
    struct itimerval itv;
    clock_t prevClock;
    int maxSigs; /* Number of signals to catch before exiting */
    int sigCnt; /* Number of signals so far caught */
    struct sigaction sa;

    if (argc > 1 && strcmp(argv[1], "--help") == 0)
        usageErr("%s [secs [usecs [int-secs [int-usecs]]]]\n", argv[0]);

    sigCnt = 0;

    sigemptyset(&sa.sa_mask);
    sa.sa_flags = 0;
    sa.sa_handler = sigalrmHandler;
    ③ if (sigaction(SIGALRM, &sa, NULL) == -1)
        errExit("sigaction");

    /* Exit after 3 signals, or on first signal if interval is 0 */

    maxSigs = (itv.it_interval.tv_sec == 0 &&
        itv.it_interval.tv_usec == 0) ? 1 : 3;

    displayTimes("START:", FALSE);

    /* Set timer from the command-line arguments */

    itv.it_value.tv_sec = (argc > 1) ? getLong(argv[1], 0, "secs") : 2;
    itv.it_value.tv_usec = (argc > 2) ? getLong(argv[2], 0, "usecs") : 0;

```

```

itv.it_interval.tv_sec = (argc > 3) ? getLong(argv[3], 0, "int-secs") : 0;
itv.it_interval.tv_usec = (argc > 4) ? getLong(argv[4], 0, "int-usecs") : 0;

④  if (setitimer(ITIMER_REAL, &itv, 0) == -1)
        errExit("setitimer");

    prevClock = clock();
    sigCnt = 0;

⑤  for (;;) {

        /* Inner loop consumes at least 0.5 seconds CPU time */

        while (((clock() - prevClock) * 10 / CLOCKS_PER_SEC) < 5) {
            ⑥  if (gotAlarm) {                                /* Did we get a signal? */
                    gotAlarm = 0;
                    displayTimes("ALARM:", TRUE);

                    sigCnt++;
                    ⑦  if (sigCnt >= maxSigs) {
                            printf("That's all folks\n");
                            exit(EXIT_SUCCESS);
                        }
                }

            prevClock = clock();
            displayTimes("Main: ", TRUE);
        }
    }
}

```

timers/real_timer.c

更为简单的定时器接口：alarm()

系统调用 `alarm()` 为创建一次性实时定时器提供了一个简单接口。（历史上，`alarm()` 曾是设置定时器的原始 UNIX API。）

```
#include <unistd.h>

unsigned int alarm(unsigned int seconds);

        Always succeeds, returning number of seconds remaining on
        any previously set timer, or 0 if no timer previously was set
```

参数 `seconds` 表示定时器到期的秒数。到期时，会向调用进程发送 `SIGALRM` 信号。调用 `alarm()` 会覆盖对定时器的前一个设置。调用 `alarm(0)` 可屏蔽现有定时器。`alarm()` 的返回值是定时器前一设置距离到期的剩余秒数，如未设置定时器则返回 0。23.3 节提供了一个使用 `alarm()` 的例子。

本书后面的一些示例会使用 `alarm()` 来启动定时器，同时不为 `SIGALRM` 信号设置处理器函数。采用该技术可以确保，即便进程没有终止，也能将其杀死。

setitimer()和 alarm()之间的交互

Linux 中，`alarm()` 和 `setitimer()` 针对同一进程（per-process）共享同一实时定时器，这也意

味着，无论调用两者之中的哪个完成了对定时器的前一设置，同样可以调用二者中的任一函数来改变这一设置。其他 UNIX 系统的情况可能会有所不同（也就是说，这两个函数可能分别控制着不同的定时器）。对于 `setitimer()`与 `alarm()`之间的交互，以及二者与 `sleep()`函数（23.4.1 节）之间的交互，SUSv3 均未加以规范。为了确保应用程序可移植性的最大化，程序设置实时定时器的函数只能在二者中选择其一。

23.2 定时器的调度及精度

取决于当前负载和对进程的调度，系统可能会在定时器到期的瞬间（通常是几分之一秒）之后才去调度其所属进程。尽管如此，由 `setitimer()`或本章后续介绍的其他接口所创建的周期性定时器，在到期后依然会恪守其规律性。例如，假设设置一个实时定时器每两秒到期一次，虽然上述延迟可能会影响每个定时器事件的送达，但系统对后续定时器到期的调度依然会严格遵循两秒的时间间隔。换言之，间隔式定时器不受潜在错误左右。

虽然 `setitimer()`使用的 `timeval` 结构提供有微秒级精度，但是传统意义上定时器精度还是受制于软件时钟（10.6 节）频率。如果定时器值未能与软件时钟间隔的倍数严格匹配，那么定时器值则会向上取整。也就是说，假如有一个间隔为 19100 微秒（刚刚超过 19 毫秒）的定时器，如果 `jiffy`（软件时钟周期）为 4 毫秒，那么定时器实际上会每隔 20 毫秒过期一次。

高分辨率定时器

对于现代 Linux 内核而言，适才关于定时器分辨率受限于软件时钟频率的论断已经不再成立。自版本 2.6.21 开始，Linux 内核可选择是否支持高分辨率定时器。如果选择支持（通过内核配置选项 `CONFIG_HIGH_RES_TIMERS`），那么本章各种定时器以及休眠接口的精度则不再受内核 `jiffy`（软件时钟周期）的影响，可以达到底层硬件所支持的精度。在现代硬件平台上，精度达到微秒级是司空见惯的事情。

23.5.1 节将介绍函数 `clock_getres()`，可以用其返回值来判断系统是否支持高分辨率定时器。

23.3 为阻塞操作设置超时

实时定时器的用途之一是为某个阻塞系统调用设置其处于阻塞状态的时间上限。例如，当用户在一段时间内没有输入整行命令时，可能希望取消对终端的 `read()`操作。处理如下。

1. 调用 `sigaction()`为 `SIGALRM` 信号创建处理器函数，排除 `SA_RESTART` 标志以确保系统调用不会重新启动（参考 21.5 节）。
2. 调用 `alarm()`或 `setitimer()`来创建一个定时器，同时设定希望系统调用阻塞的时间上限。
3. 执行阻塞式系统调用。
4. 系统调用返回后，再次调用 `alarm()`或 `setitimer()`以屏蔽定时器（以防止系统调用在定时器到期之前就已完成的情况）。
5. 检查系统调用失败时是否将 `errno` 置为 `EINTR`（系统调用遭到中断）。

程序清单 23-2 针对 `read()`调用展示了这一技术，创建定时器时使用的是 `alarm()`。

程序清单 23-2: 运行设置了超时的 read()

```
timers/timed_read.c

#include <signal.h>
#include "tspi_hdr.h"

#define BUF_SIZE 200

static void /* SIGALRM handler: interrupts blocked system call */
handler(int sig)
{
    printf("Caught signal\n");          /* UNSAFE (see Section 21.1.2) */
}

int
main(int argc, char *argv[])
{
    struct sigaction sa;
    char buf[BUF_SIZE];
    ssize_t numRead;
    int savedErrno;

    if (argc > 1 && strcmp(argv[1], "--help") == 0)
        usageErr("%s [num-secs [restart-flag]]\n", argv[0]);

    /* Set up handler for SIGALRM. Allow system calls to be interrupted,
       unless second command-line argument was supplied. */

    sa.sa_flags = (argc > 2) ? SA_RESTART : 0;
    sigemptyset(&sa.sa_mask);
    sa.sa_handler = handler;
    if (sigaction(SIGALRM, &sa, NULL) == -1)
        errExit("sigaction");

    alarm((argc > 1) ? getInt(argv[1], GN_NONNEG, "num-secs") : 10);

    numRead = read(STDIN_FILENO, buf, BUF_SIZE - 1);

    savedErrno = errno;          /* In case alarm() changes it */
    alarm(0);                  /* Ensure timer is turned off */
    errno = savedErrno;

    /* Determine result of read() */

    if (numRead == -1) {
        if (errno == EINTR)
            printf("Read timed out\n");
        else
            errMsg("read");
    } else {
        printf("Successful read (%ld bytes): %.*s",
              (long) numRead, (int) numRead, buf);
    }

    exit(EXIT_SUCCESS);
}

timers/timed_read.c
```

注意，程序清单 23-2 中程序理论上存在导致竞争条件的可能性。如果定时器到期时处于 `alarm()` 调用之后，`read()` 调用之前，那么信号处理器函数将不会中断 `read()`。由于在这种场景下设定的超时值一般相对较大（至少几秒），故而发生上述情况的概率极低，因此这种技术实际上是可行的。[Stevens & Rago, 2005] 推荐了另一种方法，使用的是 `longjmp()`。在处理 I/O 系统调用时，还有另一种备选方案，利用了系统调用 `select()` 或 `poll()`（第 63 章）的超时特性，锦上添花的是还能同时等待多路描述符的 I/O。

23.4 暂停运行（休眠）一段固定时间

有时需要将进程挂起（固定的）一段时间。将前述定时器函数与 `sigsuspend()` 相结合固然可以达到这一目的，但使用休眠函数会更为简单。

23.4.1 低分辨率休眠：sleep()

函数 `sleep()` 可以暂停调用进程的执行达数秒之久（由参数 `seconds` 设置），或者在捕获到信号（从而中断调用）后恢复进程的运行。

```
#include <unistd.h>

unsigned int sleep(unsigned int seconds);

Returns 0 on normal completion, or number of
unlept seconds if prematurely terminated
```

如果休眠正常结束，`sleep()` 返回 0。如果因信号而中断休眠，`sleep()` 将返回剩余（未休眠）的秒数。与 `alarm()` 和 `setitimer()` 所设置的定时器相同，由于系统负载的原因，内核可能会在完成 `sleep()` 的一段（通常很短）时间后才对进程重新加以调度。

对于 `sleep()` 和 `alarm()` 以及 `setitimer()` 之间的交互方式，SUSv3 并未加以规范。Linux 将 `sleep()` 实现为对 `nanosleep()`（23.4.2 节）的调用，其结果是 `sleep()` 与定时器函数之间并无交互。不过，许多其他的实现，尤其是一些老系统，会使用 `alarm()` 以及 `SIGALRM` 信号处理器函数来实现 `sleep()`。考虑到可移植性，应避免将 `sleep()` 和 `alarm()` 以及 `setitimer()` 混用。

23.4.2 高分辨率休眠：nanosleep()

函数 `nanosleep()` 的功用与 `sleep()` 类似，但更具优势，其中包括能以更高分辨率来设定休眠间隔时间。

```
#define _POSIX_C_SOURCE 199309
#include <time.h>

int nanosleep(const struct timespec *request, struct timespec *remain);

Returns 0 on successfully completed sleep,
or -1 on error or interrupted sleep
```

参数 `request` 指定了休眠的持续时间，是一个指向如下结构的指针：

```
struct timespec {
    time_t tv_sec;        /* Seconds */
    long tv_nsec;        /* Nanoseconds */
};
```

tv_nsec 字段为纳秒值，取值范围在 0~999999999 之间。

nanosleep()的更大优势在于，SUSv3 明文规定不得使用信号来实现该函数。这意味着，与 sleep()不同，即使将 nanosleep()与 alarm()或 setitimer()混用，也不会危及程序的可移植性。

尽管 nanosleep()的实现并未使用信号，但还是可以通过信号处理器函数来将其中断。这时，nanosleep()将返回-1，并将 errno 置为 EINTR。同时，若参数 remain 不为 NULL，则该指针所指向的缓冲区将返回剩余的休眠时间。可利用这一返回值重启该系统调用以完成休眠。程序清单 23-3 演示了这一用途。程序从命令行参数中获取传入 nanosleep()的秒和纳秒值，并反复循环执行 nanosleep()，直至耗尽全部的休眠间隔时间。如果信号 SIGINT（按下 Ctrl-C 产生）的处理器函数将 nanosleep()中断，那么会以参数 remain 中的返回值重新调用 nanosleep()。其运行结果如下：

```
$ ./t_nanosleep 10 0                Sleep for 10 seconds
Type Control-C
Slept for: 1.853428 secs
Remaining: 8.146617000
Type Control-C
Slept for: 4.370860 secs
Remaining: 5.629800000
Type Control-C
Slept for: 6.193325 secs
Remaining: 3.807758000
Slept for: 10.008150 secs
Sleep complete
```

虽然 nanosleep()允许设定纳秒级精度的休眠间隔值，但其精度依然受制于软件时钟的间隔大小（10.6 节）。如果指定的间隔值并非软件时钟间隔的整数倍，那么会对其向上取整。

前文曾提及，在支持高精度定时器的系统中，休眠时间间隔的精度要比软件时钟间隔精细许多。

当以高频率接收信号时，这一取整行为会给程序清单 23-3 中程序所采用的编程手法带来问题。由于返回的 remain 时间未必是软件时钟间隔的整数倍，故而 nanosleep()的每次重启都会遭遇取整错误。其结果是，nanosleep()每次重启后的休眠时间都要长于前一调用返回的 remain 值。在信号接收频率极高的情况下（与软件时钟间隔的频率一致或更高），进程的休眠可能永远也完成不了。Linux 2.6 中，使用带有 TIMER_ABSTIME 选项的 clock_nanosleep()可以避免这一问题。23.5.4 节将对 clock_nanosleep()加以讨论。

在 2.4 以及更早期的 Linux 内核版本中，nanosleep()的实现存在着一种奇怪的特性。假设正在执行 nanosleep()的进程因信号而停止，当进程于稍后截获 SIGCONT 而继续运行时，nanosleep()会如期调用失败并返回 EINTR 错误。不过，如果进程接着重启 nanosleep()调用，那么进程处于停止状态所消耗的时间将不会计入休眠间隔时间，进程的休眠时间也就比预期的要久。Linux 2.6 中去除了这一怪异特性，nanosleep()在收到 SIGCONT 信号时将自动恢复，进程处于停止状态所消耗的时间也会计入休眠间隔时间。

程序清单 23-3：使用 nanosleep()

```
-----timers/t_nanosleep.c
#define _POSIX_C_SOURCE 199309
#include <sys/time.h>
```

```

#include <time.h>
#include <signal.h>
#include "tspi_hdr.h"

static void
sigintHandler(int sig)
{
    return;                /* Just interrupt nanosleep() */
}

int
main(int argc, char *argv[])
{
    struct timeval start, finish;
    struct timespec request, remain;
    struct sigaction sa;
    int s;

    if (argc != 3 || strcmp(argv[1], "--help") == 0)
        usageErr("%s secs nanosecs\n", argv[0]);

    request.tv_sec = getLong(argv[1], 0, "secs");
    request.tv_nsec = getLong(argv[2], 0, "nanosecs");

    /* Allow SIGINT handler to interrupt nanosleep() */

    sigemptyset(&sa.sa_mask);
    sa.sa_flags = 0;
    sa.sa_handler = sigintHandler;
    if (sigaction(SIGINT, &sa, NULL) == -1)
        errExit("sigaction");

    if (gettimeofday(&start, NULL) == -1)
        errExit("gettimeofday");

    for (;;) {
        s = nanosleep(&request, &remain);
        if (s == -1 && errno != EINTR)
            errExit("nanosleep");

        if (gettimeofday(&finish, NULL) == -1)
            errExit("gettimeofday");
        printf("Slept for: %9.6f secs\n", finish.tv_sec - start.tv_sec +
              (finish.tv_nsec - start.tv_nsec) / 1000000.0);

        if (s == 0)
            break;                /* nanosleep() completed */

        printf("Remaining: %2ld.%09ld\n", (long) remain.tv_sec,
              remain.tv_nsec);
        request = remain;        /* Next sleep is with remaining time */
    }

    printf("Sleep complete\n");
    exit(EXIT_SUCCESS);
}

```

timers/t_nanosleep.c

23.5 POSIX 时钟

POSIX 时钟（原定义于 POSIX.1b）所提供的时钟访问 API 可以支持纳秒级的时间精度，其中表示纳秒级时间值的 `timespec` 结构同样也用于 `nanosleep()`（23.4.2 节）调用。

Linux 中，调用此 API 的程序必须以 `-lrt` 选项进行编译，从而与 `librt`（realtime，实时）函数库相链接。

POSIX 时钟 API 的主要系统调用包括获取时钟当前值的 `clock_gettime()`、返回时钟分辨率的 `clock_getres()`，以及更新时钟的 `clock_settime()`。

23.5.1 获取时钟的值：clock_gettime()

系统调用 `clock_gettime()` 针对参数 `clockid` 所指定的时钟返回时间。

```
#define _POSIX_C_SOURCE 199309
#include <time.h>

int clock_gettime(clockid_t clockid, struct timespec *tp);
int clock_getres(clockid_t clockid, struct timespec *res);

Both return 0 on success, or -1 on error
```

返回的时间值置于 `tp` 指针所指向的 `timespec` 结构中。虽然 `timespec` 结构提供了纳秒级精度，但 `clock_gettime()` 返回的时间值粒度可能还是要更大一点。系统调用 `clock_getres()` 在参数 `res` 中返回指向 `timespec` 结构的指针，机构中包含了由 `clockid` 所指定时钟的分辨率。

`clockid_t` 是一种由 SUSv3 定义的数据类型，用于表示时钟标识符。表 23-1 中第 1 列值即可用于设定 `clockid`。

表 23-1: POSIX.1b 时钟类型

时钟 ID	描 述
<code>CLOCK_REALTIME</code>	可设定的系统级实时时钟
<code>CLOCK_MONOTONIC</code>	不可设定的恒定态时钟
<code>CLOCK_PROCESS_CPUTIME_ID</code>	每进程 CPU 时间的时钟（自 Linux 2.6.12）
<code>CLOCK_THREAD_CPUTIME_ID</code>	每线程 CPU 时间的时钟（自 Linux 2.6.12）

`CLOCK_REALTIME` 时钟是一种系统级时钟，用于度量真实时间。与 `CLOCK_MONOTONIC` 时钟不同，它的设置是可以变更的。

SUSv3 规定，`CLOCK_MONOTONIC` 时钟对时间的度量始于“未予规范的过去某一点”，系统启动后就不会发生改变。该时钟适用于那些无法容忍系统时钟发生跳跃性变化（例如：手工改变了系统时间）的应用程序。Linux 上，这种时钟对时间的测量始于系统启动。

`CLOCK_PROCESS_CPUTIME_ID` 时钟测量调用进程所消耗的用户和系统 CPU 时间。`CLOCK_THREAD_CPUTIME_ID` 时钟的功用与之相类似，不过测量对象是进程中的单条线程。

SUSv3 规范了表 23-1 中的所有时钟，但强制要求实现的仅有 `CLOCK_REALTIME` 一种，这同时也是受到 UNIX 实现广泛支持的时钟类型。

Linux 2.6.28 增加了一种新的时钟类型：CLOCK_MONOTONIC_RAW。类似于 CLOCK_MONOTONIC，这也是一种无法设置的时钟，但是提供了对纯基于硬件时间的访问，且不受 NTP 时间调整的影响。这种非标准时钟适用于专业时钟同步应用程序。

Linux 2.6.35 又提供了两种新时钟：CLOCK_REALTIME_COARSE 和 CLOCK_MONTIC_COARSE。这些时钟类似于 CLOCK_REALTIME 和 CLOCK_MONTONIC，适用于那些希望以最小代价获取较低分辨率时间戳的程序。这些非标准时钟不会引发对硬件时钟的任何访问（访问某些硬件时钟源的代价高昂），其返回值的分辨率为 jiffy（软件时钟周期，见 10.6 节）。

23.5.2 设置时钟的值：clock_settime()

系统调用 clock_settime() 利用参数 tp 所指向缓冲区中的时间来设置由 clockid 指定的时钟。

```
#define _POSIX_C_SOURCE 199309
#include <time.h>

int clock_settime(clockid_t clockid, const struct timespec *tp);

Returns 0 on success, or -1 on error
```

如果由 tp 指定的时间并非由 clock_getres() 所返回时钟分辨率的整数倍，时间会向下取整。

特权级 (CAP_SYS_TIME) 进程可以设置 CLOCK_REALTIME 时钟。该时钟的初始值通常是自 Epoch（1970 年 1 月 1 日 0 点 0 分 0 秒）以来的时间。表 23-1 中的其他时钟均不可更改。

根据 SUSv3，系统实现可允许设置 CLOCK_PROCESS_CPUTIME_ID 和 CLOCK_THREAD_CPUTIME_ID 型时钟。撰写本书之际，这些时钟在 Linux 上依然是只读属性。

23.5.3 获取特定进程或线程的时钟 ID

要测量特定进程或线程所消耗的 CPU 时间，首先可借助本节所描述的函数来获取其时钟 ID。接着再以此返回 id 去调用 clock_gettime()，从而获得进程或线程耗费的 CPU 时间。

函数 clock_getcpuclockid() 会将隶属于 pid 进程的 CPU 时间时钟的标识符置于 clockid 指针所指向的缓冲区中。

```
#define _XOPEN_SOURCE 600
#include <time.h>

int clock_getcpuclockid(pid_t pid, clockid_t *clockid);

Returns 0 on success, or a positive error number on error
```

参数 pid 为 0 时，clock_getcpuclockid() 返回调用进程的 CPU 时间时钟 ID。

函数 pthread_getcpuclockid() 是 clock_getcpuclockid() 的 POSIX 线程版，返回的标识符所标识的时钟用于度量调用进程中指定线程消耗的 CPU 时间。

```
#define _XOPEN_SOURCE 600
#include <pthread.h>
#include <time.h>

int pthread_getcpuclockid(pthread_t thread, clockid_t *clockid);

Returns 0 on success, or a positive error number on error
```

参数 `thread` 是 POSIX 线程 ID，用于指定希望获取的 CPU 时钟 ID 所从属的线程。返回的时钟 ID 存放于 `clockid` 指针所指向的缓冲区中。

23.5.4 高分辨率休眠的改进版：`clock_nanosleep()`

类似于 `nanosleep()`，Linux 特有的 `clock_nanosleep()` 系统调用也可以暂停调用进程，直到历经一段指定的时间间隔后，亦或是收到信号才恢复运行。本节将讨论两者间的差异。

```
#define _XOPEN_SOURCE 600
#include <time.h>

int clock_nanosleep(clockid_t clockid, int flags,
    const struct timespec *request, struct timespec *remain);

Returns 0 on successfully completed sleep,
or a positive error number on error or interrupted sleep
```

参数 `request` 及 `remain` 同 `nanosleep()` 中的对应参数目的相似。

默认情况下（即 `flags` 为 0），由 `request` 指定的休眠间隔时间是相对时间（类似于 `nanosleep()`）。不过，如果在 `flags`（参考程序清单 23-4）中设定了 `TIMER_ABSTIME`，`request` 则表示 `clockid` 时钟所测量的绝对时间。这一特性对于那些需要精确休眠一段指定时间的应用程序至关重要。如果只是先获取当前时间，计算与目标时间的差距，再以相对时间进行休眠，进程可能执行到一半就被占先了¹，结果休眠时间比预期的要久。

如 23.4.2 节所述，对于那些被信号处理器函数中断并使用循环重启休眠的进程来说，“嗜睡（oversleeping）”问题尤其明显。如果以高频率接收信号，那么按相对时间休眠（`nanosleep()` 所执行的类型）的进程在休眠时间上会有较大误差。但可以通过如下方式来避免嗜睡问题：先调用 `clock_gettime()` 获取时间，加上期望休眠的时间量，再以 `TIMER_ABSTIME` 标志调用 `clock_nanosleep()` 函数（并且，如果被信号处理器中断，则会重启系统调用）。

指定 `TIMER_ABSTIME` 时，不再（且不需要）使用参数 `remain`。如果信号处理器程序中断了 `clock_nanosleep()` 调用，再次调用该函数来重启休眠时，`request` 参数不变。

将 `clock_nanosleep()` 与 `nanosleep()` 区分开来的另一特性在于，可以选择不同的时钟来测量休眠间隔时间。可在 `clockid` 中指定所期望的时钟 `CLOCK_REALTIME`、`CLOCK_MONOTONIC` 或 `CLOCK_PROCESS_CPUTIME_ID`。请参考表 23-1 对这些时钟的描述。

程序清单 23-4 演示了 `clock_nanosleep()` 的用法：针对 `CLOCK_REALTIME` 时钟，以绝对时间休眠 20 秒。

程序清单 23-4：使用 `clock_nanosleep()`

```
struct timespec request;

/* Retrieve current value of CLOCK_REALTIME clock */

if (clock_gettime(CLOCK_REALTIME, &request) == -1)
    errExit("clock_gettime");

request.tv_sec += 20;          /* Sleep for 20 seconds from now */
```

¹ 译者注：内核调度所为。

```

s = clock_nanosleep(CLOCK_REALTIME, TIMER_ABSTIME, &request, NULL);
if (s != 0) {
    if (s == EINTR)
        printf("Interrupted by signal handler\n");
    else
        errExitEN(s, "clock_nanosleep");
}

```

23.6 POSIX 间隔式定时器

使用 `setitimer()` 来设置经典 UNIX 间隔式定时器，会受到如下制约。

- 针对 `ITIMER_REAL`、`ITIMER_VIRTUAL` 和 `ITIMER_PROF` 这 3 类定时器，每种只能设置一个。
- 只能通过发送信号的方式来通知定时器到期。另外，也不能改变到期时产生的信号。
- 如果一个间隔式定时器到期多次，且相应信号遭到阻塞时，那么会只调用一次信号处理器函数。换言之，无从知晓是否出现过定时器溢出（`timer overrun`）的情况。
- 定时器的分辨率只能达到微秒级。不过，一些系统的硬件时钟提供了更为精细的时钟分辨率，软件此时应采用这一较高分辨率。

POSIX.1b 定义了一套 API 来突破这些限制，Linux 2.6 实现了这一 API。

在较老的 Linux 系统上，`glibc` 通过基于线程的实现提供了这一 API 的不完整版。不过，这种用户空间内的实现是无法提供此处描述的所有特性的。

POSIX 定时器 API 将定时器生命周期划分为如下几个阶段。

- 以系统调用 `timer_create()` 创建一个新定时器，并定义其到期时对进程的通知方法。
- 以系统调用 `timer_settime()` 来启动或停止一个定时器。
- 以系统调用 `timer_delete()` 删除不再需要的定时器。

由 `fork()` 创建的子进程不会继承 POSIX 定时器。调用 `exec()` 期间亦或进程终止时将停止并删除定时器。

Linux 上，调用 POSIX 定时器 API 的程序编译时应使用 `-lrt` 选项，从而与 `librt`（实时）函数库相链接。

23.6.1 创建定时器：`timer_create()`

函数 `timer_create()` 创建一个新定时器，并以由 `clockid` 指定的时钟来进行时间度量。

```

#define _POSIX_C_SOURCE 199309
#include <signal.h>
#include <time.h>

int timer_create(clockid_t clockid, struct sigevent *evp, timer_t *timerid);

```

Returns 0 on success, or -1 on error

设置参数 `clockid`，可以使用表 23-1 中的任意值，也可以采用 `clock_getcpuclockid()` 或 `pthread_getcpuclockid()` 返回的 `clockid` 值。函数返回时会在参数 `timerid` 所指向的缓冲区中放置定时器句柄（`handle`），供后续调用中指代该定时器之用。这一缓冲区的类型为 `timer_t`，是一

种由 SUSv3 定义的数据类型，用于标识定时器。

参数 `evp` 可决定定时器到期时对应用程序的通知方式，指向类型为 `sigevent` 的数据结构，具体定义如下：

```

union signal {
    int    sival_int;           /* Integer value for accompanying data */
    void *sival_ptr;          /* Pointer value for accompanying data */
};

struct sigevent {
    int    sigev_notify;      /* Notification method */
    int    sigev_signo;      /* Timer expiration signal */
    union signal sigev_value; /* Value accompanying signal or
                               passed to thread function */

    union {
        pid_t    _tid;      /* ID of thread to be signaled /
        struct {
            void (*_function) (union signal); /* Thread notification function */
            void *_attribute; /* Really 'pthread_attr_t *' */
        } _sigev_thread;
    } _sigev_un;
};

#define sigev_notify_function    _sigev_un._sigev_thread._function
#define sigev_notify_attributes _sigev_un._sigev_thread._attribute
#define sigev_notify_thread_id  _sigev_un._tid

```

可以表 23-2 所示值之一来设置结构中的 `sigev_notify` 字段。

表 23-2: `sigevent` 结构中 `sigev_notify` 字段的值

sigev_notify 的值	通知方法	SUSv3
SIGEV_NONE	不通知；使用 <code>timer_gettime()</code> 监测定时器	•
SIGEV_SIGNAL	发送 <code>sigev_signo</code> 信号给进程	•
SIGEV_THREAD	调用 <code>sigev_notify_function</code> 作为新线程的启动函数	•
SIGEV_THREAD_ID	发送 <code>sigev_signo</code> 信号给 <code>sigev_notify_thread_id</code> 所标识的线程	

关于 `sigev_notify` 常量值的更多细节，以及 `signal` 结构中每个常量值相关的字段，特做如下说明。

SIGEV_NONE

不提供定时器到期通知。进程可以使用 `timer_gettime()` 来监控定时器的运转情况。

SIGEV_SIGNAL

定时器到期时，为进程生成指定于 `sigev_signo` 中的信号。如果 `sigev_signal` 为实时信号，那么 `sigev_value` 字段则指定了信号的伴随数据（整型或指针）（22.8.1 节）。通过 `siginfo_t` 结构的 `si_value` 可获取这一数据，至于 `siginfo_t` 结构，既可以直接传递给该信号的处理函数，也可以由调用 `sigwaitinfo()` 或 `sigtimerdwait()` 返回。

SIGEV_THREAD

定时器到期时，会调用由 `sigev_notify_function` 字段指定的函数。调用该函数类似于调用新线程的启动函数。上述措词摘自 SUSv3，即允许系统实现以如下两种方式为周期性定时器

产生通知：要么将每个通知分别传递给一个唯一的新线程，要么将通知成系列发送给单个新线程。可将 `sigev_notify_attribytes` 字段置为 `NULL`，或是指向 `pthread_attr_t` 结构的指针，并在结构中定义线程属性。在 `sigev_value` 中设定的联合体 `sigval` 值是传递给函数的唯一参数。

SIGEV_THREAD_ID

这与 `SIGEV_SIGNAL` 相类似，只是发送信号的目标线程 ID 要与 `sigev_notify_thread_id` 相匹配。该线程应与调用线程同属一个进程。（伴随 `SIGEV_SIGNAL` 通知，会将信号置于针对整个进程的一个队列中排队，并且，如果进程包含多条线程，那么可将信号传递给进程中的任意线程。）可用 `clone()` 或 `gettid()` 的返回值对 `sigev_notify_thread_id` 赋值。设计 `SIGEV_THREAD_ID` 标志，意在供线程库使用。（要求线程实现使用 28.2.1 节描述的 `CLONE_THREAD` 选项。现代 NPTL 线程实现采用了 `CLONE_THREAD`，但较老的 LinuxThreads 线程则没有。）

除去 Linux 系统特有的 `SIGEV_THREAD_ID` 之外，SUSv3 定义了上述所有常量。

将参数 `evp` 置为 `NULL`，这相当于将 `sigev_notify` 置为 `SIGEV_SIGNAL`，同时将 `sigev_signo` 置为 `SIGALRM`（这与其他系统可能会有出入，因为 SUSv3 的措词是：一个缺省的信号值），并将 `sigev_value.sival_int` 置为定时器 ID。

在当前实现中，内核会为每个用 `timer_create()` 创建的 POSIX 定时器在队列中预分配一个实时信号结构。之所以要采取预分配，旨在确保当定时器到期时，至少有一个有效结构可服务于所产生的队列化信号。这也意味着可以创建的 POSIX 定时器数量受制于排队实时信号的数量（参考 22.8 节）。

23.6.2 配备和解除定时器：timer_settime()

一旦创建了定时器，就可以使用 `timer_settime()` 对其进行配备（启动）或解除（停止）。

```
#define _POSIX_C_SOURCE 199309
#include <time.h>

int timer_settime(timer_t timerid, int flags, const struct itimerspec *value,
                  struct itimerspec *old_value);

Returns 0 on success, or -1 on error
```

函数 `timer_settime()` 的参数 `timerid` 是一个定时器句柄（handle），由之前对 `timer_create()` 的调用返回。

参数 `value` 和 `old_value` 则类似于函数 `setitimer()` 的同名参数：`value` 中包含定时器的新设置，`old_value` 则用于返回定时器的前一设置（参考稍后对 `timer_gettime()` 的说明）。如果对定时器的前一设置不感兴趣，可将 `old_value` 设为 `NULL`。参数 `value` 和 `old_value` 都是指向结构 `itimerspec` 的指针，该结构定义如下：

```
struct itimerspec {
    struct timespec it_interval; /* Interval for periodic timer */
    struct timespec it_value; /* First expiration */
};
```

结构 `itimerspec` 中的所有字段都是 `timespec` 类型的结构，用秒和纳秒来指定时间：

```
struct timespec {
    time_t tv_sec; /* Seconds */
    long tv_nsec; /* Nanoseconds */
};
```

`it_value` 指定了定时器首次到期的时间。如果 `it_interval` 的任一子字段非 0，那么这就是一个周期性定时器，在经历了由 `it_value` 指定的初次到期后，会按这些子字段指定的频率周期性到期。如果 `it_interval` 的下属字段均为 0，那么这个定时器将只到期一次。

若将 `flags` 置为 0，则会将 `value.it_value` 视为始于 `timer_settime()`（与 `setitimer()` 类似）调用时间点的相对值。如果将 `flags` 设为 `TIMER_ABSTIME`，那么 `value.it_value` 则是一个绝对时间（从时钟值 0 开始）。一旦时钟过了这一时间，定时器会立即到期。

为了启动定时器，需要调用函数 `timer_settime()`，并将 `value.it_value` 的一个或全部下属字段设为非 0 值。如果之前曾经配备过定时器，`timer_settime()` 会将之前的设置替换掉。

如果定时器的值和间隔时间并非对应时钟分辨率（由 `clock_getres()` 返回）的整数倍，那么会对这些值做向上取整处理。

定时器每次到期时，都会按特定方式通知进程，这种方式由创建定时器的 `timer_create()` 定义。如果结构 `it_interval` 包含非 0 值，那么会用这些值来重新加载 `it_value` 结构。

要解除定时器，需要调用 `timer_settime()`，并将 `value.it_value` 的所有字段指定为 0。

23.6.3 获取定时器的当前值：timer_gettime()

系统调用 `timer_gettime()` 返回由 `timerid` 指定 POSIX 定时器的间隔以及剩余时间。

```
#define _POSIX_C_SOURCE 199309
#include <time.h>

int timer_gettime(timer_t timerid, struct itimerspec *curr_value);

Returns 0 on success, or -1 on error
```

`curr_value` 指针所指向的 `itimerspec` 结构中返回的是时间间隔以及距离下次定时器到期的时间。即使是以 `TIMER_ABSTIME` 标志创建的绝对时间定时器，在 `curr_value.it_value` 字段中返回的也是距离定时器下次到期的时间值。

如果返回结构 `curr_value.it_value` 的两个字段均为 0，那么定时器当前处于停止状态。如果返回结构 `curr_value.it_interval` 的两个字段都是 0，那么该定时器仅在 `curr_value.it_value` 给定的时间到期过一次。

23.6.4 删除定时器：timer_delete()

每个 POSIX 定时器都会消耗少量系统资源。所以，一旦使用完毕，应当用 `timer_delete()` 来移除定时器并释放这些资源。

```
#define _POSIX_C_SOURCE 199309
#include <time.h>

int timer_delete(timer_t timerid);

Returns 0 on success, or -1 on error
```

参数 `timerid` 是之前调用 `timer_create()` 时返回的句柄。对于已启动的定时器，会在移除前自动将其停止。如果因定时器到期而已经存在待定（`pending`）信号，那么信号会保持这一状态。（SUSv3 对此并未加以规范，所以其他的一些 UNIX 实现可能会有不同行为。）当进程终止时，会自动删除所有定时器。

23.6.5 通过信号发出通知

如果选择通过信号来接收定时器通知，那么处理这些信号时既可以采用信号处理器函数，也可以调用 `sigwaitinfo()` 或是 `sigtimedwait()`。接收进程借助于这两种方法可以获得一个 `siginfo_t` 结构（21.4 节），其中包含与信号相关的深入信息。（要在信号处理器函数中使用这种特性，创建信号处理器函数时需设置 `SA_SIGINFO` 标志。）在结构 `siginfo_t` 中设置如下字段。

- `si_signo`: 包含由定时器产生的信号。
- `si_code`: 置为 `SI_TIMER`，表示这是因 POSIX 定时器到期而产生的信号。
- `si_value`: 将该字段置为以 `timer_create()` 创建定时器时在 `evp.sigev_value` 中提供的值。为 `evp.sigev_value` 指定不同的值，可以将到期时发送同类信号的不同定时器区分开来。

调用 `timer_create()` 时，通常将 `evp.sigev_value.sival_ptr` 赋值为当前调用中参数 `timerid` 的地址（见程序清单 23-5）。从而允许信号处理器函数（或 `sigwaitinfo()` 调用）获得产生信号的定时器 ID。（另外，也可以将调用函数 `timer_create()` 时给定的 `timerid` 参数置于一结构中，并将结构地址赋予 `evp.sigev_value.sival_ptr`。）

Linux 还为 `siginfo_t` 结构提供了如下非标准字段。

- `si_overrun`: 包含了定时器溢出个数（在 23.6.6 节中说明）。

Linux 还支持另一个非标准字段 `si_timerid`，其中包含一个标识符，供系统内部识别定时器之用（与 `timer_create()` 返回的 ID 不同）。对于应用程序来说没什么用处。

程序清单 23-5 所演示的是使用信号作为 POSIX 定时器的通知机制。

程序清单 23-5: 使用信号进行 POSIX 定时器通知

```
-----timers/ptmr_sigev_signal.c
#define _POSIX_C_SOURCE 199309
#include <signal.h>
#include <time.h>
#include "curr_time.h"          /* Declares currTime() */
#include "itimerspec_from_str.h" /* Declares itimerspecFromStr() */
#include "tspi_hdr.h"

#define TIMER_SIG SIGRTMAX     /* Our timer notification signal */

static void
① handler(int sig, siginfo_t *si, void *uc)
{
    timer_t *tidptr;

    tidptr = si->si_value.sival_ptr;

    /* UNSAFE: This handler uses non-async-signal-safe functions
       (printf(); see Section 21.1.2) */
    printf("[%s] Got signal %d\n", currTime("%T"), sig);
    printf("    *sival_ptr      = %ld\n", (long) *tidptr);
    printf("    timer_getoverrun() = %d\n", timer_getoverrun(*tidptr));
}

int
main(int argc, char *argv[])
{
```

```

struct itimerspec ts;
struct sigaction sa;
struct sigevent sev;
timer_t *tidlist;
int j;

if (argc < 2)
    usageErr("%s secs[/nsecs][:int-secs[/int-nsecs]]...\n", argv[0]);

tidlist = calloc(argc - 1, sizeof(timer_t));
if (tidlist == NULL)
    errExit("malloc");

/* Establish handler for notification signal */

sa.sa_flags = SA_SIGINFO;
sa.sa_sigaction = handler;
sigemptyset(&sa.sa_mask);
② if (sigaction(TIMER_SIG, &sa, NULL) == -1)
    errExit("sigaction");

/* Create and start one timer for each command-line argument */

sev.sigev_notify = SIGEV_SIGNAL; /* Notify via signal */
sev.sigev_signo = TIMER_SIG; /* Notify using this signal */

for (j = 0; j < argc - 1; j++) {
③ itimerspecFromStr(argv[j + 1], &ts);

    sev.sigev_value.sival_ptr = &tidlist[j];
    /* Allows handler to get ID of this timer */

④ if (timer_create(CLOCK_REALTIME, &sev, &tidlist[j]) == -1)
    errExit("timer_create");
    printf("Timer ID: %ld (%s)\n", (long) tidlist[j], argv[j + 1]);

⑤ if (timer_settime(tidlist[j], 0, &ts, NULL) == -1)
    errExit("timer_settime");
}

⑥ for (;;) /* Wait for incoming timer signals */
    pause();
}

```

timers/ptmr_sigev_signal.c

程序清单 23-5 程序的每个命令行参数都为定时器指定了初始值及间隔时间。程序的“用法”输出中描述了这些参数的语法，并在后面的 shell 会话中做了演示。程序执行的步骤如下。

- 为用于定时器通知的信号创建处理器函数②。
- 为每一个命令行参数，创建④并配备⑤一个使用 SIGEV_SIGNAL 通知机制的 POSIX 定时器。至于将命令行参数转换③为 itimerspec 结构的函数 itimerspecFromStr(), 请参考程序清单 23-6。
- 每当一个定时器到期时，都将发送由 sev.sigev_signo 指定的信号给进程。信号处理器函数会将 sev.sigev_value.sival_ptr 中提供的值（定时器 ID, tidlist[j]）以及定时器溢出值①显示出来。
- 创建并配备定时器之后，在循环中反复调用 pause(), 以等待定时器到期⑥。

程序清单 23-6 中函数可将程序 23-5 的命令行参数转化为相应的 `itimerspec` 结构。函数可识别的字符串参数格式在源码文件开始的注释中做了说明（并在下面的 shell 会话中做了演示）。

程序清单 23-6：将“时间+间隔”的字符串转换为 `itimerspec` 的值

```
----- timers/itimerspec_from_str.c
#define POSIX_C_SOURCE 199309
#include <string.h>
#include <stdlib.h>
#include "itimerspec_from_str.h"      /* Declares function defined here */

/* Convert a string of the following form to an itimerspec structure:
   "value.sec[/value.nanosec][:interval.sec[/interval.nanosec]]".
   Optional components that are omitted cause 0 to be assigned to the
   corresponding structure fields. */

void
itimerspecFromStr(char *str, struct itimerspec *tsp)
{
    char *cptr, *sptr;

    cptr = strchr(str, ':');
    if (cptr != NULL)
        *cptr = '\0';

    sptr = strchr(str, '/');
    if (sptr != NULL)
        *sptr = '\0';

    tsp->it_value.tv_sec = atoi(str);
    tsp->it_value.tv_nsec = (sptr != NULL) ? atoi(sptr + 1) : 0;
    if (cptr == NULL) {
        tsp->it_interval.tv_sec = 0;
        tsp->it_interval.tv_nsec = 0;
    } else {
        sptr = strchr(cptr + 1, '/');
        if (sptr != NULL)
            *sptr = '\0';
        tsp->it_interval.tv_sec = atoi(cptr + 1);
        tsp->it_interval.tv_nsec = (sptr != NULL) ? atoi(sptr + 1) : 0;
    }
}
----- timers/itimerspec_from_str.c
```

如下 shell 会话演示了对程序清单 23-5 中程序的调用，创建了一个初始到期值为 2 秒，间隔时间为 5 秒的定时器。

```
$ ./ptmr_sigev_signal 2:5
Timer ID: 134524952 (2:5)
[15:54:56] Got signal 64          SIGRTMAX is signal 64 on this system
    *sival_ptr      = 134524952    sival_ptr points to the variable tid
    timer_getoverrun() = 0
[15:55:01] Got signal 64
    *sival_ptr      = 134524952
    timer_getoverrun() = 0
Type Control-Z to suspend the process
[1]+  Stopped          ./ptmr_sigev_signal 2:5
```

挂起程序后，暂停几秒钟，在恢复程序运行之前会有多个定时器到期。

```
$ fg
./ptmr_sigev_signal 2:5
[15:55:34] Got signal 64
    *sival_ptr          = 134524952
    timer_getoverrun() = 5
Type Control-C to kill the program
```

程序输出的最后一行表明发生了 5 次定时器溢出，亦即在捕获上一信号之后定时器到期了 6 次。

23.6.6 定时器溢出

假设已经选择通过信号（即 `sigev_notify` 为 `SIGEV_SIGNAL`）传递的方式来接收定时器到期通知。进一步假设，在捕获或接收相关信号之前，定时器到期多次。这可能是因为在进程再次获得调度前的延时所致。另外，不论是直接调用 `sigprocmask()`，还是在信号处理器函数里暗中处理，也都有可能堵塞相关信号的发送。如何知道发生了这些定时器溢出呢？

也许会认为使用实时信号有助于解决这个问题，因为可以对实时信号的多个实例进行排队。不过，由于对排队实时信号有数量上的限制，结果证明这种方法也无法奏效。所以 POSIX.1b 委员会选用了另一种方法：一旦选择通过信号来接收定时器通知，那么即便使用了实时信号，也绝不会对该信号的多个实例进行排队。相反，在接收信号后（无论是通过信号处理器函数还是调用 `sigwaitinfo()`），可以获取定时器溢出计数，即在信号生成与接收之间发生的定时器到期额外次数。如果上次收到信号后定时器发生了 3 次到期，那么溢出计数是 2。

接收到定时器信号之后，有两种方法可以获取定时器溢出值。

- 调用 `timer_getoverrun()`，稍后将会讨论。这是由 SUSv3 指定去获取溢出计数的方法。
- 使用随信号一同返回的结构 `siginfo_t` 中的 `si_overrun` 字段值。这种方法可以避免 `timer_getoverrun()` 的系统调用开销，但同时也是一种 Linux 扩展方法，无法移植。

每次收到定时器信号后，都会重置定时器溢出计数。若自处理或接收定时器信号之后，定时器仅到期一次，则溢出计数为 0（即无溢出）。

```
#define _POSIX_C_SOURCE 199309
#include <time.h>

int timer_getoverrun(timer_t timerid);

Returns timer overrun count on success, or -1 on error
```

函数 `timer_getoverrun()` 返回由参数 `timerid` 指定定时器的溢出值。

根据 SUSv3 规定（表 21-1），函数 `timer_getoverrun()` 是异步信号安全的函数之一，故而在信号处理器函数内部调用也是安全的。

23.6.7 通过线程来通知

`SIGEV_THREAD` 标志允许程序从一个独立的线程中调用函数来获取定时器到期通知。要理解这一标志的含义，需要具备第 29 章和第 30 章中关于 POSIX 线程的知识。如果不了解 POSIX 线程，那么在查看本节示例程序前，可能需要预先阅读一下这些章节。

程序清单 23-7 演示了 `SIGEV_THREAD` 的使用。该程序的命令行参数与程序清单 23-5 相

同。所执行的步骤如下。

- 针对每个命令行参数，程序都创建⑥并配备⑦一个使用了 SIGEV_THREAD 通知机制③的 POSIX 定时器。
- 每当定时器到期时，会在一条独立线程中调用由 sev.sigev_notify_function 指定的函数。调用函数时，使用由 sev.sigev_value.sival_ptr 指定的值作为参数。程序中会将定时器 ID (tidlist[j]) 的地址赋给该字段⑤，以便在调用通知函数时可以获得定时器 ID。
- 创建和配备所有定时器之后，主程序进入循环并等待定时器到期③。每次循环，程序都会调用 pthread_cond_wait()，等待处理定时器通知的线程就条件变量 (cond) 发出信号。
- 每次定时器到期都会调用函数 threadFunc()①。在打印消息后，增加全局变量 expireCnt 的值。考虑到定时器可能溢出，会将 timer_getoverrun() 的返回值也加入 expireCnt 变量中。(23.6.6 节解释了定时器溢出与 SIGEV_SIGNAL 通知机制之间的关系。定时器溢出还可以与 SIGEV_THREAD 机制协作使用，因为在调用通知函数前，定时器可能会多次到期。) 通知函数就条件变量 (cond) 发出信号，告知主程序定时器到期。

下面的 shell 会话日志展示了对程序清单 23-7 中程序的调用。在本例中，程序创建了两个定时器：一个定时器首次到期时间为 5 秒，并设置了 5 秒的时间间隔；另一个初次到期时间为 10 秒，并设置了 10 秒的时间间隔。

```
$ ./ptmr_sigev_thread 5:5 10:10
Timer ID: 134525024 (5:5)
Timer ID: 134525080 (10:10)
[13:06:22] Thread notify
        timer ID=134525024
        timer_getoverrun()=0
main(): count = 1
[13:06:27] Thread notify
        timer ID=134525080
        timer_getoverrun()=0
main(): count = 2
[13:06:27] Thread notify
        timer ID=134525024
        timer_getoverrun()=0
main(): count = 3
Type Control-Z to suspend the program
[1]+  Stopped          ./ptmr_sigev_thread 5:5 10:10
$ fg                               Resume execution
./ptmr_sigev_thread 5:5 10:10
[13:06:45] Thread notify
        timer ID=134525024
        timer_getoverrun()=2                There were timer overruns
main(): count = 6
[13:06:45] Thread notify
        timer ID=134525080
        timer_getoverrun()=0
main(): count = 7
Type Control-C to kill the program
```

程序清单 23-7：使用线程函数发送 POSIX 定时器通知

```
-----timers/ptmr_sigev_thread.c
#include <signal.h>
#include <time.h>
#include <pthread.h>
#include "curr_time.h"                /* Declaration of currTime() */
```

```

#include "tspi_hdr.h"
#include "itimerspec_from_str.h" /* Declares itimerspecFromStr() */

static pthread_mutex_t mtx = PTHREAD_MUTEX_INITIALIZER;
static pthread_cond_t cond = PTHREAD_COND_INITIALIZER;

static int expireCnt = 0; /* Number of expirations of all timers */

static void /* Thread notification function */
① threadFunc(union sigval sv)
{
    timer_t *tidptr;
    int s;

    tidptr = sv.sival_ptr;

    printf("[%s] Thread notify\n", currTime("%T"));
    printf(" timer ID=%ld\n", (long) *tidptr);
    printf(" timer_getoverrun()=%d\n", timer_getoverrun(*tidptr));

    /* Increment counter variable shared with main thread and signal
       condition variable to notify main thread of the change. */

    s = pthread_mutex_lock(&mtx);
    if (s != 0)
        errExitEN(s, "pthread_mutex_lock");

    expireCnt += 1 + timer_getoverrun(*tidptr);

    s = pthread_mutex_unlock(&mtx);
    if (s != 0)
        errExitEN(s, "pthread_mutex_unlock");

    ② s = pthread_cond_signal(&cond);
    if (s != 0)
        errExitEN(s, "pthread_cond_signal");
}

int
main(int argc, char *argv[])
{
    struct sigevent sev;
    struct itimerspec ts;
    timer_t *tidlist;
    int s, j;
    if (argc < 2)
        usageErr("%s secs[/nsecs][:int-secs[/int-nsecs]]...\n", argv[0]);

    tidlist = calloc(argc - 1, sizeof(timer_t));
    if (tidlist == NULL)
        errExit("malloc");

    ③ sev.sigev_notify = SIGEV_THREAD; /* Notify via thread */
    ④ sev.sigev_notify_function = threadFunc; /* Thread start function */
    sev.sigev_notify_attributes = NULL;
        /* Could be pointer to pthread_attr_t structure */

    /* Create and start one timer for each command-line argument */

```

```

for (j = 0; j < argc - 1; j++) {
    itimerspecFromStr(argv[j + 1], &ts);
⑤     sev.sigev_value.sival_ptr = &tidlist[j];
        /* Passed as argument to threadFunc() */

⑥     if (timer_create(CLOCK_REALTIME, &sev, &tidlist[j]) == -1)
        errExit("timer_create");
        printf("Timer ID: %ld (%s)\n", (long) tidlist[j], argv[j + 1]);

⑦     if (timer_settime(tidlist[j], 0, &ts, NULL) == -1)
        errExit("timer_settime");
}

/* The main thread waits on a condition variable that is signaled
   on each invocation of the thread notification function. We
   print a message so that the user can see that this occurred. */

s = pthread_mutex_lock(&mtx);
if (s != 0)
    errExitEN(s, "pthread_mutex_lock");

⑧     for (;;) {
        s = pthread_cond_wait(&cond, &mtx);
        if (s != 0)
            errExitEN(s, "pthread_cond_wait");
        printf("main(): expireCnt = %d\n", expireCnt);
    }
}

```

timers/ptmr_sigev_thread.c

23.7 利用文件描述符进行通知的定时器：timerfd API

始于版本 2.6.25，Linux 内核提供了另一种创建定时器的 API。Linux 特有的 timerfd API，可从文件描述符中读取其所创建定时器的到期通知。因为可以使用 select()、poll() 和 epoll()（将在第 63 章进行讨论）将这种文件描述符会同其他描述符一同进行监控，所以非常实用。（至于说本章讨论的其他定时器 API，想要把一个或多个定时器与一组文件描述符放在一起同时监测，可不是件容易的事。）

这组 API 中的 3 个新系统调用，其操作与 23.6 节所述的 timer_create()、timer_settime() 和 timer_gettime() 相类似。

新加入的第 1 个系统调用是 timerfd_create()，它会创建一个新的定时器对象，并返回一个指代该对象的文件描述符。

```

#include <sys/timerfd.h>

int timerfd_create(int clockid, int flags);

```

Returns file descriptor on success, or -1 on error

参数 clockid 的值可以设置为 CLOCK_REALTIME 或 CLOCK_MONOTONIC（参考表 23-1）。

timerfd_create() 的最初实现将参数 flags 预留供未来使用，必须设置为 0。不过，Linux 内核从 2.6.27 版本开始支持下面两种 flags 标志。

TFD_CLOEXEC

为新的文件描述符设置运行时关闭标志 (FD_CLOEXEC)。与 4.3.1 节介绍的 open() 标志 O_CLOEXEC 适用于相同情况。

TFD_NONBLOCK

为底层的打开文件描述符设置 O_NONBLOCK 标志，随后的读操作将是非阻塞式的。这样设置省却了对 fcntl() 的额外调用，却能达到相同效果。

timerfd_create() 创建的定时器使用完毕后，应调用 close() 关闭相应的文件描述符，以便于内核能够释放与定时器相关的资源。

系统调用 timerfd_settime() 可以配备 (启动) 或解除 (停止) 由文件描述符 fd 所指代的定时器。

```
#include <sys/timerfd.h>

int timerfd_settime(int fd, int flags, const struct itimerspec *new_value,
                   struct itimerspec *old_value);

Returns 0 on success, or -1 on error
```

参数 new_value 为定时器指定新设置。参数 old_value 用来返回定时器的前一设置 (细节请参考随后对 timerfd_gettime() 的说明)。如果不关心定时器的前一设置，可将 old_value 置为 NULL。两个参数均指向 itimerspec 结构，用法与 timer_settime() (参考 23.6.2 节) 相同。

参数 flags 与 timer_settime() 中的对应参数类似。可以是 0，此时将 new_value.it_value 的视值为相对于调用 timerfd_settime() 时间点的相对时间，也可以设为 TFD_TIMER_ABSTIME，将其视为一个绝对时间 (从时钟的 0 点开始测量)。

系统调用 timerfd_gettime() 返回文件描述符 fd 所标识定时器的间隔及剩余时间。

```
#include <sys/timerfd.h>

int timerfd_gettime(int fd, struct itimerspec *curr_value);

Returns 0 on success, or -1 on error
```

同 timer_gettime() 一样，间隔以及距离下次到期的时间均返回 curr_value 指向的结构 itimerspec 中。即使是以 TFD_TIMER_ABSTIME 标志创建的绝对时间定时器，curr_value.it_value 字段中返回值的意义也会保持不变。如果返回的结构 curr_value.it_value 中所有字段值均为 0，那么该定时器已经被解除。如果返回的结构 curr_value.it_interval 中两字段值均为 0，那么定时器只会到期一次，到期时间在 curr_value.it_value 中给出。

timerfd 与 fork() 及 exec() 之间的交互

调用 fork() 期间，子进程会继承 timerfd_create() 所创建文件描述符的拷贝。这些描述符与父进程的对应描述符均指代相同的定时器对象，任一进程都可读取定时器的到期信息。

timerfd_create() 创建的文件描述符能跨越 exec() 得以保存 (除非将描述符置为运行时关闭，如 27.4 节所述)，已配备的定时器在 exec() 之后会继续生成到期通知。

从 timerfd 文件描述符读取

一旦以 timerfd_settime() 启动了定时器，就可以从相应文件描述符中调用 read() 来读取定时

器的到期信息。出于这一目的，传给 read()的缓冲区必须足以容纳一个无符号 8 字节整型 (uint64_t) 数。

在上次使用 timerfd_settime()修改设置以后，或是最后一次执行 read()后，如果发生了一起或多起定时器到期事件，那么 read()会立即返回，且返回的缓冲区中包含了已发生的到期次数。如果并无定时器到期，read()会一直阻塞直至产生下一个到期。也可以执行 fcntl()的 F_SETFL 操作 (5.3 节) 为文件描述符设置 O_NONBLOCK 标志，这时的读动作是非阻塞式的，且如果没有定时器到期，则返回错误，并将 errno 值置为 EAGAIN。

如前所述，可以利用 select()、poll()和 epoll()对 timerfd 文件描述符进行监控。如果定时器到期，会将对应的文件描述符标记为可读。

示例程序

程序清单 23-8 演示了 timerfd API 的使用。该程序从命令行取得两个参数。第 1 个参数为必填项，用以标识定时器的初始和间隔时间。(程序清单 23-6 中的函数 itimerspecFromStr()可以用来解析这一参数。)第 2 个参数是可选项，表示程序退出之前应等待的定时器过期最大次数，其默认值为 1。

程序调用 timerfd_create()来创建一个定时器，并通过 timerfd_settime()将其启动。接着进入循环，从文件描述符中读取定时器到期通知，直至达到指定的定时器到期次数。每次 read()之后，程序都会显示定时器启动以来的逝去时间、读取到的到期次数以及至今为止的到期总数。

下面的 shell 会话日志中，通过命令行参数创建了一个初始时间为 1 秒，间隔为 1 秒，最大到期次数为 100 次的定时器。

```
$ ./demo_timerfd 1:1 100
1.000: expirations read: 1; total=1
2.000: expirations read: 1; total=2
3.000: expirations read: 1; total=3
Type Control-Z to suspend program in background for a few seconds
[1]+  Stopped                  ./demo_timerfd 1:1 100
$ fg                            Resume program in foreground
./demo_timerfd 1:1 100
14.205: expirations read: 11; total=14  Multiple expirations since last read()
15.000: expirations read: 1; total=15
16.000: expirations read: 1; total=16
Type Control-C to terminate the program
```

从以上结果可以看出，程序在后台暂停时定时器出现了多次到期。在程序恢复运行之后，第 1 次 read()调用就返回了所有这些到期。

程序清单 23-8: 使用 timerfd API

```
timers/demo_timerfd.c
#include <sys/timerfd.h>
#include <time.h>
#include <stdint.h>                /* Definition of uint64_t */
#include "itimerspec_from_str.h"  /* Declares itimerspecFromStr() */
#include "tlpi_hdr.h"

int
main(int argc, char *argv[])
{
    struct itimerspec ts;
```

```

struct timespec start, now;
int maxExp, fd, secs, nanosecs;
uint64_t numExp, totalExp;
ssize_t s;

if (argc < 2 || strcmp(argv[1], "--help") == 0)
    usageErr("%s secs[/nsecs][:int-secs[/int-nsecs]] [max-exp]\n", argv[0]);

itimerspecFromStr(argv[1], &ts);
maxExp = (argc > 2) ? getInt(argv[2], GN_GT_0, "max-exp") : 1;
fd = timerfd_create(CLOCK_REALTIME, 0);
if (fd == -1)
    errExit("timerfd_create");

if (timerfd_settime(fd, 0, &ts, NULL) == -1)
    errExit("timerfd_settime");

if (clock_gettime(CLOCK_MONOTONIC, &start) == -1)
    errExit("clock_gettime");

for (totalExp = 0; totalExp < maxExp;) {

    /* Read number of expirations on the timer, and then display
       time elapsed since timer was started, followed by number
       of expirations read and total expirations so far. */

    s = read(fd, &numExp, sizeof(uint64_t));
    if (s != sizeof(uint64_t))
        errExit("read");

    totalExp += numExp;

    if (clock_gettime(CLOCK_MONOTONIC, &now) == -1)
        errExit("clock_gettime");

    secs = now.tv_sec - start.tv_sec;
    nanosecs = now.tv_nsec - start.tv_nsec;
    if (nanosecs < 0) {
        secs--;
        nanosecs += 1000000000;
    }

    printf("%d.%03d: expirations read: %llu; total=%llu\n",
           secs, (nanosecs + 500000) / 1000000,
           (unsigned long long) numExp, (unsigned long long) totalExp);
}

exit(EXIT_SUCCESS);
}

```

timers/demo_timerfd.c

23.8 总结

进程可以使用 `setitimer()` 或 `alarm()` 来设置定时器，以便于在经历指定的一段实际（或进程）时间后收到信号通知。定时器的用途之一是为系统调用的阻塞设定时间上限。

应用程序如需暂停执行一段特定间隔的实际时间，可以使用各种合适的休眠函数。

Linux 2.6 所实现的 POSIX.1b 扩展为高精度时钟和定时器定义了一套 API。POSIX.1b 定时器比传统 (settimer()) UNIX 定时器更具优势, 可以: 创建多个定时器; 选择定时器到期时的通知信号; 获取定时器溢出计数, 以便判断自上次到期通知后定时器是否又发生了多次到期; 选择通过执行线程函数而非递送信号来获取定时器通知。

Linux 特有的 timerfd API 提供了一组创建定时器的接口, 与 POSIX 定时器 API 相类似, 但允许从文件描述符中读取定时器通知。还可使用 select()、poll() 和 epoll() 来监控这些描述符。

更多信息

在每个函数的原理 (rationael) 部分, SUSv3 就本章所述 (标准) 定时器和休眠接口一一指出其要点所在。[Callmeister, 1995] 则探讨了 POSIX.1b 时钟和定时器。

23.9 练习

- 23-1.** 尽管 Linux 将 alarm() 实现为系统调用, 但这当属蛇足。请用 setitimer() 实现 alarm()。
- 23-2.** 试着将程序清单 23-3 (t_nanosleep.c) 程序置于后台运行, 并设置 60 秒的休眠间隔, 同时使用如下命令发送尽可能多的 SIGINT 信号给后台进程:

```
$ while true; do kill -INT pid; done
```

应该能注意到程序休眠时间要长于预期。将 nanosleep() 用 clock_gettime() (使用 CLOCK_REALTIME 时钟) 和设置 TIMER_ABSTIME 标志的 clock_nanosleep() 来替换。(此练习需要 Linux 2.6 版本。) 反复测试修改后的程序, 并解释新老程序间的差别。

- 23-3.** 编写一个程序验证: 如果调用 timer_create() 时将参数 evp 置为 NULL, 那么这就等同于将 evp 设为指向 sigevent 结构的指针, 并将该结构中的 sigev_notify 置为 SIGEV_SIGNAL, 将 sigev_signo 置为 SIGALRM, 将 si_value.sival_int 置为定时器 ID。
- 23-4.** 修改程序清单 23-5 (ptmr_sigev_signal.c) 中程序, 并用 sigwaitinfo() 替换信号处理器函数。

第 24 章

进程的创建

本章以及随后的 3 章将探讨进程的创建和终止，以及进程执行新程序的过程。本章主要讨论进程的创建，不过，在切入正题之前，将首先概括一下这 4 章所涵盖的主要系统调用。

24.1 fork()、exit()、wait()以及 execve()的简介

本章以及随后几章的议题会集中在 fork()、exit()、wait()以及 execve()这几个系统调用上。上述每种系统调用都各有变体，后续会一一论及。此处将首先对这 4 个系统调用及其典型用法简单加以介绍。

- 系统调用 fork()允许一进程（父进程）创建一新进程（子进程）。具体做法是，新的子进程几近于对父进程的翻版：子进程获得父进程的栈、数据段、堆和执行文本段（6.3 节）的拷贝。可将此视为把父进程一分为二，术语 fork 也由此得名。
- 库函数 exit(status)终止一进程，将进程占用的所有资源（内存、文件描述符等）归还内核，交其进行再次分配。参数 status 为一整型变量，表示进程的退出状态。父进程可使用系统调用 wait()来获取该状态。

库函数 exit()位于系统调用 _exit()之上。第 25 章将解释二者之间的差异。这里只是强调，在调用 fork()之后，父、子进程中一般只有一个会通过调用 exit()退出，而另一进程则应使用 _exit()终止。

- 系统调用 wait(&status)的目的有二：其一，如果子进程尚未调用 exit()终止，那么 wait()会挂起父进程直至子进程终止；其二，子进程的终止状态通过 wait()的 status 参数返回。
- 系统调用 execve(pathname, argv, envp)加载一个新程序（路径名为 pathname，参数列表为 argv，环境变量列表为 envp）到当前进程的内存。这将丢弃现存的程序文本段，并为新程序重新创建栈、数据段以及堆。通常将这一动作称为执行（execing）一个新程序。稍后会介绍构建于 execve()之上的多个库函数，每种都为编程接口提供了实用的变体。在彼此差异无关宏旨的场合，循例会将此类函数统称为 exec()，尽管实际上并没有以之命名的系统调用或者库函数。

其他一些操作系统则将 `fork()` 和 `exec()` 的功能合二为一，形成单一的 `spawn` 操作——创建一个新进程并执行指定程序。比较而言，UNIX 的方案通常更为简单和优雅。两步走的策略使得 API 更为简单（系统调用 `fork()` 无需参数），程序也得以在这两步之间执行一些其他操作，因而更具弹性。另外，只执行 `fork()` 而不执行 `exec()` 的场景也颇为常见。

SUSv3 所详细规定的 `posix_spawn()` 函数，就将 `fork()` 和 `exec()` 的功能结合起来，但规范并未对实现此函数做强制要求。Linux 的 `glibc` 函数库实现了该函数以及 SUSv3 中的其他几个相关 API。将 `posix_spawn()` 纳入 SUSv3，意在为缺乏交换（`swap`）设施或内存管理单元（`memory-management units`）的硬件架构（嵌入式系统大多如此）编写具备可移植性的应用程序。在此类架构上实现传统意义的 `fork()`，即便存在可能性，难度也很大。

图 24-1 对 `fork()`、`exit()`、`wait()` 以及 `execve()` 之间的相互协同作了总结。（此图勾勒了 shell 执行一条命令所历经的步骤：shell 读取命令，进行各种处理，随之创建子进程以执行该命令，如此循环不已。）

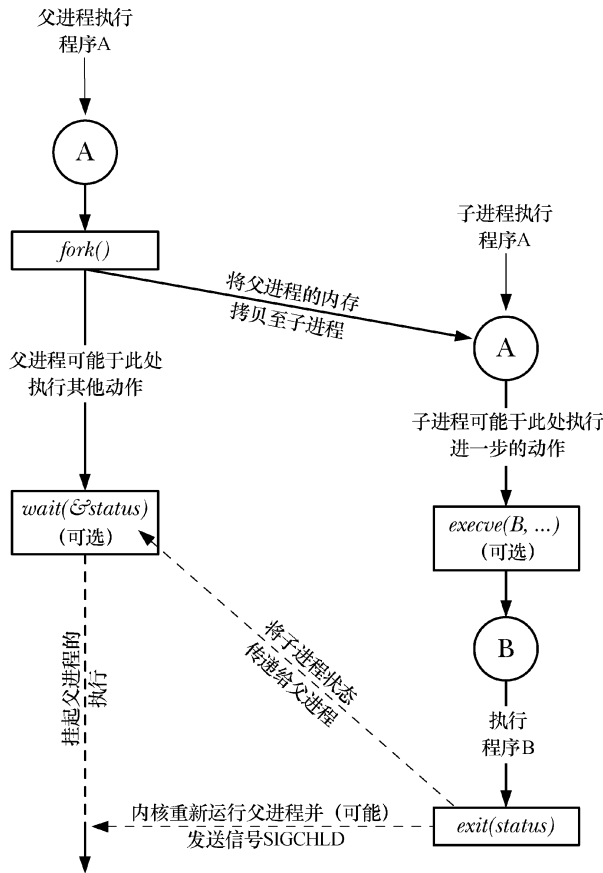


图 24-1：概述函数 `fork()`、`exit()`、`wait()` 和 `execve()` 的协同使用

图中对 `execve()` 的调用并非必须。有时，让子进程继续执行与父进程相同的程序反而会有妙用。最终，两种情况殊途同归：总是要通过调用 `exit()`（或接收一个信号）来终止子进程，而父进程可调用 `wait()` 来获取其终止状态。

同样，对 `wait()` 的调用也属于可选项。父进程可以对子进程不闻不问，继续我行我素。不

过，由后续内容可知，对 `wait()` 的使用通常也是不可或缺的，每每在 `SIGCHLD` 信号的处理程序中使用。当子进程终止时，内核会为其父进程产生此类信号（默认的处理是忽略 `SIGCHLD` 信号，下图将此标记为可选，原因正在于此）。

24.2 创建新进程：fork()

在诸多应用中，创建多个进程是任务分解时行之有效的方法。例如，某一网络服务器进程可在侦听客户端请求的同时，为处理每一请求而创建一新的子进程，与此同时，服务器进程会继续侦听更多的客户端连接请求。以此类手法分解任务，通常会简化应用程序的设计，同时提高了系统的并发性。（即，可同时处理更多的任务或请求。）

系统调用 `fork()` 创建一新进程（child），几近于对调用进程（parent）的翻版。

```
#include <unistd.h>

pid_t fork(void);

    In parent: returns process ID of child on success, or -1 on error;
    in successfully created child: always returns 0
```

理解 `fork()` 的诀窍是，要意识到，完成对其调用后将存在两个进程，且每个进程都会从 `fork()` 的返回处继续执行。

这两个进程将执行相同的程序文本段，但却各自拥有不同的栈段、数据段以及堆段拷贝。子进程的栈、数据以及栈段开始时是对父进程内存相应各部分的完全复制。执行 `fork()` 之后，每个进程均可修改各自的栈数据、以及堆段中的变量，而并不影响另一进程。

程序代码则可通过 `fork()` 的返回值来区分父、子进程。在父进程中，`fork()` 将返回新创建子进程的进程 ID。鉴于父进程可能需要创建，进而追踪多个子进程（通过 `wait()` 或类似方法），这种安排还是很实用的。而 `fork()` 在子进程中则返回 0。如有必要，子进程可调用 `getpid()` 以获取自身的进程 ID，调用 `getppid()` 以获取父进程 ID。

当无法创建子进程时，`fork()` 将返回 -1。失败的原因可能在于，进程数量要么超出了系统针对此真实用户（real user ID）在进程数量上所施加的限制（`RLIMIT_NPROC`，36.3 节将对此加以描述），要么是触及允许该系统创建的最大进程数这一系统级上限。

调用 `fork()` 时，有时会采用如下习惯用语：

```
pid_t childPid;          /* Used in parent after successful fork()
                          to record PID of child */
switch (childPid = fork()) {
case -1:                 /* fork() failed */
    /* Handle error */

case 0:                 /* Child of successful fork() comes here */
    /* Perform actions specific to child */

default:                /* Parent comes here after successful fork() */
    /* Perform actions specific to parent */
}
```

调用 `fork()` 之后，系统将率先“垂青”于哪个进程（即调度其使用 CPU），是无法确定的，意识到这一点极为重要。在设计拙劣的程序中，这种不确定性可能会导致所谓“竞争条件（race condition）”的错误，24.2 节会对此做进一步说明。

程序清单 24-1 展示了 `fork()` 的用法。该程序创建一子进程，并对继承自 `fork()` 的全局及自动变量拷贝进行修改。

使用 `sleep()`（存在于由父进程所执行的代码中），意在允许子进程先于父进程获得系统调度并使用 CPU，以便在父进程继续运行之前完成自身任务并退出。要想确保这一结果，`sleep()` 的这种用法并非万无一失，24.5 节中的方法更胜一筹。

运行程序清单 24-1 中程序，其输出如下。

```
$ ./t_fork
PID=28557 (child)  idata=333  istack=666
PID=28556 (parent) idata=111  istack=222
```

以上输出表明，子进程在 `fork()` 时拥有了自己的栈和数据段拷贝，且其对这些段中变量的修改将不会影响父进程。

程序清单 24-1：调用 `fork()`

```
----- procexec/t_fork.c
#include "tlpi_hdr.h"

static int idata = 111;          /* Allocated in data segment */

int
main(int argc, char *argv[])
{
    int istack = 222;          /* Allocated in stack segment */
    pid_t childPid;

    switch (childPid = fork()) {
    case -1:
        errExit("fork");

    case 0:
        idata *= 3;
        istack *= 3;
        break;

    default:
        sleep(3);              /* Give child a chance to execute */
        break;
    }

    /* Both parent and child come here */

    printf("PID=%ld %s idata=%d istack=%d\n", (long) getpid(),
           (childPid == 0) ? "(child) " : "(parent)", idata, istack);

    exit(EXIT_SUCCESS);
}
----- procexec/t_fork.c
```

24.2.1 父、子进程间的文件共享

执行 `fork()` 时，子进程会获得父进程所有文件描述符的副本。这些副本的创建方式类似于 `dup()`，这也意味着父、子进程中对应的描述符均指向相同的打开文件句柄（即 `open file description`，详见 5.4 节译注）。正如 5.4 节所述，打开文件句柄包含有当前文件偏移量（由 `read()`、`write()` 和 `lseek()` 修改）以及文件状态标志（由 `open()` 设置，通过 `fcntl()` 的 `F_SETFL` 操作改变）。

一个打开文件的这些属性因之而在父子进程间实现了共享。举例来说，如果子进程更新了文件偏移量，那么这种改变也会影响到父进程中相应的描述符。

程序清单 24-2 所展示的正是这样一个事实；fork()之后，这些属性将在父子进程之间共享。该程序使用 mkstemp()打开一个临时文件，接着调用 fork()以创建子进程。子进程改变文件偏移量以及文件状态标志，最后退出。父进程随即获取文件偏移量和标志，以验证其可以观察到由子进程所造成的变化。此程序运行结果如下：

```
$ ./fork_file_sharing
File offset before fork(): 0
O_APPEND flag before fork() is: off
Child has exited
File offset in parent: 1000
O_APPEND flag in parent is: on
```

关于程序清单 24-2 为何要将 lseek()的返回值强制转换为 long long，参见 5.10 节。

程序清单 24-2：在父子进程间共享文件偏移量和打开文件状态标志

```
----- procexec/fork_file_sharing.c
#include <sys/stat.h>
#include <fcntl.h>
#include <sys/wait.h>
#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    int fd, flags;
    char template[] = "/tmp/testXXXXXX";

    setbuf(stdout, NULL);                /* Disable buffering of stdout */

    fd = mkstemp(template);
    if (fd == -1)
        errExit("mkstemp");

    printf("File offset before fork(): %lld\n",
           (long long) lseek(fd, 0, SEEK_CUR));

    flags = fcntl(fd, F_GETFL);
    if (flags == -1)
        errExit("fcntl - F_GETFL");
    printf("O_APPEND flag before fork() is: %s\n",
           (flags & O_APPEND) ? "on" : "off");
    switch (fork()) {
    case -1:
        errExit("fork");

    case 0:    /* Child: change file offset and status flags */
        if (lseek(fd, 1000, SEEK_SET) == -1)
            errExit("lseek");

        flags = fcntl(fd, F_GETFL);    /* Fetch current flags */
        if (flags == -1)
            errExit("fcntl - F_GETFL");
        flags |= O_APPEND;             /* Turn O_APPEND on */
        if (fcntl(fd, F_SETFL, flags) == -1)
```

```

        errExit("fcntl - F_SETFL");
        _exit(EXIT_SUCCESS);

default: /* Parent: can see file changes made by child */
    if (wait(NULL) == -1)
        errExit("wait"); /* Wait for child exit */
    printf("Child has exited\n");

    printf("File offset in parent: %lld\n",
           (long long) lseek(fd, 0, SEEK_CUR));

    flags = fcntl(fd, F_GETFL);
    if (flags == -1)
        errExit("fcntl - F_GETFL");
    printf("O_APPEND flag in parent is: %s\n",
           (flags & O_APPEND) ? "on" : "off");
    exit(EXIT_SUCCESS);
}
}
}

```

procexec/fork_file_sharing.c

父子进程间共享打开文件属性的妙用屡见不鲜。例如，假设父子进程同时写入一文件，共享文件偏移量会确保二者不会覆盖彼此的输出内容。不过，这并不能阻止父子进程的输出随意混杂在一起。要想规避这一现象，需要进行进程间同步。比如，父进程可以使用系统调用 `wait()` 来暂停运行并等待子进程退出。`shell` 就是这么做的：只有当执行命令的子进程退出后，`shell` 才会打印出提示符（除非用户在命令行最后加上 `&` 符以显式在后台运行命令）。

如果不需要这种对文件描述符的共享方式，那么在设计应用程序时，应于 `fork()` 调用后注意两点：其一，令父、子进程使用不同的文件描述符；其二，各自立即关闭不再使用的描述符（亦即那些经由其他进程使用的描述符）。如果进程之一执行了 `exec()`，那么 27.4 节所描述的运行时关闭功能（`close-on-exec`）也会很有用处。图 24-2 展示了这些步骤。

24.2.2 `fork()` 的内存语义

从概念上说来，可以将 `fork()` 认作对父进程程序段、数据段、堆段以及栈段创建拷贝的。的确，在一些早期的 UNIX 实现中，此类复制确实是原汁原味：将父进程内存拷贝至交换空间，以此创建新进程映像（`image`），而在父进程保持自身内存的同时，将换出映像置为子进程。不过，真要是简单地将父进程虚拟内存页拷贝到新的子进程，那就太浪费了。原因有很多，其中之一是：`fork()` 之后常常伴随着 `exec()`，这会用新程序替换进程的代码段，并重新初始化其数据段、堆段和栈段。大部分现代 UNIX 实现（包括 Linux）采用两种技术来避免这种浪费。

- 内核（Kernel）将每一进程的代码段标记为只读，从而使进程无法修改自身代码。这样，父、子进程可共享同一代码段。系统调用 `fork()` 在为子进程创建代码段时，其所构建的一系列进程级页表项（`page-table entries`）均指向与父进程相同的物理内存页帧。
- 对于父进程数据段、堆段和栈段中的各页，内核采用写时复制（`copy-on-write`）技术来处理。（[Bach, 1986]和[Bovert & Cersati, 2005]描述了写时复制的实现。）最初，内核做了一些设置，令这些段的页表项指向与父进程相同的物理内存页，并将这些页面自身标记为只读。调用 `fork()` 之后，内核会捕获所有父进程或子进程针对这些页面的修改企图，并为将要修改的（`about-to-be-modified`）页面创建拷贝。系统将新的页面拷贝分配给遭内核捕获的进程，还会对子进程的相应页表项做适当调整。从这一刻起，

父、子进程可以分别修改各自的页拷贝，不再相互影响。图 24-3 展示了写时复制技术。

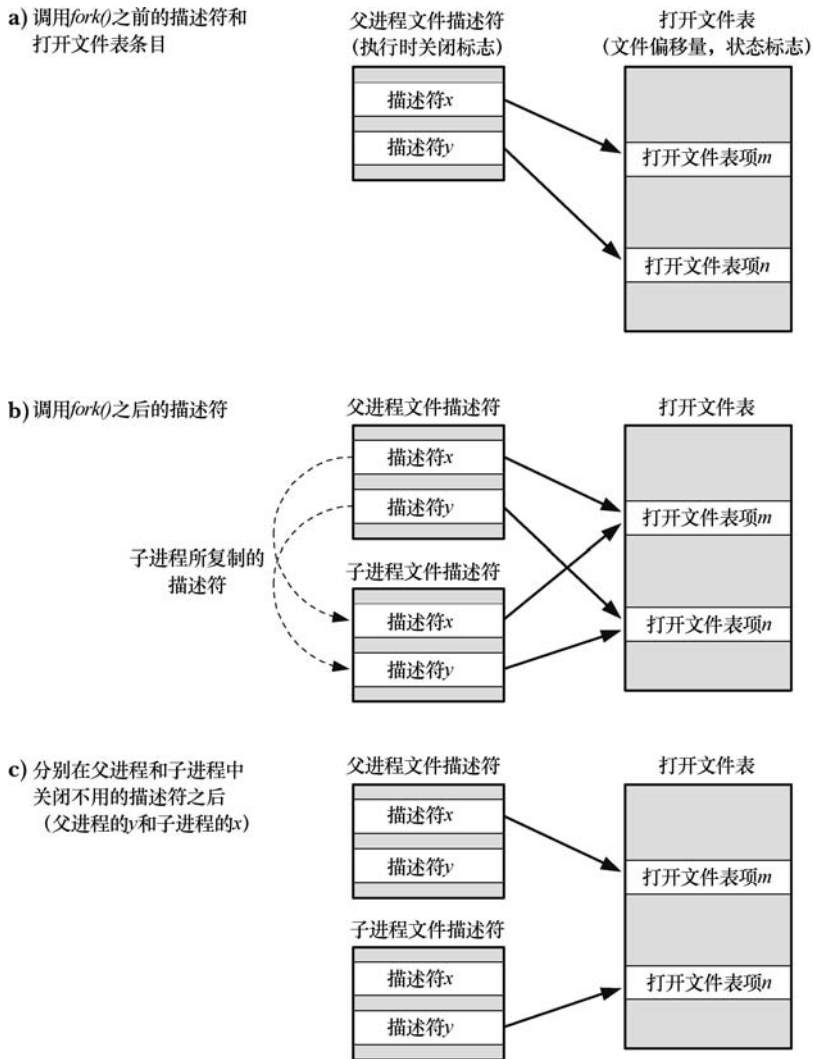


图 24-2: 执行 `fork()` 期间对文件描述符的复制，以及关闭不再使用的描述符

控制进程的内存需求

通过将 `fork()` 与 `wait()` 组合使用，可以控制一个进程的内存需求。进程的内存需求量，亦即进程所使用的虚拟内存页范围，受到多种因素的影响，例如，调用函数，或从函数返回时栈的变化情况，对 `exec()` 的调用，以及因调用 `malloc()` 和 `free()` 而对堆所做的修改——这点对这里的讨论有着特殊意义。

假设以程序清单 24-3 所示方式调用 `fork()` 和 `wait()`，且将对某函数 `func()` 的调用置于括号之中。由执行程序可知，由于所有的变化都发生于子进程，故而对 `func()` 的调用之前开始，父进程的内存使用量将保持不变。这一用法的实用性则归于如下理由。

- 若已知 `func()` 导致内存泄露，或是引发堆内存的过度碎片化，该技术则可以避免这些问题。（要是无法访问 `func()` 的源码，想要处理这些问题也就无从谈起。）

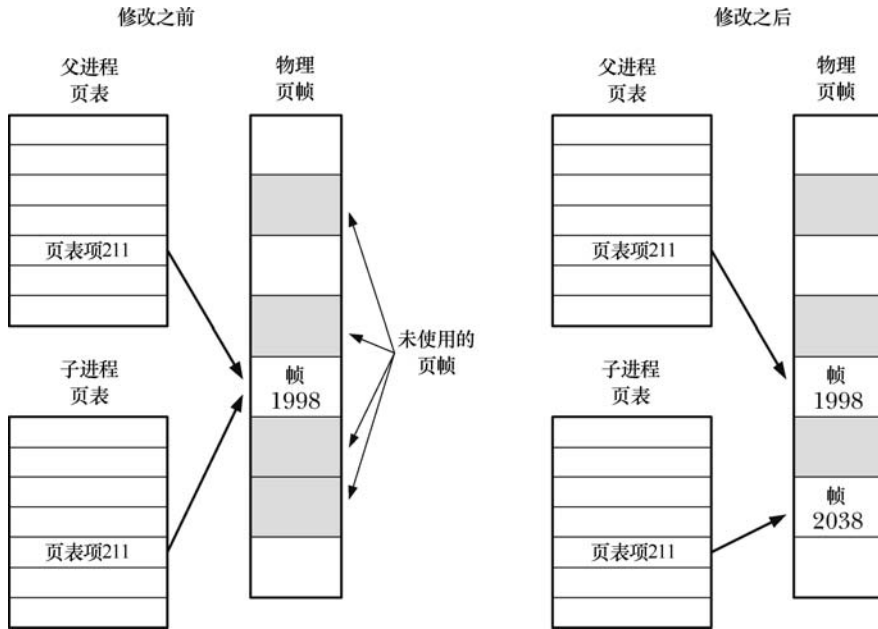


图 24-3: 对一共享写时复制页进行修改前后的页表

- 假设某一算法在做树状分析 (tree analysis) 的同时需要进行内存分配 (例如, 游戏程序需要分析一系列可能的招法以及对方的应手)。本可以调用 `free()` 来释放所有已分配的内存, 不过在某些情况下, 使用此处所描述的技术会更为简单, 返回 (父进程), 且调用者 (父进程) 的内存需求并无改变。

如程序清单 24-3 的实现所示, 必须将 `func()` 的返回结果置于 `exit()` 的 8 位传出值中, 父进程调用 `wait()` 可获得该值。不过, 也可以利用文件、管道或其他一些进程间通信技术, 使 `func()` 返回更大的结果集¹。

程序清单 24-3: 调用函数而不改变进程的内存需求量

from procexec/footprint.c

```

pid_t childPid;
int status;

childPid = fork();
if (childPid == -1)
    errExit("fork");

if (childPid == 0)           /* Child calls func() and */
    exit(func(arg));         /* uses return value as exit status */

/* Parent waits for child to terminate. It can determine the
   result of func() by inspecting 'status'. */

if (wait(&status) == -1)
    errExit("wait");

```

from procexec/footprint.c

¹ 译者注: 针对游戏程序的分析结果而言。

24.3 系统调用 vfork()

在早期的 BSD 实现中, `fork()`会对父进程的数据段、堆和栈施行严格的复制。如前所述, 这是一种浪费, 尤其是在调用 `fork()`后立即执行 `exec()`的情况下。出于这一原因, BSD 的后期版本引入了 `vmfork()`系统调用, 尽管其运作含义稍微有些不同(实则有些怪异), 但效率要远高于 BSD `fork()`。现代 UNIX 采用写时复制技术来实现 `fork()`, 其效率较之于早期的 `fork()`实现要高出许多, 进而将对 `vmfork()`的需求剔除殆尽。虽然如此, Linux (如同许多其他的 UNIX 实现一样) 还是提供了具有 BSD 语义的 `vmfork()`系统调用, 以期为程序提供尽可能快的 `fork` 功能。不过, 鉴于 `vmfork()`的怪异语义可能会导致一些难以察觉的程序缺陷(bug), 除非能给性能带来重大提升(这种情况发生的概率极小), 否则应当尽量避免使用这一调用。

类似于 `fork()`, `vmfork()`可以为调用进程创建一个新的子进程。然而, `vmfork()`是为子进程立即执行 `exec()`的程序而专门设计的。

```
#include <unistd.h>
```

```
pid_t vmfork(void);
```

In parent: returns process ID of child on success, or -1 on error;
in successfully created child: always returns 0

`vmfork()`因为如下两个特性而更具效率, 这也是其与 `fork()`的区别所在。

- 无需为子进程复制虚拟内存页或页表。相反, 子进程共享父进程的内存, 直至其成功执行了 `exec()`或是调用 `_exit()`退出。
- 在子进程调用 `exec()`或 `_exit()`之前, 将暂停执行父进程。

这两点还另有深意: 由于子进程使用父进程的内存, 因此子进程对数据段、堆或栈的任何改变将在父进程恢复执行时为其所见。此外, 如果子进程在 `vmfork()`与后续的 `exec()`或 `_exit()`之间执行了函数返回, 这同样会影响到父进程。这与 6.8 节所描述的例子(试图以 `longjmp()`进入一个已经执行了返回的函数中)相类似。同样相似的还有这一乱局的收场——以典型的段错误(SIGSEGV)而告终。

在不影响父进程的前提下, 子进程能在 `vmfork()`与 `exec()`之间所做的操作屈指可数。其中包括对打开文件描述符进行操作(但不能施之于 `stdio` 文件流)。因为系统是在内核空间为每个进程维护文件描述符表(5.4 节), 且在 `vmfork()`调用期间将复制该表, 所以子进程对文件描述符的操作不会影响到父进程。

SUSv3 指出, 在如下情况下程序行为未定义: a) 修改了除用于存储 `vmfork()`返回值的 `pid_t` 型变量之外的任何数据; b) 从调用 `vmfork()`的函数中返回; c) 在成功地调用 `_exit()`或执行 `exec()`之前, 调用了任何其他函数。

28.2 节在介绍系统调用 `clone()`时将会提及, 由 `fork()`或 `vmfork()`创建的子进程还具有少量其他进程属性的自有拷贝。

`vmfork()`的语义在于执行该调用后, 系统将保证子进程先于父进程获得调度以使用 CPU。24.2 节曾经提及 `fork()`是无法保证这一点的, 父、子进程均有可能率先获得调度。

程序清单 24-4 展示了 `vmfork()`的用法, 将其区分于 `fork()`的两种语义特性显露无遗: 子进程共享父进程的内存, 父进程会一直挂起直至子进程终止或调用 `exec()`。运行该程序, 其输出结果如下:

```

$ ./t_vfork
Child executing           Even though child slept, parent was not scheduled
Parent executing
istack=666

```

由输出的最后一行可知，子进程对变量 `istack` 的修改影响了父进程的对应变量。

程序清单 24-4：使用 `vfork()`

```

procexec/t_vfork.c

#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    int istack = 222;

    switch (vfork()) {
    case -1:
        errExit("vfork");

    case 0:
        /* Child executes first, in parent's memory space */
        sleep(3);
        /* Even if we sleep for a while,
           parent still is not scheduled */
        write(STDOUT_FILENO, "Child executing\n", 16);
        istack *= 3;
        /* This change will be seen by parent */
        _exit(EXIT_SUCCESS);

    default:
        /* Parent is blocked until child exits */
        write(STDOUT_FILENO, "Parent executing\n", 17);
        printf("istack=%d\n", istack);
        exit(EXIT_SUCCESS);
    }
}

```

procexec/t_vfork.c

除非速度绝对重要的场合，新程序应当舍 `vfork()` 而取 `fork()`。原因在于，当使用写时复制语义实现 `fork()`（大部分现代 UNIX 实现皆是如此）时，在速度几近于 `vfork()` 的同时，又避免了 `vfork()` 的上述怪异行止。（28.3 节会给出 `fork()` 与 `vfork()` 在速度方面的某些比较。）

SUSv3 将 `vfork()` 标记为已过时，SUSv4 则进一步将其从规范中删除。对于 `vfork()` 运作的诸多细节，SUSv3 颇有些语焉不详，因而可能将其实现为对 `fork()` 的调用。如此一来，那么 `vfork()` 的 BSD 语义将不复存在。一些 UNIX 系统还真就把 `vfork()` 实现为对 `fork()` 的调用，Linux 系统在内核 2.0 及其之前的版本中也是如此。

在使用时，一般应立即在 `vfork()` 之后调用 `exec()`。如果 `exec()` 执行失败，子进程应调用 `_exit()` 退出。（`vfork()` 产生的子进程不应调用 `exit()` 退出，因为这会导致对父进程 `stdio` 缓冲区的刷新和关闭。25.4 节将会详述这一点。）

`vfork()` 的其他用法，尤其当其依赖于内存共享以及进程调度方面的独特语义时，将可能破坏程序的可移植性，其中尤以将 `vfork()` 实现为简单调用 `fork()` 的情况为甚。

24.4 `fork()` 之后的竞争条件（Race Condition）

调用 `fork()` 后，无法确定父、子进程间谁将率先访问 CPU。（在多处理器系统中，它们可能会同时各自访问一个 CPU。）就应用程序而言，如果为了产生正确的结果而或明或暗

(implicitly or explicitly) 地依赖于特定的执行序列，那么将可能因竞争条件（5.1 节曾论及）而导致失败。由于此类问题的发生取决于内核根据系统当时的负载而做出的调度决定，故而往往难以发现。

可以用程序清单 24-5 中程序来验证这种不确定性。该程序循环使用 `fork()` 来创建多个子进程。在每个 `fork()` 调用后，父、子进程都会打印一条信息，其中包含循环计数器值以及标识父/子进程身份的字符串。例如，如果要求程序只产生一个子进程，其结果可能如下：

```
$ ./fork_whos_on_first 1
0 parent
0 child
```

可以使用该程序来生成大量子进程，并且分析其输出，观察父、子进程间每次到底由谁率先输出了结果。在某一 Linux/x86-32 2.2.19 系统上令此程序生成一百万个子进程，其分析结果表明，除去 332 次之外，都是由父进程先行输出结果（占总数的 99.97%）。

对程序清单 24-5 运行结果进行分析的脚本为 `procexec/fork_whos_on_first.count.swk`，在随本书发布的源代码中提供。

依据这一结果可以推测，在 Linux 2.2.19 中，`fork()` 之后总是继续执行父进程。而子进程之所以在 0.03% 的情况中首先输出结果，是因为父进程在有机会输出消息之前，其 CPU 时间片（CPU time slice）就到期了。换言之，如果该程序所代表的情况总是依赖于如下假设，即 `fork()` 之后总是调度父进程，那么程序通常可以正常运行，不过每 3000 次将会出现一次错误。当然，如果希望父进程能在调度子进程前执行大量工作，那么出错的可能性将会大增。在一个复杂程序中调试这样的错误会很困难。

程序清单 24-5: `fork()` 之后，父、子进程竞争输出信息

```
----- procexec/fork_whos_on_first.c
#include <sys/wait.h>
#include "tlpi_hdr.h"

int
main(int argc, char *argv[])
{
    int numChildren, j;
    pid_t childPid;

    if (argc > 1 && strcmp(argv[1], "--help") == 0)
        usageErr("%s [num-children]\n", argv[0]);

    numChildren = (argc > 1) ? getInt(argv[1], GN_GT_0, "num-children") : 1;

    setbuf(stdout, NULL);          /* Make stdout unbuffered */

    for (j = 0; j < numChildren; j++) {
        switch (childPid = fork()) {
            case -1:
                errExit("fork");

            case 0:
                printf("%d child\n", j);
                _exit(EXIT_SUCCESS);

            default:
```

```

        printf("%d parent\n", j);
        wait(NULL);           /* Wait for child to terminate */
        break;
    }
}

exit(EXIT_SUCCESS);
}

```

procexec/fork_whos_on_first.c

虽然 Linux 2.2.19 总是在 `fork()` 之后继续运行父进程，但在其他 UNIX 实现上，甚至不同版本的 Linux 内核之间，却不能视其为理所当然。在内核稳定版 2.4 系列中，一度曾试验性地推出了一个“`fork()` 之后由子进程先运行”的补丁，其调度结果与内核 2.2.19 完全相反。虽然这一改变之后又为 2.4 系列内核所舍弃，不过后来还是在 Linux 2.6 中采用，因此，程序假定于 2.2.19 内核的行为会在内核 2.6 中遭到推翻。

`fork()` 之后对父、子进程的调度谁先谁后？其结果孰优孰劣？最近的一些实验又推翻了内核开发者关于这一问题的评估。从 Linux 2.6.32 开始，父进程再度成为 `fork()` 之后，默认情况下率先调度的对象。将 Linux 专有文件 `/proc/sys/kernel/sched_child_runs_first` 设为非 0 值可以改变该默认设置。

要了解支持“`fork()` 之后先调度子进程”行为的理由，可考虑当 `fork()` 产生的子进程立即执行 `exec()` 时“写时复制”所发生的情况。此时，一方面父进程在 `fork()` 之后继续修改数据页和栈页，另一方面内核要为子进程复制那些“将要修改”的页。由于子进程一旦获得调度会立即执行 `exec()`，故而这一页复制动作纯属浪费。基于这一论点，先调度子进程的决策更佳。如此一来，等到下次调度到父进程时，就无需复制内存页了。在一个繁忙的 Linux/X86-32 系统上（内核版本为 2.6.30），利用程序清单 24-5 中程序创建一百万个子进程，结果表明子进程率先输出的情况占总数的 99.98%。（这一百分比的精确值取决于诸如系统负载之类的因素。）在其他 UNIX 实现中测试该程序的结果则表明，对于由哪一进程在 `fork()` 之后率先获得调度的问题，各系统的处理规则差异巨大。

Linux 2.6.32 改回“`fork()` 之后先调度父进程”，其论据则基于如下发现：`fork()` 之后，父进程在 CPU 中正处于活跃状态，并且其内存管理信息也被置于硬件内存管理单元的转译后备缓冲器（TLB, translation look-aside buffer）中。所以，先运行父进程将提高性能。在非正式场合下，针对分别采取上述两种行为的内核构建版本进行了时间度量，其结果也证实了这一点。

总之，值得强调的是：两种行为间的性能差异很小，对于大部分应用程序并无影响。

上述讨论清楚地阐明，**不对 `fork()` 之后执行父、子进程的特定顺序做任何假设**。若确需保证某一特定执行顺序，则必须采用某种同步技术。后续各章将会介绍多种同步技术，其中包括信号量（semaphore）、文件锁（file lock）以及进程间经由管道（pipe）的消息发送。接下来会描述另一种方法，那就是使用信号（signal）。

24.5 同步信号以规避竞争条件

调用 `fork()` 之后，如果进程某甲需等待进程某乙完成某一动作，那么某乙（即活动进程）可在动作完成后向某甲发送信号；某甲则等待即可。

程序清单 24-6 演示了这一技术。该程序假设父进程必须等待子进程完成某些动作。如果

是子进程反过来要等待父进程，那么将父、子进程中与信号相关的调用对掉即可。父、子进程甚至可能多次互发信号以协调彼此行为，尽管实际上更有可能采用信号量、文件锁或消息传递等技术来进行此类协调。

[Stevens & Rago, 2005] 建议将此类同步方法（阻塞信号，发送信号，捕获信号）封装为一组标准的进程同步函数。这一做法的优点在于，如果有意，后续可以其他进程间通信（IPC）机制替换信号的使用。

需要注意：程序清单 24-6 在 `fork()` 之前就阻塞了同步信号（`SIGUSR1`）。若父进程试图在 `fork()` 之后阻塞该信号，则避之唯恐不及的竞争条件恐怕将不期而遇。（此程序假设与子进程的信号掩码状态无关；如有必要，可以在 `fork()` 之后的子进程中解除对 `SIGUSR1` 的阻塞。）

如下 shell 会话日志（log）则展示了程序清单 24-6 的运行情况：

```
$ ./fork_sig_sync
[17:59:02 5173] Child started - doing some work
[17:59:02 5172] Parent about to wait for signal
[17:59:04 5173] Child about to signal parent
[17:59:04 5172] Parent got signal
```

程序清单 24-6：利用信号来同步进程间动作

```
----- procexec/fork_sig_sync.c
#include <signal.h>
#include "curr_time.h"          /* Declaration of currTime() */
#include "tspi_hdr.h"

#define SYNC_SIG SIGUSR1      /* Synchronization signal */

static void                    /* Signal handler - does nothing but return */
handler(int sig)
{
}

int
main(int argc, char *argv[])
{
    pid_t childPid;
    sigset_t blockMask, origMask, emptyMask;
    struct sigaction sa;

    setbuf(stdout, NULL);      /* Disable buffering of stdout */

    sigemptyset(&blockMask);
    sigaddset(&blockMask, SYNC_SIG); /* Block signal */
    if (sigprocmask(SIG_BLOCK, &blockMask, &origMask) == -1)
        errExit("sigprocmask");

    sigemptyset(&sa.sa_mask);
    sa.sa_flags = SA_RESTART;
    sa.sa_handler = handler;
    if (sigaction(SYNC_SIG, &sa, NULL) == -1)
        errExit("sigaction");

    switch (childPid = fork()) {
    case -1:
```

```

    errExit("fork");

case 0: /* Child */

    /* Child does some required action here... */
    printf("[%s %ld] Child started - doing some work\n",
           currTime("%T"), (long) getpid());
    sleep(2);          /* Simulate time spent doing some work */

    /* And then signals parent that it's done */

    printf("[%s %ld] Child about to signal parent\n",
           currTime("%T"), (long) getpid());
    if (kill(getppid(), SYNC_SIG) == -1)
        errExit("kill");

    /* Now child can do other things... */

    _exit(EXIT_SUCCESS);

default: /* Parent */

    /* Parent may do some work here, and then waits for child to
       complete the required action */

    printf("[%s %ld] Parent about to wait for signal\n",
           currTime("%T"), (long) getpid());
    sigemptyset(&emptyMask);
    if (sigsuspend(&emptyMask) == -1 && errno != EINTR)
        errExit("sigsuspend");
    printf("[%s %ld] Parent got signal\n", currTime("%T"), (long) getpid());

    /* If required, return signal mask to its original state */

    if (sigprocmask(SIG_SETMASK, &origMask, NULL) == -1)
        errExit("sigprocmask");

    /* Parent carries on to do other things... */

    exit(EXIT_SUCCESS);
}
}

```

procexec/fork_sig_sync.c

24.6 总结

系统调用 `fork()` 通过复制一个与调用进程（父进程）几乎完全一致的拷贝来创建一个新进程（子进程）。系统调用 `vfork()` 是一种更为高效的 `fork()` 版本，不过因为其语义独特——`vfork()` 产生的子进程将使用父进程内存，直至其调用 `exec()` 或退出；于此同时，将会挂起（suspended）父进程，**所以应尽量避免使用**。

调用 `fork()` 之后，不应対父、子进程获得调度以使用 CPU 的先后顺序有所依赖。对执行顺序做出假设的程序易于产生所谓“竞争条件”的错误。由于此类错误的产生依赖于诸如系统负载之类的外部因素，故而其发现和调试将非常困难。

更多信息

[Bach, 1986] 和 [Goodheart & Cox, 1994] 论述了 UNIX 系统中 `fork()`、`execve()`、`wait()` 以及 `exit()` 的实现细节。[Bovert & Cesati, 2005] 和 [Love, 2010] 则就进程的创建和终止提供了专属于 Linux 系统的实现细节。

24.7 练习

- 24-1. 程序在执行完如下一系列 `fork()` 调用后会产生多少新进程（假定没有调用失败）？

```
fork();
fork();
fork();
```
- 24-2. 编写一个程序以便验证调用 `vfork()` 之后，子进程可以关闭一文件描述符（例如：描述符 0）而不影响对应父进程中的文件描述符。
- 24-3. 假设可以修改程序源代码，如何在某一特定时刻生成一核心转储（core dump）文件，而同时进程得以继续执行？
- 24-4. 在其他 UNIX 实现上实验程序清单 24-5（`fork_whos_on_first.c`）中的程序，并判断在执行 `fork()` 后这些系统是如何调度父子进程的。
- 24-5. 假定在程序清单 24-6 的程序中，子进程也需要等待父进程完成某些操作。为确保达成这一目的，应如何修改程序？

第 25 章

进程的终止

本章所述为进程的退出过程。首先说明如何调用 `exit()`和`_exit()`以终止一个进程。接着讨论运用退出处理程序（`exit handler`），在进程调用 `exit()`时自动执行清理动作。最后，将探讨 `fork()`、`stdio` 缓冲区以及 `exit()`之间的某些交互。

25.1 进程的终止： `_exit()`和 `exit()`

通常，进程有两种终止方式。其一为异常（`abnormal`）终止，如 20.1 节所述，由对一信号的接收而引发，该信号的默认动作为终止当前进程，可能产生核心转储（`core dump`）。此外，进程可使用 `_exit()`系统调用正常（`normally`）终止。

```
#include <unistd.h>

void _exit(int status);
```

`_exit()`的 `status` 参数定义了进程的终止状态（`termination status`），父进程可调用 `wait()`以获取该状态。虽然将其定义为 `int` 类型，但仅有低 8 位可为父进程所用。按照惯例，终止状态为 0 表示进程“功成身退”，而非 0 值则表示进程因异常而退出。对非 0 返回值的解释则并无定例；不同的应用程序自成一派，并会在文档中加以描述。SUSv3 规定有两个常量：`EXIT_SUCCESS(0)`和 `EXIT_FAILURE(1)`，本书中大部分程序就采用了这一约定。

调用 `_exit()`的程序总会成功终止（即，`_exit()`从不返回）。

虽然可将 0~255 之间的任意值赋给 `_exit()`的 `status` 参数，并传递给父进程，不过如取值大于 128 将在 `shell` 脚本中引发混乱。原因在于，当以信号（`signal`）终止一命令时，`shell` 会将变量 `$?`置为 128 与信号值之和，以表征这一事实。如果这与进程调用 `_exit()`时所使用的相同 `status` 值混杂起来，将令 `shell` 无法区分。

程序一般不会直接调用 `_exit()`，而是调用库函数 `exit()`，它会在调用 `_exit()`前执行各种动作。

```
#include <stdlib.h>
```

```
void _exit(int status);
```

`exit()`会执行的动作如下。

- 调用退出处理程序（通过 `atexit()`和 `on_exit()`注册的函数），其执行顺序与注册顺序相反（见 25.3 节）。
- 刷新 `stdio` 流缓冲区。
- 使用由 `status` 提供的值执行 `_exit()`系统调用。

与专属于 UNIX 的 `_exit()`不同，`exit()`则属于标准 C 语言函数库，也就是说，所有的 C 语言实现都支持 `exit()`。

程序的另一种终止方法是从 `main()`函数中返回（`return`），或者或明或暗地一直执行到 `main()`函数的结尾处¹。执行 `return n` 等同于执行对 `exit(n)`的调用，因为调用 `main()`的运行函数会将 `main()`的返回值作为 `exit()`的参数。

存在一种情况，从 `main()`函数中返回与调用 `exit()`并不相同。如果在退出的处理过程中所执行的任何步骤需要访问 `main()`函数的本地变量，那么从 `main()`函数中返回会导致未定义的行为。例如，在调用 `setvbuf()`或 `setbuff()`（见 13.2 节）时引用了 `main` 函数的本地变量，就会发生这种情况。

执行未指定返回值的 `return`，或是无声无息地执行到 `main()`函数结尾，同样会导致 `main()`的调用者执行 `exit()`函数，不过，视所支持的不同 C 语言标准版本，以及所使用的不同编译器选项，其结果也有所不同。

- C89 标准未就上述情况下的行为进行定义，程序可以返回任意的 `status` 值。Linux `gcc` 的默认行为就是如此，程序的退出状态是取自于栈或特定 CPU 寄存器中的随机值。应避免以这一方式终止程序。
- C99 标准则要求，执行至 `main` 函数结尾处的情况应等同于调用 `exit(0)`。如果使用 `gcc-std=c99` 在 Linux 中编译程序，将会获得这种效果。

25.2 进程终止的细节

无论进程是否正常终止，都会发生如下动作。

- 关闭所有打开文件描述符、目录流（18.8 节）、信息目录描述符（参考手册页 `catopen(3)`和 `catgets(3)`），以及（字符集）转换描述符（见 `iconv_open(3)`手册页）。
- 作为文件描述符关闭的后果之一，将释放该进程所持有的任何文件锁（第 55 章）。
- 分离（`detach`）任何已连接的 System V 共享内存段，且对应于各段的 `shm_nattch` 计数器值将减一。（参考 48.8 节。）
- 进程为每个 System V 信号量所设置的 `semadj` 值将会被加到信号量值中（参考 47.8 节）。
- 如果该进程是一个管理终端（`terminal`）的管理进程，那么系统会向该终端前台（`foreground`）进程组中的每个进程发送 `SIGHUP` 信号，接着终端会与会话（`session`）脱离。34.6 节将就此进行深入讨论。
- 将关闭该进程打开的任何 POSIX 有名信号量，类似于调用 `sem_close()`。

¹ 译者注：即 `main()`函数尾部无 `return` 语句。

- 将关闭该进程打开的任何 POSIX 消息队列，类似于调用 `mq_close()`。
- 作为进程退出的后果之一，如果某进程组成为孤儿，且该组中存在任何已停止进程（stopped processes），则组中所有进程都将收到 `SIGHUP` 信号，随之为 `SIGCONT` 信号。34.7.4 节将深入讨论这一点。
- 移除该进程通过 `mlock()` 或 `mlockall()`（50.2 节）所建立的任何内存锁。
- 取消该进程调用 `mmap()` 所创建的任何内存映射（mapping）。

25.3 退出处理程序

有时，应用程序需要在进程终止时自动执行一些操作。试以一个应用程序库为例，如果进程使用了该程序库，那么在进程终止前该库需要自动执行一些清理动作。因为库本身对于进程何时以及如何退出并无控制权，也无法要求主程序在退出前调用库中特定的清理函数，故而也不能保证一定会执行清理动作。解决这一问题的方法之一是使用退出处理程序（exit handler）。老版 System V 手册则使用术语“程序终止过程”（program termination routine）。

退出处理程序是一个由程序设计者提供的函数，可于进程生命周期的任意时点注册，并在该进程调用 `exit()` 正常终止时自动执行。如果程序直接调用 `_exit()` 或因信号而异常终止，则不会调用退出处理程序。

当进程收到信号而终止时，将不会调用退出处理程序。这一事实一定程度上限制了它们的效用。此时最佳的应对方式莫若为可能发送给进程的信号建立信号处理程序，并于其中设置标志位，令主程序据此来调用 `exit()`。因为 `exit()` 不属于表 21-1 所列的异步信号安全（async-signal-safe）函数，所以通常不能在信号处理程序中对其发起调用。即便如此，还是无法处理 `SIGKILL` 信号，因为无法改变 `SIGKILL` 的默认行为。这也是应该避免使用 `SIGKILL` 来终止进程的另一原因（如 20.2 节所述）。建议使用信号 `SIGTERM`，这也是 `kill` 命令默认发送的信号。

注册退出处理程序

GNU C 语言函数库提供两种方式来注册退出处理程序。第一种方法是使用由 SUSv3 定义的 `atexit()` 函数。

```
#include <stdlib.h>

int atexit(void (*func)(void));
```

Returns 0 on success, or nonzero on error

函数 `atexit()` 将 `func` 加到一个函数列表中，进程终止时会调用该函数列表的所有函数。应将函数 `func` 定义为不接受任何参数，也无返回值，一般格式如下：

```
void
func(void)
{
    /* Perform some actions */
}
```

注意 `atexit()` 在出错时返回非 0 值（不一定是 -1）。

可以注册多个退出处理程序（甚至可以将同一函数注册多次）。当应用程序调用 `exit()` 时，

这些函数的执行顺序与注册顺序相反。这一设计很符合逻辑，因为，一般情况下较早注册的函数所执行的是更为基本的清理动作，可能需要在调用后续注册的函数后再执行。

本质上，可以在退出处理程序中执行任何希望的动作，包括注册附加的退出处理程序，并将其置于有待调用的剩余函数列表的头部。不过，一旦有任一退出处理程序无法返回——无论因为调用了 `_exit()` 还是进程因收到信号而终止（例如，退出处理程序调用函数 `raise()`），那么就不会再调用剩余的处理程序。此外，调用 `exit()` 时通常需要执行的剩余动作也将不再执行。

SUSv3 规定，若退出处理程序自身调用 `exit()`，其结果未定义。在 Linux 上，会照常调用剩余的退出处理程序。不过，在某些系统上，这将导致对所有退出处理程序的再次调用，并引发无限循环调用（直至栈溢出将该进程杀死）。为保障可移植性，应用程序应避免在退出处理程序内部调用 `exit()`。

SUSv3 要求系统实现应允许一个进程能够注册至少 32 个退出处理程序。使用系统调用 `sysconf(_SC_ATEXIT_MAX)`，应用程序即可确定由实现所定义的可注册退出处理程序的数量上限。（但是，并无方法获知有多少已注册的处理程序。）通过运用动态分配链表将已注册的处理程序串接起来，glibc 允许注册的退出处理程序数量近乎于无限。对于 Linux，`sysconf(_SC_ATEXIT_MAX)` 返回 2147482647（即，32 位有符号整型数的最大值）。换言之，在触及可注册函数数量的这一上限前，总会有其他原因（例如，内存不足）导致程序崩溃。

通过 `fork()` 创建的子进程会继承父进程注册的退出处理函数。而进程调用 `exec()` 时，会移除所有已注册的退出处理程序。（这是结果势所必然，因为 `exec()` 会替换包括退出处理程序在内的所有原程序代码段。）

无法取消经由 `atexit()` 或 `on_exit()`（见稍后的描述）注册的退出处理程序。不过，可以令退出处理程序在执行动作之前检查全局执行标志是否置位，或者清除该标志来屏蔽退出处理程序。

经由 `atexit()` 注册的退出处理程序会受到两种限制。其一，退出处理程序在执行时无法获知传递给 `exit()` 的状态。有时候，知道状态是必要的；例如，退出处理程序会视进程退出成功与否而执行不同的动作。其二，无法给退出处理程序指定参数。如果拥有这一特性，退出处理程序能根据传入参数的不同而执行不同动作，或使用不同参数多次注册同一个函数。

为摆脱这些限制，glibc 提供了一个（非标准的）替代方法：`on_exit()`。

```
#define _BSD_SOURCE          /* Or: #define _SVID_SOURCE */
#include <stdlib.h>

int on_exit(void (*func)(int, void *), void *arg);

Returns 0 on success, or nonzero on error
```

函数 `on_exit()` 的参数 `func` 是一个指针，指向如下类型的函数：

```
void
func(int status, void *arg)
{
    /* Perform cleanup actions */
}
```

调用时，会传递两个参数给 `func()`：提供给 `exit()` 的 `status` 参数和注册时供给 `on_exit()` 的一份 `arg` 参数拷贝。虽然定义为指针类型，参数 `arg` 的意义仍然可由设计者支配。可将其用作指向结构的指针，同样，通过审慎地强制转换，也可将其作为整型或其他标量类型使用。

类似于 `atexit()`，`on_exit()` 出错时返回非 0 值（不一定是 -1）。

如同 `atexit()` 一样，通过 `on_exit()` 可以注册多个退出处理程序。使用 `atexit()` 和 `on_exit()` 注册的函数位于同一函数列表。如果在程序中同时用到了这两种方式，则会按照使用这两个方法注册的相反顺序来执行相应的退出处理程序。

虽然比 `atexit()` 更灵活，但对于要保障可移植性的程序来说，还是应避免使用 `on_exit()`。因为并无标准涵盖到它，并且几乎也没有其他 UNIX 实现支持这一用法。

程序范例

程序清单 25-1 展示了如何利用 `atexit()` 和 `on_exit()` 注册退出处理程序的例子。运行程序会得到如下输出结果：

```
$ ./exit_handlers
on_exit function called: status=2, arg=20
atexit function 2 called
atexit function 1 called
on_exit function called: status=2, arg=10
```

程序清单 25-1：使用退出处理程序

```
----- procexec/exit_handlers.c
#define _BSD_SOURCE      /* Get on_exit() declaration from <stdlib.h> */
#include <stdlib.h>
#include "tspi_hdr.h"

static void
atexitFunc1(void)
{
    printf("atexit function 1 called\n");
}

static void
atexitFunc2(void)
{
    printf("atexit function 2 called\n");
}

static void
onexitFunc(int exitStatus, void *arg)
{
    printf("on_exit function called: status=%d, arg=%ld\n",
          exitStatus, (long) arg);
}

int
main(int argc, char *argv[])
{
    if (on_exit(onexitFunc, (void *) 10) != 0)
        fatal("on_exit 1");
    if (atexit(atexitFunc1) != 0)
        fatal("atexit 1");
    if (atexit(atexitFunc2) != 0)
        fatal("atexit 2");
    if (on_exit(onexitFunc, (void *) 20) != 0)
        fatal("on_exit 2");
}
```

```
    exit(2);
}
----- procexec/exit_handlers.c
```

25.4 fork()、stdio 缓冲区以及_exit()之间的交互

程序清单 25-2 生成的输出结果乍看颇令人费解。当程序标准输出定向到终端时，会看到预期的结果：

```
$ ./fork_stdio_buf
Hello world
Ciao
```

不过，当重定向标准输出到一个文件时，结果如下：

```
$ ./fork_stdio_buf > a
$ cat a
Ciao
Hello world
Hello world
```

以上输出中有两件怪事：`printf()`的输出行出现了两次，且 `write()`的输出先于 `printf()`。

程序清单 25-2：fork()与 stdio 缓冲区的交互

```
----- procexec/fork_stdio_buf.c
#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    printf("Hello world\n");
    write(STDOUT_FILENO, "Ciao\n", 5);

    if (fork() == -1)
        errExit("fork");

    /* Both child and parent continue execution here */

    exit(EXIT_SUCCESS);
}
----- procexec/fork_stdio_buf.c
```

要理解为什么 `printf()`的输出消息出现了两次，首先要记住，是在进程的用户空间内存中（参考 13.2 节）维护 `stdio` 缓冲区的。因此，通过 `fork()`创建子进程时会复制这些缓冲区。当标准输出定向到终端时，因为缺省为行缓冲，所以会立即显示函数 `printf()`输出的包含换行符的字符串。不过，当标准输出重定向到文件时，由于缺省为块缓冲，所以在本例中，当调用 `fork()`时，`printf()`输出的字符串仍在父进程的 `stdio` 缓冲区中，并随子进程的创建而产生一份副本。父、子进程调用 `exit()`时会刷新各自的 `stdio` 缓冲区，从而导致重复的输出结果。

可以采用以下任一方法来避免重复的输出结果。

- 作为针对 `stdio` 缓冲区问题的特定解决方案，可以在调用 `fork()`之前使用函数 `fflush()`来刷新 `stdio` 缓冲区。作为另一种选择，也可以使用 `setvbuf()`和 `setbuf()`来关闭 `stdio` 流的缓冲功能。

- 子进程可以调用 `_exit()` 而非 `exit()`，以便不再刷新 `stdio` 缓冲区。这一技术例证了一个更为通用的原则：在创建子进程的应用中，典型情况下仅有一个进程（一般为父进程）应通过调用 `exit()` 终止，而其他进程应调用 `_exit()` 终止，从而确保只有一个进程调用退出处理程序并刷新 `stdio` 缓冲区，这也算是众望所归吧。

还存在其他方法，可以（有时很有必要）允许父子进程都调用 `exit()`。例如，可以设计这样的退出处理程序，即使是从多个进程中调用，它们也能够正确地处理，或者令应用程序仅在调用 `fork()` 之后才去安装退出处理程序。此外，有时可能确实希望所有的应用程序都在 `fork()` 之后刷新 `stdio` 缓冲区。这时，可以见机行事，要么选择使用 `exit()` 来终止进程，要么在每个进程中均显式调用 `fflush()`。

程序清单 25-2 中 `write()` 的输出并未出现两次，这是因为 `write()` 会将数据直接传给内核缓冲区，`fork()` 不会复制这一缓冲区。

程序输出重定向到文件时出的第二件怪事，原因现在也清楚了。`write()` 的输出结果先于 `printf()` 而出现，是因为 `write()` 会将数据立即传给内核高速缓存，而 `printf()` 的输出则需要等到调用 `exit()` 刷新 `stdio` 缓冲区时。（如 13.7 节所述，通常，在混合使用 `stdio` 函数和系统调用对同一文件进行 I/O 处理时，需要特别谨慎。）

25.5 总结

进程的终止分为正常和异常两种。异常终止可能是由于某些信号引起，其中的一些信号还可能导致进程产生一个核心转储文件。

正常的终止可以通过调用 `_exit()` 完成，更多的情况下，则是使用 `_exit()` 的上层函数 `exit()` 完成。`_exit()` 和 `exit()` 都需要一个整型参数，其低 8 位定义了进程的终止状态。依照惯例，状态 0 用来表示进程成功完成，非 0 则表示异常退出。

不管进程正常终止与否，内核都会执行多个清理步骤。调用 `exit()` 正常终止一个进程，将会引发执行经由 `atexit()` 和 `on_exit()` 注册的退出处理程序（执行顺序与注册顺序相反），同时刷新 `stdio` 缓冲区。

更多信息

请参考 24.6 节所列的深入信息来源。

25.6 练习

如果子进程调用 `exit(-1)`，父进程将会看到何种退出状态（由 `WEXITSTATUS()` 返回）？

第 26 章

监控子进程

在很多应用程序的设计中，父进程需要知道其某个子进程于何时改变了状态——子进程终止或因收到信号而停止。本章描述两种用于监控子进程的技术：系统调用 `wait()`（及其变体）以及信号 `SIGCHLD`。

26.1 等待子进程

对于许多需要创建子进程的应用来说，父进程能够监测子进程的终止时间和过程是很有必要的。`wait()`以及若干相关的系统调用提供了这一功能。

26.1.1 系统调用 `wait()`

系统调用 `wait()`等待调用进程的任一子进程终止，同时在参数 `status` 所指向的缓冲区中返回该子进程的终止状态。

```
#include <sys/wait.h>

pid_t wait(int *status);

Returns process ID of terminated child, or -1 on error
```

系统调用 `wait()`执行如下动作。

1. 如果调用进程并无之前未被等待的子进程终止¹，调用将一直阻塞，直至某个子进程终止。如果调用时已有子进程终止，`wait()`则立即返回。
2. 如果 `status` 非空，那么关于子进程如何终止的信息则会通过 `status` 指向的整型变量返回。26.1.3 节将讨论自 `status` 返回的信息。
3. 内核将会为父进程下所有子进程的运行总量追加进程 CPU 时间（10.7 节）以及资源使用数据。

¹ 译者注：原文为 “if on (previously unwaited-for) child of the calling process has yet terminated”。

4. 将终止子进程的 ID 作为 wait()的结果返回。

出错时，wait()返回-1。可能的错误原因之一是调用进程并无之前未被等待的¹子进程，此时会将 errno 置为 ECHILD。换言之，可使用如下代码中的循环来等待调用进程的所有子进程退出。

```
while ((childPid = wait(NULL)) != -1)
    continue;
if (errno != ECHILD)
    errExit("wait"); /* An unexpected error... */
```

程序清单 26-1 演示了 wait()的用法。该程序创建多个子进程，每个子进程对应于一个（整型）命令行参数。每个子进程休眠若干秒后退出，休眠时间分别由相应各命令行参数指定。与此同时，在创建所有的子进程之后，父进程循环调用 wait()来监控这些子进程的终止。而直到 wait()返回-1 时才会退出循环。（这并非唯一的手段，另一种退出循环的方法是当记录终止子进程数量的变量 numDead 与创建的子进程数目相同时，也会退出循环。）以下 shell 会话日志显示了使用该程序创建 3 个子进程时的情况。

```
$ ./multi_wait 7 1 4
[13:41:00] child 1 started with PID 21835, sleeping 7 seconds
[13:41:00] child 2 started with PID 21836, sleeping 1 seconds
[13:41:00] child 3 started with PID 21837, sleeping 4 seconds
[13:41:01] wait() returned child PID 21836 (numDead=1)
[13:41:04] wait() returned child PID 21837 (numDead=2)
[13:41:07] wait() returned child PID 21835 (numDead=3)
No more children - bye!
```

如果于同一时点存在多个子进程退出，SUSv3 并未对 wait()处理这些子进程的顺序加以规定，换言之，该顺序取决于具体实现。即使不同的 Linux 内核版本之间，行为也有所不同。

程序清单 26-1：创建并等待多个子进程

```
----- procexec/multi_wait.c
#include <sys/wait.h>
#include <time.h>
#include "curr_time.h" /* Declaration of currTime() */
#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    int numDead; /* Number of children so far waited for */
    pid_t childPid; /* PID of waited for child */
    int j;

    if (argc < 2 || strcmp(argv[1], "--help") == 0)
        usageErr("%s sleep-time...\n", argv[0]);

    setbuf(stdout, NULL); /* Disable buffering of stdout */

    for (j = 1; j < argc; j++) { /* Create one child for each argument */
        switch (fork()) {
```

¹ 译者注：此处原文为 previously unwaited-for。

```

    case -1:
        errExit("fork");

    case 0:
        /* Child sleeps for a while then exits */
        printf("[%s] child %d started with PID %ld, sleeping %s "
            "seconds\n", currTime("%T"), j, (long) getpid(), argv[j]);
        sleep(getInt(argv[j], GN_NONNEG, "sleep-time"));
        _exit(EXIT_SUCCESS);

    default:
        /* Parent just continues around loop */
        break;
    }
}

numDead = 0;
for (;;) {
    /* Parent waits for each child to exit */
    childPid = wait(NULL);
    if (childPid == -1) {
        if (errno == ECHILD) {
            printf("No more children - bye!\n");
            exit(EXIT_SUCCESS);
        } else {
            /* Some other (unexpected) error */
            errExit("wait");
        }
    }

    numDead++;
    printf("[%s] wait() returned child PID %ld (numDead=%d)\n",
        currTime("%T"), (long) childPid, numDead);
}
}

```

procexec/multi_wait.c

26.1.2 系统调用 waitpid()

系统调用 `wait()` 存在诸多限制，而设计 `waitpid()` 则意在突破这些限制。

- 如果父进程已经创建了多个子进程，使用 `wait()` 将无法等待某个特定子进程的完成，只能按顺序等待下一个子进程的终止。
- 如果没有子进程退出，`wait()` 总是保持阻塞。有时候会希望执行非阻塞的等待：是否有子进程退出，立判可知。
- 使用 `wait()` 只能发现那些已经终止的子进程。对于子进程因某个信号（如 `SIGSTOP` 或 `SIGTTIN`）而停止，或是已停止子进程收到 `SIGCONT` 信号后恢复执行的情况就无能为力了。

```

#include <sys/wait.h>

pid_t waitpid(pid_t pid, int *status, int options);

```

Returns process ID of child, 0 (see text), or -1 on error

`waitpid()` 与 `wait()` 的返回值以及参数 `status` 的意义相同。（对 `status` 中返回值的解释请参考 26.1.3 节。）参数 `pid` 用来表示需要等待的具体子进程，意义如下：

- 如果 `pid` 大于 0，表示等待进程 ID 为 `pid` 的子进程。

- 如果 pid 等于 0，则等待与调用进程（父进程）同一个进程组（process group）的所有子进程。34.2 节将描述进程组的概念。
- 如果 pid 小于 -1，则会等待进程组标识符与 pid 绝对值相等的所有子进程。
- 如果 pid 等于 -1，则等待任意子进程。wait(&status)的调用与 waitpid(-1, &status, 0) 等价。

参数 options 是一个位掩码（bit mask），可以包含（按位或操作）0 个或多个如下标志（均在 SUSv3 中加以规范）。

WUNTRACED

除了返回终止子进程的信息外，还返回因信号而停止的子进程信息。

WCONTINUED (自 Linux2.6.10 以来)

返回那些因收到 SIGCONT 信号而恢复执行的已停止子进程的状态信息。

WNOHANG

如果参数 pid 所指定的子进程并未发生状态改变，则立即返回，而不会阻塞，亦即 poll（轮询）。在这种情况下，waitpid()返回 0。如果调用进程并无与 pid 匹配的子进程，则 waitpid()报错，将错误号置为 ECHILD。

程序清单 26-3 演示了 waitpid()的使用。

SUSv3 在其对 waitpid()的原理阐述中特别指出，WUNTRACED 的名称是源于 BSD 的历史产物。BSD 有两种停止进程的方法：作为系统调用 ptrace()追踪的结果，或者因收到一个信号而停止。当通过 ptrace()追踪一个子进程时，那么（除 SIGKILL 之外的）任何信号都会造成子进程停止，接着会将信号 SIGCHLD 发给父进程。即使子进程忽略这些信号，这一行为仍会发生。不过，如果子进程阻塞了这些信号（除非是无法阻塞的 SIGSTOP 信号），子进程就不会停止。

26.1.3 等待状态值

由 wait()和 waitpid()返回的 status 的值，可用来区分以下子进程事件。

- 子进程调用 _exit()（或 exit()）而终止，并指定一个整型值作为退出状态。
- 子进程收到未处理信号而终止。
- 子进程因为信号而停止，并以 WUNTRACED 标志调用 waitpid()。
- 子进程因收到信号 SIGCONT 而恢复，并以 WCONTINUED 标志调用 waitpid()。

此处用术语“等待状态”（wait status）来涵盖上述所有情况，而使用“终止状态”（termination status）的称谓来指代前两种情况。（在 shell 中，可通过读取 \$? 变量值来获取上次执行命令的终止状态。）

虽然将变量 status 定义为整型（int），但实际上仅使用了其最低的 2 个字节。对这 2 个字节的填充方式取决于子进程所发生的具体事件，如图 26-1 所示。

图 26-1 所示为 Linux/x86-32 下等待状态值的格式。不同的实现版本细节会有所不同。SUSv3 并未对信息格式做出具体规定，也未规定只能使用 status 变量的最低 2 个字节。要保证应用程序的可移植性，应总是使用本节介绍的宏（macro）来获取相应的值，而不应直接按位读取其内容。

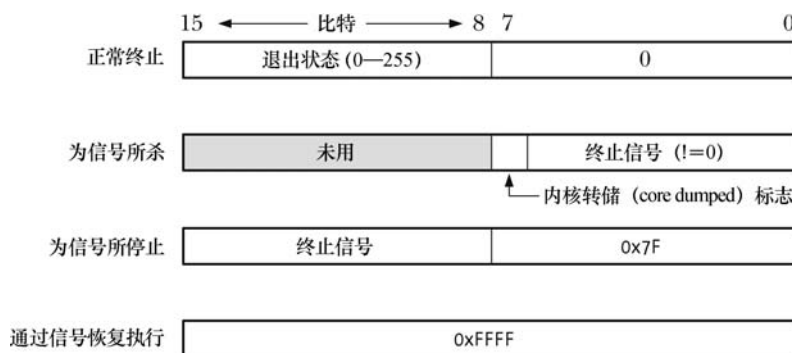


图 26-1: 自 wait()和 waitpid()的 status 参数所返回的值

头文件<sys/wait.h>定义了用于解析等待状态值的一组标准宏。对自 wait()或 waitpid()返回的 status 值进行处理时，以下列表中各宏只有一个会返回真（true）值。如列表所示，另有其他宏可对 status 值做进一步分析。

WIFEXITED (status)

若子进程正常结束则返回真（true）。此时，宏 WEXITSTATUS(status)返回子进程的退出状态。（如 25.1 节所述，父进程仅关注子进程退出状态的最低 8 位。）

WIFSIGNALED (status)

若通过信号杀掉子进程则返回真（true）。此时，宏 WTERMSIG(status)返回导致子进程终止的信号编号。若子进程产生内核转储文件，则宏 WCOREDUMP(status)返回真值（true）。SUSv3 并未规范宏 WCOREDUMP(), 不过大部分 UNIX 实现均支持该宏。

WIFSTOPPED (status)

若子进程因信号而停止，则此宏返回为真值（true）。此时，宏 WSTOPSIG(status)返回导致子进程停止的信号编号。

WIFCONTINUED (status)

若子进程收到 SIGCONT 而恢复执行，则此宏返回真值（true）。自 Linux 2.6.10 之后开始支持该宏。

注意：尽管上述宏的参数也以 status 命名，不过此处所指只是简单的整型变量，而非像 wait()和 waitpid()所要求的那样是指向整型的指针。

示例程序

程序清单 26-2 中的函数 printWaitStatus()使用了上述所有宏。此函数分析并输出了等待状态值的内容。

程序清单 26-2: 输出 wait()及相关调用返回的状态值

```

----- procexec/print_wait_status.c
#define _GNU_SOURCE /* Get strsignal() declaration from <string.h> */
#include <string.h>
#include <sys/wait.h>

```

```

#include "print_wait_status.h" /* Declaration of printWaitStatus() */
#include "tldpi_hdr.h"

/* NOTE: The following function employs printf(), which is not
   async-signal-safe (see Section 21.1.2). As such, this function is
   also not async-signal-safe (i.e., beware of calling it from a
   SIGCHLD handler). */

void          /* Examine a wait() status using the W* macros */
printWaitStatus(const char *msg, int status)
{
    if (msg != NULL)
        printf("%s", msg);

    if (WIFEXITED(status)) {
        printf("child exited, status=%d\n", WEXITSTATUS(status));

    } else if (WIFSIGNALED(status)) {
        printf("child killed by signal %d (%s)",
              WTERMSIG(status), strsignal(WTERMSIG(status)));
#ifdef WCOREDUMP /* Not in SUSv3, may be absent on some systems */
        if (WCOREDUMP(status))
            printf(" (core dumped)");
#endif
        printf("\n");

    } else if (WIFSTOPPED(status)) {
        printf("child stopped by signal %d (%s)\n",
              WSTOPSIG(status), strsignal(WSTOPSIG(status)));

#ifdef WIFCONTINUED /* SUSv3 has this, but older Linux versions and
                    some other UNIX implementations don't */
        } else if (WIFCONTINUED(status)) {
            printf("child continued\n");
        #endif

    } else { /* Should never happen */
        printf("what happened to this child? (status=%x)\n",
              (unsigned int) status);
    }
}

```

————— procexec/print_wait_status.c

程序清单 26-3 使用了 `printWaitStatus()` 函数。该程序创建了一个子进程，该子进程会循环调用 `pause()`（在此期间可以向子进程发送信号），但如果在命令行中指定了整型参数，则子进程会立即退出，并以该整型值作为退出状态。同时，父进程通过 `waitpid()` 监控子进程，打印子进程返回的状态值并将其作为参数传递给 `printWaitStatus()`。一旦发现子进程已正常退出，亦或因某一信号而终止，父进程会随即退出。

如下 shell 会话展示了执行程序清单 26-3 程序的几个例子。首先，创建一子进程并立即退出，且其状态值为 23：

```

$ ./child_status 23
Child started with PID = 15807
waitpid() returned: PID=15807; status=0x1700 (23,0)
child exited, status=23

```

接下来，在后台运行该程序，并向子进程发送 `SIGSTOP` 和 `SIGCONT` 信号。

```

$ ./child_status &
[1] 15870
$ Child started with PID = 15871
kill -STOP 15871
$ waitpid() returned: PID=15871; status=0x137f (19,127)
child stopped by signal 19 (Stopped (signal))
kill -CONT 15871
$ waitpid() returned: PID=15871; status=0xffff (255,255)
child continued

```

输出的最后两行只会在 Linux 2.6.10 及其之后的内核版本中出现，因为早期内核并不支持 `waitpid()` 的 `WCONTINUED` 选项。（由于后台运行程序的输出有时会与 `shell` 提示符混在一起，故而该 `shell` 会话稍微有些难以阅读。）

```

接着，再发送 SIGABRT 信号来终止子进程：
kill -ABRT 15871
$ waitpid() returned: PID=15871; status=0x0006 (0,6)
child killed by signal 6 (Aborted)
Press Enter, in order to see shell notification that background job has terminated
[1]+ Done ./child_status
$ ls -l core
ls: core: No such file or directory
$ ulimit -c
0

```

Display RLIMIT_CORE limit

虽然 `SIGABRT` 的默认行为是产生一个内核转储文件并终止进程，但这里并未产生转储文件。这是由于屏蔽内核转储所致，如以上命令 `ulimit` 的输出所示，将 `RLIMIT_CORE` 资源软限制（见 36.3 节）置为 0，该限制规定了转储文件大小的最大值。

再次重复同一实验，不过这次在发送信号 `SIGABRT` 给予进程之前，放开了对转储文件大小的限制。

```

$ ulimit -c unlimited
$ ./child_status &
[1] 15902
$ Child started with PID = 15903
kill -ABRT 15903
$ waitpid() returned: PID=15903; status=0x0086 (0,134)
child killed by signal 6 (Aborted) (core dumped)
Press Enter, in order to see shell notification that background job has terminated
[1]+ Done ./child_status
$ ls -l core
-rw----- 1 mtk users 65536 May 6 21:01 core

```

Allow core dumps

Send SIGABRT to child

This time we get a core dump

程序清单 26-3：使用 `waitpid()` 获取子进程状态

```

----- procexec/child_status.c
#include <sys/wait.h>
#include "print_wait_status.h" /* Declares printWaitStatus() */
#include "tlpi_hdr.h"

int
main(int argc, char *argv[])
{
    int status;
    pid_t childPid;
    if (argc > 1 && strcmp(argv[1], "--help") == 0)
        usageErr("%s [exit-status]\n", argv[0]);

    switch (fork()) {

```

```

case -1: errExit("fork");

case 0:          /* Child: either exits immediately with given
                  status or loops waiting for signals */
printf("Child started with PID = %ld\n", (long) getpid());
if (argc > 1)   /* Status supplied on command line? */
    exit(getInt(argv[1], 0, "exit-status"));
else
    /* Otherwise, wait for signals */
    for (;;)
        pause();
exit(EXIT_FAILURE);          /* Not reached, but good practice */

default:        /* Parent: repeatedly wait on child until it
                  either exits or is terminated by a signal */
    for (;;) {
        childPid = waitpid(-1, &status, WUNTRACED
#ifdef WCONTINUED          /* Not present on older versions of Linux */
            | WCONTINUED
#endif
        );
        if (childPid == -1)
            errExit("waitpid");

        /* Print status in hex, and as separate decimal bytes */

        printf("waitpid() returned: PID=%ld; status=0x%04x (%d,%d)\n",
                (long) childPid,
                (unsigned int) status, status >> 8, status & 0xff);
        printWaitStatus(NULL, status);

        if (WIFEXITED(status) || WIFSIGNALED(status))
            exit(EXIT_SUCCESS);
    }
}
}

```

procexec/child_status.c

26.1.4 从信号处理程序中终止进程

如表 20-1 所列，默认情况下某些信号会终止进程。有时，可能希望在进程终止之前执行一些清理步骤。为此，可以设置一个处理程序（handler）来捕获这些信号，随即执行清理步骤，再终止进程。如果这么做，需要牢记的是：通过 wait()和 waitpid()调用，父进程依然可以获取子进程的终止状态。例如，如果在信号处理程序中调用_exit(EXIT_SUCCESS)，父进程会认为子进程是正常终止。

如果需要通知父进程自己因某个信号而终止，那么子进程的信号处理程序应首先将自己废除，然后再次发出相同信号，该信号这次将终止这一子进程。信号处理程序需包含如下代码：

```

void
handler(int sig)
{
    /* Perform cleanup steps */

    signal(sig, SIG_DFL);          /* Disestablish handler */
}

```

```
    raise(sig);                /* Raise signal again */
}
```

26.1.5 系统调用 waitid()

与 waitpid()类似，waitid()返回子进程的状态。不过，waitid()提供了 waitpid()所没有的扩展功能。该系统调用源于系统 V (System V)，不过现在已获 SUSv3 采用，并从版本 2.6.9 开始，将其加入 Linux 内核。

在 Linux 2.6.9 之前，通过 glibc 实现提供了一版 waitid()。然而，由于完全实现该接口需要内核的支持，因此 glibc 版实现并未提供比 waitpid()更多的功能。

```
#include <sys/wait.h>

int waitid(idtype_t idtype, id_t id, siginfo_t *infop, int options);

Returns 0 on success or if WNOHANG was specified and
there were no children to wait for, or -1 on error
```

参数 idtype 和 id 指定需要等待哪些子进程，如下所示。

- 如果 idtype 为 P_ALL，则等待任何子进程，同时忽略 id 值。
- 如果 idtype 为 P_PID，则等待进程 ID 为 id 进程的子进程。
- 如果 idtype 为 P_PGID，则等待进程组 ID 为 id 各进程的所有子进程。

请注意，与 waitpid()不同，不能靠指定 id 为 0 来表示与调用者属于同一进程组的所有进程。相反，必须以 getpgid()的返回值来显式指定调用者的进程组 ID。

waitpid()与 waitid()最显著的区别在于，对于应该等待的子进程事件，waitid()可以更为精确地控制。可通过在 options 中指定一个或多个如下标识（按位或运算）来实现这种控制。

WEXITED

等待已终止的子进程，而无论其是否正常返回。

WSTOPPED

等待已通过信号而停止的子进程。

WCONTINUED

等待经由信号 SIGCONT 而恢复的子进程。

以下附加标识也可以通过按位或运算加入 options 中。

WNOHANG

与其在 waitpid()中的意义相同。如果匹配 id 值的子进程中并无状态信息需要返回，则立即返回（一个轮询）。此时，waitid()返回 0。如果调用进程并无子进程与 id 的值相匹配，则 waitid 调用失败，且错误号为 ECHILD。

WNOWAIT

通常，一旦通过 waitid()来等待子进程，那么必然会去处理所谓“状态事件”。不过，如

果指定了 `WNOHANG`，则会返回子进程状态，但子进程依然处于可等待的（`waitable`）状态，稍后可再次等待并获取相同信息。

执行成功，`waitid()`返回 0，且会更新指针 `infop` 所指向的 `siginfo_t` 结构，以包含子进程的相关信息。以下是结构 `siginfo_t` 的字段情况。

si_code

该字段包含以下值之一：`CLD_EXITED`，表示子进程已通过调用 `_exit()` 而终止；`CLD_KILLED`，表示子进程为某个信号所杀；`CLD_STOPPED`，表示子进程因某个信号而停止；`CLD_CONTINUED`，表示（之前停止的）子进程因接收到（`SIGCONT`）信号而恢复执行。

si_pid

该字段包含状态发生变化子进程的进程 ID。

si_signo

总是将该字段置为 `SIGCHLD`。

si_status

该字段要么包含传递给 `_exit()` 的子进程退出状态，要么包含导致子进程停止、继续或终止的信号值。可以通过读取 `si_code` 值来判定具体包含的是哪一种类型的信息。

si_uid

该字段包含子进程的真正用户 ID。大部分其他 UNIX 实现不会设置该字段。

在 Solaris 系统中，此结构还包含两个附加字段：`si_stime` 和 `si_utime`，分别包含子进程使用的系统和用户 CPU 时间。SUSv3 并不要求 `waitid()` 处理这两个字段。

`waitid()` 操作的一处细节需要进一步澄清。如果在 `options` 中指定了 `WNOHANG`，那么 `waitid()` 返回 0 意味着以下两种情况之一：在调用时子进程的状态已经改变（关于子进程的相关信息保存在 `infop` 指针所指向的结构 `siginfo_t` 中），或者没有任何子进程的状态有所改变。对于没有任何子进程改变状态的情况，一些 UNIX 实现（包括 Linux）会将 `siginfo_t` 结构内容清 0。这也是区分两种情况的方法之一：检查 `si_pid` 的值是否为 0。不幸的是，SUSv3 并未规范这一行为，一些 UNIX 实现此时会保持结构 `siginfo_t` 原封不动。（未来针对 SUSv4 的勘误表可能会增加在这种情况下将 `si_pid` 和 `si_signo` 置 0 的要求。）区分这两种情况唯一可移植的方法是：在调用 `waitid()` 之前就将结构 `siginfo_t` 的内容置为 0，正如以下代码所示：

```
siginfo_t info;
...
memset(&info, 0, sizeof(siginfo_t));
if (waitid(idtype, id, &info, options | WNOHANG) == -1)
    errExit("waitid");
if (info.si_pid == 0) {
    /* No children changed state */
} else {
    /* A child changed state; details are provided in 'info' */
}
```

26.1.6 系统调用 `wait3()` 和 `wait4()`

系统调用 `wait3()` 和 `wait4()` 执行与 `waitpid()` 类似的工作。主要的语义差别在于，`wait3()` 和 `wait4()`

在参数 `rusage` 所指向的结构中返回终止子进程的资源使用情况。其中包括进程使用的 CPU 时间总量以及内存管理的统计数据。36.1 节将在介绍系统调用 `getrusage()` 时详细讨论 `rusage` 结构。

```
#define _BSD_SOURCE      /* Or #define _XOPEN_SOURCE 500 for wait3() */
#include <sys/resource.h>
#include <sys/wait.h>

pid_t wait3(int *status, int options, struct rusage *rusage);
pid_t wait4(pid_t pid, int *status, int options, struct rusage *rusage);

Both return process ID of child, or -1 on error
```

除了对参数 `rusage` 的使用之外，调用 `wait3()` 等同于以如下方式调用 `waitpid()`：

```
waitpid(-1, &status, options);
```

与之相类似，对 `wait4()` 的调用等同于对 `waitpid()` 的如下调用：

```
waitpid(pid, &status, options);
```

换言之，`wait3()` 等待的是任意子进程，而 `wait4()` 则可以用于等待选定的一个或多个子进程。

在一些 UNIX 实现中，`wait3()` 和 `wait4()` 仅返回已终止子进程的资源使用情况。而对于 Linux 系统，如果在 `options` 中指定了 `WUNTRACED` 选项，则还可以获取到停止子进程的资源使用信息。

这两个系统调用的名称来自于它们所使用参数的个数。虽然源自 BSD 系统，不过现在大部分的 UNIX 实现都支持它们。这两个系统调用均未获得 SUSv3 标准的接纳。（SUSv2 标准纳入了 `wait3()`，但将其标记为“已过时”。）

本书一般会避免使用 `wait3()` 和 `wait4()`。通常情况下，此类调用所返回的额外信息没有什么价值。此外，未获业界标准的接纳也会限制其可移植性。

26.2 孤儿进程与僵尸进程

父进程与子进程的生命周期一般都不相同，父、子进程间互有长短。这就引出了下面两个问题。

- 谁会是孤儿 (orphan) 子进程的父进程？进程 ID 为 1 的众进程之祖——`init` 会接管孤儿进程。换言之，某一子进程的父进程终止后，对 `getppid()` 的调用将返回 1。这是判定某一子进程之“生父”是否“在世”的方法之一（前提是假设该子进程由 `init` 之外的进程创建）。

使用参数 `PR_SET_PDEATHSIG` 调用 Linux 特有的系统调用 `prctl()`，将有可能导致某一进程在成为孤儿时收到特定信号。

- 在父进程执行 `wait()` 之前，其子进程就已经终止，这将会发生什么？此处的要点在于，即使子进程已经结束，系统仍然允许其父进程在之后的某一时刻去执行 `wait()`，以确定该子进程是如何终止的。内核通过将子进程转为僵尸进程 (zombie) 来处理这种情况。这也意味着将释放子进程所把持的大部分资源，以便供其他进程重新使用。该进程所唯一保留的是内核进程表中的一条记录，其中包含了子进程 ID、终止状态、资源使用数据 (36.1 节) 等信息。

至于僵尸进程名称的由来，则源于 UNIX 系统对电影情节的效仿——无法通过信号来杀死僵尸进程，即便是 (银弹) `SIGKILL`。这就确保了父进程总是可以执行 `wait()` 方法。

当父进程执行 `wait()` 后，由于不再需要子进程所剩余的最后信息，故而内核将删除僵尸进程。另一方面，如果父进程未执行 `wait()` 随即退出，那么 `init` 进程将接管子进程并自动调用 `wait()`，

从而从系统中移除僵尸进程。

如果父进程创建了某一子进程，但并未执行 `wait()`，那么在内核的进程表中将为该子进程永久保留一条记录。如果存在大量此类僵尸进程，它们势必将填满内核进程表，从而阻碍新进程的创建。既然无法用信号杀死僵尸进程，那么从系统中将其移除的唯一方法就是杀掉它们的父进程（或等待其父进程终止），此时 `init` 进程将接管和等待这些僵尸进程，从而从系统中将它们清理掉。

在设计长生命周期的父进程（例如：会创建众多子进程的网络服务器和 `Shell`）时，这些语义具有重要意义。换句话说，在此类应用中，父进程应执行 `wait()` 方法，以确保系统总是能够清理那些死去的子进程，避免使其成为长寿僵尸。如 26.3.1 节所述，父进程在处理 `SIGCHLD` 信号时，对 `wait()` 的调用既可同步，也可异步。

程序清单 26-4 展示了一个僵尸进程的创建，以及发送 `SIGKILL` 信号无法杀死僵尸进程的例子。运行这一程序的输出如下：

```
$ ./make_zombie
Parent PID=1013
Child (PID=1014) exiting
 1013 pts/4    00:00:00 make_zombie                Output from ps(1)
 1014 pts/4    00:00:00 make_zombie <defunct>
After sending SIGKILL to make_zombie (PID=1014):
 1013 pts/4    00:00:00 make_zombie                Output from ps(1)
 1014 pts/4    00:00:00 make_zombie <defunct>
```

以上输出中，`ps(1)`所输出的字符串`<defunct>`表示进程处于僵尸状态。

程序清单 26-4 使用 `system()` 函数来执行通过字符串参数传入的 `shell` 命令。27.6 节将会详细描述 `system()` 函数。

程序清单 26-4：创建一个僵尸子进程

```
----- procexec/make_zombie.c
#include <signal.h>
#include <libgen.h>          /* For basename() declaration */
#include "tlpi_hdr.h"

#define CMD_SIZE 200

int
main(int argc, char *argv[])
{
    char cmd[CMD_SIZE];
    pid_t childPid;

    setbuf(stdout, NULL);    /* Disable buffering of stdout */

    printf("Parent PID=%ld\n", (long) getpid());

    switch (childPid = fork()) {
    case -1:
        errExit("fork");

    case 0:    /* Child: immediately exits to become zombie */
        printf("Child (PID=%ld) exiting\n", (long) getpid());
        exit(EXIT_SUCCESS);
    default:    /* Parent */
        sleep(3);    /* Give child a chance to start and exit */
        snprintf(cmd, CMD_SIZE, "ps | grep %s", basename(argv[0]));
```

```

cmd[CMD_SIZE - 1] = '\0';      /* Ensure string is null-terminated */
system(cmd);                  /* View zombie child */

/* Now send the "sure kill" signal to the zombie */

if (kill(childPid, SIGKILL) == -1)
    errMsg("kill");
sleep(3);                     /* Give child a chance to react to signal */
printf("After sending SIGKILL to zombie (PID=%ld):\n", (long) childPid);
system(cmd);                  /* View zombie child again */

exit(EXIT_SUCCESS);
}
}

```

procexec/make_zombie.c

26.3 SIGCHLD 信号

子进程的终止属异步事件。父进程无法预知其子进程何时终止。（即使父进程向子进程发送 SIGKILL 信号，子进程终止的确切时间还依赖于系统的调度：子进程下一次在何时使用 CPU。）之前已经论及，父进程应使用 wait()（或类似调用）来防止僵尸子进程的累积，以及采用如下两种方法来避免这一问题。

- 父进程调用不带 WNOHANG 标志的 wait()，或 waitpid()方法，此时如果尚无已经终止的子进程，那么调用将会阻塞。
- 父进程周期性地调用带有 WNOHANG 标志的 waitpid()，执行针对已终止子进程的非阻塞式检查（轮询）。

这两种方法使用起来都有所不便。一方面，可能并不希望父进程以阻塞的方式来等待子进程的终止。另一方面，反复调用非阻塞的 waitpid()会造成 CPU 资源的浪费，并增加应用程序设计的复杂度。为了规避这些问题，**可以采用针对 SIGCHLD 信号的处理程序**。

26.3.1 为 SIGCHLD 建立信号处理程序

无论一个子进程于何时终止，系统都会向其父进程发送 SIGCHLD 信号。对该信号的默认处理是将其忽略，不过也可以安装信号处理程序来捕获它。在处理程序中，可以使用 wait()（或类似方法）来收拾僵尸进程。不过，使用这一方法时需要掌握一些窍门。

由 20.10 节和 20.12 节可知，当调用信号处理程序时，会暂时将引发调用的信号阻塞起来（除非为 sigaction()指定了 SA_NODEFER 标志），且不会对 SIGCHLD 之流的标准信号进行排队处理。这样一来，当 SIGCHLD 信号处理程序正在为一个终止的子进程运行时，如果相继有两个子进程终止，即使产生了两次 SIGCHLD 信号，父进程也只能捕获到一个。结果是，如果父进程的 SIGCHLD 信号处理程序每次只调用一次 wait()，那么一些僵尸子进程可能会成为“漏网之鱼”。

解决方案是：在 SIGCHLD 处理程序内部循环以 WNOHANG 标志来调用 waitpid()，直至再无其他终止的子进程需要处理为止。通常 SIGCHLD 处理程序都简单地由以下代码组成，仅仅捕获已终止子进程而不关心其退出状态。

```

while (waitpid(-1, NULL, WNOHANG) > 0)
    continue;

```

上述循环会一直持续下去，直至 waitpid()返回 0，表明再无僵尸子进程存在，或-1，表示有

错误发生（可能是 ECHILD，意即再无更多的子进程）。

SIGCHLD 处理程序的设计问题

假设创建 SIGCHLD 处理程序的时候，该进程已经有子进程终止。那么内核会立即为父进程产生 SIGCHLD 信号吗？SUSv3 对这一点并未规定。一些源自系统 V (System V) 的实现在这种情况下会产生 SIGCHLD 信号；而另一些系统，包括 Linux，则不这么做。为保障可移植性，应用应在创建任何子进程之前就设置好 SIGCHLD 处理程序，将这一隐患消解于无形。（无疑，这也是顺其自然的处事之道。）

需要更深入考虑的问题是可重入性 (reentrancy)。21.1.2 节特别指出，在信号处理程序中使用系统调用（如 waitpid()）可能会改变全局变量 errno 的值。当主程序企图显式设置 errno（参考 35.1 节对 getpriority() 的讨论）或是在系统调用失败后检查 errno 值时，这一变化会与之发生冲突。出于这一原因，有时在编写 SIGCHLD 信号处理程序时，需要在一进入处理程序时就使用本地变量来保存 errno 值，而在返回前加以恢复。请参考程序清单 26-5。

范例程序

程序清单 26-5 提供了一个更为复杂的 SIGCHLD 信号处理程序示例。该处理程序为所捕获的每个子进程输出进程号及其等待状态①。为了模拟调用处理程序期间产生多个 SIGCHLD 信号而无法排队的效果，利用 sleep() 调用②人为地拉长了处理程序的执行时间。主程序为每个（整型）命令行参数创建一个子进程④。每个子进程持续休眠其对应命令行参数所指定的秒数，随即退出⑤。从程序下面的执行例子可以看出，尽管有 3 个子进程退出，而父进程只捕获到两次 SIGCHLD 信号。

```
$ ./multi_SIGCHLD 1 2 4
16:45:18 Child 1 (PID=17767) exiting
16:45:18 handler: Caught SIGCHLD           First invocation of handler
16:45:18 handler: Reaped child 17767 - child exited, status=0
16:45:19 Child 2 (PID=17768) exiting       These children terminate during...
16:45:21 Child 3 (PID=17769) exiting       first invocation of handler
16:45:23 handler: returning                End of first invocation of handler
16:45:23 handler: Caught SIGCHLD           Second invocation of handler
16:45:23 handler: Reaped child 17768 - child exited, status=0
16:45:23 handler: Reaped child 17769 - child exited, status=0
16:45:28 handler: returning
16:45:28 All 3 children have terminated; SIGCHLD was caught 2 times
```

请注意，在程序清单 26-5 中，在创建子进程之前使用了 sigprocmask() 来阻塞 SIGCHLD 信号③。这一做法确保了父进程中 sigsuspend() 循环的正确操作。如果以此方式未能阻塞 SIGCHLD 信号，而某一子进程又在对 numLiveChildren 的检查和执行 sigsuspend() 调用（也可以是 pause() 调用）之间终止，那么 sigsuspend() 调用会永远阻塞，等待一个早已捕获过的信号⑥。22.9 节详细描述了处理此类竞争条件的要求。

程序清单 26-5：通过 SIGCHLD 信号处理程序捕获已终止的子进程

```
----- procexec/multi_SIGCHLD.c
#include <signal.h>
#include <sys/wait.h>
#include "print_wait_status.h"
#include "curr_time.h"
#include "tspi_hdr.h"

static volatile int numLiveChildren = 0;
/* Number of children started but not yet waited on */
```

```

static void
sigchldHandler(int sig)
{
    int status, savedErrno;
    pid_t childPid;

    /* UNSAFE: This handler uses non-async-signal-safe functions
       (printf(), printWaitStatus(), currTime()); see Section 21.1.2) */

    savedErrno = errno;    /* In case we modify 'errno' */

    printf("%s handler: Caught SIGCHLD\n", currTime("%T"));

    while ((childPid = waitpid(-1, &status, WNOHANG)) > 0) {
        ①    printf("%s handler: Reaped child %ld - ", currTime("%T"),
                (long) childPid);
            printWaitStatus(NULL, status);
            numLiveChildren--;
        }

        if (childPid == -1 && errno != ECHILD)
            errMsg("waitpid");
        ②    sleep(5);    /* Artificially lengthen execution of handler */
            printf("%s handler: returning\n", currTime("%T"));

        errno = savedErrno;
    }

    int
    main(int argc, char *argv[])
    {
        int j, sigCnt;
        sigset_t blockMask, emptyMask;
        struct sigaction sa;

        if (argc < 2 || strcmp(argv[1], "--help") == 0)
            usageErr("%s child-sleep-time...\n", argv[0]);

        setbuf(stdout, NULL);    /* Disable buffering of stdout */

        sigCnt = 0;
        numLiveChildren = argc - 1;

        sigemptyset(&sa.sa_mask);
        sa.sa_flags = 0;
        sa.sa_handler = sigchldHandler;
        if (sigaction(SIGCHLD, &sa, NULL) == -1)
            errMsg("sigaction");

        /* Block SIGCHLD to prevent its delivery if a child terminates
           before the parent commences the sigsuspend() loop below */

        sigemptyset(&blockMask);
        sigaddset(&blockMask, SIGCHLD);
        ③    if (sigprocmask(SIG_SETMASK, &blockMask, NULL) == -1)
            errMsg("sigprocmask");

        ④    for (j = 1; j < argc; j++) {

```

```

switch (fork()) {
case -1:
    errExit("fork");

case 0:          /* Child - sleeps and then exits */
    ⑤ sleep(getInt(argv[j], GN_NONNEG, "child-sleep-time"));
    printf("%s Child %d (PID=%ld) exiting\n", currTime("%T"),
           j, (long) getpid());
    _exit(EXIT_SUCCESS);

default:        /* Parent - loops to create next child */
    break;
}
}
/* Parent comes here: wait for SIGCHLD until all children are dead */

sigemptyset(&emptyMask);
while (numLiveChildren > 0) {
    ⑥ if (sigsuspend(&emptyMask) == -1 && errno != EINTR)
        errExit("sigsuspend");
    sigCnt++;
}

printf("%s All %d children have terminated; SIGCHLD was caught "
       "%d times\n", currTime("%T"), argc - 1, sigCnt);

exit(EXIT_SUCCESS);
}

```

procexec/multi_SIGCHLD.c

26.3.2 向已停止的子进程发送 SIGCHLD 信号

正如可以使用 `waitpid()` 来监测已停止的子进程一样，当信号导致子进程停止时，父进程也就有可能收到 `SIGCHLD` 信号。调用 `sigaction()` 设置 `SIGCHLD` 信号处理程序时，如传入 `SA_NOCLDSTOP` 标志即可控制这一行为。若未使用该标志，系统会在子进程停止时向父进程发送 `SIGCHLD` 信号；反之，如果使用了这一标志，那么就不会因子进程的停止而发出 `SIGCHLD` 信号。（22.7 节中对 `signal()` 的实现就未指定 `SA_NOCLDSTOP`。）

因为默认情况下会忽略信号 `SIGCHLD`，`SA_NOCLDSTOP` 标志仅在设置 `SIGCHLD` 信号处理程序时才有意义。而且，`SA_NOCLDSTOP` 只对 `SIGCHLD` 信号起作用。

`SUSv3` 也允许，当信号 `SIGCONT` 导致已停止的子进程恢复执行时，向其父进程发送 `SIGCHLD` 信号。（相当于 `waitpid()` 的 `WCONTINUED` 标志。）始于版本 2.6.9，Linux 内核实现了这一特性。

26.3.3 忽略终止的子进程

更有可能像这样处理终止子进程：将对 `SIGCHLD` 的处置 (disposition) 显式置为 `SIG_IGN`，系统从而会将其后终止的子进程立即删除，毋庸转为僵尸进程。这时，会将子进程的状态弃之不问，故而所有后续的 `wait()`（或类似）调用不会返回子进程的任何信息。

注意，虽然对信号 `SIGCHLD` 的默认处置就是将其忽略，但显式设置对 `SIG_IGN` 标志的处置还是会导致这里所描述的行为差异。在这方面，对信号 `SIGCHLD` 的处理非常独特，不同于其他信号。

如同许多 UNIX 实现一样，在 Linux 系统中将对 SIGCHLD 信号的处处置为 SIG_IGN 并不会影响任何既有僵尸进程的状态，对它们的等待仍然要照常进行。在其他一些 UNIX 实现中（例如 Solaris 8），将对 SIGCHLD 的处处置为 SIG_IGN 确实会删除所有已有的僵尸进程。

信号 SIGCHLD 的 SIG_IGN 语义由来已久，源于系统 V（System V）。SUSv3 也规定了此处所描述的行为，不过原始的 POSIX.1 标准对此则未作表述。因此，在一些较老的 UNIX 实现中，忽略 SIGCHLD 并不影响僵尸进程的创建。要防止产生僵尸进程，唯一完全可移植的方法就是（可能是从 SIGCHLD 信号处理程序的内部）调用 wait() 或者 waitpid()。

老版本 Linux 内核实现与 SUSv3 标准的差异

SUSv3 规定，如果将对 SIGCHLD 的处处置为 SIG_IGN，那么将丢弃子进程的资源使用信息，且若指定 RUSAGE_CHILDREN 标志调用 getrusage() 函数，其返回总量中也将不包含该项信息（36.1 节）。然而，在版本 2.6.9 之前的 Linux 内核中，还是会记录子进程的 CPU 使用时间以及资源的使用情况，并可通过 getrusage() 调用获取。这一违规行为直至 Linux 2.6.9 才得以修正。

将对 SIGCHLD 的处处置为 SIG_IGN 还会阻止 times()（10.7 节）返回的结构中包含子进程的 CPU 使用时间。不过，在 Linux 2.6.9 之前，times() 所返回的信息同样存在违规行为。

SUSv3 规定，如果将对 SIGCHLD 的处处置为 SIG_IGN，同时，父进程已终止的子进程中并无处于僵尸状态且未被等待的情况，那么 wait()（或 waitpid()）调用将一直阻塞，直至所有子进程都终止，届时该调用将返回错误 ECHILD。Linux 2.6 符合这一要求。不过在 Linux 2.4 以及更早期的版本中，wait() 只会阻塞到下一个子进程终止的时刻，随即返回该子进程的进程 ID 及状态（亦即，此行为与未将对 SIGCHLD 信号的处处置为 SIG_IGN 时一样）。

sigaction() 的 SA_NOCLDWAIT 标志

SUSv3 规定了 SA_NOCLDWAIT 标志，可在调用 sigaction() 对 SIGCHLD 信号的处处置时使用此标志。设置该标志的作用类似于将对 SIGCHLD 的处处置为 SIG_IGN 时的效果。Linux 2.4 及其早期版本并未实现该标志，直至 Linux 2.6 才实现对其支持。

将对 SIGCHLD 的处处置为 SIG_IGN 与采用 SA_NOCLDWAIT 之间最主要的区别在于，当以 SA_NOCLDWAIT 设置信号处理程序时，SUSv3 并未规定系统在子进程终止时是否向其父进程发送 SIGCHLD 信号。换言之，当指定 SA_NOCLDWAIT 时允许系统发送 SIGCHLD 信号，则应用程序即可捕捉这一信号（尽管由于内核已经丢弃了僵尸进程，造成 SIGCHLD 处理程序无法用 wait() 来获得子进程状态）。在包括 Linux 在内的一些 UNIX 实现中，内核确实会为父进程产生 SIGCHLD 信号。而在另一些 UNIX 实现中，则不会。

当为 SIGCHLD 信号设置 SA_NOCLDWAIT 标志时，老版本 Linux 内核的行为细节同样与 SUSv3 不符，正如之前在将对 SIGCHLD 的处处置为 SIG_IGN 处所讨论的那样。

系统 V 的 SIGCLD 信号

Linux 系统中，信号 SIGCLD 与信号 SIGCHLD 意义相同。之所以两个名称并存，是由历史原因造成的。SIGCHLD 信号源自 BSD，POSIX 标准采用了这一名称，同时对 BSD 信号模型做了大量标准化工作。系统 V 则提供了相应的 SIGCLD 信号，在语义上稍许有些不同。

BSD SIGCHLD 信号与系统 V SIGCLD 间的主要差别在于，将信号处处置为 SIG_IGN 时不

同的处理方式。

- 在历史（和一些现代）的 BSD 实现中，即使忽略了 SIGCHLD 信号，系统仍会继续将无人等待的子进程变为僵尸进程。
- 在系统 V 上，使用 signal()（而非 sigaction()）忽略 SIGCHLD 信号将导致子进程在终止时不会转为僵尸进程。

如前所述，原始的 POSIX.1 标准对于忽略 SIGCHLD 的后果未作规定，从而也认可了系统 V 的行为。而今，系统 V 的行为已成为 SUSv3 标准的一部分（不过将仍然使用 SIGCHLD 的名称）。衍生自系统 V 的现代系统实现中对该信号使用了 SIGCHLD 这一标准名称，同时继续提供具有相同含义的 SIGCLD 信号。关于 SIGCLD 的更多信息可参考[Stevens & Rago, 2005]。

26.4 总结

使用 wait()和 waitpid()（以及其他相关函数），父进程可以得到其终止或停止子进程的状态。该状态表明子进程是正常终止（带有表示成功或失败的退出状态），还是异常中止，因收到某个信号而停止，还是因收到 SIGCONT 信号而恢复执行。

如果子进程的父进程终止，那么子进程将变为孤儿进程，并为进程 ID 为 1 的 init 进程接管。

子进程终止后会变为僵尸进程，仅当其父进程调用 wait()（或类似函数）获取子进程退出状态时，才能将其从系统中删除。在设计长时间运行的程序，诸如 shell 程序以及守护进程（daemon）时，应总是捕获其所创建子进程的状态，因为系统无法杀死僵尸进程，而未处理的僵尸进程最终将塞满内核进程表。

捕获终止子进程的一般方法是为信号 SIGCHLD 设置信号处理程序。当子进程终止时（也可选择子进程因信号而停止时），其父进程会收到 SIGCHLD 信号。还有另一种移植性稍差的处理方法，进程可选择将对 SIGCHLD 信号的处处置为忽略（SIG_IGN），这时将立即丢弃终止子进程的状态（因此其父进程从此也无法获取到这些信息），子进程也不会成为僵尸进程。

更多信息

请参考列于 24.6 节中的更多信息来源。

26.5 练习

- 26-1.** 编写一程序以验证当一子进程的父进程终止时，调用 getppid()将返回 1（进程 init 的进程 ID）。
- 26-2.** 假设存在 3 个相互关联的进程（祖父、父及子进程），祖父进程没有在父进程退出之后立即执行 wait()，所以父进程变成僵尸进程。那么请指出孙进程何时被 init 进程收养（即孙进程调用 getppid()将返回 1），是在父进程终止后，还是祖父进程调用 wait()后？请编写程序验证结果。
- 26-3.** 使用 waitid()替换程序清单 26-3（child_status.c）中的 waitpid()。需要将对函数 print WaitStatus()的调用替换为打印 waitid()所返回 siginfo_t 结构中相关字段的代码。
- 26-4.** 程序清单 26-4（make_zombie.c）调用了 sleep()，以便允许子进程在父进程执行函数 system()前得到机会去运行并终止。这一方法理论上存在产生竞争条件的可能。修改此程序，使用信号来同步父子进程以消除该竞争条件。

第 27 章

程序的执行

承接前几章对于进程创建和终止的探讨，本章首先将介绍系统调用 `execve()`，通过该调用，进程能以全新程序来替换当前运行的程序；接下来会讨论函数 `system()` 的实现，该函数可允许调用者执行任意 shell 命令。

27.1 执行新程序：execve()

系统调用 `execve()` 可以将新程序加载到某一进程的内存空间。在这一操作过程中，将丢弃旧有程序，而进程的栈、数据以及堆段会被新程序的相应部件所替换。在执行了各种 C 语言函数库的运行时启动代码以及程序的初始化代码后，例如，C++ 静态构造函数，或者以 `gcc constructor` 属性（见 42.4 节）声明的 C 语言函数，新程序会从 `main()` 函数处开始执行。

由 `fork()` 生成的子进程对 `execve()` 的调用最为频繁，不以 `fork()` 调用为先导而单独调用 `execve()` 的做法在应用中实属罕见。

基于系统调用 `execve()`，还提供了一系列冠以 `exec` 来命名的上层库函数，虽然接口方式各异，但功能相同。通常将调用这些函数加载一个新程序的过程称作 `exec` 操作，或是简单地以 `exec()` 来表示。下面将先描述 `execve()`，然后再对相关库函数进行说明。

```
#include <unistd.h>

int execve(const char *pathname, char *const argv[], char *const envp[]);

Never returns on success; returns -1 on error
```

参数 `pathname` 包含准备载入当前进程空间的新程序的路径名，既可以是绝对路径（冠之以 /），也可以是相对于调用进程当前工作目录（`current working directory`）的相对路径。

参数 `argv` 则指定了传递给新进程的命令行参数。该数组对应于 C 语言 `main()` 函数的第 2 个参数（`argv`），且格式也与之相同：是由字符串指针所组成的列表，以 `NULL` 结束。`argv[0]` 的值则对应于命令名。通常情况下，该值与 `pathname` 中的 `basename`（路径名的最后部分）相同。

最后一个参数 `envp` 指定了新程序的环境列表。参数 `envp` 对应于新程序的 `environ` 数组：

也是由字符串指针组成的列表，以 NULL 结束，所指向的字符串格式为 name=value（6.7 节）。

Linux 所特有的/proc/PID/exe 文件是一个符号链接，包含 PID 对应进程中正在运行可执行文件的绝对路径名。

调用 `execve()` 之后，因为同一进程依然存在，所以进程 ID 仍保持不变。如 28.4 节所述，还有少量其他的进程属性也未发生变化。

如果对 `pathname` 所指定的程序文件设置了 `set-user-ID`（`set-group-ID`）权限位，那么系统调用会在执行此文件时将进程的有效（effective）用户（组）ID 置为程序文件的属主（组）ID。利用这一机制，可令用户在运行特定程序时临时获取特权。（参考 9.3 节）。

无论是否更改了有效 ID，也不管这一变化是否生效，`execve()` 都会以进程的有效用户 ID 去覆盖已保存的（saved）`set-user-ID`，以进程的有效组 ID 去覆盖已保存的（saved）`set-group-ID`。

由于是将调用程序取而代之，对 `execve()` 的成功调用将永不返回，而且也无需检查 `execve()` 的返回值，因为该值总是雷打不动地等于 -1。实际上，一旦函数返回，就表明发生了错误。通常，可以通过 `errno` 来判断出错原因。可能自 `errno` 返回的错误如下：

EACCES

参数 `pathname` 没有指向一个常规（regular）文件，未对该文件赋予可执行权限，或者因为 `pathname` 中某一级目录不可搜索（not searchable）（即，关闭了该目录的可执行权限）。还有一种可能，是以 `MS_NOEXEC` 标志（14.8.1 节）来挂载（mount）文件所在的文件系统，从而导致这一错误。

ENOENT

`pathname` 所指代的文件并不存在。

ENOEXEC

尽管对 `pathname` 所指代文件赋予了可执行权限，但系统却无法识别其文件格式。一个脚本文件，如果没有包含用于指定脚本解释器（interpreter）（以字符 #! 开头）的起始行，就可能导致这一错误。

ETXTBSY

存在一个或多个进程已经以写入方式打开 `pathname` 所指代的文件（4.3.2 节）。

E2BIG

参数列表和环境列表所需空间总和超出了允许的最大值。

当上述任一条件作用于执行脚本的脚本解释器，或是执行程序的 ELF 解释器时，同样会产生相应错误。

ELF（Executable and Linking Format）是一种广为实现的标准，描述了可执行文件的布局。在执行期间，进程映像（image）通常是由可执行文件的各段（segment）构造而成（6.3 节）。不过，ELF 规格也允许定义一个解释器（ELF 程序头部的 `PT_INTERP` 元素）来运行程序。如果定义了解释器，内核则基于指定解释器可执行文件的各段来构建进程映像，转由解释器负责加载和执行程序。第 41 章会对 ELF 解释器做进一步描述，并给出对深层信息的一些指引。

示例程序

程序清单 27-1 展示了 `execve()` 的用法。该程序首先为新程序创建参数列表和环境列表，接着调用 `execve()` 来执行由命令行参数 (`argv[1]`) 所指定的程序路径名。

程序清单 27-2 中所展示的程序，是设计专供程序清单 27-1 中程序来执行的。该程序只是简单显示一下自身的命令行参数以及环境列表（对后者的访问使用了全局变量 `environ`，如 6.7 节所述）。

如下 shell 会话 (session) 演示了对程序清单 27-1 和程序清单 27-2 的使用（本例在指定执行程序时使用的是相对路径名）：

```
$ ./t_execve ./envargs
argv[0] = envargs                All of the output is printed by envargs
argv[1] = hello world
argv[2] = goodbye
environ: GREET=salut
environ: BYE=adieu
```

程序清单 27-1：调用函数 `execve()` 来执行新程序

```
----- procexec/t_execve.c
#include "t1pi_hdr.h"

int
main(int argc, char *argv[])
{
    char *argVec[10];             /* Larger than required */
    char *envVec[] = { "GREET=salut", "BYE=adieu", NULL };

    if (argc != 2 || strcmp(argv[1], "--help") == 0)
        usageErr("%s pathname\n", argv[0]);

    argVec[0] = strrchr(argv[1], '/');    /* Get basename from argv[1] */
    if (argVec[0] != NULL)
        argVec[0]++;
    else
        argVec[0] = argv[1];
    argVec[1] = "hello world";
    argVec[2] = "goodbye";
    argVec[3] = NULL;             /* List must be NULL-terminated */

    execve(argv[1], argVec, envVec);
    errExit("execve");          /* If we get here, something went wrong */
}
----- procexec/t_execve.c
```

程序清单 27-2：显示参数列表和环境列表

```
----- procexec/envargs.c
#include "t1pi_hdr.h"

extern char **environ;

int
main(int argc, char *argv[])
{
    int j;
    char **ep;
```

```

    for (j = 0; j < argc; j++)
        printf("argv[%d] = %s\n", j, argv[j]);

    for (ep = environ; *ep != NULL; ep++)
        printf("environ: %s\n", *ep);

    exit(EXIT_SUCCESS);
}

```

procexec/envargs.c

27.2 exec()库函数

本节所讨论的库函数为执行 `exec()` 提供了多种 API 选择。所有这些函数均构建于 `execve()` 调用之上，只是在为新程序指定程序名、参数列表以及环境变量的方式上有所不同。

```

#include <unistd.h>

int execl(const char *pathname, const char *arg, ...
          /* , (char *) NULL, char *const envp[] */ );
int execlp(const char *filename, const char *arg, ...
          /* , (char *) NULL */);
int execvp(const char *filename, char *const argv[]);
int execv(const char *pathname, char *const argv[]);
int execl(const char *pathname, const char *arg, ...
          /* , (char *) NULL */);

    None of the above returns on success; all return -1 on error

```

各函数名称的最后一个字母为区分这些函数提供了线索。表 27-1 总结了这些差异，下面则是详细说明。

- 大部分 `exec()` 函数要求提供欲加载新程序的路径名。而 `execlp()` 和 `execvp()` 则允许只提供程序的文件名。系统会在由环境变量 `PATH` 所指定的目录列表中寻找相应的执行文件（稍后将详细解释）。这与 `shell` 对键入命令的搜索方式一致。这些函数名都包含字母 `p`（表示 `PATH`），以示在操作上有所不同。如果文件名中包含“/”，则将其视为相对或绝对路径名，不再使用变量 `PATH` 来搜索文件。
- 函数 `execl()`、`execlp()` 和 `execl()` 要求开发者在调用中以字符串列表形式来指定参数，而不使用数组来描述 `argv` 列表。首个参数对应于新程序 `main()` 函数的 `argv[0]`，因而通常与参数 `filename` 或 `pathname` 的 `basename` 部分相同。必须以 `NULL` 指针来终止参数列表，以便于各调用定位列表的尾部。（上述各原型注释中的 `(char*)NULL` 部分透露了这一要求。至于为何需要对 `NULL` 进行强制类型转换，请参考附录 C。）这些函数的名称都包含字母 `l`（表示 `list`），以示与那些将以 `NULL` 结尾的数组作为参数列表的函数有所区别。后者（`execve()`、`execvp()` 和 `execv()`）名称中则包含字母 `v`（表示 `vector`）。
- 函数 `execve()` 和 `execl()` 则允许开发者通过 `envp` 为新程序显式指定环境变量，其中 `envp` 是一个以 `NULL` 结束的字符串指针数组。这些函数命名均以字母 `e`（`environment`）结尾。其他 `exec()` 函数将使用调用者的当前环境（即 `environ` 中内容）作为新程序的环境。

`glibc 2.11` 曾加入一个非标准函数 `execve(file, argv, envp)`。该函数与 `execvp()` 类似，不过并非通过 `environ` 来取得新程序的环境，而是通过参数 `envp`（类似于函数 `execve()` 和 `execl()`）来指定新环境。

后面几页会演示部分 `exec()` 函数变体的使用。

表 27-1: `exec()` 函数间的差异总结

函 数	对程序文件的描述 (-, p)	对参数的描述 (v, l)	环境变量来源 (e, -)
<code>execve()</code>	路径名	数组	<code>envp</code> 参数
<code>execl()</code>	路径名	列表	<code>envp</code> 参数
<code>execlp()</code>	文件名+PATH	列表	调用者的 <code>environ</code>
<code>execvp()</code>	文件名+PATH	数组	调用者的 <code>environ</code>
<code>execv()</code>	路径名	数组	调用者的 <code>environ</code>
<code>execel()</code>	路径名	列表	调用者的 <code>environ</code>

27.2.1 环境变量 PATH

函数 `execvp()` 和 `execlp()` 允许调用者只提供欲执行程序的文件名。二者均使用环境变量 `PATH` 来搜索文件。`PATH` 的值是一个以冒号 (:) 分隔, 由多个目录名, 也将其称为路径前缀 (path prefixes) 组成的字符串。下例中的 `PATH` 包含 5 个目录:

```
$ echo $PATH
/home/mtk/bin:/usr/local/bin:/usr/bin:/bin:.
```

对于一个登录 shell 而言, 其 `PATH` 值将由系统级和特定用户的 shell 启动脚本来设置。由于子进程继承其父进程的环境变量, shell 执行每个命令时所创建的进程也就继承了 shell 的 `PATH`。

`PATH` 中指定的路径名既可以是绝对路径名 (以/开始), 也可以是相对路径名。对相对路径名的诠释是基于调用进程的当前工作目录 (current working directory)。正如前面例子中所示, 可以以. (点) 来表示当前工作目录。

在 `PATH` 中包含一个长度为 0 的前缀, 也可以用来指定当前工作目录。表示方式有: 连续的冒号、起始冒号或尾部冒号 (例如, `/usr/bin:/bin:.`)。SUSv3 废止了这一技术; 当前工作目录应该用. (点) 来显式指定。

如果没有定义变量 `PATH`, 那么 `execvp()` 和 `execlp()` 会采用默认的路径列表: `./usr/bin:/bin`。

出于安全方面的考虑, 通常会将当前工作目录排除在超级用户 (root) 的 `PATH` 之外。这是为了防止 root 用户发生如下意外情况: 执行当前工作目录下与标准命令同名的程序 (事先由恶意用户故意放置), 或者将常用命令拼错而执行了当前工作目录下的其他程序 (例如, 输入 `sl` 而非 `ls`)。一些 Linux 发行版还将当前工作目录排除在非特权用户的 `PATH` 缺省值之外。这里假定, 在本书展示的所有 shell 会话日志中, 对 `PATH` 的定义均不包含当前工作目录, 而书中示例在执行当前工作目录下的程序时都冠以前缀 `./`, 原因也正在于此。(同时还有一重妙用: 在本书的 shell 会话日志中, 从表现形式上将示例程序与标准命令区分开来。)

函数 `execvp()` 和 `execlp()` 会在 `PATH` 包含的每个目录中搜索文件, 从列表开头的目录开始, 直至成功执行了既定文件。如果不清楚可执行程序的具体位置, 或是不想因硬编码 (hard-code) 而对具体位置产生依赖, 对 `PATH` 环境变量的这种使用方式是非常有效的。

应该避免在设置了 `set-user-ID` 或 `set-group-ID` 的程序中调用 `execvp()` 和 `execlp()`, 至少应当慎用。需要特别谨慎地控制 `PATH` 环境变量, 以防运行恶意程序。在实际操作中, 这意味着应用程序应该使用已知安全的目录列表来覆盖之前定义的任何 `PATH` 值。

程序清单 27-3 提供了一个使用 `execlp()` 的例子。下面的 shell 会话日志则演示了如何通过该程序来调用 `echo` 命令 (`/bin/echo`):

```
$ which echo
/bin/echo
$ ls -l /bin/echo
-rwxr-xr-x  1 root      15428 Mar 19 21:28 /bin/echo
$ echo $PATH
/home/mtk/bin:/usr/local/bin:/usr/bin:/bin    /bin is in PATH
$ ./t_execlp echo
hello world
```

在上例中，程序清单 27-3 程序将字符串 `hello world` 作为第 3 个参数传递给 `execlp()` 调用。接下来，重新对 `PATH` 进行定义，从中移去包含程序 `echo` 的目录 `/bin`:

```
$ PATH=/home/mtk/bin:/usr/local/bin:/usr/bin
$ ./t_execlp echo
ERROR [ENOENT No such file or directory] execlp
$ ./t_execlp /bin/echo
hello world
```

如你所见，当仅向 `execlp()` 提供文件名（即，字符串中不包含斜杠“/”）时，调用会失败。这是因为在 `PATH` 包含的目录列表中无法找到名为 `echo` 的文件。另一方面，当提供了包含一个或多个斜杠的路径名时，`execlp()` 则会忽略 `PATH` 的内容。

程序清单 27-3: 使用 `execlp()` 在 `PATH` 中搜索文件

```
----- procexec/t_execlp.c
#include "tlpi_hdr.h"

int
main(int argc, char *argv[])
{
    if (argc != 2 || strcmp(argv[1], "--help") == 0)
        usageErr("%s pathname\n", argv[0]);

    execlp(argv[1], argv[1], "hello world", (char *) NULL);
    errExit("execlp");    /* If we get here, something went wrong */
}
----- procexec/t_execlp.c
```

27.2.2 将程序参数指定为列表

如果在编程时已知某个 `exec()` 的参数个数，调用 `execle()`、`execlp()` 或者 `execl()` 时就可以将参数作为列表传入。较之于将参数装配于一个 `argv` 向量中，代码要少一些，便于使用。程序清单 27-4 中程序收效与程序清单 27-1 相同，只是调用了 `execle()` 而非 `execve()`。

程序清单 27-4: 使用 `execle()`，将程序参数指定为列表

```
----- procexec/t_execle.c
#include "tlpi_hdr.h"

int
main(int argc, char *argv[])
{
    char *envVec[] = { "GREET=salut", "BYE=adieu", NULL };
    char *filename;
```

```

if (argc != 2 || strcmp(argv[1], "--help") == 0)
    usageErr("%s pathname\n", argv[0]);

filename = strrchr(argv[1], '/');    /* Get basename from argv[1] */
if (filename != NULL)
    filename++;
else
    filename = argv[1];

execl(argv[1], filename, "hello world", (char *) NULL, envVec);
errExit("execl");    /* If we get here, something went wrong */
}

```

procxec/t_execl.c

27.2.3 将调用者的环境传递给新程序

函数 `execlp()`、`execvp()`、`execl()` 和 `execv()` 不允许开发者显式指定环境列表，新程序的环境继承自调用进程（6.7 节）。这一举措的后果可谓是喜忧参半。出于安全方面的考虑，有时希望确保程序在一个已知（安全）的环境列表下运行。38.8 节将对此做深入讨论。

程序清单 27-5 演示了如何运用函数 `execl()` 使新程序继承调用者的环境。对于通过 `fork()` 从 `shell` 处所继承的环境，程序首先用函数 `putenv()` 进行了修改，接着执行 `printenv` 程序来显示环境变量 `USER` 和 `SHELL` 的值。运行程序的输出如下：

```

$ echo $USER $SHELL           Display some of the shell's environment variables
blv /bin/bash
$ ./t_execl
Initial value of USER: blv    Copy of environment was inherited from the shell
britta                       These two lines are displayed by exced printenv
/bin/bash

```

程序清单 27-5：调用函数 `execl()`，将调用者的环境传递给新程序

```

#include <stdlib.h>
#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    printf("Initial value of USER: %s\n", getenv("USER"));
    if (putenv("USER=britta") != 0)
        errExit("putenv");

    execl("/usr/bin/printenv", "printenv", "USER", "SHELL", (char *) NULL);
    errExit("execl");    /* If we get here, something went wrong */
}

```

procxec/t_execl.c

27.2.4 执行由文件描述符指代的程序：fexecve()

`glibc` 自版本 2.3.2 开始提供函数 `fexecve()`，其行为与 `execve()` 类似，只是指定将要执行的程序是以打开文件描述符 `fd` 的方式，而非通过路径名。有些应用程序需要打开某个程序文件，通过执行校验和（`checksum`）来验证文件内容，然后再运行该程序，这一场景就较为适宜使用函数 `fexecve()`。

```
#define _GNU_SOURCE
```



```
#include <unistd.h>

int fexecve(int fd, char *const argv[], char *const envp[]);

Doesn't return on success; returns -1 on error
```

当然，即便没有 `fexecve()` 函数，也可以调用 `open()` 来打开文件，读取并验证其内容，并最终运行。然而，在打开与执行文件之间，存在将该文件替换的可能性（持有打开文件描述符并不能阻止创建同名新文件），最终造成验证者并非执行者的情况。

27.3 解释器脚本

所谓解释器（interpreter），就是能够读取并执行文本格式命令的程序。（相形之下，编译器则是将输入源代码译为可在真实或虚拟机上执行的机器语言。）各种 UNIX shell，以及诸如 `awk`、`sed`、`perl`、`python` 和 `ruby` 之类的程序都属于解释器。除了能够交互式地读取和执行命令之外，解释器通常还具备这样一种能力：从被称为脚本（script）的文本文件中读取和执行命令。

UNIX 内核运行解释器脚本的方式与二进制（binary）程序无异，前提是脚本必须满足下面两点要求：首先，必须赋予脚本文件可执行权限；其次，文件的起始行（initial line）必须指定运行脚本解释器的路径名。格式如下：

```
#! interpreter-path [ optional-arg ]
```

字符 `#!` 必须置于该行起始处，这两个字符串与解释器路径名之间可以以空格分隔。在解释该路径名时不会使用环境变量 `PATH`，因而一般应采用绝对路径。使用相对路径固然可行，但很少见。对其解释则相对于启动解释器进程的当前工作目录。解释器路径名后还可跟随可选参数（稍后将解释其目的），二者之间以空格分隔。可选参数中不应包含空格。

作为例子，UNIX shell 脚本通常以下面这行开始，指定运行该脚本的 shell：

```
#!/bin/sh
```

解释器脚本文件首行中的可选参数不应包含空格，因为空格此处所起的作用完全取决于实现。Linux 系统不会对可选参数（`optional-arg`）中的空格做特殊解释，将从参数起始直至行尾的所有文本视为一个单词（正如后面所述，再将其作为一个参数传递给解释器）。注意，对空格的这种处理方式与 shell 的做法形成鲜明对比，后者总是将其视为命令行中各单词的界定符。

其他 UNIX 实现在处理可选参数中的空格时，其做法与 Linux 有同有异。在 6.0 版本之前的 FreeBSD 上，可在解释器路径（`interpreter-path`）之后跟随多个以空格分隔的可选参数（并作为多个独立的单词传递给解释器）；而到了 6.0 版本，其行为又转而与 Linux 趋同。而 Solaris 8 则使用空格来表征可选参数的结束，同时忽略 `#!` 行中之后的任何剩余文本。

Linux 内核要求脚本的 `#!` 起始行不得超过 127 个字节，其中不包括行尾的换行符（`newline`）。超出部分会被悄无声息地略去。

SUSv3 并未对脚本解释器的 `#!` 行技术加以规范，不过大多数 UNIX 实现都支持这一特性。

不同的 UNIX 实现对于 `#!` 行的长度限制有所不同。例如，OpenBSD 3.1 的限制为 64 个字节，而 Tru64 5.1 则为 1024 字节。在一些早期的实现（例如 SunOS 4）中，这一限制甚至低至 32 字节。

解释器脚本的执行

因为脚本并不包含二进制机器码，所以当调用 `execve()` 来运行脚本时，显然发生了一些不同寻常的事件。`execve()` 如果检测到传入的文件以两字节序列 “#!” 开始，就会析取该行的剩余部分（路径名以及参数），然后按如下参数列表来执行解释器程序：

```
interpreter-path [ optional-arg ] script-path arg...
```

这里，`interpreter-path`（解释器路径）和 `optional-arg`（可选参数）都取自脚本的#!行，`script-path`（脚本路径）是传递给 `execve()` 的路径名，`arg ...` 则是通过变量 `argv` 传递给 `execve()` 的参数列表（不过将 `argv[0]` 排除在外）。图 27-1 对每个脚本参数的起源做了总结。

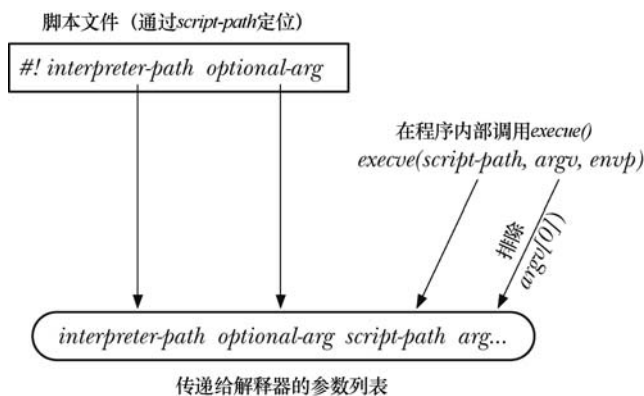


图 27-1：提供给可执行脚本的参数列表

编写一个脚本，用程序清单 6-2 (`necho.c`) 程序作为解释器，用于说明解释器参数的来源。该程序只是简单地输出所有的命令行参数。接着，再使用 27-1 中程序来执行该脚本：

```
$ cat > necho.script          Create script
#!/home/mtk/bin/necho some argument
Some junk
Type Control-D
$ chmod +x necho.script      Make script executable
$ ./t_execve necho.script    And exec the script
argv[0] = /home/mtk/bin/necho First 3 arguments are generated by kernel
argv[1] = some argument      Script argument is treated as a single word
argv[2] = necho.script       This is the script path
argv[3] = hello world        This was argv[1] given to execve()
argv[4] = goodbye            And this was argv[2]
```

在本例中，“解释器（`necho`）”并不关心脚本的内容（`necho.script`），脚本的第 2 行（`Some junk`）在执行时不起作用。

2.2 内核在执行脚本时将只传递 `interpreter-path`（解释器路径）的 `basename` 部分，以作为调用脚本的首个参数。所以，对于 Linux 2.2 来说，`argv[0]` 的输出行会只显示值 `necho`。

大多数 UNIX shell 和解释器会视字符#为注释的开始。因此，这些解释器在解释脚本时会忽略带有#的初始行。

使用脚本的 `optional-arg`（可选参数）

在脚本的#!起始行中，`optional-arg` 的用途之一是为解释器指定命令行参数。对于 `awk` 之

类的解释器而言，这是非常实用的特性。

自 20 世纪 70 年代末期开始，awk 解释器业已成为 UNIX 系统的一部分。在介绍 awk 语言的诸多书籍之中，就有一本[Aho 等，1988]是由该语言的 3 位发明者所著，而该语言的命名也源于 3 人名字的首字母。Awk 的长处在于，能快速为文本处理程序创建原型。作为一门弱类型语言，其设计中富含多种文本处理原素，语法结构则以 C 语言为基础。对于时下风光无限的诸多脚本语言（诸如 JavaScript 和 PHP）而言，awk 的始祖地位毋庸置疑。

向 awk 提供脚本有两种不同方式。默认方式是将脚本作为 awk 的首个命令行参数：

```
$ awk 'script' input-file...
```

也可以将 awk 脚本保存于文件之中，正如下面显示最长输入行长度的例子那样：

```
$ cat longest_line.awk
#!/usr/bin/awk
length > max { max = length; }
END          { print max; }
```

假设使用如下 C 代码来执行这一脚本：

```
execl("longest_line.awk", "longest_line.awk", "input.txt", (char *) NULL);
```

execl()转而调用 execve()，以如下参数列表来运行 awk：

```
/usr/bin/awk longest_line.awk input.txt
```

由于 awk 会把字符串 longest_line.awk 解释为一个包含无效 awk 命令的脚本，故而 execve()调用将以失败告终。这就需要有一种方法来通知 awk：该参数实际上是包含脚本的文件名称。在脚本的 #!起始行中加入 -f 可选参数，就可达到这一目的。这等于告诉 awk，后面的参数是一个脚本文件：

```
#!/usr/bin/awk -f
length > max { max = length; }
END          { print max; }
```

现在，新的 execl()调用会使用如下参数列表：

```
/usr/bin/awk -f longest_line.awk input.txt
```

这样，awk 就可以成功地执行 longest_line.awk 脚本来处理输入文件 input.txt。

使用 execlp()和 execvp()执行脚本

通常，脚本缺少 #!起始行将导致 exec()函数执行失败。不过，execlp()和 execvp()的行事方式多少有些不同。前面提到，这些函数会通过环境变量 PATH 来获取目录列表，并在其中搜索将要执行的文件。两个函数无论谁找到该文件，如果既具有可执行权限，又并非二进制格式，且起始行也不以 #!开始，那么就会使用 shell 来解释这一文件。Linux 中，会将这类文件视同于包含 #!/bin/sh 起始行的文件来进行处理。

27.4 文件描述符与 exec()

默认情况下，由 exec()的调用程序所打开的所有文件描述符在 exec()的执行过程中会保持打开状态，且在新程序中依然有效。这通常很实用，因为调用程序可能会以特定的描述符来打开文件，而在新程序中这些文件将自动有效，无需再去了解文件名或是把它们重新打开。

shell 利用这一特性为其所执行的程序处理 I/O 重定向。例如，假设键入如下的 shell 命令：

```
$ ls /tmp > dir.txt
```

shell 运行该命令时，执行了以下步骤。

1. 调用 `fork()` 创建子进程，子进程也会运行 shell 的一份拷贝（因此命令行也有一份拷贝）。
2. 子 shell 以描述符 1（标准输出）打开文件 `dir.txt` 用于输出。可能会采取以下任一方式。
 - a) 子 shell 关闭描述符 1（`STDOUT_FILENO`）后，随即打开文件 `dir.txt`。因为 `open()` 在为描述符取值时总是取最小值，而标准输入（描述符 0）又仍处于打开状态，所以会以描述符 1 来打开文件。
 - b) shell 打开文件 `dir.txt`，获取一个新的文件描述符。之后，如果该文件描述符不是标准输出，那么 shell 会使用 `dup2()` 强制将标准输出复制为新描述符的副本，并将此时已然无用的新描述符关闭。（这种方法较之前者更为安全，因为它并不依赖于打开文件描述符的低值取数原则。）源代码顺序大体如下：

```
fd = open("dir.txt", O_WRONLY | O_CREAT,  
          S_IRUSR | S_IWUSR | S_IRGRP | S_IWGRP | S_IROTH | S_IWOTH);  
/* rw-rw-rw- */  
if (fd != STDOUT_FILENO) {  
    dup2(fd, STDOUT_FILENO);  
    close(fd);  
}
```
3. 子 shell 执行程序 `ls`。`ls` 将其结果输出到标准输出，亦即文件 `dir.txt` 中。

此处对 shell 处理 I/O 重定向的解释有所简化。特别是，某些命令，即所谓 shell 内建命令，是由 shell 直接运行的，并未调用 `fork()` 或者 `exec()`。在处理 I/O 重定向时，针对这样的命令必须进行特殊处理。

将某一 shell 命令实现为内建命令，不外乎如下两个目的：效率以及会对 shell 产生副作用（side effect）。一些频繁使用的命令（如 `pwd`、`echo` 和 `test`）逻辑都很简单，放在 shell 内部实现效率会更高。将其他命令内置于 shell 实现，则是希望命令对 shell 本身能产生副作用：更改 shell 所存储的信息，修改 shell 进程的属性，亦或是影响 shell 进程的运行。例如，`cd` 命令必须改变 shell 自身的工作目录，故而不应在一个独立进程中执行。产生副作用的内建命令还包括 `exec`、`exit`、`read`、`set`、`source`、`ulimit`、`umask`、`wait` 以及 shell 的作业控制（job-control）命令（`jobs`、`fg` 和 `bg`）。想了解 shell 支持的全套内建命令，可参考 shell 手册页（manual page）文档。

执行时关闭（close-on-exec）标志（`FD_CLOEXEC`）

在执行 `exec()` 之前，程序有时需要确保关闭某些特定的文件描述符。尤其是在特权进程中调用 `exec()` 来启动一个未知程序时（并非自己编写），亦或是启动程序并不需要使用这些已打开的文件描述符时，从安全编程的角度出发，应当在加载新程序之前确保关闭那些不必要的文件描述符。对所有此类描述符施以 `close()` 调用就可达到这一目的，然而这一做法存在如下局限性。

- 某些描述符可能是由库函数打开的。但库函数无法使主程序在执行 `exec()` 之前关闭相应的文件描述符。作为基本原则，库函数应总是为其打开的文件设置执行时关闭（close-on-exec）标志，稍后将介绍所使用的技术。
- 如果 `exec()` 因某种原因而调用失败，可能还需要使描述符保持打开状态。如果这些描述符已然关闭，将它们重新打开并指向相同文件的难度很大，基本上不太可能。

为此，内核为每个文件描述符提供了执行时关闭标志。如果设置了这一标志，那么在成功执行 `exec()` 时，会自动关闭该文件描述符，如果调用 `exec()` 失败，文件描述符则会保持打开状态。可以通过系统调用 `fcntl()`（5.2 节）来访问执行时关闭标志。`fcntl()` 的 `F_GETFD` 操作可

以获取文件描述符标志的一份拷贝：

```
int flags;

flags = fcntl(fd, F_GETFD);
if (flags == -1)
    errExit("fcntl");
```

获取这些标志后，可以对 `FD_CLOEXEC` 位进行修改，再调用 `fcntl()` 的 `F_SETFD` 操作令其生效：

```
flags |= FD_CLOEXEC;
if (fcntl(fd, F_SETFD, flags) == -1)
    errExit("fcntl");
```

实际上 `FD_CLOEXEC` 是文件描述符标志中唯一可以操作的一位。该位对应值为 1。在较老一些的程序中，有时可以看到以 `fcntl(fd, F_SETFD, 1)` 的调用方式来设置执行时关闭标志，其事实依据是这种操作不可能影响到其他位。但理论上，情况不会总是一成不变（未来，一些 UNIX 系统可能会实现其他的标志位），所以还是使用正文所示的技术较为稳妥。

包括 Linux 在内的许多 UNIX 实现，还允许以另外两种非标准的 `ioctl()` 调用来修改执行时关闭标志：以 `ioctl(fd, FIOCLEX)` 为 `fd` 设置此标志，以 `ioctl(fd, FIONCLEX)` 来清除此标志。

当使用 `dup()`、`dup2()` 或 `fcntl()` 为一文件描述符创建副本时，总是会清除副本描述符的执行时关闭标志。（这一现象既有其历史渊源，也顺应了 SUSv3 的要求。）

程序清单 27-6 展示了对执行时关闭标志的操作。如果运行时带有命令行参数（可为任意字符串），该程序首先为标准输出设置执行时关闭标志，随后执行 `ls` 命令。程序运行的结果如下：

```
$ ./closeonexec                               Exec ls without closing standard output
-rwxr-xr-x  1 mtk  users  28098 Jun 15 13:59 closeonexec
$ ./closeonexec n                             Sets close-on-exec flag for standard output
ls: write error: Bad file descriptor
```

上例中第 2 次运行该程序时，`ls` 检测出其标准输出已然关闭，故向标准错误（`stderr`）输出了一条错误信息。

程序清单 27-6：为一文件描述符设置执行时关闭标志

```
----- procexec/closeonexec.c

#include <fcntl.h>
#include "tlpi_hdr.h"

int
main(int argc, char *argv[])
{
    int flags;

    if (argc > 1) {
        flags = fcntl(STDOUT_FILENO, F_GETFD);          /* Fetch flags */
        if (flags == -1)
            errExit("fcntl - F_GETFD");

        flags |= FD_CLOEXEC;                            /* Turn on FD_CLOEXEC */

        if (fcntl(STDOUT_FILENO, F_SETFD, flags) == -1) /* Update flags */
            errExit("fcntl - F_SETFD");
    }
}
```

```
    execlp("ls", "ls", "-l", argv[0], (char *) NULL);
    errExit("execlp");
}
```

procexec/closeonexec.c

27.5 信号与 exec()

exec()在执行时会把现有进程的文本段丢弃。该文本段可能包含了由调用进程创建的信号处理器程序。既然处理器已经不知所踪，内核就会将对所有已设信号的处置重置为 SIG_DFL。而对所有其他信号（即将处置置为 SIG_IGN 或 SIG_DFL 的信号）的处置则保持不变。这也符合 SUSv3 的要求。

不过，遭忽略的 SIGCHLD 信号属于 SUSv3 中的特例。（之前曾在 26.3.3 节提及，忽略 SIGCHLD 能够阻止僵尸进程的产生）。至于调用 exec()之后，是继续让遭忽略的 SIGCHLD 信号保持被忽略状态，还是将对其处置重置为 SIG_DFL，SUSv3 对此不置可否。Linux 的操作取其前者，而其他一些 UNIX 实现（如：Solaris）则采用后者。这就意味着，对于忽略 SIGCHLD 的程序而言，要最大限度的保证可移植性，就应该在调用 exec()之前执行 signal (SIGCHLD, SIG_DFL)。此外，程序也不应当假设对 SIGCHLD 处置的初始设置是 SIG_DFL 之外的其他值。

老程序的数据段、堆以及栈悉数被毁，这也意味着通过 sigaltstack() (21.3 节) 所创建的任何备选信号栈都会丢失。由于 exec()在调用期间不会保护备选信号栈，故而也会将所有信号的 SA_ONSTACK 位清除掉。

在调用 exec()期间，进程信号掩码以及挂起 (pending) 信号的设置均得以保存。这一特性允许对新程序的信号进行阻塞和排队处理。不过，SUSv3 指出，许多现有应用程序的编写都基于如下的错误假设：程序启动时将对某些特定信号的处置置为 SIG_DFL，又或者并未阻塞这些信号。（特别是，C 语言标准对信号的规范很弱，对信号阻塞也未置一词，所以为非 UNIX 系统所编写的 C 程序也不可能去解除对信号的阻塞。）为此，SUSv3 建议，在调用 exec()执行任何程序的过程中，不应当阻塞或忽略信号。这里的“任何程序”是指并非由 exec()的调用者所编写的程序。至于说如果执行和被执行的程序均出自一人之手，又或者对运行程序处理信号的手法知根知底，那自然又另当别论。

27.6 执行 shell 命令：system()

程序可通过调用 system()函数来执行任意的 shell 命令。本节将讨论 system()的操作，下一节将介绍如何运用 fork()、exec()、wait()和 exit()来实现 system()。

44.5 节所介绍的 popen()和 pclose()函数同样可以用来执行 shell 命令，而且还允许调用程序向命令发送输入信息，或是读取命令的输出。

```
#include <stdlib.h>
```

```
int system(const char *command);
```

See main text for a description of return value

函数 system()创建一个子进程来运行 shell，并以之执行命令 command。其调用示例如下：

```
system("ls | wc");
```

system()的主要优点在于简便。

- 无需处理对 fork()、exec()、wait()和 exit()的调用细节。
- system()会代为处理错误和信号。
- 因为 system()使用 shell 来执行命令 (command)，所以会在执行 command 之前对其进行所有的常规 shell 处理、替换以及重定向操作。为应用增加“执行一条 shell 命令”的功能不过是举手之劳。(许多交互式应用程序以“! command”的形式提供了这一功能。)

但这些优点是以低效率为代价的。使用 system()运行命令需要创建至少两个进程。一个用于运行 shell，另外一个或多个则用于 shell 所执行的命令 (执行每个命令都会调用一次 exec())。

如果对效率或者速度有所要求，最好还是直接调用 fork()和 exec()来执行既定程序。

system()的返回值如下。

- 当 command 为 NULL 指针时，如果 shell 可用则 system()返回非 0 值，若不可用则返回 0。这种返回值方式源于 C 语言标准，因为并未与任何操作系统绑定，所以如果 system()运行在非 UNIX 系统上，那么该系统可能是没有 shell 的。此外，即便所有 UNIX 实现都有 shell，如果程序在调用 system()之前又调用了 chroot()，那么 shell 依然可能无效。若 command 不为 NULL，则 system()的返回值由本列表中的余下规则决定。
- 如果无法创建子进程或是无法获取其终止状态，那么 system()返回-1。
- 若子进程不能执行 shell，则 system()的返回值会与子 shell 调用 _exit(127)终止时一样。
- 如果所有的系统调用都成功，system()会返回执行 command 的子 shell 的终止状态。shell 的终止状态是其执行最后一条命令时的退出状态；如果命令为信号所杀，大多数 shell 会以值 128+n 退出，其中 n 为信号编号 (如果是子 shell 为信号所杀，那么其终止状态如 26.1.3 节所述)。

至于调用失败是由于 system()无法执行 shell，还是 shell 以状态 127 退出 (若 shell 未能发现并执行既定名称的程序，就会导致后一种情况的发生)，(通过 system()的返回值)是无法区分的。

在最后两种情况中，system()的返回值与 waitpid()所返回的等待状态 (wait status) 形式相同。因此，可以使用 26.1.3 节所述函数来分析返回值，并以 printWaitStatus()函数 (见程序清单 26-2) 加以显示。

示例程序

程序清单 27-7 演示了 system()的用法。程序循环读取命令字符串，再使用 system()来执行命令，并对 system()的返回值进行分析和展示。下面是一个运行的例子：

```
$ ./t_system
Command: whoami
mtk
system() returned: status=0x0000 (0,0)
child exited, status=0
Command: ls | grep XYZ
system() returned: status=0x0100 (1,0)
child exited, status=1
Command: exit 127
system() returned: status=0x7f00 (127,0)
(Probably) could not invoke shell
Command: sleep 100
Type Control-Z to suspend foreground process group
```

*Shell terminates with the status of...
its last command (grep), which...
found no match, and so did an exit(1)*

Actually, not true in this case

```

[1]+  Stopped                  ./t_system
$ ps | grep sleep              Find PID of sleep
29361 pts/6    00:00:00 sleep
$ kill 29361                  And send a signal to terminate it
$ fg                           Bring t_system back into foreground
./t_system
system() returned: status=0x000f (0,15)
child killed by signal 15 (Terminated)
Command: ^D$                  Type Control-D to terminate program

```

程序清单 27-7：通过 system() 执行 shell 命令

```

procexec/t_system.c

#include <sys/wait.h>
#include "print_wait_status.h"
#include "tlpi_hdr.h"

#define MAX_CMD_LEN 200

int
main(int argc, char *argv[])
{
    char str[MAX_CMD_LEN];      /* Command to be executed by system() */
    int status;                /* Status return from system() */

    for (;;) {                 /* Read and execute a shell command */
        printf("Command: ");
        fflush(stdout);
        if (fgets(str, MAX_CMD_LEN, stdin) == NULL)
            break;             /* end-of-file */

        status = system(str);
        printf("system() returned: status=0x%04x (%d,%d)\n",
              (unsigned int) status, status >> 8, status & 0xff);

        if (status == -1) {
            errExit("system");
        } else {
            if (WIFEXITED(status) && WEXITSTATUS(status) == 127)
                printf("(Probably) could not invoke shell\n");
            else
                /* Shell successfully executed command */
                printWaitStatus(NULL, status);
        }
    }

    exit(EXIT_SUCCESS);
}

procexec/t_system.c

```

在设置用户 ID (set-user-ID) 和组 ID (set-group-ID) 程序中避免使用 system()

当设置了用户 ID 和组 ID 的程序在特权模式下运行时，绝不能调用 system()。即便此类程序并未允许用户指定需要执行的命令文本，鉴于 shell 对操作的控制有赖于各种环境变量，故而使用 system() 会不可避免地给系统带来安全隐患。

例如，在较老的 Bourne shell 中，环境变量 IFS（定义了用于将命令行拆分为独立单词的内部字段分隔符）就引发了若干起对系统入侵的成功案例。如果将 IFS 定义为 a，那么 shell 会将字

字符串 `shar` 视为带有参数 `r` 的单词 `sh`，并启动另一 `shell` 进程来执行当前工作目录下名为 `r` 的脚本，这就一改命令的原意（执行名为 `shar` 的命令）。对这一安全漏洞的修复举措是，将 `IFS` 只应用于 `shell` 扩展所产生的单词。此外，现代 `shell` 会在启动时重置 `IFS`（为由空格、`Tab` 以及换行 3 个字符组成的字符串），以确保即使 `IFS` 的继承值很奇怪，脚本的行为也会保持一致。作为进一步的安全举措，当从设置用户（组）`ID` 程序中调用时，`bash` 会回转为实际用户（组）`ID` 的“真身”。

应用需要加载其他程序时，为确保安全过关，应当直接调用 `fork()` 和 `exec()` 系函数（`execlp()` 和 `execvp()` 除外）之一。

27.7 system()的实现

本节将说明如何实现 `system()` 的功能。首先给出一个简化版实现，接着指出这一实现的缺失所在，最后再展示了一个完整的实现。

对 system()的简化实现

命令 `sh` 的参数 `-c` 提供了一种简单的方法，可以执行包含任意命令的字符串：

```
$ sh -c "ls | wc"
    38    38    444
```

因此，为了实现 `system()`，需要使用 `fork()` 来创建一个子进程，并以对应于上例 `sh` 命令的参数来调用 `execl()`：

```
execl("/bin/sh", "sh", "-c", command, (char *) NULL);
```

为了收集 `system()` 所创建的子进程状态，还以指定的子进程 `ID` 调用了 `waitpid()`。（使用 `wait()` 并不合适，因为 `wait()` 等待的是任一子进程，因而无意间所获取的子进程状态可能属于其他子进程。）程序清单 27-8 是对 `system()` 的简化实现。

程序清单 27-8：一个缺乏信号处理的 `system()` 实现

```
----- procexec/simple_system.c
#include <unistd.h>
#include <sys/wait.h>
#include <sys/types.h>

int
system(char *command)
{
    int status;
    pid_t childPid;

    switch (childPid = fork()) {
    case -1: /* Error */
        return -1;
    case 0: /* Child */
        execl("/bin/sh", "sh", "-c", command, (char *) NULL);
        _exit(127); /* Failed exec */
    default: /* Parent */
        if (waitpid(childPid, &status, 0) == -1)
            return -1;
        else
            return status;
    }
}
```

}

在 system() 内部正确处理信号

给 system() 的实现带来复杂性的是对信号的正确处理。

首先需要考虑的信号是 SIGCHLD。假设调用 system() 的程序还直接创建其他子进程，对 SIGCHLD 的信号处理器自身也执行了 wait()。在这种情况下，当由 system() 所创建的子进程退出并产生 SIGCHLD 信号时，在 system() 有机会调用 waitpid() 之前，主程序的信号处理器程序可能会率先得以执行（收集子进程的状态）。这是竞争条件（race condition）的又一例证。这会产生两种不良后果。

- 调用程序会误以为其所创建的某个子进程终止了。
- system() 函数却无法获取其所创建子进程的终止状态。

所以，system() 在运行期间必须阻塞 SIGCHLD 信号。

其他需要关注的信号则是分别由终端的中断（interrupt）（通常为 Ctrl-C）和退出（quit）（通常为 Ctrl-\）符所产生的 SIGINT 和 SIGQUIT 信号。考虑执行如下调用的后果：

```
system("sleep 20");
```

此时此刻，会有 3 个进程在运行：执行调用程序的进程、一个 shell 进程，以及 sleep 进程。如图 27-2 所示。

为提高效率，如果赋予 -c 选项的是一条简单命令，较之于管道（pipeline）或序列（sequence），一些 shell（包括 bash）会直接执行该命令，而不会再去创建一个子 shell。对于采用此类优化的 shell 而言，因为只有两个进程（调用进程和 sleep 进程），所以图 27-2 有失准确。不过，本节关于 system() 如何处理信号的论述仍然适用。

图 27-2 中所示的所有进程构成终端前台进程组的一部分。（34.2 节将详细讨论进程组。）所以，在输入中断或退出符时，会将相应信号发送给所有 3 个进程。shell 在等待子进程期间会忽略 SIGINT 和 SIGQUIT 信号。不过，默认情况下这些信号会杀死调用程序与 sleep 进程。

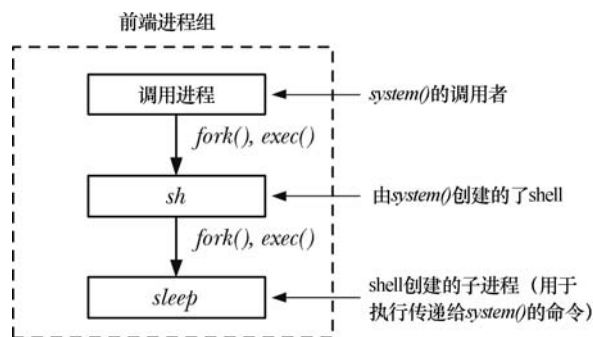


图 27-2: 执行 system("sleep 20") 期间的进程情况

调用进程和所执行的命令应当如何应对这些信号呢？SUSv3 的规定如下。

- 调用进程在执行命令期间应忽略 SIGINT 和 SIGQUIT 信号。
- 子进程对上述两信号的处理，如同调用进程调用了 fork() 和 exec() 一般，也就是说，将对已处理信号的处置重置为默认值，而对其他信号的处置则保持不变。

按照 SUSv3 所规范的方式来处理信号是最为合理的，其原因如下。

- 让所有的进程都对这些信号做出响应是没有意义的，用户会对应用的行为困惑不已。
- 与上述相类似，一面在执行命令的进程中忽略这些信号，而同时又在调用进程中按对它们的缺省处置来行事，这同样也说不通。用户借此可以将调用进程杀掉，同时放任其所执行的命令继续运行。这与实际情况也并不相符：当命令传递给 `system()` 执行时，调用进程实际上已经放弃了控制权（即阻塞于 `waitpid()` 调用中）。
- `system()` 运行的可能是一个交互式应用，让此类应用响应终端产生的信号是合理的。

SUSv3 要求按上述方式来处理 `SIGINT` 和 `SIGQUIT`，但同时指出，对于暗中调用 `system()` 来执行任务的程序，这一做法可能会产生不良后果。执行命令时如按下 `Ctrl-C` 或 `Ctrl-\`，将只会杀掉 `system()` 的子进程，而应用程序会继续运行（用户并不希望如此）。以此方式调用 `system()` 的程序应当检查 `system()` 所返回的终止状态，一旦发现命令因信号而终止，应采取相应措施。

system()实现的改进版

程序清单 27-9 所示为遵循上述规则的 `system()` 实现。关于该实现，需注意以下几点。

- 如前所述，当 `command` 为 `NULL` 指针时，若 `shell` 可用，则 `system()` 应返回非 0 值；如不可用，则返回 0。要得出结论，唯一可靠的办法就是尝试运行 `shell`。程序这里的做法是：递归调用 `system()` 去运行 `shell` 命令“:”并检查该递归调用的返回状态是否为 0 ①。“:”是一个 `shell` 内建命令，无所作为却总是返回成功。执行命令 `exit 0` 也可获得相同效果。（仅仅通过 `access()` 来判断文件 `/bin/sh` 存在与否，是否具有可执行权限的做法存在局限性。在 `chroot()` 环境中，即使具有可执行权限的 `shell` 文件存在，如果与其进行动态链接的共享库无效，依然无法执行 `shell`。）
- 只有父进程（`system()` 的调用者）才需要阻塞 `SIGCHLD` ②，同时还需要忽略 `SIGINT` 和 `SIGQUIT` ③。不过，必须在调用 `fork()` 之前执行这些动作，因为如果在父进程的 `fork()` 之后执行，将出现竞争条件。（假设：如果在父进程有机会阻塞 `SIGCHLD` 之前，子进程就退出了。）结果是，如同稍后所述，子进程必须取消对信号属性的这些变更。
- 父进程并未对 `sigaction()` 以及 `sigprocmask()` 调用进行错误检查 ②③⑨，这两者分别用于操作对信号的处置和信号掩码。这样做原因有二。其一，这些调用失手的可能性不大。实际上，只有指定参数有误时才会失败，而只要一开始调试就能搞定此类问题。其二，这里假定，较之于此类信号操控函数，调用者更关注 `fork()` 或 `waitpid()` 的成败与否。同理，在 `system()` 尾部的信号处理操作前后，分别有代码来保存和恢复 `errno`，以便一旦 `fork()` 或 `waitpid()` 失败，调用者能查明原因。如果因信号操作失败而返回 -1，那么调用者会误认为是 `system()` 执行 `command` 失败所致。

SUSv3 仅仅指出在创建子进程失败或无法获取子进程状态时，`system()` 返回 -1。并未提及 `system()` 在处理信号失败时也会返回 -1。

- 子进程中对于信号相关的系统调用也未执行错误检查 ④⑤。一方面，无法报告此类错误（`_exit(127)` 是预留给执行 `shell` 时报告错误之用）；另一方面，这一失败也不会殃及 `system()` 的调用者，二者分属于不同进程。
- 子进程刚从 `fork()` 返回时，会将 `SIGINT` 和 `SIGQUIT` 的处处置为 `SIG_IGN`（继承自父进程）。不过，如前所述，子进程处理这些信号时就如同 `system()` 的调用者执行了 `fork()` 和 `exec()`。`fork()` 不会改变子进程对这些信号的处理方式。而 `exec()` 则会将对

已处理信号的处置重置为默认值，但不改变对其他信号的处置（27.5 节）。因此，如果调用者对 SIGINT 和 SIGQUIT 的处置设置并非 SIG_IGN，那么子进程会将其置为 SIG_DFL^④。

一些 system()实现反而会将于进程对 SIGINT 和 SIGQUIT 的处置重置为在调用者中生效的设置。这一做法的依据是，后续对 execl()的调用会自动将对这些已处理信号的处置重置为默认值。不过，如果调用者正在处理两个信号之一时，这可能会导致不希望的行为发生。在这种情况下，如果在调用 execl()之前的瞬间有信号送达子进程，那么在信号经由 sigprocmask()解除阻塞后，子进程还是会调用信号处理器程序。

- 子进程如果调用 execl()失败，就会以 _exit()，而非 exit()来终止进程^⑤。这是为了防止对子进程 stdio 缓冲区中的任何未写入数据进行刷新。
- 父进程必须使用 waitpid()来专候其所创建的特定子进程^⑦。如果使用 wait()，不经意间可能会捕获到其他子进程的状态。
- 虽然 system()实现并未强制要求使用信号处理器程序，但调用程序还是可能会去创建它们，从而中断对 waitpid()的阻塞调用。SUSv3 明确要求在这种情况下必须重新等待。所以，如果发生 EINTR 错误^⑦，则循环调用 waitpid()以期成功重启。如果是其他错误，则退出 waitpid()循环。

程序清单 27-9: system()的实现

```
procexec/system.c

#include <unistd.h>
#include <signal.h>
#include <sys/wait.h>
#include <sys/types.h>
#include <errno.h>

int
system(const char *command)
{
    sigset_t blockMask, origMask;
    struct sigaction saIgnore, saOrigQuit, saOrigInt, saDefault;
    pid_t childPid;
    int status, savedErrno;

    ① if (command == NULL)                /* Is a shell available? */
        return system(":") == 0;
    sigemptyset(&blockMask);            /* Block SIGCHLD */
    sigaddset(&blockMask, SIGCHLD);
    ② sigprocmask(SIG_BLOCK, &blockMask, &origMask);

    saIgnore.sa_handler = SIG_IGN;      /* Ignore SIGINT and SIGQUIT */
    saIgnore.sa_flags = 0;
    sigemptyset(&saIgnore.sa_mask);
    ③ sigaction(SIGINT, &saIgnore, &saOrigInt);
    sigaction(SIGQUIT, &saIgnore, &saOrigQuit);

    switch (childPid = fork()) {
    case -1: /* fork() failed */
        status = -1;
        break;                            /* Carry on to reset signal attributes */
```

```

case 0: /* Child: exec command */
    saDefault.sa_handler = SIG_DFL;
    saDefault.sa_flags = 0;
    sigemptyset(&saDefault.sa_mask);

④    if (saOrigInt.sa_handler != SIG_IGN)
        sigaction(SIGINT, &saDefault, NULL);
    if (saOrigQuit.sa_handler != SIG_IGN)
        sigaction(SIGQUIT, &saDefault, NULL);

⑤    sigprocmask(SIG_SETMASK, &origMask, NULL);

    execl("/bin/sh", "sh", "-c", command, (char *) NULL);
⑥    _exit(127); /* We could not exec the shell */

default: /* Parent: wait for our child to terminate */
⑦    while (waitpid(childPid, &status, 0) == -1) {
        if (errno != EINTR) { /* Error other than EINTR */
            status = -1;
            break; /* So exit loop */
        }
    }
    break;
}

/* Unblock SIGCHLD, restore dispositions of SIGINT and SIGQUIT */

⑧    savedErrno = errno; /* The following may change 'errno' */

⑨    sigprocmask(SIG_SETMASK, &origMask, NULL);
    sigaction(SIGINT, &saOrigInt, NULL);
    sigaction(SIGQUIT, &saOrigQuit, NULL);

⑩    errno = savedErrno;

    return status;
}

```

procexec/system.c

关于 system()的更多细节

为保障应用程序的可移植性，应确保在将对 SIGCHLD 的处置置为 SIG_IGN 的情况下不去调用 system()，因为此时 waitpid()将无法获取子进程的状态。（忽略 SIGCHLD 会导致立即丢弃子进程状态，如 26.3.3 节所述。）

在一些 UNIX 实现中，如果在将对 SIGCHLD 的处置置为 SIG_IGN 的情况下调用 system()，system()的应对策略是：临时将其改为 SIG_DFL。只要在把对 SIGCHLD 的处置重置为 SIG_IGN 时，UNIX 实现能够处理僵尸子进程（Linux 不在此列），这种方法就是可行的。（如果系统做不到这一点，按此方式实现 system()将产生如下不良后果：在调用者执行 system()期间，如果另一子进程终止了，那么该子进程将成为僵尸进程，且无法回收。）

对于一些 UNIX 实现（尤其是 Solaris），/bin/sh 并非标准的 POSIX shell。若希望确保执行标准 shell，则必须使用库函数 confstr()来获取配置变量 _CS_PATH 的值。该值的风格与 PATH 相同，包含了标准系统工具的目录列表。可以将该列表赋给变量 PATH，随即调用 execlp()以执行标准 shell，具体如下：

```

char path[PATH_MAX];

if (confstr(_CS_PATH, path, PATH_MAX) == 0)
    _exit(127);
if (setenv("PATH", path, 1) == -1)
    _exit(127);
execlp("sh", "sh", "-c", command, (char *) NULL);
_exit(127);

```

27.8 总结

进程可使用 `execve()` 以一新程序替换当前增长运行的程序。`execve()` 的参数允许为新程序指定参数列表 (`argv`) 和环境列表。构建于 `execve()` 之上, 存在多种命名相似的函数, 功能相同, 但提供的接口不同。

所有的 `exec()` 函数均可用于加载二进制的可执行文件或是执行解释器脚本。当进程执行脚本时, 脚本解释器程序将替换进程当前执行的程序。脚本的起始行 (以 `#!` 开头) 指定了解释器的路径名, 供识别解释器之用。如果没有这一起始行, 那么只能通过 `execlp()` 或 `execvp()` 来执行脚本, 并默认把 `shell` 作为脚本解释器。

本章还展示了如何组合使用 `fork()`、`exec()`、`exit()` 和 `wait()` 来实现 `system()` 函数, 该函数可用于执行任意 `shell` 命令。

更多的信息

请参考 24.6 节所列的更多信息来源。

27.9 练习

27-1. 如下 `shell` 会话的最后一条命令使用程序清单 27-3 程序来执行程序 `xyz`。结果如何?

```

$ echo $PATH
/usr/local/bin:/usr/bin:/bin:./dir1:./dir2
$ ls -l dir1
total 8
-rw-r--r--  1 mtk      users      7860 Jun 13 11:55 xyz
$ ls -l dir2
total 28
-rwxr-xr-x  1 mtk      users      27452 Jun 13 11:55 xyz
$ ./t_execlp xyz

```

27-2. 试用 `execve()` 实现 `execlp()`。需使用 `stdarg(3)` API 来处理 `execlp()` 所提供的变长参数列表。还需要使用 `malloc` 函数库中函数为参数以及环境向量分配空间。最后, 请注意, 要检查特定目录下某个文件是否存在且可以执行, 有一种简单方法: 尝试执行该文件即可。

27-3. 如果赋予如下脚本可执行权限并以 `exec()` 运行, 输出结果如何?

```

#!/bin/cat -n
Hello world

```

27-4. 下列代码会有什么效果? 在何种情况下会起作用?

```

childPid = fork();
if (childPid == -1)
    errExit("fork1");

```

```

if (childPid == 0) { /* Child */
    switch (fork()) {
        case -1: errExit("fork2");

        case 0: /* Grandchild */
            /* ----- Do real work here ----- */
            exit(EXIT_SUCCESS); /* After doing real work */

        default:
            exit(EXIT_SUCCESS); /* Make grandchild an orphan */
    }
}

/* Parent falls through to here */

if (waitpid(childPid, &status, 0) == -1)
    errExit("waitpid");

/* Parent carries on to do other things */

```

27-5. 运行如下程序时无输出。试问原因何在？

```

#include "tspi_hdr.h"

int
main(int argc, char *argv[])
{
    printf("Hello world");
    execlp("sleep", "sleep", "0", (char *) NULL);
}

```

27-6. 假设父进程为信号 SIGCHLD 创建了一处理器程序，同时阻塞该信号。随后，其某一子进程退出，父进程接着执行 wait() 以获取该子进程的状态。当父进程解除对 SIGCHLD 的阻塞时，会发生什么？编写一个程序来验证答案。这一结果与调用 system() 函数的程序之间有什么关联？

第 28 章

详述进程创建和程序执行

本章对第 24 章到第 27 章的内容进行了拓展，涵盖进程创建和程序执行的多个主题。首先是进程记账（process accounting），这一内核特性会使系统在每个进程结束后记录一条账单信息。接着，会讨论 Linux 特有的系统调用 `clone()`，Linux 系统创建线程就有赖于这一底层 API。然后对 `fork()`、`vfork()` 和 `clone()` 的性能进行了比较。最后，本章就 `fork()` 和 `exec()` 对进程属性的影响做了总结。

28.1 进程记账

打开进程记账功能后，内核会在每个进程终止时将一条记账信息写入系统级的进程记账文件。这条账单记录包含了内核为该进程所维护的多种信息，包括终止状态以及进程消耗的 CPU 时间。借助于标准工具（`sa(8)` 对账单文件进行汇总，`lastcomm(1)` 则就先前执行的命令列出相关信息）或是定制应用，可对记账文件进行分析。

内核 2.6.10 之前，内核会为基于 NPTL 线程实现所创建的每个线程单独记录一条进程记账信息。自内核 2.6.10 开始，只有当最后一个线程退出时才会为整个进程保存一条账单记录。至于更老的 LinuxThread 线程实现，则会为每个线程单独记录一条进程记账信息。

从历史上看，进程记账主要用于在多用户 UNIX 系统上针对用户所消耗的系统资源进行计费。不过，如果进程的信息并未由其父进程进行监控和报告，那么就可以使用进程记账来获取。

虽然大部分 UNIX 实现都支持进程记账功能，但 SUSv3 并未对其进行规范。账单记录的格式、记账文件的位置也随系统实现的不同而多少存在差别。本节所述是针对 Linux 系统的细节，但会在论述过程中点出其与其他 Unix 系统的差异。

Linux 系统的进程记账功能属于可选内核组件，可以通过 `CONFIG_BSD_PROCESS_ACCT` 选项进行配置。

打开和关闭进程记账功能

特权进程可利用系统调用 `acct()` 来打开和关闭进程记账功能。应用程序很少使用这一系统调用。一般会将相应命令置于系统启动脚本中，在系统每次重启时开启进程记账功能。

```
#define _BSD_SOURCE
#include <unistd.h>

int acct(const char *acctfile);

Returns 0 on success, or -1 on error
```

为了打开进程账单功能，需要在参数 `acctfile` 中指定一个现有常规文件的路径名。记账文件通常的路径名是 `/var/log/pacct` 或 `/usr/account/pacct`。若想关闭进程记账功能，则指定 `acctfile` 为 `NULL` 即可。

程序清单 28-1 中程序使用 `acct()` 来开关进程的记账功能。该程序的作用类似于 `shell` 命令 `accton(8)`。

程序清单 28-1：打开和关闭进程记账功能

```
----- procexec/acct_on.c
#define _BSD_SOURCE
#include <unistd.h>
#include "tlpi_hdr.h"

int
main(int argc, char *argv[])
{
    if (argc > 2 || (argc > 1 && strcmp(argv[1], "--help") == 0))
        usageErr("%s [file]\n");
    if (acct(argv[1]) == -1)
        errExit("acct");

    printf("Process accounting %s\n",
           (argv[1] == NULL) ? "disabled" : "enabled");
    exit(EXIT_SUCCESS);
}
----- procexec/acct_on.c
```

进程账单记录

一旦打开进程记账功能，每当一进程终止时，就会有一条 `acct` 记录写入记账文件。`acct` 结构定义于头文件 `<sys/acct.h>` 中，具体如下：

```
typedef u_int16_t comp_t; /* See text */

struct acct {
    char    ac_flag; /* Accounting flags (see text) */
    u_int16_t ac_uid; /* User ID of process */
    u_int16_t ac_gid; /* Group ID of process */
    u_int16_t ac_tty; /* Controlling terminal for process (may be
                     0 if none, e.g., for a daemon) */
    u_int32_t ac_btime; /* Start time (time_t; seconds since the Epoch) */
    comp_t    ac_utime; /* User CPU time (clock ticks) */
    comp_t    ac_stime; /* System CPU time (clock ticks) */
};
```

```

comp_t  ac_etime;    /* Elapsed (real) time (clock ticks) */
comp_t  ac_mem;     /* Average memory usage (kilobytes) */
comp_t  ac_io;      /* Bytes transferred by read(2) and write(2)
                    (unused) */
comp_t  ac_rw;      /* Blocks read/written (unused) */
comp_t  ac_minflt; /* Minor page faults (Linux-specific) */
comp_t  ac_majflt; /* Major page faults (Linux-specific) */
comp_t  ac_swaps;   /* Number of swaps (unused; Linux-specific) */
u_int32_t ac_exitcode; /* Process termination status */
#define ACCT_COMM 16
char     ac_comm[ACCT_COMM+1];
                    /* (Null-terminated) command name
                    (basename of last execed file) */
char     ac_pad[10]; /* Padding (reserved for future use) */
};

```

关于 acct 结构需要注意以下几点。

- 数据类型 `u_int16_t` 和 `u_int32_t` 分别是 16 位和 32 位的无符号整数类型。
- `ac_flag` 字段 (field) 是为进程记录多种事件 (event) 的位掩码 (bit mask)。表 28-1 展示了在该字段中可能出现的位。如表中所示,并非所有的 UNIX 实现都支持这些位。另有少数实现为该字段还提供了一些附加的位。
- `ac_comm` 字段记录了该进程最后执行的命令 (程序文件) 名称。内核会在每次调用 `execve()` 时记录该值。一些 UNIX 实现将该字段的大小限制在 8 个字节以内。
- 类型 `comp_t` 是一种浮点型 (floating-point) 数字。有时也将该类型值称为压缩时钟周期 (compressed clock tick)。该浮点值由 3 位 (bit) 以 8 为底的指数以及 13 位 (bit) 小数组成,指数用来表示值范围在 $8^0=1\sim 8^7$ (2097152) 之间的因子。例如,尾数为 125,指数部分为 1 就表示值为 1000。程序清单 28-2 中定义的函数 (`comptToLL()`) 可以将该类型转换为 `long long`。因为在 x86-32 架构下的系统中,用于表示无符号长整型的 32 位数并不足以保存 `comp_t` 型的最大值: $(2^{13} - 1) \times 8^7$ 。
- 3 个定义为 `comp_t` 型的时间字段其度量单位为系统时钟周期。要将它们转换成秒,必须除以 `sysconf(_SC_CLK_TCK)` 的返回值。
- `ac_exitcode` 字段保存着进程的退出状态 (如 26.1.3 节所述)。其他大多数 UNIX 实现则提供了一个名为 `ac_stat` 的单字节字段来代替 `ac_exitcode`,其中仅记录了杀死进程的信号值 (如果进程为信号所杀) 和一个标志位,用于标识是否因该信号而导致进程转储核心 (dump core)。二者在源于 BSD 的实现中均未提供。

表 28-1: 进程账单记录中 `ac_flag` 字段各位的值

位	说 明
AFORK	由 <code>fork()</code> 创建的进程,终止前并未调用 <code>exec()</code>
ASU	拥有超级用户特权的进程
AXSIG	进程因信号而终止 (有些实现未支持)
ACORE	进程产生了核心转储 (有些实现未支持)

因为只在进程终止时才记录账单信息,所以对这些记录的排序也是按照进程的终止时间

(并未写入记录), 而非启动时间 (ac_btime)。

如果系统崩溃, 也不会为当前运行的进程记录任何记账信息。

由于向记账文件中写入信息可能会加速对磁盘空间的消耗, 为了对进程记账行为加以控制, Linux 系统提供了名为 /proc/sys/kernel/acct 的虚拟文件。此文件包含 3 个值, 按顺序分别定义了如下参数: 高水位 (high-water)、低水位 (low-water) 和频率 (frequency)。3 个参数通常的默认值为 4、2 和 30。如果开启进程记账特性且磁盘空闲空间低于低水位 (low-water) 百分比, 将暂停记账。如果磁盘空闲空间升高水位百分比之上, 则恢复记账。频率值则规定了两次检查空闲磁盘空间占比之间的间隔时间 (以秒为单位)。

示例程序

程序清单 28-2 中程序显示了某进程记账文件记录中特定字段的信息。以下 shell 会话演示了对该程序的使用。首先新建一个空的进程记账文件, 同时开启进程记账功能。

```
$ su Need privilege to enable process accounting
Password:
# touch pacct
# ./acct_on pacct This process will be first entry in accounting file
Process accounting enabled
# exit Cease being superuser
```

从开启进程记账功能到现在, 已经有 3 个进程退出, 分别执行了 acct_on、su 和 bash 程序。进程 bash 由 su 启动, 负责运行特权级 shell 会话。

接着运行一系列命令, 借此向记账文件加入更多记录:

```
$ sleep 15 &
[1] 18063
$ ulimit -c unlimited Allow core dumps (shell built-in)
$ cat Create a process
Type Control-\ (generates SIGQUIT, signal 3) to kill cat process
Quit (core dumped)
$
Press Enter to see shell notification of completion of sleep before next shell prompt
[1]+ Done sleep 15
$ grep xxx badfile grep fails with status of 2
grep: badfile: No such file or directory
$ echo $? The shell obtained status of grep (shell built-in)
2
```

下面两个命令执行的是前面章节中展示过的两个程序 (程序清单 27-1 和程序清单 24-1)。第一条命令运行的程序执行了 /bin/echo, 因此, 写入账单记录中的命令名是 echo。第二条命令创建了一个子进程, 该子进程并未调用 exec()。

```
$ ./t_execve /bin/echo
hello world goodbye
$ ./t_fork
PID=18350 (child) idata=333 istack=666
PID=18349 (parent) idata=111 istack=222
```

最后, 运行程序清单 28-2 中程序来查看记账文件的内容。

```
$ ./acct_view pacct
command flags term. user start time CPU elapsed
          status
acct_on -S-- 0 root 2010-07-23 17:19:05 0.00 0.00
```

bash	----	0	root	2010-07-23	17:18:55	0.02	21.10
su	-S--	0	root	2010-07-23	17:18:51	0.01	24.94
cat	--XC	0x83	mtk	2010-07-23	17:19:55	0.00	1.72
sleep	----	0	mtk	2010-07-23	17:19:42	0.00	15.01
grep	----	0x200	mtk	2010-07-23	17:20:12	0.00	0.00
echo	----	0	mtk	2010-07-23	17:21:15	0.01	0.01
t_fork	F---	0	mtk	2010-07-23	17:21:36	0.00	0.00
t_fork	----	0	mtk	2010-07-23	17:21:36	0.00	3.01

输出中的每行都对应于 shell 会话所创建的一个进程。ulimit 和 echo 都是 shell 的内建命令，所以并不会创建新进程。注意，记账文件中 sleep 之所以出现在 cat 之后，是因为 sleep 在 cat 之后才终止。

大部分输出的含义均不言而喻。flags 列中的各个字母表示每条记录对哪些 ac_flag 位进行了置位（参考表 28-1）。至于应如何解释 term.status 列中的终止状态，26.1.3 节有相关描述。

程序清单 28-2：显示进程记账文件中的数据

```

-----procexec/acct_view.c
#include <fcntl.h>
#include <time.h>
#include <sys/stat.h>
#include <sys/acct.h>
#include <limits.h>
#include "ugid_functions.h"          /* Declaration of userNameFromId() */
#include "tlpi_hdr.h"

#define TIME_BUF_SIZE 100

static long long          /* Convert comp_t value into long long */
comptToLL(comp_t ct)
{
    const int EXP_SIZE = 3;          /* 3-bit, base-8 exponent */
    const int MANTISSA_SIZE = 13;    /* Followed by 13-bit mantissa */
    const int MANTISSA_MASK = (1 << MANTISSA_SIZE) - 1;
    long long mantissa, exp;

    mantissa = ct & MANTISSA_MASK;
    exp = (ct >> MANTISSA_SIZE) & ((1 << EXP_SIZE) - 1);
    return mantissa << (exp * 3);    /* Power of 8 = left shift 3 bits */
}

int
main(int argc, char *argv[])
{
    int acctFile;
    struct acct ac;
    ssize_t numRead;
    char *s;
    char timeBuf[TIME_BUF_SIZE];
    struct tm *loc;
    time_t t;

    if (argc != 2 || strcmp(argv[1], "--help") == 0)
        usageErr("%s file\n", argv[0]);

    acctFile = open(argv[1], O_RDONLY);
    if (acctFile == -1)
        errExit("open");

```

```

printf("command flags  term. user  "
      "start time      CPU  elapsed\n");
printf("
      status          "
      "time    time\n");

while ((numRead = read(acctFile, &ac, sizeof(struct acct))) > 0) {
    if (numRead != sizeof(struct acct))
        fatal("partial read");

    printf("%-8.8s ", ac.ac_comm);

    printf("%c", (ac.ac_flag & AFORK) ? 'F' : '-');
    printf("%c", (ac.ac_flag & ASU)  ? 'S' : '-');
    printf("%c", (ac.ac_flag & AXSIG) ? 'X' : '-');
    printf("%c", (ac.ac_flag & ACORE) ? 'C' : '-');

#ifdef __linux__
    printf(" %#6lx ", (unsigned long) ac.ac_exitcode);
#else
    /* Many other implementations provide ac_stat instead */
    printf(" %#6lx ", (unsigned long) ac.ac_stat);
#endif

    s = userNameFromId(ac.ac_uid);
    printf("%-8.8s ", (s == NULL) ? "???" : s);

    t = ac.ac_btime;
    loc = localtime(&t);
    if (loc == NULL) {
        printf("???Unknown time??? ");
    } else {
        strftime(timeBuf, TIME_BUF_SIZE, "%V-%m-%d %T ", loc);
        printf("%s ", timeBuf);
    }

    printf("%5.2f %7.2f ", (double) (comptToLL(ac.ac_utime) +
                                     comptToLL(ac.ac_stime)) / sysconf(_SC_CLK_TCK),
           (double) comptToLL(ac.ac_etime) / sysconf(_SC_CLK_TCK));
    printf("\n");
}

if (numRead == -1)
    errExit("read");

exit(EXIT_SUCCESS);
}

```

procexec/acct_view.c

进程记账文件格式（版本 3）

从内核 2.6.8 开始，Linux 引入了另一版本的进程记账文件以备选用，意在突破传统记账文件的一些限制。若有意使用这种被称为版本 3 的备选格式，需要在编译内核前打开内核配置选项 `CONFIG_BSD_PROCESS_ACCT_V3`。

使用版本 3，操作进程记账时唯一的差别在于，写入记账文件的记录格式不同。新格式的定义如下：

```

struct acct_v3 {
    char    ac_flag;        /* Accounting flags */

```

```

char    ac_version;    /* Accounting version (3) */
u_int16_t ac_tty;    /* Controlling terminal for process */
u_int32_t ac_exitcode; /* Process termination status */
u_int32_t ac_uid;    /* 32-bit user ID of process */
u_int32_t ac_gid;    /* 32-bit group ID of process */
u_int32_t ac_pid;    /* Process ID */
u_int32_t ac_ppid;    /* Parent process ID */
u_int32_t ac_btime;    /* Start time (time_t) */
float    ac_etime;    /* Elapsed (real) time (clock ticks) */
comp_t    ac_utime;    /* User CPU time (clock ticks) */
comp_t    ac_stime;    /* System CPU time (clock ticks) */
comp_t    ac_mem;    /* Average memory usage (kilobytes) */
comp_t    ac_io;    /* Bytes read/written (unused) */
comp_t    ac_rw;    /* Blocks read/written (unused) */
comp_t    ac_minflt;    /* Minor page faults */
comp_t    ac_majflt;    /* Major page faults */
comp_t    ac_swaps;    /* Number of swaps (unused; Linux-specific) */
#define ACCT_COMM 16
char    ac_comm[ACCT_COMM]; /* Command name */
};

```

以下是 `acct_v3` 结构与传统 Linux `acct` 结构的主要差别。

- 增加 `ac_version` 字段。该字段包含本类型账单记录的版本号。对于 `acct_v3` 来说，总是等于 3。
- 增加 `ac_pid` 和 `ac_ppid` 字段，分别包含终止进程的进程 ID 及其父进程 ID。
- 字段 `ac_uid` 和 `ac_gid` 从 16 位扩展至 32 位，旨在容纳 Linux 2.4 所引入的 32 位用户 ID 和组 ID。（传统 `acct` 文件无法正确记录大数值的用户和组 ID。）
- 将 `ac_etime` 字段类型从 `comp_t` 改为 `float`，意在能够记录更长的逝去时间。

随本书发布的源代码文件 `procexec/acct_v3_view.c` 中提供了程序清单 28-2 中程序基于 `v3` 格式的新版本。

28.2 系统调用 `clone()`

类似于 `fork()` 和 `vfork()`，Linux 特有的系统调用 `clone()` 也能创建一个新进程。与前两者不同的是，后者在进程创建期间对步骤的控制更为精准。`clone()` 主要用于线程库的实现。由于 `clone()` 有损于程序的可移植性，故而应避免在应用程序中直接使用。之所以在这里讨论 `clone()`，意在通过对第 29 章至第 33 章所论述的 POSIX 线程有所铺垫，同时也利于进一步阐明 `fork()` 和 `vfork()` 的操作。

```

#define _GNU_SOURCE
#include <sched.h>

int clone(int (*func) (void *), void *child_stack, int flags, void *func_arg, ...
        /* pid_t *ptid, struct user_desc *tls, pid_t *ctid */ );

        Returns process ID of child on success, or -1 on error

```

如同 `fork()`，由 `clone()` 创建的新进程几近于父进程的翻版。

但与 `fork()` 不同的是，克隆生成的子进程继续运行时不以调用处为起点，转而去调用以参数 `func` 所指定的函数，`func` 又称为子函数（child function）。调用子函数时的参数由 `func_arg`

指定。经过适当转换，子函数可对该参数的含义自由解读，例如，可以作为整型值（int），也可视为指向结构的指针。（之所以可以作为指针处理，是因为克隆产生的子进程对调用进程的内存既可获取，也可共享。）

对于内核而言，fork()、vfork()以及 clone()最终均由同一函数实现（kernel/fork.c 中的 do_fork()）。在这一层次上，clone 与 fork 更为接近：sys_clone()并没有 func 和 func_arg 参数，且调用后 sys_clone()在子进程中返回的方式也与 fork()相同。正文所述的 clone()是由 glibc 为 sys_clone()提供的封装函数。（对该函数的定义位于 glibc 针对特定架构的汇编源码中，例如 sysdeps/unix/sysv/linux/i386/clone.S.）sys_clone()在子进程中返回之后，由 clone()发起对 func 函数的调用。

当函数 func 返回（此时其返回值即为进程的退出状态）或是调用 exit()（或 _exit()）之后，克隆产生的子进程就会终止。照例，父进程可以通过 wait()一类函数来等待克隆子进程。

因为克隆产生的子进程可能（类似 vfork()）共享父进程的内存，所以它不能使用父进程的栈。相反，调用者必须分配一块大小适中的内存空间供子进程的栈使用，同时将这块内存的指针置于参数 child_stack 中。在大多数硬件架构中，栈空间的生长方向是向下的，所以参数 child_stack 应当指向所分配内存块的高端。

栈增长方向对架构的依赖是 clone()设计的一处缺陷。Interl IA-64 架构就提供了一款经过改善的克隆 API，称为 clone2()。该系统调用对子进程栈范围的定义方式不依赖于栈的生长方向，只需要提供栈的起始地址以及大小即可。详情请参阅手册页。

函数 clone()的参数 flags 服务于双重目的。首先，其低字节中存放着子进程的终止信号（terminateion signal），子进程退出时其父进程将收到这一信号。（如果克隆产生的子进程因信号而终止，父进程依然会收到 SIGCHLD 信号。）该字节也可能为 0，这时将不会产生任何信号。（借助于 Linux 特有的/proc/PID/stat 文件，可以判定任何进程的终止信号，详情请参阅 proc(5)手册页。）

对于 fork()和 vfork()而言，就无从选择终止信号，只能是 SIGCHLD。

参数 flags 的剩余字节则存放了位掩码，用于控制 clone()的操作。表 28-2 对这些位掩码值进行了总结，28.2.1 节会进一步加以说明。

表 28-2: clone()参数 flags 的位掩码值

标 志	设置后的效果
CLONE_CHILD_CLEARTID	当子进程调用 exec()或 _exit()时，清除 ctid（从版本 2.6 开始）
CLONE_CHILD_SETTID	将子进程的线程 ID 写入 ctid（从 2.6 版本开始）
CLONE_FILES	父、子进程共享打开文件描述符表
CLONE_FS	父、子进程共享与文件系统相关的属性
CLONE_IO	子进程共享父进程的 I/O 上下文环境（从 2.6.25 版本开始）
CLONE_NEWIPC	子进程获得新的 System V IPC 命名空间（从 2.6.19 开始）
CLONE_NEWNET	子进程获得新的网络命名空间（从 2.4.24 版本开始）

标 志	设置后的效果
CLONE_NEWNS	子进程获得父进程挂载（mount）命名空间的副本（从 2.4.19 版本开始）
CLONE_NEWPID	子进程获得新的进程 ID 命名空间（从 2.6.23 版本开始）
CLONE_NEWUSER	子进程获得新的用户 ID 命名空间（从 2.6.23 版本开始）
CLONE_NEWUTS	子进程获得新的 UTS（utsname()）命名空间（从 2.6.19 版本开始）
CLONE_PARENT	将子进程的父进程置为调用者的父进程（从 2.4 版本开始）
CLONE_PARENT_SETTID	将子进程的线程 ID 写入 ptid（从 2.6 版本开始）
CLONE_PID	标志已废止，仅用于系统启动进程（直至 2.4 版本为止）
CLONE_PTRACE	如果正在跟踪父进程，那么子进程也照此办理
CLONE_SETTLS	tls 描述子进程的线程本地存储（从 2.6 开始）
CLONE_SIGHAND	父、子进程共享对信号的处置设置
CLONE_SYSVSEM	父、子进程共享信号量还原（undo）值（从 2.6 版本开始）
CLONE_THREAD	将子进程置于父进程所属的线程组中（从 2.4 开始）
CLONE_UNTRACED	不强制对子进程设置 CLONE_PTRACE（从 2.6 版本开始）
CLONE_VFORK	挂起父进程直至子进程调用 exec()或_exit()
CLONE_VM	父、子进程共享虚拟内存

clone()的余下参数分别是：ptid、tls 和 ctid。这些参数与线程的实现相关，尤其是在针对线程 ID 以及线程本地存储的使用方面。28.2.1 节在说明 flags 位掩码值时，会论及这些参数的使用。（在 Linux 2.4 及其之前的版本中，clone()尚未提供上述 3 个参数。直到 Linux 2.6，为了支持 NPTL POSIX 的线程实现，才特意加入了这些参数。）

示例程序

程序清单 28-3 是使用 clone()创建子进程的一个简单例子。主程序所做工作如下。

- 打开一个文件描述符（打开设备/dev/null），在子进程中将其关闭②。
- 若提供有命令行参数，则将 clone()的 flags 参数置为 CLONE_FILES③，以便父、子进程共享同一文件描述符表。若没有提供命令行参数，则将 flags 置 0。
- 分配一个栈供子进程使用④。
- 若 CHILD_SIG 非 0 且不等于 SIGCHLD⑤，则忽略之，以防该信号将子进程终止。之所以未忽略 SIGCHLD，是因为那将导致无法收集子进程的退出状态。
- 调用 clone()创建子进程⑥。第三个参数（位掩码）包含了终止信号。第四个参数（func_arg）指定了之前打开的文件描述符（在②处）。
- 等待子进程终止⑦。
- 尝试调用 write()，以检查文件描述符（在②处打开）是否仍处于打开状态⑧。程序报告 write()操作是否成功。

克隆产生的子进程从 childFunc()处开始执行，该函数（利用参数 arg）接收由主程序打开的文件描述符（在②处）。子进程关闭文件描述符并调用 return 以终止①。


```

#define _GNU_SOURCE
#include <signal.h>
#include <sys/wait.h>
#include <fcntl.h>
#include <sched.h>
#include "tlpi_hdr.h"

#ifdef CHILD_SIG
#define CHILD_SIG SIGUSR1      /* Signal to be generated on termination
                               of cloned child */
#endif

static int childFunc(void *arg) /* Startup function for cloned child */
{
①   if (close(*(int *) arg) == -1)
        errExit("close");

        return 0;              /* Child terminates now */
}

int main(int argc, char *argv[])
{
    const int STACK_SIZE = 65536; /* Stack size for cloned child */
    char *stack;                 /* Start of stack buffer */
    char *stackTop;              /* End of stack buffer */
    int s, fd, flags;

②   fd = open("/dev/null", O_RDWR); /* Child will close this fd */
    if (fd == -1)
        errExit("open");
    /* If argc > 1, child shares file descriptor table with parent */

③   flags = (argc > 1) ? CLONE_FILES : 0;

    /* Allocate stack for child */

④   stack = malloc(STACK_SIZE);
    if (stack == NULL)
        errExit("malloc");
    stackTop = stack + STACK_SIZE; /* Assume stack grows downward */

    /* Ignore CHILD_SIG, in case it is a signal whose default is to
       terminate the process; but don't ignore SIGCHLD (which is ignored
       by default), since that would prevent the creation of a zombie. */

⑤   if (CHILD_SIG != 0 && CHILD_SIG != SIGCHLD)
        if (signal(CHILD_SIG, SIG_IGN) == SIG_ERR)
            errExit("signal");

    /* Create child; child commences execution in childFunc() */

⑥   if (clone(childFunc, stackTop, flags | CHILD_SIG, (void *) &fd) == -1)
        errExit("clone");

```

```

/* Parent falls through to here. Wait for child; __WCLONE is
   needed for child notifying with signal other than SIGCHLD. */
⑦ if (waitpid(-1, NULL, (CHILD_SIG != SIGCHLD) ? __WCLONE : 0) == -1)
    errExit("waitpid");
printf("child has terminated\n");

/* Did close() of file descriptor in child affect parent? */

⑧ s = write(fd, "x", 1);
if (s == -1 && errno == EBADF)
    printf("file descriptor %d has been closed\n", fd);
else if (s == -1)
    printf("write() on file descriptor %d failed "
           "unexpectedly (%s)\n", fd, strerror(errno));
else
    printf("write() on file descriptor %d succeeded\n", fd);

exit(EXIT_SUCCESS);
}

```

procexec/t_clone.c

运行程序清单 28-3 中程序，没有命令行参数时输出如下：

```

$ ./t_clone                               Doesn't use CLONE_FILES
child has terminated
write() on file descriptor 3 succeeded      Child's close() did not affect parent

```

带有命令行参数运行程序时，两个进程将共享文件描述符表：

```

$ ./t_clone x                               Uses CLONE_FILES
child has terminated
file descriptor 3 has been closed          Child's close() affected parent

```

随本书发布的源代码文件 `procexec/demo_clone.c` 提供一个更为复杂的 `clone()` 用例。

28.2.1 `clone()` 的 flags 参数

`clone()` 的 flags 参数是各种位掩码的组合（“或”操作），下面将对它们一一说明。讲述时并未按字母顺序展开，而是着眼于促进对概念理解，从实现 POSIX 线程所使用的标志开始。从线程实现的角度来看，下文多次出现的“进程”一词都可用“线程”替代。

这里需要指出，某种意义上，对术语“线程”和“进程”的区分不过是在玩弄文字游戏而已。引入术语“内核调度实体（KSE, kernel scheduling entity）”（某些教科书以之来指代内核调度器所处理的对象）的概念对解释这一点会有所助益。实际上，线程和进程都是 KSE，只是与其他 KSE 之间对属性（虚拟内存、打开文件描述符、对信号的处置、进程 ID 等）的共享程度不同。针对线程间属性共享的方案不少，POSIX 线程规范只是其中之一。

在下面的说明中，有时会提及 Linux 平台对 POSIX 线程的两种主要实现：年长的 LinuxThreads，以及较为年轻的 NPTL。关于这两种实现的更多细节可以在 33.5 节找到。

从内核 2.6.16 开始，Linux 提供了新的系统调用 `unshare()`，由 `clone()`（或 `fork()`、`vfork()`）创建的子进程利用该调用可以撤销对某些属性的共享（即反转一些 `clone()` flags 位的效果）。详细情况请参考 `unshare(2)` 手册页。

共享文件描述符表：CLONE_FILES

如果指定了 `CLONE_FILES` 标志，父、子进程会共享同一个打开文件描述符表。也就是说，无论哪个进程对文件描述符的分配和释放（`open()`、`close()`、`dup()`、`pipe()`、`socket()`等），都会影响到另一进程。如果未设置 `CLONE_FILES`，那么也就不会共享文件描述符表，子进程获取的是父进程调用 `clone()`时文件描述符表的一份拷贝。这些描述符副本与其父进程中的相应描述符均指向相同的打开文件（和 `fork()`和 `vfork()`的情况一样）。

POSIX 线程规范要求进程中的所有线程共享相同的打开文件描述符。

共享与文件系统相关的信息：CLONE_FS

如果指定了 `CLONE_FS` 标志，那么父、子进程将共享与文件系统相关的信息（file system-related information）：权限掩码（`umask`）、根目录以及当前工作目录。也就是说，无论在哪个进程中调用 `umask()`、`chdir()`或者 `chroot()`，都将影响到另一个进程。如果未设置 `CLONE_FS`，那么父、子进程对此类信息则会各持一份（与 `fork()`和 `vfork()`的情况相同）。

POSIX 线程规范要求实现 `CLONE_FS` 标志所提供的属性共享。

共享对信号的处置设置：CLONE_SIGHAND

如果设置了 `CLONE_SIGHAND`，那么父、子进程将共享同一个信号处置表。无论在哪个进程中调用 `sigaction()`或 `signal()`来改变对信号处置的设置，都会影响其他进程对信号的处置。若未设置 `CLONE_SIGHAND`，则不共享对信号的处置设置，子进程只是获取父进程信号处置表的一份副本（如同 `fork()`和 `vfork()`）。`CLONE_SIGHAND` 不会影响到进程的信号掩码以及对挂起（pending）信号的设置，父子进程的此类设置是绝不相同的。从 Linux 2.6 开始，如果设置了 `CLONE_SIGHAND`，就必须同时设置 `CLONE_VM`。

POSIX 线程规范要求共享对信号的处置设置。

共享父进程的虚拟内存：CLONE_VM

如果设置了 `CLONE_VM` 标志，父、子进程会共享同一份虚拟内存页（如同 `vfork()`）。无论哪个进程更新了内存，或是调用了 `mmap()`、`munmap()`，另一进程同样会观察到这些变化。如果未设置 `CLONE_VM`，那么子进程得到的是对父进程虚拟内存的拷贝（如同 `fork()`）。

共享同一虚拟内存是线程的关键属性之一，POSIX 线程标准对此也有要求。

线程组：CLONE_THREAD

若设置了 `CLONE_THREAD`，则会将子进程置于父进程的线程组中。如果未设置该标志，那么会将子进程置于新的线程组中。

POSIX 标准规定，进程的所有线程共享同一进程 ID（即每个线程调用 `getpid()`都应返回相同值），Linux 从 2.4 版本开始引入了线程组（threads group），以满足这一需求。如图 28-1 所示，线程组就是共享同一线程组标识（TGID）（thread group identifier）的一组 KSE。在对 `CLONE_THREAD` 的后续讨论中，会将 KSE 视同线程看待。

始于 Linux 2.4，`getpid()`所返回的就是调用者的 TGID。换言之，TGID 和进程 ID 是一回事。

在 2.2 以及更早的 Linux 系统中，对 `clone()`的实现并不支持 `CLONE_THREAD`。相反，

LinuxThreads 曾将 POSIX 线程实现为共享了多种属性（例如，虚拟内存）、进程 ID 又各不相同的进程。考虑到兼容性因素，即便是在当前的 Linux 内核中，LinuxThreads 实现也未提供 CLONE_THREAD，因为按此方式实现的线程就可以继续拥有不同的进程 ID。

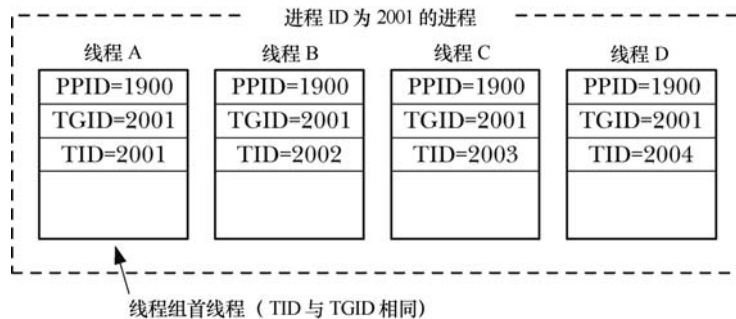


图 28-1: 包含 4 个线程的线程组

一个线程组内的每个线程都拥有一个唯一的线程标识符（thread identifier, TID），用以标识自身。Linux 2.4 提供了一个新的系统调用 `gettid()`，线程可通过该调用来获取自己的线程 ID（与线程调用 `clone()` 时的返回值相同）。线程 ID 与进程 ID 都使用相同的数据类型 `pid_t` 来表示。线程 ID 在整个系统中是唯一的，且除了线程担当进程中线程组首线程的情况之外，内核能够保证系统中不会出现线程 ID 与进程 ID 相同的情况。

线程组中首个线程的线程 ID 与其线程组 ID 相同，也将该线程称之为线程组首线程（thread group leader）。

此处讨论的线程 ID 与 POSIX 线程所使用的线程 ID（以数据类型 `pthread_t` 表示）不同。后者由 POSIX 线程实现（在用户空间）自行生成并维护。

线程组中的所有线程拥有同一父进程 ID，即与线程组首线程 ID 相同。仅当线程组中的所有线程都终止后，其父进程才会收到 `SIGCHLD` 信号（或其他终止信号）。这些行为符合 POSIX 线程规范的要求。

当一个设置了 `CLONE_THREAD` 的线程终止时，并没有信号会发送给该线程的创建者（即调用 `clone()` 创建终止线程的线程）。相应的，也不可能调用 `wait()`（或类似函数）来等待一个以 `CLONE_THREAD` 标志创建的线程。这与 POSIX 的要求一致。POSIX 线程与进程不同，不能使用 `wait()` 等待，相反，必须调用 `pthread_join()` 来加入。为检测以 `CLONE_THREAD` 标志创建的线程是否终止，需要使用一种特殊的同步原语——`futex`（参考下文对 `CLONE_PARENT_SETTID` 标志的讨论）。

如果一个线程组中的任一线程调用了 `exec()`，那么除了首线程之外的其他线程都会终止（这一行为也符合 POSIX 线程规范的要求），新进程将在首线程中执行。换言之，新程序中的 `gettid()` 调用将会返回首线程的线程 ID。调用 `exec()` 期间，会将该进程发送给其父进程的终止信号重置为 `SIGCHLD`。

如果线程组中的某个线程调用 `fork()` 或 `vfork()` 创建了子进程，那么组中的任何线程都可使用 `wait()` 或类似函数来监控该子进程。

从 Linux 2.6 开始，如果设置了 `CLONE_THREAD`，同时也必须设置 `CLONE_SIGHAND`。这也与 POSIX 线程标准的深入要求相契合，详细内容可参考 33.2 节关于 POSIX 线程与信号交互的相关讨论。（内核针对 `CLONE_THREAD` 线程组的信号处理对应于 POSIX 标准对进程中线程如何处理信号的规范。）

线程库支持: CLONE_PARENT_SETTID、CLONE_CHILD_SETTID 和 CLONE_CHILD_CLEARTID

为实现 POSIX 线程, Linux 2.6 提供了对 CLONE_PARENT_SETTID、CLONE_CHILD_SETTID 和 CLONE_CHILD_CLEARTID 的支持。这些标志会影响 clone() 对参数 ptid 和 ctid 的处理。NPTL 的线程实现使用了 CLONE_CHILD_SETTID 和 CLONE_CHILD_CLEARTID。

如果设置了 CLONE_PARENT_SETTID, 内核会将子线程的线程 ID 写入 ptid 所指向的位置。在对父进程的内存进行复制之前, 会将线程 ID 复制到 ptid 所指位置。这也意味着, 即使没有设置 CLONE_VM, 父、子进程均能在此位置获得子进程的线程 ID。(如上所述, 创建 POSIX 线程时总是指定了 CLONE_VM 标志。)

CLONE_PARENT_SETTID 之所以存在, 意在为线程实现获取新线程 ID 提供一种可靠的手段。注意, 通过 clone() 的返回值并不足以获取新线程的线程 ID。

```
tid = clone(...);
```

问题在于, 因为赋值操作只能在 clone() 返回后才会发生, 所以上面代码会导致各种竞争条件。例如, 假设新线程终止, 而在完成对 tid 的赋值前就调用了终止信号的处理程序。此时, 处理程序无法有效访问 tid。(在线程库内部, 可能会将 tid 置于一个用以跟踪所有线程状态的全局结构中。) 程序通常可以通过直接调用 clone() 来规避这种竞争条件。不过, 线程库无法控制其调用者程序的行为。使用 CLONE_PARENT_SETTID 可以保证在 clone() 返回之前就将新线程的 ID 赋值给 ptid 指针, 从而使线程库避免了这种竞争条件。

如果设置了 CLONE_CHILD_SETTID, 那么 clone() 会将子线程的线程 ID 写入指针 ctid 所指向的位置。对 ctid 的设置只会发生在子进程的内存中, 不过如果设置了 CLONE_VM, 还是会影响到父进程。虽然 NPTL 并不需要 CLONE_CHILD_SETTID, 但这一标识还是能给其他的线程库实现带来灵活性。

如果设置了 CLONE_CHILD_CLEARTID 标志, 那么 clone() 会在子进程终止时将 ctid 所指向的内存内容清零。

借助于参数 ctid 所提供的机制(稍后描述), NPTL 线程实现可以获得线程终止的通知。函数 pthread_join() 正需要这样的通知, POSIX 线程利用该函数来等待另一线程的终止。

使用 pthread_create() 创建线程时, NPTL 会调用 clone(), 其 ptid 和 ctid 均指向同一位置。(这正是 NPTL 不需要 CLONE_CHILD_SETTID 的原因所在。) 设置了 CLONE_PARENT_SETTID 标志, 就会以新的线程 ID 对该位置进行初始化。当子进程终止, ctid 遭清除时, 进程中的所有线程都会目睹这一变化(因为设置了 CLONE_VM)。

内核将 ctid 指向的位置视同 futex——一种有效的同步机制来处理。(关于 futex 的更多内容请参考 futex(2) 手册页。) 执行系统调用 futex() 来监测 ctid 所指位置的内容变化, 就可获得线程终止的通知。(这正是 pthread_join() 所做的幕后工作。) 内核在清除 ctid 的同时, 也会唤醒那些调用了 futex() 来监控该地址内容变化的任一内核调度实体(即线程)。(在 POSIX 线程的层面上, 这会导致 pthread_join() 调用去解除阻塞。)

线程本地存储: CLONE_SETTLS

如果设置了 CLONE_SETTLS, 那么参数 tls 所指向的 user_desc 结构会对该线程所使用的线程本地存储缓冲区加以描述。为了支持 NPTL 对线程本地存储的实现, Linux 2.6 开始加入这一标志(31.4 节)。关于 user_desc 结构的详情, 可参考 2.6 内核代码中对该结构的定义和使用, 以及 set_thread_area(2) 手册页。

共享 System V 信号量的撤销值：CLONE_SYSVSEM

如果设置了 CLONE_SYSVSEM，父、子进程将共享同一个 System V 信号量撤销值列表（47.8 节）。如果未设置该标志，父、子进程各自持有取消列表，且子进程的列表初始为空。

内核从 2.6 版本开始支持 CLONE_SYSVSEM，提供 POSIX 线程规范所要求的共享语义。

每进程挂载命名空间：CLONE_NEWNS

Linux 从内核 2.4.19 开始支持每进程挂载（mount）命名空间的概念。挂载命名空间是由对 mount()和 umount()的调用来维护的一组挂载点。挂载命名空间会影响将路径名解析为真实文件的过程，也会波及诸如 chdir()和 chroot()之类的系统调用。

默认情况下，父、子进程共享同一挂载命名空间，一个进程调用 mount()或 umount()对命名空间所做的改变，也会为其他进程所见（如同 fork()和 vfork()）。特权级（CAP_SYS_ADMIN）进程可以指定 CLONE_NEWNS 标志，以便子进程去获取对父进程挂载命名空间的一份拷贝。这样一来，进程对命名空间的修改就不会为其他进程所见。（早期的 2.4.x 内核以及更老的版本认为，系统的所有进程共享同一个系统级挂载命名空间。）

可以利用每进程挂载命名空间来创建类似于 chroot()监禁区（jail）的环境，而且更加安全、灵活，例如，可以向遭到监禁的进程提供一个挂载点，而该点对于其他进程是不可见的。设置虚拟服务器环境时也会用到挂载命名空间。

在同一 clone()调用中同时指定 CLONE_NEWNS 和 CLONE_FS 纯属无聊，也不允许这样做。

将子进程的父进程置为调用者的父进程：CLONE_PARENT

默认情况下，当调用 clone()创建新进程时，新进程的父进程（由 getppid()返回）就是调用 clone()的进程（同 fork()和 vfork()）。如果设置了 CLONE_PARENT，那么调用者的父进程就成为子进程的父进程。换言之，CLONE_PARENT 等同于这样的设置：子进程.PPID = 调用者.PPID。（未设置 CLONE_PARENT 的默认情况是：子进程.PPID = 调用者.PID。）子进程终止时会向父进程（子进程.PPID）发出信号。

Linux 从版本 2.4 之后开始支持 CLONE_PARENT。其设计初衷意图是对 POSIX 线程的实现提供支持，不过内核 2.6 找出一种无需此标志而支持线程（之前所述的 CLONE_THREAD）的新方法。

将子进程的进程 ID 置为与父进程相同：CLONE_PID（已废止）

如果设置了 CLONE_PID，那么子进程就拥有与父进程相同的进程 ID。若未设置此标志，那么父、子进程的进程 ID 则不同（如同 fork()和 vfork()）。只有系统引导进程（进程 ID 为 0）可能会使用该标志，用于初始化多处理器系统。

CLONE_PID 的设计初衷并非供用户级应用使用。Linux 2.6 已将其移除，并以 CLONE_IDLETASK 取而代之，将新进程的 ID 置为 0。CLONE_IDLETASK 仅供内核内部使用（即使在 clone()的参数中指定，系统也会对其视而不见）。使用此标志可为每颗 CPU 创建隐身的空闲进程（idle process），在多处理器系统中可能存在有多个实例。

进程跟踪：CLONE_PTRACE 和 CLONE_UNTRACED

如果设置了 CLONE_PTRACE 且正在跟踪调用进程，那么也会对子进程进行跟踪。关于进程跟踪（由调试器和 strace 命令使用）的细节，请参考 ptrace(2)手册页。

从内核 2.6 开始，即可设置 `CLONE_UNTRACED` 标志，这也意味着跟踪进程不能强制将其子进程设置为 `CLONE_PTRACE`。`CLONE_UNTRACED` 标志供内核创建内核线程时内部使用。

挂起（suspending）父进程直至子进程退出或调用 `exec()`：`CLONE_VFORK`

如果设置了 `CLONE_VFORK`，父进程将一直挂起，直至子进程调用 `exec()` 或 `_exit()` 来释放虚拟内存资源（如同 `vfork()`）为止。

支持容器（container）的 `clone()` 新标志

Linux 从 2.6.19 版本开始给 `clone()` 加入了一些新标志：`CLONE_IO`、`CLONE_NEWIPC`、`CLONET_NEWNET`、`CLONE_NEWPID`、`CLONE_NEWUSER` 和 `CLONE_NEWUTS`。（参考 `clone(2)` 手册页可获得有关这些标志的详细说明。）

这些标志中的大部分都是为容器（container）的实现提供支持（[Bhattiprolu et al., 2008]）。容器是轻量级虚拟化的一种形式，将运行于同一内核的进程组从环境上彼此隔离，如同运行在不同机器上一样。容器可以嵌套，一个容器可以包含另一个容器。与完全虚拟化将每个虚拟环境运行于不同内核的手法相比，容器的运作方式可谓是大相径庭。

为实现容器，内核开发者不得不为内核中的各种全局系统资源提供一个间接层，以便每个容器能为这些资源提供各自的实例。这些资源包括：进程 ID、网络协议栈、`uname()` 返回的 ID、System V IPC 对象、用户和组 ID 命名空间……

容器的用途很多，如下所示。

- 控制系统的资源分配，诸如网络带宽或 CPU 时间（例如，授予容器某甲 75% 的 CPU 时间，某乙则获取 25%）。
- 在单台主机上提供多个轻量级虚拟服务器。
- 冻结某个容器，以此来挂起该容器中所有进程的执行，并于稍后重启，可能是在迁移到另一台机器之后。
- 允许转储应用程序的状态信息，记录于检查点（checkpointed），并于之后再行恢复（或许在应用程序崩溃之后，亦或是计划内、外的系统停机后），从检查点开始继续运行。

`clone()` 标志的使用

大体上说来，`fork()` 相当于仅设置 `flags` 为 `SIGCHLD` 的 `clone()` 调用，而 `vfork()` 则对应于设置如下 `flags` 的 `clone()`：

```
CLONE_VM | CLONE_VFORK | SIGCHLD
```

自 2.3.3 版本以来，作为 NPTL 线程实现的一部分，glibc 所提供的封装函数 `fork()` 绕开了内核的 `fork()` 系统调用，转而调用了 `clone()`。该封装函数会去调用任何由调用者通过 `pthread_atfork()`（参考 33.3 节）所设置的 `fork` 处理器程序。

LinuxThreads 线程实现使用 `clone()`（仅用到前 4 个参数）来创建线程，对 `flags` 的设置如下：
`CLONE_VM | CLONE_FILES | CLONE_FS | CLONE_SIGHAND`

NPTL 线程实现则使用 `clone()`（使用了所有 7 个参数）来创建线程，对 `flags` 的设置如下：

```
CLONE_VM | CLONE_FILES | CLONE_FS | CLONE_SIGHAND | CLONE_THREAD |  
CLONE_SETTLS | CLONE_PARENT_SETTID | CLONE_CHILD_CLEARID | CLONE_SYSVSEM
```

28.2.2 因克隆生成的子进程而对 waitpid()进行的扩展

为等待由 clone()产生的子进程，waitpid()、wait3()和 wait4()的位掩码参数 options 可以包含如下附加（Linux 特有）值。

__WCLONE

一经设置，只会等待克隆子进程。如未设置，只会等待非克隆子进程。在这种情况下，克隆子进程终止时发送给其父进程的信号并非 SIGCHLD。如果同时还指定了 __WALL，那么将忽略 __WCLONE。

__WALL（自 Linux2.4 之后）

等待所有子进程，无论类型（克隆、非克隆通吃）。

__WNOTHREAD（自 Linux2.4 之后）

默认情况下，等待（wait）类调用所等待的子进程，其父进程的范围遍及与调用者隶属同一线程组的任何进程。指定 __WNOTHREAD 标志则限制调用者只能等待自己的子进程。

waitid()不能使用上述标志。

28.3 进程的创建速度

表 28-3 对采用不同方法创建进程的速度进行了比较。测试程序在循环中反复创建子进程并等待子进程终止，从而获得了这一结果。比较过程中使用了 3 种不同大小的进程内存，如表中虚拟内存总量（total virtual memory）值所示。对不同大小内存的模拟，依赖于程序计时之前在堆中分配（malloc()）的额外内存。

表 28-3 中的进程大小（虚拟内存总量）取自命令 ps -o “pid vsz cmd” 输出的 VSZ 值。

表 28-3: 使用 fork()、vfork()和 clone()创建 10 万个进程所需的时间

进程的创建方法	虚拟内存总量					
	1.70 MB		2.70 MB		11.70 MB	
	时间 (秒)	速率	时间 (秒)	速率	时间 (秒)	速率
fork()	22.27 (7.99)	4544	26.38 (8.98)	4135	126.93 (52.55)	1276
vfork()	3.52 (2.49)	28955	3.55 (2.50)	28621	3.53 (2.51)	28810
clone()	2.97 (2.14)	34333	2.98 (2.13)	34217	2.93 (2.10)	34688
fork() + exec()	135.72 (12.39)	764	146.15 (16.69)	719	260.34 (61.86)	435
vfork() + exec()	107.36 (6.27)	969	107.81 (6.35)	964	107.97 (6.38)	960

表 28-3 对于每种进程大小都提供了两类统计数据。

- 第 1 项统计包含两种度量时间。以执行 10 万次进程创建期间所逝去的（实际）时间为主（较大值），以父进程所消耗的 CPU 时间（括号内的值）为辅。由于测试环境并无其他负载，两者之差应是测试期间创建子进程所消耗的时间总量。
- 第 2 项数据显示每（实际）秒创建的进程数，即创建速率，取各种情况下运行 20 次的平均值。实验基于 x86-32 系统，内核版本为 2.6.27。

前 3 行针对的是简单的进程创建（子进程不运行新程序）。子进程在创建后立即退出，父进程等待子进程终止后再去创建下一个子进程。

第 1 行取自系统调用 `fork()`。由数据可知，进程所占内存越大，`fork()`所需时间也就越长。额外时间花在了为子进程复制那些逐渐变大的页表，以及将数据段、堆段以及栈段的页记录标记为只读的工作上。因为子进程并未修改数据段或栈段，所以也没有对页（page）复制。

第 2 行取自 `vfork()`。可以看出，尽管进程大小在增加，但所用时间保持不变，因为调用 `vfork()`时并未复制页表或页，调用进程的虚拟内存大小并未造成影响。`fork()`和 `vfork()`在时间统计上的差值就是复制进程页表所需的时间总量。

表 28-3 中 `vfork()`和 `clone()`的各自数据在不同的进程内存大小下。之所以存在微小的差异，要归因于采样误差以及调度的变化。即使创建 300MB 大小的进程，两个系统调用的时间仍将保持不变。

第 3 行数据的统计信息来自对 `clone()`的调用，所使用的标志如下：

`CLONE_VM | CLONE_VFORK | CLONE_FS | CLONE_SIGHAND | CLONE_FILES`

前两个标志模拟 `vfork()`的行为。剩余的标志则要求父、子进程应当共享文件系统属性文件权限掩码（`umask`）、根目录和当前工作目录，信号处置表以及打开文件描述符表。`clone()`和 `vfork()`之间的数据差值则代表了 `vfork()`将这些信息拷贝到子进程的少量额外工作。拷贝文件系统属性和信号处置表的成本是固定的。不过，拷贝打开文件描述符表的开销则取决于描述符数量。例如：父进程打开 100 个文件，`vfork()`的实际时间（表中第 1 列）会从 3.52 秒增至 5.04 秒，但不会影响 `clone()`所需要的时间。

对 `clone()`的计时针对的是 `glibc` 库的封装函数 `clone()`，而非直接调用 `sys_clone()`。另有测试（在此恕不一一列出）对 `sys_clone()`和 `clone()`（以子函数调用并立即退出）做了比较，实验结果表明，时间上的差异可以忽略不计。

`fork()`和 `vfork()`之间的差别非常明显，但仍需要注意以下几点。

- 最后一列数据表明，在大进程情况下，`vfork()`要比 `fork()`快逾 30 倍。而针对普通进程，则会近乎于表中前两列的数据。
- 因为进程的创建时间往往比 `exec()`的执行时间要少得多，所以如果随后接着执行 `exec()`，那么两者间的差异也就不再明显。表 28-3 的最后两行数据说明了这一点，其中的每个子进程都去调用 `exec()`，而非直接退出。程序执行的是 `true` 命令（`/bin/true`，选择该程序的原因是因为它不产生任何输出）。这时，`fork()`和 `vfork()`之间的相对差距就小了许多。

事实上，表 28-3 中所示数据并未揭示 `exec()`的全部开销，因为测试程序的每个循环中子进程均执行同一程序。根本就未计入把程序文件读入内存的磁盘 I/O 开销，因为第一次运行 `exec()`时就会将程序读入内核缓冲区，并一直保存在那里。如果测试每次循环执行的程序不同（例如，复制同一程序，并以不同文件名命名），那么应该可以观察到 `exec()`的开销要大出许多。

28.4 exec()和 fork()对进程属性的影响

进程有多种属性，其中一部分已经在前面几章有所说明，后续章节将讨论其他一些属性。关于这些属性，存在两个问题。

- 当进程执行 exec()时，这些属性将发生怎样的变化？
- 当执行 fork()时，子进程会继承哪些属性？

表 28-4 是对这些问题的回答。exec()列注明，调用 exec()期间哪些属性得以保存。fork()列则表明调用 fork()之后子进程继承，或（有时是）共享了哪些属性。除了标注为 Linux 特有的属性之外，列出的所有属性均获得了标准 UNIX 实现的支持，调用 exec()和 fork()期间对它们的处理也都符合 SUSv3 规范。

表 28-4: exec()和 fork()对进程属性的影响

进程属性	exec()	fork()	影响属性的接口；额外说明
进程地址空间			
文本段	否	共享	子进程与父进程共享文本段
栈段	否	是	函数入口/出口；alloca()、longjmp()、siglongjmp()
数据段和堆段	否	是	brk()、sbrk()
环境变量	见注释	是	putenv()、setenv()；直接修改 environ。execle()和 execve()会对其改写，其他 exec()调用则会加以保护
内存映射	否	是；见注释	mmap()、munmap()。跨越 fork()进程，映射的 MAP_NOERREVE 标志得以继承。带有 madvise(MADV_DONTFORK)标志的映射则不会跨 fork()继承
内存锁	否	否	mlock()、munlock()
进程标识符和凭证			
进程 ID	是	否	
父进程 ID	是	否	
进程组 ID	是	是	setpgid()
会话 ID	是	是	setsid()
实际 ID	是	是	setuid()、setgid()，以及相关调用
有效和保存设置 (saved set) ID	见注释	是	setuid()、setgid()，以及相关调用。第 9 章解释了 exec()是如何影响这些 ID 的
补充组 ID	是	是	setgroups()、initgroups()
文件、文件 IO 和目录			
打开文件描述符	见注释	是	open()、close()、dup()、pipe()、socket()等。文件描述符在跨越 exec()调用的过程中得以保存，除非对其设置了运行时关闭 (close-on-exec) 标志。父、子进程中的描述符指向相同的打开文件描述，参考 5.4 节
运行时关闭 (close-on-exec) 标志	是 (如果关闭)	是	fcntl (F_SETFD)
文件偏移	是	共用	lseek()、read()、write()、readv()、writev()。父、子进程共享文件偏移

续表

进 程 属 性	exec()	fork()	影响属性的接口；额外说明
文件、文件 IO 和目录			
打开文件状态标志	是	共用	open()、fcntl(F_SETFL)。父、子进程共享打开文件状态标志
异步 I/O 操作	见注释	否	aio_read()、aio_write()以及相关调用。调用 exec()期间，会取消尚未完成的操作
目录流	否	是：见注释	opendir()、readdir()。SUSv3 规定，子进程获得父进程目录流的一份副本，不过这些副本可以（也可以不）共享目录流的位置。Linux 系统不共享目录流的位置
文件系统			
当前工作目录	是	是	chdir()
根目录	是	是	chroot()
文件模式创建掩码	是	是	umask()
信号			
信号处置	见注释	是	signal()、sigaction()。将处置设置成默认或忽略的信号在执行 exec()期间保持不变；已捕获的信号会恢复为默认处置。参考 27.5 节
信号掩码	是	是	信号传递；sigprocmask()、sigaction()
挂起（pending）信号集合	是	否	信号传递；raise()、kill()、sigqueue()
备选信号栈	否	是	sigaltstack()
定时器			
间隔定时器	是	否	setitimer()
由 alarm()设置的定时器	是	否	alarm()
POSIX 定时器	否	否	timer_create()及其相关调用
POSIX 线程			
线程	否	见注释	fork()调用期间，子进程只会复制调用线程
线程可撤销状态与类型	否	是	exec()之后，可将撤销类型和状态分别重置为 PTHREAD_CANCEL_ENABLE 和 PTHREAD_CANCEL_DEFERRED
互斥量与条件变量	否	是	关于调用 fork()期间对互斥量以及其他线程资源的处理细节可参考 33.3 节
优先级与调度			
nice 值	是	是	nice()、setpriority()
调度策略及优先级	是	是	sched_setscheduler()、sched_setparam()
资源与 CPU 时间			
资源限制	是	是	setrlimit()
进程和子进程的 CPU 时间	是	否	由 times()返回
资源使用量	是	否	由 getrusage()返回

进 程 属 性	exec()	fork()	影响属性的接口；额外说明
进程间通信			
System V 共享内存段	否	是	shmat()、shmdt()
POSIX 共享内存	否	是	shm_open()及其相关调用
POSIX 消息队列	否	是	mq_open()及其相关调用。父、子进程的描述符都指向同一打开消息队列描述。子进程并不继承父进程的消息通知注册信息
POSIX 命名信号量	否	共用	sem_open()及其相关调用。子进程与父进程共享对相同信号量的引用
POSIX 未命名信号量	否	见注释	sem_init()及其相关调用。如果信号量位于共享内存区域，那么子进程与父进程共享信号量；否则，子进程拥有属于自己的信号量拷贝
System V 信号量调整	是	否	参考 47.8 节
文件锁	是	见注释	flock()。子进程自父进程处继承对同一锁的引用
记录锁	见注释	否	fcntl(F_SETLK)。除非将指代文件的文件描述符标记为执行时关闭，否则会跨越 exec()对锁加以保护
杂项			
地区设置	否	是	setlocale()。作为 C 运行时初始化的一部分，执行新程序后会调用 setlocale(LC_ALL, "C")的等效函数
浮点环境	否	是	运行新程序时，将浮点环境状态重置为默认值，参考 fenv(3)
控制终端	是	是	
退出处理器程序	否	是	atexit()、on_exit()
Linux 特有			
文件系统 ID	见注释	是	setfsuid()、setfsgid()。一旦相应的有效 ID 发生变化，那么这些 ID 也会随之改变
timerfd 定时器	是	见注释	timerfd_create()，子进程继承的文件描述符与父进程指向相同的定时器
能力	见注释	是	capset()。执行 exec()期间对能力的处理一如 39.5 节所述
功能外延集合	是	是	
能力安全位 (securebits) 标志	见注释	是	执行 exec()期间，会保全所有的安全位标志，SECBIT_KEEP_CAPS 除外，总是会清除该标志
CPU 黏性 (affinity)	是	是	sched_setaffinity()
SCHED_RESET_ON_FORK	是	否	参考 35.3.2 节
允许的 CPU	是	是	参考 cpuset(7)手册页
允许的内存节点	是	是	参考 cpuset(7)手册页
内存策略	是	是	参考 set_mempolicy(2)手册页
文件租约	是	见注释	fcntl(F_SETLEASE)。子进程从父进程处继承对相同租约的引用

进 程 属 性	exec()	fork()	影响属性的接口；额外说明
Linux 特有			
目录变更通知	是	否	dnotify API, 通过 fcntl (F_NOTIFY) 来实现支持
prctl (PR_SET_DUMPABLE)	见注释	是	exec()执行期间会设置 PR_SET_DUMPABLE 标志, 执行设置用户或组 ID 程序的情况除外, 此时将清除该标志
prctl (PR_SET_PDEATHSIG)	是	否	
prctl (PR_SET_NAME)	否	是	
oom_adj	是	是	参考 49.9 节
coredump_filter	是	是	参考 22.1 节

28.5 总结

当打开进程记账功能时, 内核会在系统中每一进程终止时将其账单记录写入一个文件。该记录包含进程使用资源的统计数据。

如同函数 fork(), Linux 特有的 clone()系统调用也会创建一个新进程, 但其对父子间的共享属性有更为精确的控制。该系统调用主要用于线程库的实现。

本章对 fork()、vfork()和 clone()的进程创建速度做了比较。尽管 vfork()要快于 fork(), 但较之于子进程随后调用 exec()所耗费的时间, 二者间的时间差异也就微不足道了。

fork()创建的子进程会从其父进程处继承 (有时是共享) 某些进程属性的副本, 而对其他进程属性则不做继承。例如, 子进程继承了父进程文件描述符表和信号处置的副本, 但并不继承父进程的间隔定时器、记录锁或是挂起信号集合。相应地, 进程执行 exec()时, 某些进程属性保持不变, 而会将其他属性重置为缺省值。例如, 进程 ID 保持不变, 文件描述符保持打开 (除非设置了执行时关闭标志), 间隔定时器得以保存, 挂起信号依然挂起, 不过会将已处理信号的处置重置为默认设置, 同时与共享内存段“脱钩”。

更多信息

请参考 24.6 节列出的更多信息来源。[Frisch, 2002]第 17 章描述了对进程记账的管理, 以及不同 UNIX 实现之间的差异。[Bovert & Vesati, 2005]介绍了系统调用 clone()的实现。

28.6 练习

- 28-1.** 编写一程序, 观察 fork()和 vfork()系统调用在读者系统中的速度。要求每个子进程必须立即退出, 而父进程应在创建下一个子进程之前调用 wait(), 等待当前子进程退出。将两个系统调用之间的差别与表 28-3 相比较。shell 内建命令 time 可用来测量程序的执行时间。

第 29 章

线程：介绍

本章及随后几章将讨论 POSIX 线程（thread），亦即 Pthreads。鉴于 Pthreads 范围极广，本书无意涵盖其所有 API。关于线程的深入信息，本章会在结尾处列出其来源。

后续各章将主要描述 Pthreads API 所规定的标准行为。33.5 节则会探讨 Linux 的两种主流线程实现——LinuxThreads 和 Native POSIX Threads Library (NPTL)——与线程标准间的出入。

本章会对线程操作加以概述，随之将述及线程的创建和销毁过程。此外，在应用程序设计中，是选择多线程还是多进程方式？本章将在最后讨论影响这一选择的因素。

29.1 概述

与进程（process）类似¹，线程（thread）是允许应用程序并发执行多个任务的一种机制。如图 29-1 所示，一个进程可以包含多个线程。同一程序中的所有线程均会独立执行相同程序，且共享同一份全局内存区域，其中包括初始化数据段（initialized data）、未初始化数据段（uninitialized data），以及堆内存段（heap segment）。（传统意义上的 UNIX 进程只是多线程程序的一个特例，该进程只包含一个线程。）

图 29-1 其实做了一些简化。特别是，线程栈（thread stack）的位置可能会与共享库和共享内存区域混杂在一起，这取决于创建线程、加载共享库，以及映射共享内存的具体顺序。而且，对于不同的 Linux 发行版，线程栈地址也会有所不同。

同一进程中的多个线程可以并发执行。在多处理器环境下，多个线程可以同时并行。如果一线程因等待 I/O 操作而遭阻塞，那么其他线程依然可以继续运行。（虽然有时单独创建一个专门执行 I/O 操作的线程颇为有用，但采用另一种 I/O 模型则更为可取，第 63 章会对此加以描述。）

对于某些应用而言，线程要优于进程。传统 UNIX 通过创建多个进程来实现并行任务。以网络服务器的设计为例，服务器进程（父进程）在接受客户端的连接后，会调用 fork() 来创建一个单独的子进程，以处理与客户端的通信（可参考 60.3 节）。采用这种设计，服务器就能

¹ 译者注：此处指多进程并发。

同时为多个客户端提供服务。虽然这种方法在很多情境下都屡试不爽，但对于某些应用来说也确实存在如下一些限制。

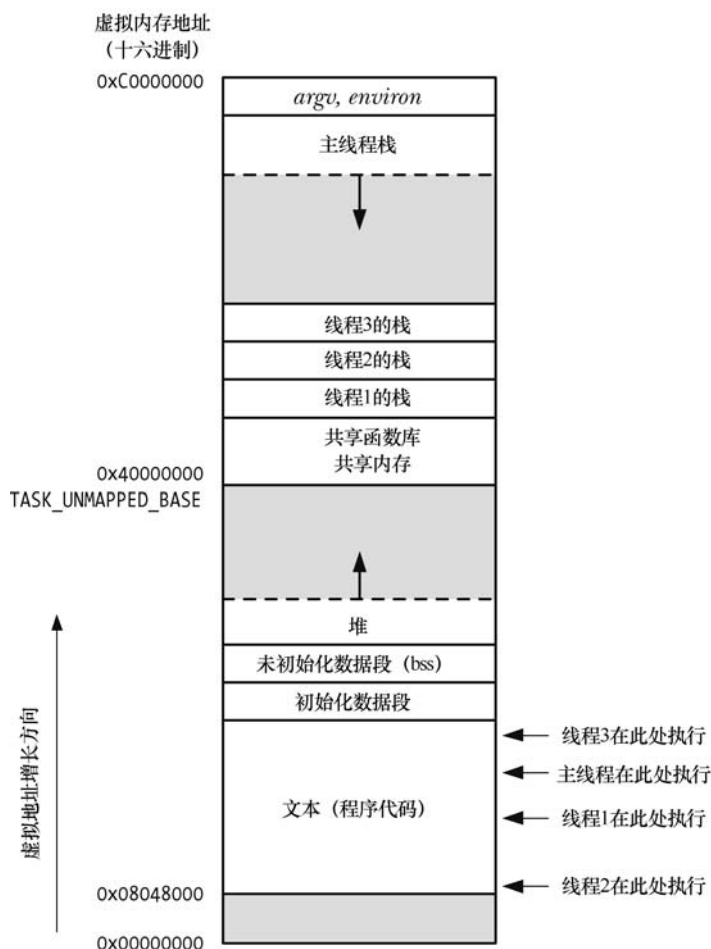


图 29-1: 同时执行 4 个线程的进程 (Linux/x86-32)

- 进程间的信息难以共享。由于除去只读代码段外，父子进程并未共享内存，因此必须采用一些进程间通信 (inter-process communication, 简称 IPC) 方式，在进程间进行信息交换。
- 调用 `fork()` 来创建进程的代价相对较高。即便利用 24.2.2 节所描述的写时复制 (copy-on-write) 技术，仍然需要复制诸如内存页表 (page table) 和文件描述符表 (file descriptor table) 之类的多种进程属性，这意味着 `fork()` 调用在时间上的开销依然不菲。线程解决了上述两个问题。
- 线程之间能够方便、快速地共享信息。只需将数据复制到共享 (全局或堆) 变量中即可。不过，要避免出现多个线程试图同时修改同一份信息的情况，这需要使用第 30 章描述的同步技术。
- 创建线程比创建进程通常要快 10 倍甚至更多。(在 Linux 中，是通过系统调用 `clone()` 来实现线程的，表 28-3 展示了 `fork()` 和 `clone()` 在速度上的差异。) 线程的创建之所以较快，是因为调用 `fork()` 创建子进程时所需复制的诸多属性，在线程间本来就是共享的。特别是，既无需采用写时复制来复制内存页，也无需复制页表。

除了全局内存之外，线程还共享了一干其他属性（这些属性对于进程而言是全局性的，而并非针对某个特定线程），包括以下内容。

- 进程 ID（process ID）和父进程 ID。
- 进程组 ID 与会话 ID（session ID）。
- 控制终端。
- 进程凭证（process credential）（用户 ID 和组 ID）。
- 打开的文件描述符。
- 由 `fcntl()` 创建的记录锁（record lock）。
- 信号（signal）处置。
- 文件系统的相关信息：文件权限掩码（umask）、当前工作目录和根目录。
- 间隔定时器（`setitimer()`）和 POSIX 定时器（`timer_create()`）。
- 系统 V（system V）信号量撤销（undo, `semadj`）值（47.8 节）。
- 资源限制（resource limit）。
- CPU 时间消耗（由 `times()` 返回）。
- 资源消耗（由 `getrusage()` 返回）。
- nice 值（由 `setpriority()` 和 `nice()` 设置）。

各线程所独有的属性，如下列出了其中一部分。

- 线程 ID（thread ID, 29.5 节）。
- 信号掩码（signal mask）。
- 线程特有数据（31.3 节）。
- 备选信号栈（`sigaltstack()`）。
- `errno` 变量。
- 浮点型（floating-point）环境（见 `fenv(3)`）。
- 实时调度策略（real-time scheduling policy）和优先级（35.2 节和 35.3 节）。
- CPU 亲和力（affinity, Linux 所特有, 35.4 节将加以描述）。
- 能力（capability, Linux 所特有, 第 39 章将加以描述）。
- 栈，本地变量和函数的调用链接（linkage）信息。

如图 29-1 所示，所有的线程栈均驻留于同一虚拟地址空间。这也意味着，利用一个合适的指针，各线程可以在对方栈中相互共享数据。这种方法偶尔也能派上用场，但由于局部变量的状态有效与否依赖于其所驻留栈帧的生命周期，故而需要在编程中谨慎处理这一问题。（当函数返回时，该函数栈帧所占用的内存区域有可能为后续的函数调用所重新使用。如果线程终止，那么新线程有可能会对已终止线程的栈所占用的内存空间重新加以利用）。若无法正确处理这一依赖关系，由此而产生的程序 bug 将难以捕获。

29.2 Pthreads API 的详细背景

20 世纪 80 年代末、90 年代初，存在着数种不同的线程接口。1995 年 POSIX.1c 对 POSIX 线程 API 进行了标准化，该标准后来为 SUSv3 所接纳。

有几个概念贯穿整个 Pthreads API，在深入探讨 API 之前，将简单予以介绍。

线程数据类型 (Pthreads data type)

Pthreads API 定义了一干数据类型，表 29-1 列出了其中的一部分。后续章节会对这些数据类型中的绝大部分加以描述。

表 29-1: Pthreads 数据类型

数据类型	描述
pthread_t	线程 ID
pthread_mutex_t	互斥对象 (Mutex)
pthread_mutexattr_t	互斥属性对象
pthread_cond_t	条件变量 (condition variable)
pthread_condattr_t	条件变量的属性对象
pthread_key_t	线程特有数据的键 (Key)
pthread_once_t	一次性初始化控制上下文 (control context)
pthread_attr_t	线程的属性对象

SUSv3 并未规定如何实现这些数据类型，可移植的程序应将其视为“不透明”数据。亦即，程序应避免对此类数据类型变量的结构或内容产生任何依赖。尤其是，不能使用 C 语言的比较操作符 (==) 去比较这些类型的变量。

线程和 errno

在传统 UNIX API 中，errno 是一全局整型变量。然而，这无法满足多线程程序的需要。如果线程调用的函数通过全局 errno 返回错误时，会与其他发起函数调用并检查 errno 的线程混淆在一起。换言之，这将引发竞争条件 (race condition)。因此，在多线程程序中，每个线程都有属于自己的 errno。在 Linux 中，线程特有 errno 的实现方式与大多数 UNIX 实现相类似：将 errno 定义为一个宏，可展开为函数调用，该函数返回一个可修改的左值 (lvalue)，且为每个线程所独有。(因为左值可以修改，多线程程序依然能以 errno=value 的方式对 errno 赋值。)

一言以蔽之，errno 机制在保留传统 UNIX API 报错方式的同时，也适应了多线程环境。

最初的 POSIX.1 标准沿袭 K&R 的 C 语言用法，允许程序将 errno 声明为 extern int errno。SUSv3 却不允许这一做法 (这一变化实际发生于 1995 年的 POSIX.1c 标准之中)。如今，需要声明 errno 的程序必须包含 <errno.h>，以启用对 errno 的线程级实现。

Pthreads 函数返回值

从系统调用和库函数中返回状态，传统的做法是：返回 0 表示成功，返回 -1 表示失败，并设置 errno 以标识错误原因。Pthreads API 则反其道而行之。所有 Pthreads 函数均以返回 0 表示成功，返回一正值表示失败。这一失败时的返回值，与传统 UNIX 系统调用置于 errno 中的值含义相同。

由于多线程程序对 errno 的每次引用都会带来函数调用的开销，因此，本书示例并不会直接将 Pthreads 函数的返回值赋给 errno，而是使用一个中间变量，并利用自己实现的诊断函数 errExitEN() (3.5.2 节)，如下所示：

```
pthread_t *thread;
int s;

s = pthread_create(&thread, NULL, func, &arg);
if (s != 0)
    errExitEN(s, "pthread_create");
```

编译 Pthreads 程序

在 Linux 平台上, 在编译调用了 Pthreads API 的程序时, 需要设置 `cc -pthread` 的编译选项。使用该选项的效果如下。

- 定义 `_REENTRANT` 预处理宏。这会公开对少数可重入 (reentrant) 函数的声明。
- 程序会与库 `libpthread` 进行链接 (等价于 `-lpthread`)。

编译多线程程序时的具体编译选项会因实现及编译器的不同而不同。其他一些实现 (例如 Tru64) 使用 `cc -pthread`, 而 Solaris 和 HP-UX 则使用 `cc -mt`。

29.3 创建线程

启动程序时, 产生的进程只有单条线程, 称之为初始 (initial) 或主 (main) 线程。本节将讨论其他线程的创建过程。

函数 `pthread_create()` 负责创建一条新线程。

```
#include <pthread.h>

int pthread_create(pthread_t *thread, const pthread_attr_t *attr,
                  void *(*start)(void *), void *arg);

Returns 0 on success, or a positive error number on error
```

新线程通过调用带有参数 `arg` 的函数 `start` (即 `start(arg)`) 而开始执行。调用 `pthread_create()` 的线程会继续执行该调用之后的语句。(如 28.2 节所述, 这一行为与 `glibc` 库对系统调用 `clone()` 的包装函数行为相同。)

将参数 `arg` 声明为 `void*` 类型, 意味着可以将指向任意对象的指针传递给 `start()` 函数。一般情况下, `arg` 指向一个全局或堆变量, 也可将其置为 `NULL`。如果需要向 `start()` 传递多个参数, 可以将 `arg` 指向一个结构, 该结构的各个字段则对应于待传递的参数。通过审慎的类型强制转换, `arg` 甚至可以传递 `int` 类型的值。

严格说来, 对于 `int` 与 `void*` 之间相互强制转换的后果, C 语言标准并未加以定义。不过, 大部分 C 语言编译器允许这样的操作, 并且也能达成预期的目的, 即 `int j == (int)((void*)j)`。

`start()` 的返回值类型为 `void*`, 对其使用方式与参数 `arg` 相同。对后续 `pthread_join()` 函数的描述中, 将论及对该返回值的使用方式。

将经强制转换的整型数作为线程 `start` 函数的返回值时, 必须小心谨慎。原因在于, 取消线程 (见第 32 章) 时的返回值 `PTHREAD_CANCELED`, 通常是由实现所定义的整型值, 再经强制转换为 `void*`。若线程某甲的 `start` 函数将此整型值返回给正在执行 `pthread_join()` 操作的线程某乙, 某乙会误认为某甲遭到了取消。应用如果采用了线程取消技术并选择将 `start` 函数的返回值强制转换为整型, 那么就必须确保线程正常结束时的返回值与当前 Pthreads 实现中的 `PTHREAD_CANCELED` 不同。如欲保证程序的可移植性, 则在任何将要运行该应用的实现中,

正常退出线程的返回值应不同于相应的 `PTHREAD_CANCELED` 值。

参数 `thread` 指向 `pthread_t` 类型的缓冲区，在 `pthread_create()` 返回前，会在此保存一个该线程的唯一标识。后续的 Pthreads 函数将使用该标识来引用此线程。

SUSv3 明确指出，在新线程开始执行之前，实现无需对 `thread` 参数所指向的缓冲区进行初始化，即新线程可能会在 `pthread_create()` 返回给调用者之前已经开始运行。如新线程需要获取自己的线程 ID，则只能使用 `pthread_self()`（29.5 节描述）方法。

参数 `attr` 是指向 `pthread_attr_t` 对象的指针，该对象指定了新线程的各种属性。29.8 节将述及其中的部分属性。如果将 `attr` 设置为 `NULL`，那么创建新线程时将使用各种默认属性，本书的大部分示例程序都采用这一做法。

调用 `pthread_create()` 后，应用程序无从确定系统接着会调度哪一个线程来使用 CPU 资源（在多处理器系统中，多个线程可能会在不同 CPU 上同时执行）。程序如隐含了对特定调度顺序的依赖，则无疑会对 24.4 节所述的竞争条件打开方便之门。如果对执行顺序确有强制要求，那么就必须采用第 30 章所描述的同步技术。

29.4 终止线程

可以如下方式终止线程的运行。

- 线程 `start` 函数执行 `return` 语句并返回指定值。
- 线程调用 `pthread_exit()`（详见后述）。
- 调用 `pthread_cancel()` 取消线程（在 32.1 节讨论）。
- 任意线程调用了 `exit()`，或者主线程执行了 `return` 语句（在 `main()` 函数中），都会导致进程中的所有线程立即终止。

`pthread_exit()` 函数将终止调用线程，且其返回值可由另一线程通过调用 `pthread_join()` 来获取。

```
include <pthread.h>

void pthread_exit(void *retval);
```

调用 `pthread_exit()` 相当于在线程的 `start` 函数中执行 `return`，不同之处在于，可在线程 `start` 函数所调用的任意函数中调用 `pthread_exit()`。

参数 `retval` 指定了线程的返回值。`Retval` 所指向的内容不应分配于线程栈中，因为线程终止后，将无法确定线程栈的内容是否有效。（例如，系统可能会立刻将该进程虚拟内存的这片区域重新分配，供一个新的线程栈使用。）出于同样的理由，也不应在线程栈中分配线程 `start` 函数的返回值。

如果主线程调用了 `pthread_exit()`，而非调用 `exit()` 或是执行 `return` 语句，那么其他线程将继续运行。

29.5 线程 ID（Thread ID）

进程内部的每个线程都有一个唯一标识，称为线程 ID。线程 ID 会返回给 `pthread_create()` 的调用者，一个线程可以通过 `pthread_self()` 来获取自己的线程 ID。

```
include <pthread.h>

pthread_t pthread_self(void);
```

Returns the thread ID of the calling thread

线程 ID 在应用程序中非常有用，原因如下。

- 不同的 Pthreads 函数利用线程 ID 来标识要操作的目标线程。这些函数包括 `pthread_join()`、`pthread_detach()`、`pthread_cancel()`和 `pthread_kill()`等，后续章节将会加以讨论。
- 在一些应用程序中，以特定线程的线程 ID 作为动态数据结构的标签，这颇有用处，既可用于识别某个数据结构的创建者或属主线程，又可以确定随后对该数据结构执行操作的具体线程。

函数 `pthread_equal()`可检查两个线程的 ID 是否相同。

```
include <pthread.h>

int pthread_equal(pthread_t t1, pthread_t t2);

Returns nonzero value if t1 and t2 are equal, otherwise 0
```

例如，为了检查调用线程的线程 ID 与保存于变量 `t1` 中的线程 ID 是否一致，可以编写如下代码：

```
if (pthread_equal(tid, pthread_self())
    printf("tid matches self\n");
```

因为必须将 `pthread_t` 作为一种不透明的数据类型加以对待，所以函数 `pthread_equal()`是必须的。Linux 将 `pthread_t` 定义为无符号长整型（`unsigned long`），但在其他实现中，则有可能是一个指针或结构。

在 NPTL 中，`pthread_t` 实际上是一个经强制转化而为无符号长整型的指针。

SUSv3 并未要求将 `pthread_t` 实现为一个标量（`scalar`）类型，该类型也可以是一个结构。因此，下列显示线程 ID 的代码实例并不具有可移植性（尽管该实例在包括 Linux 在内的许多实现上均可正常运行，而且有时在调试程序时还很实用）。

```
pthread_t thr;
```

```
printf("Thread ID = %ld\n", (long) thr); /* WRONG! */
```

在 Linux 的线程实现中，线程 ID 在所有进程中都是唯一的。不过在其他实现中则未必如此，SUSv3 特别指出，应用程序若使用线程 ID 来标识其他进程的线程，其可移植性将无法得到保证。此外，在对已终止线程施以 `pthread_join()`，或者在已分离（`detached`）线程退出后，实现可以复用该线程的线程 ID。（下一节和 29.7 节将分别解释 `pthread_join()`和线程的分离。）

POSIX 线程 ID 与 Linux 专有的系统调用 `gettid()`所返回的线程 ID 并不相同。POSIX 线程 ID 由线程库实现来负责分配和维护。`gettid()`返回的线程 ID 是一个由内核（`Kernel`）分配的数值，类似于进程 ID（`process ID`）。虽然在 Linux NPTL 线程实现中，每个 POSIX 线程都对应一个唯一的内核线程 ID，但应用程序一般无需了解内核线程 ID（况且，如果程序依赖于这一信息，也将无法移植）。

29.6 连接（`joining`）已终止的线程

函数 `pthread_join()`等待由 `thread` 标识的线程终止。（如果线程已经终止，`pthread_join()`会

立即返回)。这种操作被称为连接(joining)。

```
include <pthread.h>

int pthread_join(pthread_t thread, void **retval);
Returns 0 on success, or a positive error number on error
```

若 `retval` 为一非空指针，将会保存线程终止时返回值的拷贝，该返回值亦即线程调用 `return` 或 `pthread_exit()` 时所指定的值。

如向 `pthread_join()` 传入一个之前已然连接过的线程 ID，将会导致无法预知的行为。例如，相同的线程 ID 在参与一次连接后恰好为另一新建线程所重用，再度连接的可能就是这个新线程。

若线程并未分离 (`detached`，见 29.7 节)，则必须使用 `pthread_join()` 来进行连接。如果未能连接，那么线程终止时将产生僵尸线程，与僵尸进程 (`zombie process`) 的概念相类似 (参考 26.2 节)。除了浪费系统资源以外，僵尸线程若累积过多，应用将再也无法创建新的线程。

`pthread_join()` 执行的功能类似于针对进程的 `waitpid()` 调用，不过二者之间存在一些显著差别。

- 线程之间的关系是对等的 (`peers`)。进程中的任意线程均可以调用 `pthread_join()` 与该进程的任何其他线程连接起来。例如，如果线程 A 创建线程 B，线程 B 再创建线程 C，那么线程 A 可以连接线程 C，线程 C 也可以连接线程 A。这与进程间的层次关系不同，父进程如果使用 `fork()` 创建了子进程，那么它也是唯一能够对子进程调用 `wait()` 的进程。调用 `pthread_create()` 创建的新线程与发起调用的线程之间，就没有这样的关系。
- 无法“连接任意线程” (对于进程，则可以通过调用 `waitpid(-1, &status, options)` 做到这一点)，也不能以非阻塞 (`nonblocking`) 方式进行连接 (类似于设置 `WHOHANG` 标志的 `waitpid()`)。使用条件 (`condition`) 变量可以实现类似的功能，30.2.4 节会给出示例。

限制 `pthread_join()` 只能连接特定线程 ID，这样做是“别有用心”的。其用意在于，程序应只能连接它所“知道的”线程。线程之间并无层次关系，如果听任“与任意线程连接”的操作发生，那么所谓“任意”线程就可以包括由库函数私自创建的线程，从而带来问题。(30.2.4 所展示的条件变量技术也只允许线程连接它“知道的”其他线程。) 结果是，函数库在获取线程返回状态时将不再能与该线程连接¹，只会一错再错，试图连接一个已然连接过的线程 ID。换言之，“连接任意线程”的操作与模块化的程序设计理念背道而驰。

示例程序

程序清单 29-1 中的程序创建了一个线程，并与之连接。

程序清单 29-1：一个使用 Pthreads 的简单程序

```
----- threads/simple_thread.c

#include <pthread.h>
#include "tspi_hdr.h"

static void *
threadFunc(void *arg)
{
    char *s = (char *) arg;
```

¹ 译者注：本句的两处线程均指前文由库函数私自创建的线程。

```

    printf("%s", s);

    return (void *) strlen(s);
}
int
main(int argc, char *argv[])
{
    pthread_t t1;
    void *res;
    int s;

    s = pthread_create(&t1, NULL, threadFunc, "Hello world\n");
    if (s != 0)
        errExitEN(s, "pthread_create");

    printf("Message from main()\n");
    s = pthread_join(t1, &res);
    if (s != 0)
        errExitEN(s, "pthread_join");

    printf("Thread returned %ld\n", (long) res);

    exit(EXIT_SUCCESS);
}

```

threads/simple_thread.c

当运行程序清单 29-1 的程序时，可以看到如下输出：

```

$ ./simple_thread
Message from main()
Hello world
Thread returned 12

```

依赖于系统对两个线程的具体调度，第 1 行与第 2 行的输出顺序可能会颠倒过来。

29.7 线程的分离

默认情况下，线程是可连接的(joinable)，也就是说，当线程退出时，其他线程可以通过调用 `pthread_join()` 获取其返回状态。有时，程序员并不关心线程的返回状态，只是希望系统在线程终止时能够自动清理并移除之。在这种情况下，可以调用 `pthread_detach()` 并向 `thread` 参数传入指定线程的标识符，将该线程标记为处于分离 (detached) 状态。

```

#include <pthread.h>

int pthread_detach(pthread_t thread);

Returns 0 on success, or a positive error number on error

```

如下例所示，使用 `pthread_detach()`，线程可以自行分离：

```
pthread_detach(pthread_self());
```

一旦线程处于分离状态，就不能再使用 `pthread_join()` 来获取其状态，也无法使其重返“可连接”状态。

其他线程调用了 `exit()`，或是主线程执行 `return` 语句时，即便遭到分离的线程也还是会受到影响。此时，不管线程处于可连接状态还是已分离状态，进程的所有线程会立即终止。换

言之，`pthread_detach()`只是控制线程终止之后所发生的事情，而非何时或如何终止线程。

29.8 线程属性

前面已然提及 `pthread_create()` 中类型为 `pthread_attr_t` 的 `attr` 参数，可利用其在创建线程时指定新线程的属性。本书无意深入这些属性的细节（关于此类细节，可参考本章结尾处的参考资料列表），也不会将操作 `pthread_attr_t` 对象的各种 Pthreads 函数原型一一列出，只会点出如下之类的一些属性：线程栈的位置和大小、线程调度策略和优先级（类似于 35.2 节和 35.3 节所描述的进程实时调度策略和优先级），以及线程是否处于可连接或分离状态。

作为线程属性的使用示例，程序清单 29-2 中的代码创建了一个新线程，该线程刚一创建即遭分离（而非之后再调用 `pthread_detach()`）。这段代码首先以缺省值对线程属性结构进行初始化，接着为创建分离线程而设置属性，最后再以此线程属性结构来创建新线程。线程一旦创建，就无需再保留该属性对象，故而程序将其销毁。

程序清单 29-2：使用分离属性创建线程

```
----- from threads/detached_attrib.c

pthread_t thr;
pthread_attr_t attr;
int s;

s = pthread_attr_init(&attr);           /* Assigns default values */
if (s != 0)
    errExitEN(s, "pthread_attr_init");

s = pthread_attr_setdetachstate(&attr, PTHREAD_CREATE_DETACHED);
if (s != 0)
    errExitEN(s, "pthread_attr_setdetachstate");

s = pthread_create(&thr, &attr, threadFunc, (void *) 1);
if (s != 0)
    errExitEN(s, "pthread_create");

s = pthread_attr_destroy(&attr);       /* No longer needed */
if (s != 0)
    errExitEN(s, "pthread_attr_destroy");

----- from threads/detached_attrib.c
```

29.9 线程 VS 进程

将应用程序实现为一组线程还是进程？本节将简单考虑一下可能影响这一决定的部分因素。先从多线程方法的优点开始。

- 线程间的数据共享很简单。相形之下，进程间的数据共享需要更多的投入。（例如，创建共享内存段或者使用管道 `pipe`）。
- 创建线程要快于创建进程。线程间的上下文切换（`context-switch`），其消耗时间一般也比进程要短。

线程相对于进程的一些缺点如下所示。

- 多线程编程时，需要**确保调用线程安全（thread-safe）的函数**，或者以线程安全的方式

式来调用函数。(31.1 节将讨论线程安全的概念。)多进程应用则无需关注这些。

- 某个线程中的 bug (例如, 通过一个错误的指针来修改内存) 可能会危及该进程的所有线程, 因为它们共享着相同的地址空间和其他属性。相比之下, 进程间的隔离更彻底。
- 每个线程都在争用宿主进程 (host process) 中有限的虚拟地址空间。特别是, 一旦每个线程栈以及线程特有数据 (或线程本地存储) 消耗掉进程虚拟地址空间的一部分, 则后续线程将无缘使用这些区域。虽然有效地址空间很大 (例如, 在 x86-32 平台上通常有 3GB), 但当进程分配大量线程, 亦或线程使用大量内存时, 这一因素的限制作用也就突显出来。与之相反, 每个进程都可以使用全部的有效虚拟内存, 仅受制于实际内存和交换 (swap) 空间。

影响选择的还有如下几点。

- 在多线程应用中处理信号, 需要小心设计。(作为通则, 一般建议在三线程程序中避免使用信号。)关于线程与信号, 33.2 节会做深入讨论。
- 在多线程应用中, 所有线程必须运行同一个程序 (尽管可能是位于不同函数中)。对于多进程应用, 不同的进程可以运行不同的程序。
- 除了数据, 线程还可以共享某些其他信息 (例如, 文件描述符、信号处置、当前工作目录, 以及用户 ID 和组 ID)。优劣之判, 视应用而定。

29.10 总结

在多线程程序中, 多个线程并发执行同一程序。所有线程共享相同的全局和堆变量, 但每个线程都配有用来存放局部变量的私有栈。同一进程中的线程还共享若干其他属性, 包括进程 ID、打开的文件描述符、信号处置、当前工作目录以及资源限制。

线程与进程间的关键区别在于, 线程比进程更易于共享信息, 这也是许多应用程序舍进程而取线程的主要原因。对于某些操作来说 (例如, 创建线程比创建进程快), 线程还可以提供更好的性能。但是, 在程序设计的进程/线程之争中, 这往往不会是决定性因素。

可使用 `pthread_create()` 来创建线程。每个线程随后可调用 `pthread_exit()` 独立退出。(如有任一线程调用了 `exit()`, 那么所有线程将立即终止。)除非将线程标记为分离状态 (例如通过调用 `pthread_detach()`), 其他线程要连接该线程, 则必须使用 `pthread_join()`, 由其返回遭连接线程的退出状态。

进阶信息

[Butenhof, 1996]对 Pthreads 进行了透彻而清晰的阐述。[Robbins & Robbins, 2003]对 Pthreads 的各方面都有涉及。[Tanen-baum, 2007]对线程概念的介绍更具理论化, 涵盖主题包括互斥量 (mutex)、临界区 (critical region)、条件变量以及死锁 (deadlock) 检测及规避。[Vahalia, 1996]提供了线程实现的背景知识。

29.11 练习

29-1. 若一线程执行了如下代码, 可能会产生什么结果?

```
pthread_join(pthread_self(), NULL);
```

在 Linux 上编写一个程序, 观察一下实际会发生什么情况。假设代码中有一变量 `tid`, 其中包含了某个线程 ID, 在自身发起 `pthread_join(tid, NULL)` 调用时, 要避免造成

与上述语句相同的后果，该线程应采取何种措施？

29-2. 除了缺少错误检查，以及对各种变量和结构的声明外，下列程序还有什么问题？

```
static void *
threadFunc(void *arg)
{
    struct someStruct *pbuf = (struct someStruct *) arg;

    /* Do some work with structure pointed to by 'pbuf' */
}

int
main(int argc, char *argv[])
{
    struct someStruct buf;

    pthread_create(&thr, NULL, threadFunc, (void *) &buf);
    pthread_exit(NULL);
}
```

第 30 章

线程：线程同步

本章介绍线程用来同步彼此行为的两个工具：互斥量（`mutex`）和条件变量（`condition variable`）。互斥量可以帮助线程同步对共享资源的使用，以防如下情况发生：线程某甲试图访问一共享变量时，线程某乙正在对其进行修改。条件变量则是在此之外的拾遗补缺，允许线程相互通知共享变量（或其他共享资源）的状态发生了变化。

30.1 保护对共享变量的访问：互斥量

线程的主要优势在于，能够通过全局变量来共享信息。不过，这种便捷的共享是有代价的：必须确保多个线程不会同时修改同一变量，或者某一线程不会读取正由其他线程修改的变量。术语临界区（`critical section`）是指访问某一共享资源的代码片段，并且这段代码的执行应为原子（`atomic`）操作，亦即，同时访问同一共享资源的其他线程不应中断该片段的执行。

程序清单 30-1 中的简单示例，展示了以非原子方式访问共享资源时所发生的问题。该程序创建了两个线程，且均执行同一函数。该函数执行一个循环，重复以下步骤：将 `glob` 复制到本地变量 `loc` 中，然后递增 `loc`，再把 `loc` 复制回 `glob`，以此不断增加全局变量 `glob` 的值。因为 `loc` 是分配于线程栈中的自动变量（`automatic variable`），所以每个线程都有一份。循环重复的次数要么由命令行参数指定，要么取默认值。

程序清单 30-1：两线程以错误方式递增全局变量的值

```
threads/thread_incr.c
#include <pthread.h>
#include "tspi_hdr.h"

static int glob = 0;

static void *
threadFunc(void *arg)
{
    int loops = *((int *) arg);
    int loc, j;
```

```

    for (j = 0; j < loops; j++) {
        loc = glob;
        loc++;
        glob = loc;
    }

    return NULL;
}

int
main(int argc, char *argv[])
{
    pthread_t t1, t2;
    int loops, s;

    loops = (argc > 1) ? getInt(argv[1], GN_GT_0, "num-loops") : 10000000;

    s = pthread_create(&t1, NULL, threadFunc, &loops);
    if (s != 0)
        errExitEN(s, "pthread_create");
    s = pthread_create(&t2, NULL, threadFunc, &loops);
    if (s != 0)
        errExitEN(s, "pthread_create");

    s = pthread_join(t1, NULL);
    if (s != 0)
        errExitEN(s, "pthread_join");
    s = pthread_join(t2, NULL);
    if (s != 0)
        errExitEN(s, "pthread_join");

    printf("glob = %d\n", glob);
    exit(EXIT_SUCCESS);
}

```

threads/thread_incr.c

运行程序清单 30-1 中的示例，并指定每个线程均对该变量递增 1000 次，看起来一切正常。

```

$ ./thread_incr 1000
glob = 2000

```

不过，很有可能会发生如下情况：在线程某乙尚未得以运行时，线程某甲已经执行完毕并且退出了。如果加大每个线程的工作量，结果将完全不同。

```

$ ./thread_incr 10000000
glob = 16517656

```

执行到最后，glob 的值本应为 2000 万。问题的原因是由于如下的执行序列（参见图 30-1）。

1. 线程 1 将 glob 值赋给局部变量 loc。假设 blob 的当前值为 2000。
2. 线程 1 的时间片期满，线程 2 开始执行。
3. 线程 2 执行多次循环：将全局变量 glob 的值置于局部变量 loc，递增 loc，再将结果写回变量 glob。第 1 次循环时，glob 的值为 2000。假设线程 2 的时间片到期时，glob 的值已经增至 3000。
4. 线程 1 获得另一时间片，并从上次停止处恢复执行。线程 1 在上次运行时，已将 glob 的值（2000）赋给 loc，现在递增 loc，再将 loc 的值 2001 写回 glob。此时，线程 2 此前递增操作的结果遭到覆盖。

如果使用同样的命令行参数将该程序运行多次，glob 的值会波动很大：

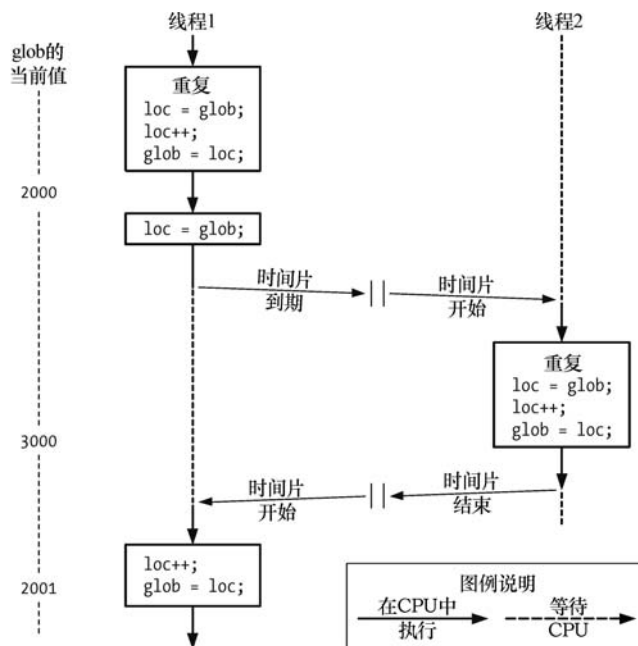


图 30-1: 两个线程不使用同步技术递增全局变量的值

```
$ ./thread_incr 1000000
glob = 10880429
$ ./thread_incr 10000000
glob = 13493953
```

这一行为的不确定性，实应归咎于内核 CPU 调度决定的难以预见。若在复杂程序中发生这一不确定行为，则意味着此类错误将偶尔发作，难以重现，因此也很难发现。

使用如下语句，将程序清单 30-1 中函数 `threadFunc()` 内 `for` 循环中的 3 条语句加以替换，似乎可以解决这一问题：

```
glob++; /* or: ++glob; */
```

不过，在很多硬件架构上（例如，RISC 系统），编译器依然会将这条语句转换成机器码，其执行步骤仍旧等同于 `threadFunc` 循环内的 3 条语句。换言之，尽管 C 语言的递增符看似简单，其操作也未必就属于原子操作，依然可能发生上述行为。

为避免线程更新共享变量时所出现问题，必须使用互斥量（`mutex` 是 `mutual exclusion` 的缩写）来确保同时仅有一个线程可以访问某项共享资源。更为全面的说法是，可以使用互斥量来保证对任意共享资源的原子访问，而保护共享变量是其最常见的用法。

互斥量有两种状态：已锁定（`locked`）和未锁定（`unlocked`）。任何时候，至多只有一个线程可以锁定该互斥量。试图对已经锁定的某一互斥量再次加锁，将可能阻塞线程或者报错失败，具体取决于加锁时使用的方法。

一旦线程锁定互斥量，随即成为该互斥量的所有者。只有所有者才能给互斥量解锁。这一属性改善了使用互斥量的代码结构，也顾及到对互斥量实现的优化。因为所有权的关系，有时会使用术语获取（`acquire`）和释放（`release`）来替代加锁和解锁。

一般情况下，对每一共享资源（可能由多个相关变量组成）会使用不同的互斥量，每一线程在访问同一资源时将采用如下协议。

- 针对共享资源锁定互斥量。

- 访问共享资源。
- 对互斥量解锁。

如果多个线程试图执行这一代码块（一个临界区），事实上只有一个线程能够持有该互斥量（其他线程将遭到阻塞），即同时只有一个线程能够进入这段代码区域，如图 30-2 所示。

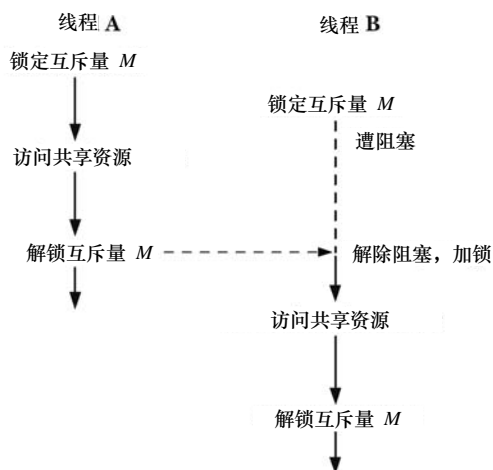


图 30-2: 使用互斥量来保护临界区

最后请注意，使用互斥锁仅是一种建议，而非强制。亦即，线程可以考虑不使用互斥量而仅访问相应的共享变量。为了安全地处理共享变量，所有线程在使用互斥量时必须互相协调，遵守既定的锁定规则。

30.1.1 静态分配的互斥量

互斥量既可以像静态变量那样分配，也可以在运行时动态创建（例如，通过 `malloc()` 在一块内存中分配）。动态互斥量的创建稍微有些复杂，将延后至 30.1.5 节再做讨论。

互斥量是属于 `pthread_mutex_t` 类型的变量。在使用之前必须对其初始化。对于静态分配的互斥量而言，可如下例所示，将 `PTHREAD_MUTEX_INITIALIZER` 赋给互斥量。

```
pthread_mutex_t mtx = PTHREAD_MUTEX_INITIALIZER;
```

依照 SUSv3 的规定，对某一互斥量的副本（copy）执行本节（30.1 节）后续所述的操作将导致未定义的结果。此类操作只能施之于如下两类互斥量的“真身”，经由 `PTHREAD_MUTEX_INITIALIZER` 初始化的静态互斥量或者经由 `pthread_mutex_init()`（在 30.1.5 节讨论）初始化的动态互斥量。

30.1.2 加锁和解锁互斥量

初始化之后，互斥量处于未锁定状态。函数 `pthread_mutex_lock()` 可以锁定某一互斥量，而函数 `pthread_mutex_unlock()` 则可以将一个互斥量解锁。

```
#include <pthread.h>
```

```
int pthread_mutex_lock(pthread_mutex_t *mutex);
int pthread_mutex_unlock(pthread_mutex_t *mutex);
```

Both return 0 on success, or a positive error number on error

要锁定互斥量，在调用 `pthread_mutex_lock()` 时需要指定互斥量。如果互斥量当前处于未锁定状态，该调用将锁定互斥量并立即返回。如果其他线程已经锁定了这一互斥量，那么 `pthread_mutex_lock()` 调用会一直堵塞，直至该互斥量被解锁，到那时，调用将锁定互斥量并返回。

如果发起 `pthread_mutex_lock()` 调用的线程自身之前已然将目标互斥量锁定，对于互斥量的默认类型而言，可能会产生两种后果——视具体实现而定：线程陷入死锁（`deadlock`），因试图锁定已为自己所持有的互斥量而遭到阻塞；或者调用失败，返回 `EDEADLK` 错误。在 Linux 上，默认情况下线程会发生死锁。（30.1.7 节在讨论互斥量类型时会述及一些其他的可能行为。）

函数 `pthread_mutex_unlock()` 将解锁之前已遭调用线程锁定的互斥量。以下行为均属错误：对处于未锁定状态的互斥量进行解锁，或者解锁由其他线程锁定的互斥量。

如果有不止一个线程在等待获取由函数 `pthread_mutex_unlock()` 解锁的互斥量，则无法判断究竟哪个线程将如愿以偿。

示例程序

程序清单 30-2 是对程序清单 30-1 的修改，使用了一个互斥量来保护对全局变量 `glob` 的访问。使用与之前类似的命令行来运行这个改版程序，可以看到对 `glob` 的累加总是能够保持正确。

```
$ ./thread_incr_mutex 10000000
glob = 20000000
```

程序清单 30-2：使用互斥量保护对全局变量的访问

```
----- threads/thread_incr_mutex.c
#include <pthread.h>
#include "tspi_hdr.h"

static int glob = 0;
static pthread_mutex_t mtx = PTHREAD_MUTEX_INITIALIZER;

static void *
threadFunc(void *arg)
{
    /* Loop 'arg' times incrementing 'glob' */
    int loops = *((int *) arg);
    int loc, j, s;

    for (j = 0; j < loops; j++) {
        s = pthread_mutex_lock(&mtx);
        if (s != 0)
            errExitEN(s, "pthread_mutex_lock");
        loc = glob;
        loc++;
        glob = loc;

        s = pthread_mutex_unlock(&mtx);
        if (s != 0)
            errExitEN(s, "pthread_mutex_unlock");
    }

    return NULL;
}

int
main(int argc, char *argv[])
{
```

```

pthread_t t1, t2;
int loops, s;

loops = (argc > 1) ? getInt(argv[1], GN_GT_0, "num-loops") : 1000000;

s = pthread_create(&t1, NULL, threadFunc, &loops);
if (s != 0)
    errExitEN(s, "pthread_create");
s = pthread_create(&t2, NULL, threadFunc, &loops);
if (s != 0)
    errExitEN(s, "pthread_create");

s = pthread_join(t1, NULL);
if (s != 0)
    errExitEN(s, "pthread_join");
s = pthread_join(t2, NULL);
if (s != 0)
    errExitEN(s, "pthread_join");

printf("glob = %d\n", glob);
exit(EXIT_SUCCESS);
}

```

threads/thread_incr_mutex.c

pthread_mutex_trylock()和 pthread_mutex_timedlock()

Pthreads API 提供了 `pthread_mutex_lock()`函数的两个变体：`pthread_mutex_trylock()`和 `pthread_mutex_timedlock()`。可参考手册页（manual page）获取这些函数的原型。

如果信号量已然锁定，对其执行函数 `pthread_mutex_trylock()`会失败并返回 `EBUSY` 错误，除此之外，该函数与 `pthread_mutex_lock()`行为相同。

除了调用者可以指定一个附加参数 `abstime`（设置线程等待获取互斥量时休眠的时间限制）外，函数 `pthread_mutex_timedlock()`与 `pthread_mutex_lock()`没有差别。如果参数 `abstime` 指定的时间间隔期满，而调用线程又没有获得对互斥量的所有权，那么函数 `pthread_mutex_timedlock()`返回 `ETIMEDOUT` 错误。

函数 `pthread_mutex_trylock()` 和 `pthread_mutex_timedlock()`比 `pthread_mutex_lock()`的使用频率要低很多。在大多数经过良好设计的应用程序中，线程对互斥量的持有时间应尽可能短，以避免妨碍其他线程的并发执行。这也保证了遭堵塞的其他线程可以很快获取对互斥量的锁定。若某一线程使用 `pthread_mutex_trylock()`周期性地轮询是否可以对互斥量加锁，则有可能要承担这样的风险：当队列中的其他线程通过调用 `pthread_mutex_lock()`相继获得对互斥量的访问时，该线程将始终与此互斥量无缘。

30.1.3 互斥量的性能

使用互斥量的开销有多大？前面已经展示了递增共享变量程序的两个不同版本：没有使用互斥量的程序清单 30-1 和使用互斥量的程序清单 30-2。在 x86-32 架构的 Linux 2.6.31（含 NPTL）系统下运行这两个程序，如令单一线程循环 1000 万次，前者共花费了 0.35 秒（并产生错误结果），而后者则需要 3.1 秒。

乍看起来，代价极高。不过，考虑一下前者（程序清单 30-1）执行的主循环。在该版本中，函数 `threadFunc()`于 `for` 循环中，先递增循环控制变量，再将其与另一变量进行比较，随后执行两个复制操作和一个递增操作，最后返回循环起始处开始下一次循环。而后者——使用互斥量

的版本（程序清单 30-2）执行了相同步骤，不过在每次循环的前后多了加锁和解锁互斥量的工作。换言之，对互斥量加锁和解锁的开销略低于第 1 个程序的 10 次循环操作。成本相对比较低廉。此外，在通常情况下，线程会花费更多时间去做其他工作，对互斥量的加锁和解锁操作相对要少得多，因此使用互斥量对于大部分应用程序的性能并无显著影响。

进而言之，在相同系统上运行一些简单的测试程序，结果显示，如将使用函数 `fcntl()`（见 55.3 节）加锁、解锁一片文件区域的代码循环 2000 万次，需耗时 44 秒，而将对系统 V 信号量（`semaphore`）（见 47 章）的递增和递减代码循环 2000 万次，则需要 28 秒。文件锁和信号量的问题在于，其锁定和解锁总是需要发起系统调用（`system call`），而每个系统调用的开销虽小，但颇为可观（见 3.1 节）。与之相反，互斥量的实现采用了机器语言级的原子操作（在内存中执行，对所有线程可见），只有发生锁的争用时才会执行系统调用。

Linux 上，互斥量的实现采用了 `futex`（源自“快速用户空间互斥量”[`fast user space mutex`]的首字母缩写），而对锁争用的处理则使用了 `futex()` 系统调用。本书无意描述 `futex`，其设计意图也并非供用户空间（`user space`）应用程序直接使用，不过[Drepper, 2004(a)]给出了详细描述，还讨论了如何使用 `futexes` 来实现互斥量。[Franke et al., 2002]是一篇由 `futex` 开发人员所撰写的论文（已经过时），介绍了 `futex` 的早期实现以及因其所带来的性能提升。

30.1.4 互斥量的死锁

有时，一个线程需要同时访问两个或更多不同的共享资源，而每个资源又都由不同的互斥量管理。当超过一个线程加锁同一组互斥量时，就有可能发生死锁。图 30-3 展示了一个死锁的例子，其中每个线程都成功地锁住一个互斥量，接着试图对已为另一线程锁定的互斥量加锁。两个线程将无限期地等待下去。

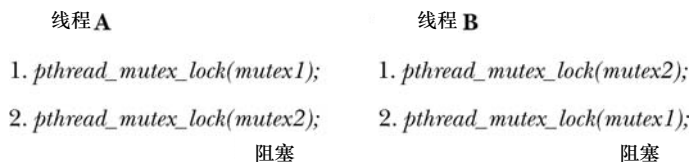


图 30-3: 两个线程分别锁定两个互斥量所导致的死锁

要避免此类死锁问题，最简单的方法是定义互斥量的层级关系。当多个线程对一组互斥量操作时，总是应该以相同顺序对该组互斥量进行锁定。例如，在图 30-3 所示场景中，如果两个线程总是先锁定 `mutex1` 再锁定 `mutex2`，死锁就不会出现。有时，互斥量间的层级关系逻辑清晰。不过，即便没有，依然可以设计出所有线程都必须遵循的强制层级顺序。

另一种方案的使用频率较低，就是“尝试一下，然后恢复”。在这种方案中，线程先使用函数 `pthread_mutex_lock()` 锁定第 1 个互斥量，然后使用函数 `pthread_mutex_trylock()` 来锁定其余互斥量。如果任一 `pthread_mutex_trylock()` 调用失败（返回 `EBUSY`），那么该线程将释放所有互斥量，也许经过一段时间间隔，从头再试。较之于按锁的层级关系来规避死锁，这种方法效率要低一些，因为可能需要历经多次循环。另一方面，由于无需受制于严格的互斥量层级关系，该方法也更为灵活。[Butenhof, 1996]中载有这一方案的范例。

30.1.5 动态初始化互斥量

静态初始值 `PTHREAD_MUTEX_INITIALIZER`，只能用于对如下互斥量进行初始化：经由静态分配且携带默认属性。其他情况下，必须调用 `pthread_mutex_init()` 对互斥量进行动态初始化。


```
#include <pthread.h>

int pthread_mutex_init(pthread_mutex_t *mutex, const pthread_mutexattr_t *attr);
Returns 0 on success, or a positive error number on error
```

参数 `mutex` 指定函数执行初始化操作的目标互斥量。参数 `attr` 是指向 `pthread_mutexattr_t` 类型对象的指针，该对象在函数调用之前已经过了初始化处理，用于定义互斥量的属性。（下节会介绍更多互斥量属性。）若将 `attr` 参数置为 `NULL`，则该互斥量的各种属性会取默认值。

SUSv3 规定，初始化一个业已初始化的互斥量将导致未定义的行为，应当避免这一行为。

在如下情况下，必须使用函数 `pthread_mutex_init()`，而非静态初始化互斥量。

- 动态分配于堆中的互斥量。例如，动态创建针对某一结构的链表，表中每个结构都包含一个 `pthread_mutex_t` 类型的字段来存放互斥量，借以保护对该结构的访问。
- 互斥量是在栈中分配的自动变量。
- 初始化经由静态分配，且不使用默认属性的互斥量。

当不再需要经由自动或动态分配的互斥量时，应使用 `pthread_mutex_destroy()` 将其销毁。（对于使用 `PTHREAD_MUTEX_INITIALIZER` 静态初始化的互斥量，无需调用 `pthread_mutex_destroy()`。）

```
#include <pthread.h>

int pthread_mutex_destroy(pthread_mutex_t *mutex);
Returns 0 on success, or a positive error number on error
```

只有当互斥量处于未锁定状态，且后续也无任何线程企图锁定它时，将其销毁才是安全的。若互斥量驻留于动态分配的一片内存区域中，应在释放（`free`）此内存区域前将其销毁。对于自动分配的互斥量，也应在宿主函数返回前将其销毁。

经由 `pthread_mutex_destroy()` 销毁的互斥量，可调用 `pthread_mutex_init()` 对其重新初始化。

30.1.6 互斥量的属性

如前所述，可以在 `pthread_mutex_init()` 函数的 `arg` 参数中指定 `pthread_mutexattr_t` 类型对象，对互斥量的属性进行定义。通过 `pthread_mutexattr_t` 类型对象对互斥量属性进行初始化和读取操作的 `Pthreads` 函数有多个。本书不打算深入讨论互斥量属性的细节，也不会将初始化 `pthread_mutexattr_t` 对象内属性的各种函数原型一一列出。不过，下一节会讨论互斥量的属性之一：类型。

30.1.7 互斥量类型

前面几页对互斥量的行为做了若干论述。

- 同一线程不应为同一互斥量加锁两次。
- 线程不应为不为自己所拥有的互斥量解锁（亦即，尚未锁定互斥量）。
- 线程不应为一尚未锁定的互斥量做解锁动作。

准确地说，上述情况的结果将取决于互斥量类型（`type`）。SUSv3 定义了以下互斥量类型。

PTHREAD_MUTEX_NORMAL

该类型的互斥量不具有死锁检测（自检）功能。如线程试图对已由自己锁定的互斥量加

锁，则发生死锁。互斥量处于未锁定状态，或者已由其他线程锁定，对其解锁会导致不确定的结果。（在 Linux 上，对这类互斥量的上述两种操作都会成功。）

PTHREAD_MUTEX_ERRORCHECK

对此类互斥量的所有操作都会执行错误检查。所有上述 3 种情况都会导致相关 Pthreads 函数返回错误。这类互斥量运行起来比一般类型要慢，不过可将其作为调试工具，以发现程序在哪里违反了互斥量使用的基本原则。

PTHREAD_MUTEX_RECURSIVE

递归互斥量维护有一个锁计数器。当线程第 1 次取得互斥量时，会将锁计数器置 1。后续由同一线程执行的每次加锁操作会递增锁计数器的数值，而解锁操作则递减计数器计数。只有当锁计数器值降至 0 时，才会释放（release，亦即可为其他线程所用）该互斥量。解锁时如目标互斥量处于未锁定状态，或是已由其他线程锁定，操作都会失败。

Linux 的线程实现针对以上各种类型的互斥量提供了非标准的静态初始值（例如，PTHREAD_RECURSIVE_MUTEX_INITIALIZER_NP），以便对那些通过静态分配的互斥量进行初始化，而无需使用 pthread_mutex_init() 函数。不过，为保证程序的可移植性，应该避免使用这些初始值。

除了上述类型，SUSv3 还定义了 PTHREAD_MUTEX_DEFAULT 类型。使用 PTHREAD_MUTEX_INITIALIZER 初始化的互斥量，或是经调用参数 attr 为 NULL 的 pthread_mutex_init() 函数所创建的互斥量，都属于此类型。至于该类型互斥量在本节开始处 3 个场景中的行为，规范有意未作定义，意在为互斥量的高效实现保留最大的灵活性。Linux 上，PTHREAD_MUTEX_DEFAULT 类型互斥量的行为与 PTHREAD_MUTEX_NORMAL 类型相仿。

程序清单 30-3 演示了如何设置互斥量类型，本例创建了一个带有错误检查属性（error-checking）的互斥量。

程序清单 30-3：设置互斥量类型

```
pthread_mutex_t mtx;
pthread_mutexattr_t mtxAttr;
int s, type;

s = pthread_mutexattr_init(&mtxAttr);
if (s != 0)
    errExitEN(s, "pthread_mutexattr_init");
s = pthread_mutexattr_settype(&mtxAttr, PTHREAD_MUTEX_ERRORCHECK);
if (s != 0)
    errExitEN(s, "pthread_mutexattr_settype");

s = pthread_mutex_init(mtx, &mtxAttr);
if (s != 0)
    errExitEN(s, "pthread_mutex_init");

s = pthread_mutexattr_destroy(&mtxAttr);      /* No longer needed */
if (s != 0)
    errExitEN(s, "pthread_mutexattr_destroy");
```

30.2 通知状态的改变：条件变量（Condition Variable）

互斥量防止多个线程同时访问同一共享变量。条件变量允许一个线程就某个共享变量（或

其他共享资源)的状态变化通知其他线程,并让其他线程等待(堵塞于)这一通知。

一个未使用条件变量的简单例子有助于展示条件变量的重要性。假设由若干线程生成一些“产品单元(result unit)”供主线程消费。还使用了一个由互斥量保护的变量 `avail` 来代表待消费产品的数量:

```
static pthread_mutex_t mtx = PTHREAD_MUTEX_INITIALIZER;

static int avail = 0;
本节引用的代码片段摘自于随本书发布的源代码文件 threads/prod_no_condvar.c。
生产者线程的源代码如下:
/* Code to produce a unit omitted */

s = pthread_mutex_lock(&mtx);
if (s != 0)
    errExitEN(s, "pthread_mutex_lock");

avail++; /* Let consumer know another unit is available */

s = pthread_mutex_unlock(&mtx);
if (s != 0)
    errExitEN(s, "pthread_mutex_unlock");
主线程(消费者)的代码如下:
for (;;) {
    s = pthread_mutex_lock(&mtx);
    if (s != 0)
        errExitEN(s, "pthread_mutex_lock");

    while (avail > 0) { /* Consume all available units */
        /* Do something with produced unit */
        avail--;
    }

    s = pthread_mutex_unlock(&mtx);
    if (s != 0)
        errExitEN(s, "pthread_mutex_unlock");
}
```

上述代码虽然可行,但由于主线程不停地循环检查变量 `avail` 的状态,故而造成 CPU 资源的浪费。采用了条件变量(condition variable),这一问题就迎刃而解:允许一个线程休眠(等待)直至接获另一线程的通知(收到信号)去执行某些操作(例如,出现一些“情况”后,等待者必须立即做出响应)。

条件变量总是结合互斥量使用。条件变量就共享变量的状态改变发出通知,而互斥量则提供对该共享变量访问的互斥(mutual exclusion)。这里使用的术语“信号”(signal),与第 20 章至第 22 章所述信号(signal)无关,而是发出信号的意思。

30.2.1 由静态分配的条件变量

如同互斥量一样,条件变量的分配,有静态和动态之分。条件变量的动态创建延后到 30.2.5 节再行描述,这里先讨论一下静态分配。

条件变量的数据类型是 `pthread_count_t`。类似于互斥量,使用条件变量前必须对其初始化。对于经由静态分配的条件变量,将其赋值为 `PTHREAD_COND_INITIALIZER` 即完成初始化操作。可参考下面的例子:

```
pthread_cond_t cond = PTHREAD_COND_INITIALIZER;
```

依据 SUSv3 规定，将本节后续所描述的操作施之于一个条件变量的副本（copy）时，其结果未定义。所有操作仅能针对条件变量的原本执行，要么经由 PTHREAD_COND_INITIALIZER 进行了静态初始化，要么使用 pthread_cond_init() 做了动态初始化（30.2.5 节描述）处理。

30.2.2 通知和等待条件变量

条件变量的主要操作是发送信号（signal）和等待（wait）。发送信号操作即通知一个或多个处于等待状态的线程，某个共享变量的状态已经改变。等待操作是指在收到一个通知前一直处于阻塞状态。

函数 pthread_cond_signal() 和 pthread_cond_broadcast() 均可针对由参数 cond 所指定的条件变量而发送信号。pthread_cond_wait() 函数将阻塞一线程，直至收到条件变量 cond 的通知。

```
#include <pthread.h>

int pthread_cond_signal(pthread_cond_t *cond);
int pthread_cond_broadcast(pthread_cond_t *cond);
int pthread_cond_wait(pthread_cond_t *cond, pthread_mutex_t *mutex);

All return 0 on success, or a positive error number on error
```

函数 pthread_cond_signal() 和 pthread_cond_broadcast() 之间的差别在于，二者对阻塞于 pthread_cond_wait() 的多个线程处理方式不同。pthread_cond_signal() 函数只保证唤醒至少一条遭到阻塞的线程，而 pthread_cond_broadcast() 则会唤醒所有遭阻塞的线程。

使用函数 pthread_cond_broadcast() 总能产生正确结果（因为所有线程应都能处理多余和虚假的唤醒动作），但函数 pthread_cond_signal() 会更为高效。不过，只有当仅需唤醒一条（且无论是其中哪条）等待线程来处理共享变量的状态变化时，才应使用 pthread_cond_signal()。应用这种方式的典型情况是，所有等待线程都在执行完全相同的任务。基于这些假设，函数 pthread_cond_signal() 会比 pthread_cond_broadcast() 更具效率，因为这可以避免发生如下情况。

1. 同时唤醒所有等待线程。
2. 某一线程首先获得调度。此线程检查了共享变量的状态（在相关互斥量的保护之下），发现还有任务需要完成。该线程执行了所需工作，并改变共享变量状态，以表明任务完成，最后释放对相关互斥量的锁定。
3. 剩余的每个线程轮流锁定互斥量并检测共享变量的状态。不过，由于第一个线程所做的工作，余下的线程发现无事可做，随即解锁互斥量转而休眠（即再次调用 pthread_cond_wait()）。

相形之下，函数 pthread_cond_broadcast() 所处理的情况是：处于等待状态的所有线程执行的任务不同（即各线程关联于条件变量的判定条件不同）。

条件变量并不保存状态信息，只是传递应用程序状态信息的一种通讯机制。发送信号时若无任何线程在等待该条件变量，这个信号也就会不了了之。线程如在此后等待该条件变量，只有当再次收到此变量的下一信号时，方可解除阻塞状态。

函数 pthread_cond_timedwait() 与函数 pthread_cond_wait() 几近相同，唯一的区别在于，由参数 abstime 来指定一个线程等待条件变量通知时休眠时间的上限。

```
#include <pthread.h>
```

```
int pthread_cond_timedwait(pthread_cond_t *cond, pthread_mutex_t *mutex,
                           const struct timespec *abstime);
```

Returns 0 on success, or a positive error number on error

参数 `abstime` 是一个 `timespec` 类型的结构（见 23.4.2 节），用以指定自 Epoch（参考 10.1 节）以来以秒和纳秒（nanosecond）为单位表示的绝对（absolute）时间。如果 `abstime` 指定的时间间隔到期且无相关条件变量的通知，则返回 `ETIMEOUT` 错误。

在生产者-消费者（producer-consumer）示例中使用条件变量

下面对前面的示例作出修改，引入条件变量。对全局变量、相关互斥量以及条件变量的声明代码如下：

```
static pthread_mutex_t mtx = PTHREAD_MUTEX_INITIALIZER;
static pthread_cond_t cond = PTHREAD_COND_INITIALIZER;

static int avail = 0;
```

本节中的代码片段摘自随本书发布的源代码文件 `threads/prod_condvar.c`。

除了增加对函数 `pthread_cond_signal()` 的调用外，生产者线程的代码与之前并无变化：

```
s = pthread_mutex_lock(&mtx);
if (s != 0)
    errExitEN(s, "pthread_mutex_lock");

avail++;
/* Let consumer know another unit is available */

s = pthread_mutex_unlock(&mtx);
if (s != 0)
    errExitEN(s, "pthread_mutex_unlock");

s = pthread_cond_signal(&cond);
if (s != 0)
    errExitEN(s, "pthread_cond_signal");
```

在分析消费者代码之前，需要对 `pthread_cond_wait()` 函数做更为详细的解释。前文已经指出，条件变量总是要与一个互斥量相关。将这些对象通过函数参数传递给 `pthread_cond_wait()`，后者执行如下操作步骤。

- 解锁互斥量 `mutex`。
- 堵塞调用线程，直至另一线程就条件变量 `cond` 发出信号。
- 重新锁定 `mutex`。

设计 `pthread_cond_wait()` 执行上述步骤，是因为通常情况下代码会以如下方式访问共享变量：

```
s = pthread_mutex_lock(&mtx);
if (s != 0)
    errExitEN(s, "pthread_mutex_lock");

while (/* Check that shared variable is not in state we want */)
    pthread_cond_wait(&cond, &mtx);

/* Now shared variable is in desired state; do some work */

s = pthread_mutex_unlock(&mtx);
if (s != 0)
    errExitEN(s, "pthread_mutex_unlock");
```

下一节将会介绍为何将 `pthread_cond_wait()` 调用置于 `while` 循环中，而非 `if` 语句中。

在以上代码中，两处对共享变量的访问都必须置于互斥量的保护之下，其原因之前已做了解释。换言之，条件变量与互斥量之间存在着天然的关联关系。

1. 线程在准备检查共享变量状态时锁定互斥量。
2. 检查共享变量的状态。
3. 如果共享变量未处于预期状态，线程应在等待条件变量并进入休眠前解锁互斥量（以便其他线程能访问该共享变量）。
4. 当线程因为条件变量的通知而被再度唤醒时，必须对互斥量再次加锁，因为在典型情况下，线程会立即访问共享变量。

函数 `pthread_cond_wait()` 会自动执行最后两步中对互斥量的解锁和加锁动作。第 3 步中互斥量的释放与陷入对条件变量的等待同属于一个原子操作。换句话说，在函数 `pthread_cond_wait()` 的调用线程陷入对条件变量的等待之前，其他线程不可能获取到该互斥量，也不可能就该条件变量发出信号。

通过观察得出推论：条件变量与互斥量之间存在天然关系，同时等待相同条件变量的所有线程在调用 `pthread_cond_wait()` 或 `pthread_cond_timedwait()` 时必须指定同一互斥量。实际上，`pthread_cond_wait()` 在调用期间能将条件变量与一个唯一的互斥量做动态绑定。SUSv3 规定，在针对同一条件变量并发调用 `pthread_cond_wait()` 时，若使用多个互斥量会导致未定义的结果。

结合以上所有细节，使用 `pthread_cond_wait()` 修改主（消费者）线程的代码如下：

```
for (;;) {
    s = pthread_mutex_lock(&mtx);

    if (s != 0)
        errExitEN(s, "pthread_mutex_lock");

    while (avail == 0) {                /* Wait for something to consume */
        s = pthread_cond_wait(&cond, &mtx);
        if (s != 0)
            errExitEN(s, "pthread_cond_wait");
    }

    while (avail > 0) {                /* Consume all available units */
        /* Do something with produced unit */
        avail--;
    }

    s = pthread_mutex_unlock(&mtx);
    if (s != 0)
        errExitEN(s, "pthread_mutex_unlock");

    /* Perhaps do other work here that doesn't require mutex lock */
}
```

最后，再看一下 `pthread_cond_signal()` 和 `pthread_cond_broadcast()` 的使用。前面展示的生产者代码先调用了 `pthread_mutex_unlock()`，接着调用了 `pthread_cond_signal()`；换言之，先解锁与共享变量相关的互斥量，再就对应的条件变量发出信号。也可以将这两步颠倒执行，SUSv3 允许以任意顺序执行这两个调用。

[Butenhof, 1996]指出，在某些实现中，先解锁互斥量再通知条件变量可能比反序执行效率要高。如果仅在发出条件变量信号后才解锁互斥量，执行 `pthread_cond_wait()` 调用的线程可能会在互斥量仍处于加锁状态时就醒来，当其发现互斥量仍未解锁，会立即再次休眠。这会导致两个多余的上下文切换（context switch）。有些实现运用等待变形（wait morphing）技术解决了这一问题：将等待接收信号的线程从条件变量的等待队列转移至互斥量等待队列。这样，即便互斥量处于加锁状态，也无需切换上下文。

30.2.3 测试条件变量的判断条件（predicate）

每个条件变量都有与之相关的判断条件，涉及一个或多个共享变量。例如，在上一节的代码中，与 `cond` 相关的判断是 `(avail == 0)`。这段代码展示了一个通用的设计原则：必须由一个 `while` 循环，而不是 `if` 语句，来控制对 `pthread_cond_wait()` 的调用。这是因为，当代码从 `pthread_cond_wait()` 返回时，并不能确定判断条件的状态，所以应该立即重新检查判断条件，在条件不满足的情况下继续休眠等待。

从 `pthread_cond_wait()` 返回时，之所以不能对判断条件的状态做任何假设，其理由如下。

- 其他线程可能会率先醒来。也许有多个线程在等待获取与条件变量相关的互斥量。即使就互斥量发出通知的线程将判断条件置为预期状态，其他线程依然有可能率先获取互斥量并改变相关共享变量的状态，进而改变判断条件的状态。
- 设计时设置“宽松的”判断条件或许更为简单。有时，用条件变量来表征可能性而非确定性，在设计应用程序时会更为简单。换言之，就条件变量发送信号意味着“可能有些事情”需要接收信号的线程去响应，而不是“一定有一些事情”要做。使用这种方法，可以基于判断条件的近似情况来发送条件变量通知，接收信号的线程可以通过再次检查判断条件来确定是否真的需要做些什么。
- 可能会发生虚假唤醒的情况。在一些实现中，即使没有任何其他线程真地就条件变量发出信号，等待此条件变量的线程仍有可能醒来。在一些多处理器系统上，为确保高效实现而采用的技术会导致此类（不常见的）虚假唤醒。SUSv3 对此予以明确认可。

30.2.4 示例程序：连接任意已终止线程

前面已然提及，使用 `pthread_join()` 只能连接一个指定线程。且该函数也未提供任何机制去连接任意的已终止线程。本节展示如何使用条件变量绕过这一限制。

程序清单 30-4 为其每个命令行参数创建一个线程。每个线程休眠一段时间后随即退出，休眠时间由相应命令行参数所指定的秒数决定。这里用休眠间隔来模拟线程工作了一段时间。

该程序维护有一组全局变量，记录了所有已创建线程的信息。对于每个线程，全局数组 `thread` 中都含有一元素记录其线程 ID（字段 `tid`）以及当前状态（字段 `state`）。状态字段 `state` 可设置为以下值：`TS_ALIVE`，表示线程是活动的；`TS_TERMINATED`，代表线程已经终结但尚未被连接；`TS_JOINED`，表示线程终止且已被连接。

当线程终止时，将 `TS_TERMINATED` 赋给数组 `thread` 中对应元素的 `state` 字段，对表征已终止但尚未连接线程的全局计数器（`numUnjoined`）加一，并就条件变量 `threadDied` 发出信号。

主线程使用循环不断等待条件变量 `threadDied`。当收到 `threadDied` 信号，且存在已终止线程尚未被连接时，主线程将扫描 `thread` 数组，寻找 `state` 为 `TS_TERMINATED` 的数组元素。对

处于该状态的每个线程，以数组 `thread` 中的对应 `tid` 字段调用 `pthread_join()` 函数，并将 `state` 置为 `TS_JOINED`。当由主线程创建的所有线程终止时，即全局变量 `numLive` 值为 0 时，主循环结束。

以下 shell 会话日志展示了对程序清单 30-4 中程序的调用：

```
$ ./thread_multijoin 1 1 2 3 3          Create 5 threads
Thread 0 terminating
Thread 1 terminating
Reaped thread 0 (numLive=4)
Reaped thread 1 (numLive=3)
Thread 2 terminating
Reaped thread 2 (numLive=2)
Thread 3 terminating
Thread 4 terminating
Reaped thread 3 (numLive=1)
Reaped thread 4 (numLive=0)
```

最后要指出，虽然示例中的线程都被创建为处于可连接状态，且终止后立即由 `pthread_join()` 予以捕获，其实无需采用这一方法来发现线程的终止。可以将线程置为分离态 (`detached`)，无需使用 `pthread_join()`，简单地利用 `thread` 数组（及其他相关全局变量）作为记录每个线程终结的手段。

程序清单 30-4：可以连接任意已终止线程的主线程

```
----- threads/thread_multijoin.c
#include <pthread.h>
#include "tspi_hdr.h"

static pthread_cond_t threadDied = PTHREAD_COND_INITIALIZER;
static pthread_mutex_t threadMutex = PTHREAD_MUTEX_INITIALIZER;
/* Protects all of the following global variables */

static int totThreads = 0; /* Total number of threads created */
static int numLive = 0; /* Total number of threads still alive or
                        terminated but not yet joined */
static int numUnjoined = 0; /* Number of terminated threads that
                            have not yet been joined */

enum tstate {
    TS_ALIVE, /* Thread is alive */
    TS_TERMINATED, /* Thread terminated, not yet joined */
    TS_JOINED /* Thread terminated, and joined */
};

static struct {
    pthread_t tid; /* Info about each thread */
    enum tstate state; /* ID of this thread */
    int sleepTime; /* Thread state (TS_* constants above) */
} *thread; /* Number seconds to live before terminating */

static void *
threadFunc(void *arg) /* Start function for thread */
{
    int idx = *((int *) arg);
    int s;
    sleep(thread[idx].sleepTime); /* Simulate doing some work */
    printf("Thread %d terminating\n", idx);

    s = pthread_mutex_lock(&threadMutex);
    if (s != 0)
```



```

        errExitEN(s, "pthread_mutex_lock");

numUnjoined++;
thread[idx].state = TS_TERMINATED;

s = pthread_mutex_unlock(&threadMutex);
if (s != 0)
    errExitEN(s, "pthread_mutex_unlock");
s = pthread_cond_signal(&threadDied);
if (s != 0)
    errExitEN(s, "pthread_cond_signal");

return NULL;
}

int
main(int argc, char *argv[])
{
    int s, idx;

    if (argc < 2 || strcmp(argv[1], "--help") == 0)
        usageErr("%s nsecs...\n", argv[0]);

    thread = calloc(argc - 1, sizeof(*thread));
    if (thread == NULL)
        errExit("calloc");

    /* Create all threads */

    for (idx = 0; idx < argc - 1; idx++) {
        thread[idx].sleepTime = getInt(argv[idx + 1], GN_NONNEG, NULL);
        thread[idx].state = TS_ALIVE;
        s = pthread_create(&thread[idx].tid, NULL, threadFunc, &idx);
        if (s != 0)
            errExitEN(s, "pthread_create");
    }

    totThreads = argc - 1;
    numLive = totThreads;

    /* Join with terminated threads */

    while (numLive > 0) {
        s = pthread_mutex_lock(&threadMutex);
        if (s != 0)
            errExitEN(s, "pthread_mutex_lock");
        while (numUnjoined == 0) {
            s = pthread_cond_wait(&threadDied, &threadMutex);
            if (s != 0)
                errExitEN(s, "pthread_cond_wait");
        }

        for (idx = 0; idx < totThreads; idx++) {
            if (thread[idx].state == TS_TERMINATED){
                s = pthread_join(thread[idx].tid, NULL);
                if (s != 0)
                    errExitEN(s, "pthread_join");

                thread[idx].state = TS_JOINED;
            }
        }
    }
}

```

```

        numLive--;
        numUnjoined--;

        printf("Reaped thread %d (numLive=%d)\n", idx, numLive);
    }
}

s = pthread_mutex_unlock(&threadMutex);
if (s != 0)
    errExitEN(s, "pthread_mutex_unlock");
}

exit(EXIT_SUCCESS);
}

```

threads/thread_multijoin.c

30.2.5 经由动态分配的条件变量

使用函数 `pthread_cond_init()` 对条件变量进行动态初始化。需要使用 `pthread_cond_init()` 的情形类似于使用 `pthread_mutex_init()` 来动态初始化互斥量的情况。亦即，对自动或动态分配的条件变量进行初始化时，或是对未采用默认属性经由静态分配的条件变量进行初始化时，必须使用 `pthread_cond_init()`。

```

#include <pthread.h>

int pthread_cond_init(pthread_cond_t *cond, const pthread_condattr_t *attr);

```

Returns 0 on success, or a positive error number on error

参数 `cond` 表示将要初始化的目标条件变量。类似于互斥量，可以指定之前经由初始化处理的 `attr` 参数来判定条件变量的属性。对于 `attr` 所指向的 `pthread_condattr_t` 类型对象，可使用多个 `Pthreads` 函数对其中属性进行初始化。若将 `attr` 置为 `NULL`，则使用一组缺省属性来设置条件变量。

`SUSv3` 规定，对业已初始化的条件变量进行再次初始化，将导致未定义的行为。应当避免这一做法。

当不再需要一个经由自动或动态分配的条件变量时，应调用 `pthread_cond_destroy()` 函数予以销毁。对于使用 `PTHREAD_COND_INITIALIZER` 进行静态初始化的条件变量，无需调用 `pthread_cond_destroy()`。

```

#include <pthread.h>

int pthread_cond_destroy(pthread_cond_t *cond);

```

Returns 0 on success, or a positive error number on error

对某个条件变量而言，仅当没有任何线程在等待它时，将其销毁才是安全的。如果条件变量驻留于某片动态创建的内存区域，那么应在释放该内存区域前就将其销毁。经由自动分配的条件变量应在宿主函数返回前予以销毁。

经 `pthread_cond_destroy()` 销毁的条件变量，之后可以调用 `pthread_cond_init()` 对其进行重新初始化。

30.3 总结

线程提供的强大共享是有代价的。多线程应用程序必须使用互斥量和条件变量等同步原语来协调对共享变量的访问。互斥量提供了对共享变量的独占式访问。条件变量允许一个或多个线程等候通知：其他线程改变了共享变量的状态。

更多信息

请参考 29.10 节所列的更多信息来源。

30.4 练习

- 30-1.** 修改程序清单 30-1 (`thread_incr.c`) 中的程序，以便线程起始函数在每次循环中都能输出 `glob` 的当前值以及能对线程做唯一标识的标识符。可将线程的这一唯一标识指定为创建线程的函数 `pthread_create()` 的调用参数。对于这一程序，需要将线程起始函数的参数改为指针，指向包含线程唯一标识和循环次数限制的数据结构。运行该程序，将输出重定向至一文件，查看内核在调度两线程交替执行时 `glob` 的变化情况。
- 30-2.** 实现一组线程安全的函数，以更新和搜索一个不平衡二叉树。此函数库应该包含如下形式的函数（目的明显）：

```
initialize(tree);
add(tree, char *key, void *value);
delete(tree, char *key)
Boolean lookup(char *key, void **value)
```

上述函数原型中，`tree` 是一个指向根节点的结构（为此需要定义一个合适的结构）。树的每个节点保存有一个键-值对。还需为树中每个节点定义一数据结构，其中应包含互斥量，以确保同时仅有一个线程可以访问该节点。`initialize()`、`add()`和 `lookup()`函数的实现相对简单。`delete()`的实现需要较为深入的考虑。

无需维护平衡二叉树，这极大简化了实现对锁的需求，但同时也带来了风险，特定模式的输入会导致树的执行效率低下。要维护平衡二叉树，则在执行 `add()`和 `delete()`操作时必然要在子树间移动节点，这就需要更为复杂的锁定策略。

第 31 章

线程：线程安全和每线程存储

本章将拓展对 POSIX 线程 API 的探讨，描述线程安全（thread-safe）函数以及一次性初始化。同时讨论在不改变函数接口定义的前提下，如何通过线程特有数据（thread-specific data）或线程局部存储（thread-local storage）实现已有函数的线程安全。

31.1 线程安全（再论可重入性）

若函数可同时供多个线程安全调用，则称之为线程安全函数；反之，如果函数不是线程安全的，则不能并发调用。例如，如下函数（30.1 节也有类似代码）就不是线程安全的：

```
static int glob = 0;

static void
incr(int loops)
{
    int loc, j;
    for (j = 0; j < loops; j++) {
        loc = glob;
        loc++;
        glob = loc;
    }
}
```

如果多个线程并发调用该函数，glob 的最终值将不得而知。本例展示了导致线程不安全的典型原因：使用了在所有线程之间共享的全局或静态变量。

实现线程安全有多种方式。其一是将函数与互斥量关联使用（如果函数库中的所有函数都共享同样的全局变量，那么或许应将所有函数都与该互斥量相关联），在调用函数时将其锁定，在函数返回时解锁。这一方法的优点在于简单。另一方面，这也意味着同时只能有一个线程执行该函数，亦即，对该函数的访问是串行的（serialized）。如果各线程在执行此函数时都耗费了相当多的时间，那么串行化会导致并发能力的丧失，所有线程将不再并发执行。

另一种更为复杂的解决方案是：将共享变量与互斥量关联起来。这需要程序员们确认函数

的哪些部分是使用了共享变量的临界区，且仅在执行到临界区时去获取和释放互斥量。这将允许许多线程同时执行一个函数并实现并行，除非出现多个线程需要同时执行同一临界区的情况。

非线性安全的函数

为便于开发多线程应用程序，除了表 31-1 所列函数以外（其中大部分并未在本书中提及），SUSv3 中的所有函数都需实现线程安全。

除了表 31-1 中所列函数，SUSv3 还做了如下规定。

- 如传参为 NULL 时，函数 `ctermid()` 和 `tmpnam()` 无需是线程安全的。
- 如果函数 `wcrtomb()` 和 `wcsrtombs()` 的最后一个参数（ps）为 NULL，那么这两个函数也无需是线程安全的。

SUSv4 对表 31-1 中的函数做了以下修改。

- 移除函数 `ecvt()`、`fcvt()`、`gcvt()`、`gethostbyname()` 以及 `gethostbyaddr()`，因为已从标准中删除了这些函数。
- 增加函数 `strsignal()` 和 `system()`。由于 `system()` 函数就信号处置所做的操作将影响整个进程，故而是不可重入的。

标准并未禁止将表 31-1 中的函数实现为线程安全。不过，即使在某些实现中有些函数是线程安全的，为确保应用程序的可移植性，也不应该假设这些函数在所有实现中都是如此。

表 31-1: SUSv3 不要求这些函数是线程安全的

<code>asctime()</code>	<code>fcvt()</code>	<code>getpwnam()</code>	<code>nl_langinfo()</code>
<code>basename()</code>	<code>ftw()</code>	<code>getpwuid()</code>	<code>ptsname()</code>
<code>catgets()</code>	<code>gcvt()</code>	<code>getservbyname()</code>	<code>putc_unlocked()</code>
<code>crypt()</code>	<code>getc_unlocked()</code>	<code>getservbyport()</code>	<code>putchar_unlocked()</code>
<code>ctime()</code>	<code>getchar_unlocked()</code>	<code>getservent()</code>	<code>putenv()</code>
<code>dbm_clearerr()</code>	<code>getdate()</code>	<code>getutxent()</code>	<code>pututxline()</code>
<code>dbm_close()</code>	<code>getenv()</code>	<code>getutxid()</code>	<code>rand()</code>
<code>dbm_delete()</code>	<code>getgrent()</code>	<code>getutxline()</code>	<code>readdir()</code>
<code>dbm_error()</code>	<code>getgrgid()</code>	<code>gmtime()</code>	<code>setenv()</code>
<code>dbm_fetch()</code>	<code>getgrnam()</code>	<code>hcreate()</code>	<code>setgrent()</code>
<code>dbm_firstkey()</code>	<code>gethostbyaddr()</code>	<code>hdestroy()</code>	<code>setkey()</code>
<code>dbm_nextkey()</code>	<code>gethostbyname()</code>	<code>hsearch()</code>	<code>setpwent()</code>
<code>dbm_open()</code>	<code>gethostent()</code>	<code>inet_ntoa()</code>	<code>setutxent()</code>
<code>dbm_store()</code>	<code>getlogin()</code>	<code>l64a()</code>	<code>strerror()</code>
<code>dirname()</code>	<code>getnetbyaddr()</code>	<code>lgamma()</code>	<code>strtok()</code>
<code>dlerror()</code>	<code>getnetbyname()</code>	<code>lgammaf()</code>	<code>ttyname()</code>
<code>drand48()</code>	<code>getnetent()</code>	<code>lgammal()</code>	<code>unsetenv()</code>
<code>ecvt()</code>	<code>getopt()</code>	<code>localeconv()</code>	<code>wcstombs()</code>
<code>encrypt()</code>	<code>getprotobyname()</code>	<code>localtime()</code>	<code>wctomb()</code>
<code>endgrent()</code>	<code>getprotobynumber()</code>	<code>lrand48()</code>	
<code>endpwent()</code>	<code>getprotoent()</code>	<code>mrnd48()</code>	
<code>endutxent()</code>	<code>getpwent()</code>	<code>nftw()</code>	

可重入和不可重入函数

较之于对整个函数使用互斥量，使用临界区实现线程安全虽然有明显改进，但由于存在对互斥量的加锁和解锁开销，所以多少还是有些低效。可重入函数则无需使用互斥量即可实

现线程安全。其要诀在于避免对全局和静态变量的使用。需要返回给调用者的任何信息，亦或是在对函数的历次调用间加以维护的信息，都存储于由调用者分配的缓冲区内。（初次碰到可重入问题，是在 21.1.2 节讨论信号处理器中的全局变量时。）不过，并非所有函数都可以实现为可重入，通常的原因如下。

- 根据其性质，有些函数必须访问全局数据结构。`malloc` 函数库中的函数就是这方面的典范。这些函数为堆中的空闲块维护有一个全局链表。**malloc 库函数的线程安全是通过使用互斥量来实现的。**
- 一些函数（在发明线程之前就已问世）的接口本身就定义为不可重入，要么返回指针，指向由函数自身静态分配的存储空间，要么利用静态存储对该函数（或相关函数）历次调用间的信息加以维护。表 31-1 所列函数大多属于此类。例如，函数 `asctime()`（10.2.3 节）就返回一个指针，指向经由静态分配的缓冲区，其内容为日期和时间字符串。

对于一些接口不可重入的函数，SUSv3 为其定义了以后缀 `_r` 结尾的可重入“替身”。这些“替身”函数要求由调用者来分配缓冲区，并将缓存区地址传给函数用以返回结果。这使得调用线程可以使用局部（栈）变量来存放函数结果。出于这一目的，SUSv3 定义了如下函数：`asctime_r()`、`ctime_r()`、`getgrgid_r()`、`getgrnam_r()`、`getlogin_r()`、`getpwnam_r()`、`getpwuid_r()`、`gmtime_r()`、`localtime_r()`、`rand_r()`、`readdir_r()`、`strerror_r()`、`strtok_r()`和 `ttyname_r()`。

有些系统实现为一些传统的不可重入函数也提供了附加的可重入“替身”。例如，`glibc` 就提供了函数 `crypt_r()`、`gethostbyname_r()`、`getservbyname_r()`、`getutent_r()`、`getutid_r()`、`getutline_r()`和 `ptsname_r()`。不过，为确保应用程序的可移植性，不应假设这些函数在其他实现中也存在。某些情况下，SUSv3 并未规定这些等价的可重入函数，因为功能更强、又可重入的替代函数已然存在。例如，函数 `getaddrinfo()`就更新且可重入，可用来替代函数 `gethostbyname()`和 `getservbyname()`。

31.2 一次性初始化

多线程程序有时有这样的需求：不管创建了多少线程，有些初始化动作只能发生一次。例如，可能需要执行 `pthread_mutex_init()`对带有特殊属性的互斥量进行初始化，而且必须只能初始化一次。如果由主线程来创建新线程，那么这一点易如反掌：可以在创建依赖于该初始化的线程之前进行初始化。不过，对于库函数而言，这样处理就不可行，因为调用者在初次调用库函数之前可能已经创建了这些线程。故而需要这样的库函数：无论首次为任何线程所调用，都会执行初始化动作。

库函数可以通过函数 `pthread_once()`实现一次性初始化。

```
#include <pthread.h>

int pthread_once(pthread_once_t *once_control, void (*init)(void));

Returns 0 on success, or a positive error number on error
```

利用参数 `once_control` 的状态，函数 `pthread_once()`可以确保无论有多少线程对 `pthread_once()`调用了多少次，也只会执行一次由 `init` 指向的调用者定义函数。

`init` 函数没有任何参数，形式如下：

```

void
init(void)
{
    /* Function body */
}

```

另外，参数 `once_control` 必须是一指针，指向初始化为 `PTHREAD_ONCE_INIT` 的静态变量。

```
pthread_once_t once_var = PTHREAD_ONCE_INIT;
```

调用函数 `pthread_once()`时要指定一个指针，指向类型为 `pthread_once_t` 的特定变量，对该函数的首次调用将修改 `once_control` 所指向的内容，以便对其后续调用不会再次执行 `init`。

常常将 `Pthread_once()`和线程特有数据结合使用，相关内容会在下一节描述。

`Pthreads` 的早期版本不能对互斥量进行静态初始化，只能使用 `pthread_mutex_init()` ([Butenbof, 1996])，这也是函数 `pthread_once()`存在的主要原因。随着静态分配互斥量功能的问世，库函数可以使用一个经静态分配的互斥量和一个静态布尔型 (`Boolean`) 变量来实现一次性初始化。虽然如此，出于方便的考虑，函数 `pthread_once()`得以保留。

31.3 线程特有数据

实现函数线程安全最为有效的方式就是使其可重入，应以这种方式来实现所有新的函数库。不过，对于已有的不可重入函数库（可能问世于线程流行之前）来说，采用这种方法通常需要修改函数接口，这也意味着，需要修改所有使用此类函数的应用程序。

使用线程特有数据技术，可以无需修改函数接口而实现已有函数的线程安全。较之于可重入函数，采用线程特有数据的函数效率可能要略低一些，不过对于使用了这些调用的程序而言，则省去了修改程序之劳。

如图 31-1 所示，线程特有数据使函数得以为每个调用线程分别维护一份变量的副本 (`copy`)。线程特有数据是长期存在的。在同一线程对相同函数的历次调用间，每个线程的变量会持续存在，函数可以向每个调用线程返回各自的结果缓冲区（如果需要的话）。

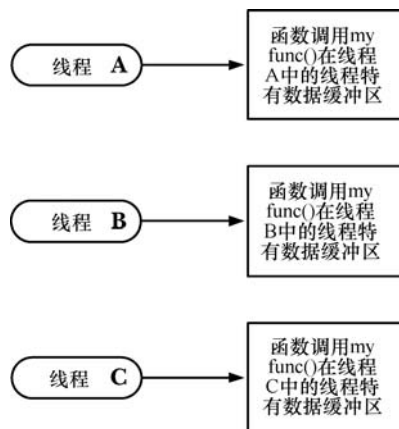


图 31-1: 线程特有数据 (TSD) 为函数提供线程内存储

31.3.1 库函数视角下的线程特有数据

要了解线程特有数据相关 API 的使用，需要从使用这一技术的库函数角度来考虑如下问题。

- 该函数必须为每个调用者线程分配单独的存储，且只需在线程初次调用此函数时分配一次即可。
- 在同一线程对此函数的后续所有调用中，该函数都需要获取初次调用时线程分配的存

存储块地址。由于函数调用结束时会释放自动变量，故而函数不应利用自动变量存放存储块指针，也不能将指针存放于静态变量中，因为静态变量在进程中只有一个实例。Pthreads API 提供了函数来处理这一情况。

- 不同（无相互依赖关系）函数各自可能都需要使用线程特有数据。每个函数都需要方法来标识其自身的线程特有数据（键），以便与其他函数所使用的线程特有数据有所区分。
- 当线程退出时，函数无法控制将要发生的情况。这时，线程可能会执行该函数之外的代码。不过，一定存在某些机制（解构器），在线程退出时会自动释放为该线程所分配的存储。若非如此，随着持续不断地创建线程，调用函数和终止线程，将会引发内存泄露。

31.3.2 线程特有数据 API 概述

要使用线程特有数据，库函数执行的一般步骤如下。

1. 函数创建一个键（key），用以将不同函数使用的线程特有数据项区分开来。调用函数 `pthread_key_create()` 可创建此“键”，且只需在首个调用该函数的线程中创建一次，函数 `pthread_once()` 的使用正是出于这一目的。键在创建时并未分配任何线程特有数据块。
2. 调用 `pthread_key_create()` 还有另一个目的，即允许调用者指定一个自定义解构函数，用于释放为该键所分配的各个存储块（参见下一步）。当使用线程特有数据的线程终止时，Pthreads API 会自动调用此解构函数，同时将该线程的数据块指针作为参数传入。
3. 函数会为每个调用者线程创建线程特有数据块。这一分配通过调用 `malloc()`（或类似函数）完成，每个线程只分配一次，且只会在线程初次调用此函数时分配。
4. 为了保存上一步所分配存储块的地址，函数会使用两个 Pthreads 函数：`pthread_setspecific()` 和 `pthread_getspecific()`。调用函数 `pthread_setspecific()` 实际上是对 Pthreads 实现发起这样的请求：保存该指针，并记录其与特定键（该函数的键）以及特定线程（调用者线程）的关联性。调用 `pthread_getspecific()` 所执行的是互补操作：返回之前所保存的、与给定键以及调用线程相关联的指针。如果还没有指针与特定的键及线程相关联，那么 `pthread_getspecific()` 返回 `NULL`。函数可以利用这一点来判断自身是否是初次为某个线程所调用，若为初次，则必须为该线程分配空间。

31.3.3 线程特有数据 API 详述

本节将详述上节所提及的各个函数，并通过线程特有数据的典型实现来说明其操作方法。下一节会演示如何使用线程特有数据来实现线程安全的标准 C 语言库函数 `stderr()`。

调用 `pthread_key_create()` 函数为线程特有数据创建一个新键，并通过 `key` 所指向的缓冲区返回给调用者。

因为进程中的所有线程都可使用返回的键，所以参数 `key` 应指向一个全局变量。

```
#include <pthread.h>

int pthread_key_create(pthread_key_t *key, void (*destructor)(void *));
    Returns 0 on success, or a positive error number on error
```

参数 `destructor` 指向一个自定义函数，其格式如下：


```

void
dest(void *value)
{
    /* Release storage pointed to by 'value' */
}

```

只要线程终止时与 `key` 的关联值不为 `NULL`，`Pthreads API` 会自动执行解构函数，并将与 `key` 的关联值作为参数传入解构函数。传入的值通常是与该键关联，且指向线程特有数据块的指针。如果无需解构，那么可将 `destructor` 设置为 `NULL`。

如果一个线程有多个线程特有数据块，那么对各个解构函数的调用顺序是不确定的。对每个解构函数的设计应相互独立。

观察线程特有数据的实现有助于理解它们的使用方法。典型的实现（`NPTL` 即在此列）会包含以下数组。

- 一个全局（进程范围）数组，存放线程特有数据的键信息。
- 每个线程包含一个数组，存有为每个线程分配的线程特有数据块的指针（通过调用 `pthread_setspecific()` 来存储指针）。



图 31-2: 线程特有数据键的实现

在这一实现中，`pthread_key_create()` 返回的 `pthread_key_t` 类型值只是对全局数组的索引（`index`），标记为 `pthread_keys`，其格式如图 31-2 所示。数组的每个元素都是一个包含两个字段（`field`）的结构。第一个字段标记该数组元素是否在用（即已由之前对 `pthread_key_create()` 的调用分配）。第二个字段用于存放针对此键、线程特有数据块的解构函数指针（是函数 `pthread_key_create()` 中参数 `destructor` 的一份拷贝）。

函数 `pthread_setspecific()` 要求 `Pthreads API` 将 `value` 的副本存储于一数据结构中，并将 `value` 与调用线程以及 `key` 相关联（`key` 由之前对 `pthread_key_create()` 的调用返回）。`pthread_getspecific()` 函数执行的操作与之相反，返回之前与本线程及给定 `key` 相关的值（`value`）。

```

#include <pthread.h>

int pthread_setspecific(pthread_key_t key, const void *value);
    Returns 0 on success, or a positive error number on error

void *pthread_getspecific(pthread_key_t key);
    Returns pointer, or NULL if no thread-specific data is associated with key

```

函数 `pthread_setspecific()` 的参数 `value` 通常是一指针，指向由调用者分配的一块内存。当线程终止时，会将该指针作为参数传递给与 `key` 对应的解构函数。

参数 `value` 也可以不是一个指向内存区域的指针，而是任何可以赋值（通过强制转换）给 `void*` 的标量值。在这种情况下，先前对 `pthread_key_create()` 函数的调用应将 `destructor` 指定为 `NULL`。

图 31-3 展示了用于存储 `value` 的数据结构的常见实现。图中假设将 `pthread_keys[1]` 分配给函数 `myfunc()`。`Pthreads API` 为每个函数维护指向线程特有数据块的一个指针数组。其中每个数

组元素都与图 31-2 中全局 `pthread_keys` 数组的元素一一对应。函数 `pthread_setspecific()` 在指针数组中为每个调用线程设置与 `key` 对应的元素。

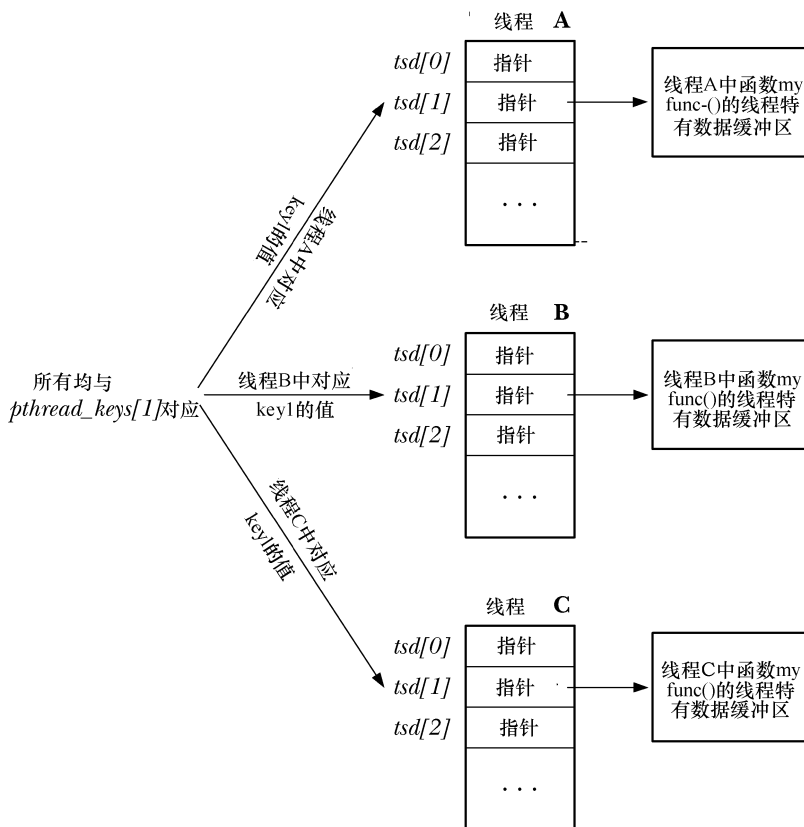


图 31-3: 用于实现线程特有数据 (TSD) 指针的数据结构

当线程刚刚创建时，会将所有线程特有数据的指针都初始化为 `NULL`。这意味着当线程初次调用库函数时，必须使用 `pthread_getspecific()` 函数来检查该线程是否已有与 `key` 对应的关联值。如果没有，那么此函数会分配一块内存并通过 `pthread_setspecific()` 保存指向该内存块的指针。在下一节实现线程安全版的 `stderr()` 函数时，将给出示例。

31.3.4 使用线程特有数据 API

3.4 节在首度论及标准 `stderr()` 函数时曾指出，可能会返回一个指向静态分配字符串的指针作为函数结果。这意味着 `stderr()` 可能不是线程安全的。后面将以数页篇幅讨论一下非线程安全的 `stderr()` 实现，接着说明如何使用线程特有数据来实现该函数的线程安全。

在包括 Linux 在内的许多 UNIX 实现中，由标准 C 语言函数库提供的 `stderr()` 函数都是线程安全的。不过，由于 SUSv3 并未规定该函数必须是线程安全的，而且这一 `stderr()` 的实现又为使用线程特有数据提供了一个简单范例，故而在此将其作为示例。

程序清单 31-1 演示了非线程安全版 `stderr()` 函数的一个简单实现。该函数利用了由 glibc 定义的一对全局变量：`_sys_errlist` 是一个指针数组，其每个元素指向一个与 `errno` 错误号相匹配的字符串（因此，例如，`_sys_errlist[EINVAL]` 即指向字符串 `Invalid operation`）；`_sys_nerr` 表

示 `_sys_errlist` 中的元素个数。

程序清单 31-1: 非线程安全版 `strerror()` 函数的一种实现

```
----- threads/strerror.c
#define _GNU_SOURCE          /* Get '_sys_nerr' and '_sys_errlist'
                             declarations from <stdio.h> */

#include <stdio.h>
#include <string.h>          /* Get declaration of strerror() */

#define MAX_ERROR_LEN 256   /* Maximum length of string
                             returned by strerror() */

static char buf[MAX_ERROR_LEN]; /* Statically allocated return buffer */

char *
strerror(int err)
{
    if (err < 0 || err >= _sys_nerr || _sys_errlist[err] == NULL) {
        snprintf(buf, MAX_ERROR_LEN, "Unknown error %d", err);
    } else {
        strncpy(buf, _sys_errlist[err], MAX_ERROR_LEN - 1);
        buf[MAX_ERROR_LEN - 1] = '\0'; /* Ensure null termination */
    }

    return buf;
}
----- threads/strerror.c
```

可以利用程序清单 31-2 中程序来展示程序清单 31-1 中非线程安全版的 `strerror()` 实现所造成的后果。该程序分别从两个不同线程中调用 `strerror()`，并且均在两个线程调用 `strerror()` 之后才显示返回结果。虽然两个线程为 `strerror()` 指定的参数值不同 (`EINVAL` 和 `EPERM`)，在与程序清单 31-1 版的 `strerror()` 链接、编译后，运行该程序将产生如下结果：

```
$ ./strerror_test
Main thread has called strerror()
Other thread about to call strerror()
Other thread: str (0x804a7c0) = Operation not permitted
Main thread: str (0x804a7c0) = Operation not permitted
```

两个线程都显示与 `EPERM` 对应的 `errno` 字符串，因为第二个线程对 `strerror()` 的调用（在函数 `threadFunc()` 中）覆盖了主线程调用 `strerror()` 时写入缓冲区的内容。检查输出结果可以发现，两个线程的局部变量 `str` 均指向同一内存地址。

程序清单 31-2: 从两个不同线程调用 `strerror()`

```
----- threads/strerror_test.c

#include <stdio.h>
#include <string.h>          /* Get declaration of strerror() */
#include <pthread.h>
#include "tlpi_hdr.h"

static void *
threadFunc(void *arg)
{
    char *str;

    printf("Other thread about to call strerror()\n");
    str = strerror(EPERM);
}
```

```

    printf("Other thread: str (%p) = %s\n", str, str);

    return NULL;
}

int
main(int argc, char *argv[])
{
    pthread_t t;
    int s;
    char *str;

    str = strerror(EINVAL);
    printf("Main thread has called strerror()\n");

    s = pthread_create(&t, NULL, threadFunc, NULL);
    if (s != 0)
        errExitEN(s, "pthread_create");

    s = pthread_join(t, NULL);
    if (s != 0)
        errExitEN(s, "pthread_join");

    printf("Main thread: str (%p) = %s\n", str, str);

    exit(EXIT_SUCCESS);
}

```

threads/strerror_test.c

程序清单 31-3 是对函数 `strerror()` 的全新实现，使用了线程特有数据来确保线程安全。

程序清单 31-3：使用线程特有数据以实现线程安全的 `strerror()` 函数

```

threads/strerror_tsd.c

#define _GNU_SOURCE          /* Get '_sys_nerr' and '_sys_errlist'
                             declarations from <stdio.h> */

#include <stdio.h>
#include <string.h>          /* Get declaration of strerror() */
#include <pthread.h>
#include "tspi_hdr.h"

static pthread_once_t once = PTHREAD_ONCE_INIT;
static pthread_key_t strerrorKey;

#define MAX_ERROR_LEN 256   /* Maximum length of string in per-thread
                             buffer returned by strerror() */

static void                 /* Free thread-specific data buffer */
① destructor(void *buf)
{
    free(buf);
}

static void                 /* One-time key creation function */
② createKey(void)
{
    int s;

    /* Allocate a unique thread-specific data key and save the address

```

```

        of the destructor for thread-specific data buffers */
③    s = pthread_key_create(&strerrorKey, destructor);
        if (s != 0)
            errExitEN(s, "pthread_key_create");
    }
    char *
    strerror(int err)
    {
        int s;
        char *buf;

        /* Make first caller allocate key for thread-specific data */

④    s = pthread_once(&once, createKey);
        if (s != 0)
            errExitEN(s, "pthread_once");

⑤    buf = pthread_getspecific(strerrorKey);
        if (buf == NULL) {           /* If first call from this thread, allocate
                                     buffer for thread, and save its location */
⑥        buf = malloc(MAX_ERROR_LEN);
            if (buf == NULL)
                errExit("malloc");

⑦        s = pthread_setspecific(strerrorKey, buf);
            if (s != 0)
                errExitEN(s, "pthread_setspecific");
        }

        if (err < 0 || err >= _sys_nerr || _sys_errlist[err] == NULL) {
            snprintf(buf, MAX_ERROR_LEN, "Unknown error %d", err);
        } else {
            strncpy(buf, _sys_errlist[err], MAX_ERROR_LEN - 1);
            buf[MAX_ERROR_LEN - 1] = '\0';           /* Ensure null termination */
        }

        return buf;
    }

```

threads/strerror_tsd.c

改进版 `strerror()` 所做的第一步是调用 `pthread_once()`④，以确保（从任何线程）对该函数的首次调用将执行 `createKey()`②。函数 `createKey()` 会调用 `pthread_key_create()` 来分配一个线程特有数据的键（key），并将其存储于全局变量 `strerrorKey`③中。对 `pthread_key_create()` 的调用同时也会记录解构函数①的地址，将使用该解构函数来释放与键对应的线程特有数据缓冲区。

接着，函数 `strerror()` 调用 `pthread_getspecific()`⑤以获取该线程中对应于 `strerrorKey` 的唯一缓冲区地址。如果 `pthread_getspecific()` 返回 `NULL`，这表明该线程是首次调用 `strerror()` 函数，因此函数会调用 `malloc()`⑥分配一个新缓冲区，并使用 `pthread_setspecific()`⑦来保存该缓冲区的地址。如果 `pthread_getspecific()` 的返回值非 `NULL`，那么该值指向业已存在的缓冲区，此缓冲区由之前对 `strerror()` 的调用所分配。

这一 `strerror()` 函数实现的剩余部分与非线程安全版的前述实现相类似，唯一的区别在于，`buf` 是线程特有数据的缓冲区地址，而非静态变量。

如果使用新版 `strerror()`（程序清单 31-3）编译链接测试程序（程序清单 31-2）`strerror_test_tsd`，程序运行会有如下结果：

```
$ ./strerror_test_tsd
Main thread has called strerror()
Other thread about to call strerror()
Other thread: str (0x804b158) = Operation not permitted
Main thread: str (0x804b008) = Invalid argument
```

根据这一输出，可以看出新版 `strerro()` 是线程安全的：两个线程中局部变量 `str` 所指向的地址是不同的。

31.3.5 线程特有数据的实现限制

正如对线程特有数据典型实现过程的描述所揭示的，实现可能要对其所支持的线程特有数据键的数量加以限制。SUSv3 要求至少支持 128 (`_POSIX_THREAD_KEYS_MAX`) 个键。应用程序要么通过对 `PTHREAD_KEY_MAX`（定义于 `<limits.h>`）的定义，要么通过调用 `sysconf(_SC_THREAD_KEYS_MAX)`，来确定实际支持的键数量。Linux 支持多达 1024 个键。

即使 128 个键对于大多数应用来说也已经绰绰有余。这是因为，每个库函数应该只会使用到少量的键，通常会只用一个。如果一个函数需要多个线程特有数据的值，通常可将这些值置于一个结构中，并将该结构仅与一个线程特有数据的键关联。

31.4 线程局部存储

类似于线程特有数据，线程局部存储提供了持久的每线程存储。作为非标准特性，诸多其他的 UNIX 实现（例如 Solaris 和 FreeBSD）为其提供了相同，或类似的接口形式。

线程局部存储的主要优点在于，比线程特有数据的使用要简单。要创建线程局部变量，只需简单地在全局或静态变量的声明中包含 `__thread` 说明符即可。

```
static __thread buf[MAX_ERROR_LEN];
```

但凡带有这种说明符的变量，每个线程都拥有一份对变量的拷贝。线程局部存储中的变量将一直存在，直至线程终止，届时会自动释放这一存储。

关于线程局部变量的声明和使用，需要注意如下几点。

- 如果变量声明中使用了关键字 `static` 或 `extern`，那么关键字 `__thread` 必须紧随其后。
- 与一般的全局或静态变量声明一样，线程局部变量在声明时可设置一个初始值。
- 可以使用 C 语言取址操作符 (`&`) 来获取线程局部变量的地址。

线程局部存储需要内核（由 Linux 2.6 提供）、Pthreads 实现（由 NPTL 提供）以及 C 编译器（在 x86-32 平台上由 gcc 3.3 或后续版本提供）的支持。

程序清单 31-4 提供了使用线程局部存储实现线程安全版 `strerror()` 函数的例子。如果用该版 `strerror()` 与测试程序（程序清单 31-2）编译、链接、生成 `strerror_test_tls`，那么运行时将产生如下结果：

```
$ ./strerror_test_tls
Main thread has called strerror()
Other thread about to call strerror()
Other thread: str (0x40376ab0) = Operation not permitted
Main thread: str (0x40175080) = Invalid argument
```

程序清单 31-4：使用线程局部存储实现线程安全版的 strerror()函数

```
----- threads/strerror_tls.c
#define _GNU_SOURCE          /* Get '_sys_nerr' and '_sys_errlist'
                             declarations from <stdio.h> */

#include <stdio.h>
#include <string.h>          /* Get declaration of strerror() */
#include <pthread.h>

#define MAX_ERROR_LEN 256   /* Maximum length of string in per-thread
                             buffer returned by strerror() */

static __thread char buf[MAX_ERROR_LEN];
                             /* Thread-local return buffer */

char *
strerror(int err)
{
    if (err < 0 || err >= _sys_nerr || _sys_errlist[err] == NULL) {
        snprintf(buf, MAX_ERROR_LEN, "Unknown error %d", err);
    } else {
        strncpy(buf, _sys_errlist[err], MAX_ERROR_LEN - 1);
        buf[MAX_ERROR_LEN - 1] = '\0';      /* Ensure null termination */
    }

    return buf;
}
----- threads/strerror_tls.c
```

31.5 总结

若一函数可由多个线程同时安全调用，则称之为线程安全的函数。使用全局或静态变量是导致函数非线程安全的通常原因。在多线程应用中，保障非线程安全函数安全的手段之一是运用互斥锁来防护对该函数的所有调用。这种方法带来了并发性能的下降，因为同一时刻只能有一个线程运行该函数。提升并发性能的另一方法是：仅在函数中操作共享变量（临界区）的代码前后加入互斥锁。

使用互斥量可以实现大部分函数的线程安全，不过由于互斥量的加、解锁开销，故而也带来了性能的下降。如能避免使用全局或静态变量，可重入函数则无需使用互斥量即可实现线程安全。

SUSv3 所规范的大部分函数都需实现线程安全。SUSv3 同时也列出了小部分无需实现线程安全的函数。一般情况下，这些函数将静态存储返回给调用者，或者在对函数的连续调用间进行信息维护。根据定义，这些函数是不可重入的，也不能使用互斥量来确保其线程安全。本章讨论了两种大致相当的编程技术——线程特有数据和线程局部存储——可在无需改变函数接口定义的情况下保障不安全函数的线程安全。这两种技术均允许函数分配持久的、基于线程的存储。

更多信息

请参考 29.10 节所列的更多信息来源。

31.6 练习

- 31-1.** 试实现函数 `one_time_init(control, init)`，要求与函数 `pthread_once()` 执行等同操作。参数 `control` 应为一指针，指向经静态分配的结构，其中包含一个布尔型变量和一个互斥量。布尔型变量用以标识函数 `init` 是否曾被调用过，而由互斥量来控制对变量的访问。为简化函数实现，可以忽略诸如 `init()` 调用失败或者在由线程初次调用时被取消的情况（亦即，无需为此做特别设计，如果真发生了此类事件，那么下一个调用 `one_time_init()` 的线程会重新调用 `init()`）。
- 31-2.** 使用线程特有数据重新实现线程安全版的函数 `dirname()` 和 `basename()`（18.14 节）。

第 32 章

线程：线程取消

在通常情况下，程序中的多个线程会并发执行，每个线程各司其职，直至其决意退出，随即会调用函数 `pthread_exit()` 或者从线程启动函数中返回。

有时候，需要将一个线程取消（cancel）。亦即，向线程发送一个请求，要求其立即退出。比如，一组线程正在执行一个运算，一旦某个线程检测到错误发生，需要其他线程退出，取消线程的功能这时就派上用场。还有一种情况，一个由图形用户界面（GUI）驱动的应用程序可能会提供一个“取消”按钮，以便用户可以终止后台某一线程正在执行的任务。这种情况下，主线程（控制图形用户界面）需要请求后台线程退出。

本章就来讨论 POSIX 线程的取消机制。

32.1 取消一个线程

函数 `pthread_cancel()` 向由 `thread` 指定的线程发送一个取消请求。

```
#include <pthread.h>

int pthread_cancel(pthread_t thread);

Returns 0 on success, or a positive error number on error
```

发出取消请求后，函数 `pthread_cancel()` 当即返回，不会等待目标线程的退出。

准确地说，目标线程会发生什么？何时发生？这些都取决于下节将要述及的线程取消状态（state）和类型（type）。

32.2 取消状态及类型

函数 `pthread_setcancelstate()` 和 `pthread_setcanceltype()` 会设定标志，允许线程对取消请求的响应过程加以控制。

```
#include <pthread.h>
```

```
int pthread_setcancelstate(int state, int *oldstate);  
int pthread_setcanceltype(int type, int *oldtype);
```

Both return 0 on success, or a positive error number on error

函数 `pthread_setcancelstate()` 会将调用线程的取消性状态置为参数 `state` 所给定的值。该参数的值如下。

PTHREAD_CANCEL_DISABLE

线程不可取消。如果此类线程收到取消请求，则会将请求挂起，直至将线程的取消状态置为启用。

PTHREAD_CANCEL_ENABLE

线程可以取消。这是新建线程取消性状态的默认值。

线程的前一取消性状态将返回至参数 `oldstate` 所指向的位置。

如果对前一状态没有兴趣，Linux 允许将 `oldstate` 置为 `NULL`。在很多其他的系统实现中，情况也是如此。不过，SUSv3 并没有规范这一特性，所以要保证应用的可移植性，就不能依赖这一特性。应该总是为 `oldstate` 设置一个非 `NULL` 的值。

如果线程执行的代码片段需要不间断地一气呵成，那么临时屏闭线程的取消性状态（`PTHREAD_CANCEL_DISABLE`）就变得很有必要。

如果线程的取消性状态为“启用”（`PTHREAD_CANCEL_ENABLE`），那么对取消请求的处理则取决于线程的取消性类型，该类型可以通过调用函数 `pthread_setcanceltype()` 时的参数 `type` 给定。参数 `type` 有如下值：

PTHREAD_CANCEL_ASYNCHRONOUS

可能会在任何时点（也许是立即取消，但不一定）取消线程。异步取消的应用场景很少，将延后至 32.6 节再做讨论。

PTHREAD_CANCEL_DEFERRED

取消请求保持挂起状态，直至到达取消点（`cancellation point`，见下节）。这也是新建线程的缺省类型。后续各节将介绍延迟取消（`deferred cancelability`）的更多细节。

线程原有的取消类型将返回至参数 `oldtype` 所指向的位置。

与函数 `pthread_setcancelstate()` 的参数 `oldstate` 类似，如果不关心原有取消类型，许多系统实现（包括 Linux）允许将 `oldtype` 置为 `NULL`。同样，SUSv3 也没有规范这一行为，所以需要保障可移植性的应用不应使用这一特性，应该总是为 `oldtype` 设置一个非 `NULL` 值。

当某线程调用 `fork()` 时，子进程会继承调用线程的取消性类型及状态。而当某线程调用 `exec()` 时，会将新程序主线程的取消性类型及状态分别重置为 `PTHREAD_CANCEL_NABLE` 和 `PTHREAD_CANCEL_DEFERRED`。

32.3 取消点

若将线程的取消性状态和类型分别置为启用和延迟，仅当线程抵达某个取消点（`cancellation`

point) 时, 取消请求才会起作用。取消点即是对由实现定义的一组函数之一加以调用。

SUSv3 规定, 实现若提供了表 32-1 中所列的函数, 则这些函数必须是取消点。其中的大部分函数都有能力将线程无限期地堵塞起来。

表 32-1: SUSv3 规定必须是取消点的函数

accept()	nanosleep()	sem_timedwait()
aio_suspend()	open()	sem_wait()
clock_nanosleep()	pause()	send()
close()	poll()	sendmsg()
connect()	pread()	sendto()
creat()	pselect()	sigpause()
fcntl(F_SETLK)	pthread_cond_timedwait()	sigsuspend()
fsync()	pthread_cond_wait()	sigtimedwait()
fdatasync()	pthread_join()	sigwait()
getmsg()	pthread_testcancel()	sigwaitinfo()
getpmsg()	putmsg()	sleep()
lockf(F_LOCK)	putpmsg()	system()
mq_receive()	pwrite()	tcdrain()
mq_send()	read()	usleep()
mq_timedreceive()	readv()	wait()
mq_timedsend()	recv()	waitid()
msgrcv()	recvfrom()	waitpid()
msgsnd()	recvmsg()	write()
msync()	select()	writv()

除表 32-1 所列函数之外, SUSv3 还指定了大量函数, 系统实现可以将其定义为取消点。其中包括 `stdio` 函数、`dlopen` API、`syslog` API、`nftw()`、`popen()`、`semop()`、`unlink()`, 以及从诸如 `utmp` 之类的系统文件中获取信息的各种函数。可移植应用程序必须正确处理这一情况: 线程在调用这些函数时有可能遭到取消。

SUSv3 规定, 除了上述两组必须或可能是可取消点的函数之外, 不得将标准中的任何其他函数视为取消点 (亦即, 调用这些函数不会招致线程取消, 可移植程序无需加以处理)。

SUSv4 在必须的可取消点函数列表中增加了 `openat()`, 并移除了函数 `sigpause()` (将其移至“可能的”取消点函数列表中) 和函数 `usleep()` (已从标准中删除)。

系统实现可随意将标准并未规范的其他函数标记为取消点。任何可能造成堵塞的函数 (有可能是因为需要访问文件) 都是取消点的理想候选对象。出于这一理由, `glibc` 将其中的许多非标准函数标记为取消点。

线程一旦收到取消请求, 且启用了取消性状态并将类型置为延迟, 则其会在下次抵达取消点时终止。如果该线程尚未分离 (`not detached`), 那么为防止其变为僵尸线程, 必须由其他线程对其进行连接 (`join`)。连接之后, 返回至函数 `pthread_join()` 中第二个参数的将是一个特殊值: `PTHREAD_CANCELED`。

示例程序

程序清单 32-1 是一个使用 `pthread_cancel()` 的简单例子。主程序创建一个线程来执行无限循

环，每次都在休眠一秒后打印循环计数器的值。（仅当向其发送取消请求或者进程退出时，该线程才会终止。）同时，主程序将休眠 3 秒，随即向新创建的线程发送取消请求。程序运行结果如下：

```
$ ./t_pthread_cancel
New thread started
Loop 1
Loop 2
Loop 3
Thread was canceled
```

程序清单 32-1：调用 pthread_cancel()取消线程

```
threads/thread_cancel.c

#include <pthread.h>
#include "tspi_hdr.h"

static void *
threadFunc(void *arg)
{
    int j;
    printf("New thread started\n");    /* May be a cancellation point */
    for (j = 1; ; j++) {
        printf("Loop %d\n", j);        /* May be a cancellation point */
        sleep(1);                      /* A cancellation point */
    }

    /* NOTREACHED */
    return NULL;
}

int
main(int argc, char *argv[])
{
    pthread_t thr;
    int s;
    void *res;

    s = pthread_create(&thr, NULL, threadFunc, NULL);
    if (s != 0)
        errExitEN(s, "pthread_create");

    sleep(3);                          /* Allow new thread to run a while */

    s = pthread_cancel(thr);
    if (s != 0)
        errExitEN(s, "pthread_cancel");

    s = pthread_join(thr, &res);
    if (s != 0)
        errExitEN(s, "pthread_join");

    if (res == PTHREAD_CANCELED)
        printf("Thread was canceled\n");
    else
        printf("Thread was not canceled (should not happen!)\n");

    exit(EXIT_SUCCESS);
}

threads/thread_cancel.c
```

32.4 线程可取消性的检测

在程序清单 32-1 中，由 `main()` 创建的线程会执行到属于取消点的函数（`sleep()` 属于取消点，`printf()` 可能也是），因而会接受取消请求。不过，假设线程执行的是一个不含取消点的循环（计算密集型 [compute-bound] 循环），这时，线程永远也不会响应取消请求。

函数 `pthread_testcancel()` 的目的很简单，就是产生一个取消点。线程如果已有处于挂起状态的取消请求，那么只要调用该函数，线程就会随之终止。

```
#include <pthread.h>

void pthread_testcancel(void);
```

当线程执行的代码未包含取消点时，可以周期性地调用 `pthread_testcancel()`，以确保对其他线程向其发送的取消请求做出及时响应。

32.5 清理函数（cleanup handler）

一旦有处于挂起状态的取消请求，线程在执行到取消点时如果只是草草收场，这会将共享变量以及 Pthreads 对象（例如互斥量）置于一种不一致状态，可能导致进程中其他线程产生错误结果、死锁，甚至造成程序崩溃。为规避这一问题，线程可以设置一个或多个清理函数，当线程遭取消时会自动运行这些函数，在线程终止之前可执行诸如修改全局变量，解锁互斥量等动作。

每个线程都可以拥有一个清理函数栈。当线程遭取消时，会沿该栈自顶向下依次执行清理函数，首先会执行最近设置的函数，接着是次新的函数，以此类推。当执行完所有清理函数后，线程终止。

函数 `pthread_cleanup_push()` 和 `pthread_cleanup_pop()` 分别负责向调用线程的清理函数栈添加和移除清理函数。

```
#include <pthread.h>

void pthread_cleanup_push(void (*routine)(void*), void *arg);
void pthread_cleanup_pop(int execute);
```

`pthread_cleanup_push()` 会将参数 `routine` 所含的函数地址添加到调用线程的清理函数栈顶。参数 `routine` 是一个函数指针，格式如下：

```
void
routine(void *arg)
{
    /* Code to perform cleanup */
}
```

执行 `pthread_cleanup_push()` 时给定的 `arg` 值，会作为调用清理函数时的参数。其参数类型为 `void*`，如果强制装换使用得当，那么通过该参数可以传入各种类型的数据。

通常，线程如在执行一段特殊代码时遭到取消，才需要执行清理动作。如果线程顺利执行完这段代码而未遭取消，那么就不再需要清理。所以，每个对 `pthread_cleanup_push()` 的调用都会伴随着对 `pthread_cleanup_pop()` 的调用。此函数从清理函数栈中移除最顶层的函数。如

果参数 `execute` 非零，那么无论如何都会执行清理函数。在函数未遭取消而又希望执行清理动作的情况下，这会非常方便。

尽管这里把 `pthread_cleanup_push()` 和 `pthread_cleanup_pop()` 描述为函数，SUSv3 却允许将它们实现为宏（macro），可展开为分别由 `{` 和 `}` 所包裹的语句序列。并非所有的 UNIX 都这样做，不过包括 Linux 在内的很多系统都是使用宏来实现的。这意味着，`pthread_cleanup_push()` 和与其配对的 `pthread_cleanup_pop()` 属于同一个语法块，必须一一对应。（一旦以此方式来实现 `pthread_cleanup_push()` 和 `pthread_cleanup_pop()`，在对两者的调用间所声明的变量，其作用域将受限于这一语法块。）例如，以下代码就不正确：

```
pthread_cleanup_push(func, arg);
...
if (cond) {
    pthread_cleanup_pop(0);
}
```

为便于编码，若线程因调用 `pthread_exit()` 而终止，则也会自动执行尚未从清理函数栈中弹出（pop）的清理函数。线程正常返回（return）时不会执行清理函数。

示例程序

程序清单 32-2 提供了一个使用清理函数的简单例子。主程序创建线程⑧，线程首先分配一块内存③，并将其地址存储于 `buf` 中，接着锁定互斥量 `mtx`④。因为线程可能会遭到取消，所以调用 `pthread_cleanup_push()`⑤ 设置清理函数，并将存储于 `buf` 中的地址作为参数传入。如果执行到清理函数，那么清理函数会释放内存①并解锁互斥量②。

线程接着进入循环，等待对条件变量 `cond` 的通知⑥。取决于可执行程序是否带有命令行参数，此循环会以以下两种方式结束。

- 若无命令行参数，则由 `main()`⑨ 函数取消线程。此时，取消操作发生在对 `pthread_cond_wait()`⑥ 的调用中，此函数可见于程序清单 32-1 中，属于取消点。作为取消动作的一部分，会自动调用由 `pthread_cleanup_push()` 设置的清理函数。
- 如果指定了命令行参数，那么在将全局变量 `glob` 设置为非零后，通知条件变量⑩。此时，线程会一直执行到 `pthread_cleanup_pop()`⑦，因为向此函数传入了非零参数，所以依然会调用清理函数。

程序清单 32-2：使用清理函数

```
----- threads/thread_cleanup.c
#include <pthread.h>
#include "tspi_hdr.h"

static pthread_cond_t cond = PTHREAD_COND_INITIALIZER;
static pthread_mutex_t mtx = PTHREAD_MUTEX_INITIALIZER;
static int glob = 0;          /* Predicate variable */

static void      /* Free memory pointed to by 'arg' and unlock mutex */
cleanupHandler(void *arg)
{
    int s;

    printf("cleanup: freeing block at %p\n", arg);
    ① free(arg);

    printf("cleanup: unlocking mutex\n");
}
```

```

②  s = pthread_mutex_unlock(&mtx);
    if (s != 0)
        errExitEN(s, "pthread_mutex_unlock");
}

static void *
threadFunc(void *arg)
{
    int s;
    void *buf = NULL;          /* Buffer allocated by thread */

③  buf = malloc(0x10000);      /* Not a cancellation point */
    printf("thread: allocated memory at %p\n", buf);

④  s = pthread_mutex_lock(&mtx); /* Not a cancellation point */
    if (s != 0)
        errExitEN(s, "pthread_mutex_lock");

⑤  pthread_cleanup_push(cleanupHandler, buf);

    while (glob == 0) {
⑥      s = pthread_cond_wait(&cond, &mtx); /* A cancellation point */
        if (s != 0)
            errExitEN(s, "pthread_cond_wait");
    }

    printf("thread: condition wait loop completed\n");
⑦  pthread_cleanup_pop(1);      /* Executes cleanup handler */
    return NULL;
}

int
main(int argc, char *argv[])
{
    pthread_t thr;
    void *res;
    int s;

⑧  s = pthread_create(&thr, NULL, threadFunc, NULL);
    if (s != 0)
        errExitEN(s, "pthread_create");

    sleep(2);                  /* Give thread a chance to get started */

    if (argc == 1) {           /* Cancel thread */
        printf("main: about to cancel thread\n");
⑨      s = pthread_cancel(thr);
        if (s != 0)
            errExitEN(s, "pthread_cancel");
    } else {                  /* Signal condition variable */
        printf("main: about to signal condition variable\n");
        glob = 1;
⑩      s = pthread_cond_signal(&cond);
        if (s != 0)
            errExitEN(s, "pthread_cond_signal");
    }

⑪  s = pthread_join(thr, &res);
    if (s != 0)
        errExitEN(s, "pthread_join");
}

```

```

    if (res == PTHREAD_CANCELED)
        printf("main:   thread was canceled\n");
    else
        printf("main:   thread terminated normally\n");

    exit(EXIT_SUCCESS);
}

```

threads/thread_cleanup.c

主程序与遭终止线程建立连接①，并报告线程是遭到取消还是正常终止。

如果执行程序清单 32-2 中程序且不带任何命令行参数，那么 main() 函数会调用 pthread_cancel()，清理函数也会得以自动执行。输出如下：

```

$ ./thread_cleanup
thread: allocated memory at 0x804b050
main:   about to cancel thread
cleanup: freeing block at 0x804b050
cleanup: unlocking mutex
main:   thread was canceled

```

如果运行该程序且带有命令行参数，那么 main() 将 glob 设置为 1 并通知条件变量，清理函数则通过 pthread_cleanup_pop() 的调用执行，可以看到如下结果：

```

$ ./thread_cleanup s
thread: allocated memory at 0x804b050
main:   about to signal condition variable
thread: condition wait loop completed
cleanup: freeing block at 0x804b050
cleanup: unlocking mutex
main:   thread terminated normally

```

32.6 异步取消

如果设定线程为可异步取消时（取消性类型为 PTHREAD_CANCEL_ASYNCHRONOUS），可以在任何时点将其取消（亦即，执行任何机器指令时），取消动作不会拖延到下一个取消点才执行。

异步取消的问题在于，尽管清理函数依然会得以执行，但处理函数却无从得知线程的具体状态。程序清单 32-2 采用了延时取消类型，只有在执行到 pthread_cond_wait() 这一唯一的取消点时，线程才会遭到取消。此时可知，已将 buf 初始化为指向新分配的内存块，并且锁定了互斥量 mtx。不过，要是采用异步取消，就可以在任意点取消线程（例如，调用 malloc() 之前，调用 malloc() 与锁定互斥量之间，或者锁定互斥量之后）。清理函数无法知道将在哪里发生取消动作，或者准确地来说，清理函数不清楚需要执行哪些清理步骤。此外，线程也很可能在对 malloc() 的调用期间被取消，这极有可能造成后续的混乱（见 7.1.3 节）。

作为一般性原则，可异步取消的线程不应该分配任何资源，也不能获取互斥量或锁。这导致大量库函数无法使用，其中就包括 Pthreads 函数的大部分。（SUSv3 中有 3 处例外 pthread_cancel()、pthread_setcancelstate() 以及 pthread_setcanceltype()，规范明确要求将它们实现为“异步取消安全（async-cancel-safe）”，亦即，实现必须确保在可异步取消的线程中可以安全调用它们。）换言之，异步取消功能鲜有应用场景，其中之一就是：取消在执行计算密集型循环的线程。

32.7 总结

函数 `pthread_cancel()` 允许某线程向另一个线程发送取消请求，要求目标线程终止。

目标线程如何响应，取决于其取消性状态和类型。如果禁用线程的取消性状态，那么请求会保持挂起（`pending`）状态，直至将线程的取消性状态置为启用。如果启用取消性状态，那么线程何时响应请求则依赖于取消性类型。若类型为延迟取消，则在线程下一次调用某个取消点（由 `SUSv3` 标准所规定的一系列函数之一）时，取消发生。如果为异步取消类型，取消动作随时可能发生（鲜有使用）。

线程可以设置一个清理函数栈，其中的清理函数属于由开发人员定义的函数，当线程遭到取消时，会自动调用这些函数以执行清理工作（例如，恢复共享变量状态，或解锁互斥量）。

更多信息

请参考列于 29.10 节的深入信息来源。

第 33 章

线程：更多细节

本章将就 POSIX 线程库各方面的细节做深入探讨，涉及线程与传统 UNIX API——尤其是信号以及进程控制原语（`fork()`、`exec()`和`_exit()`）——之间的交互，同时对 Linux 上的两个 POSIX 线程实现（LinuxThreads 和 NPTL）加以概括，并特别指出了这些实现与 SUSv3 Pthreads 标准间的偏差所在。

33.1 线程栈

创建线程时，每个线程都有一个属于自己的线程栈，且大小固定。在 Linux/x86-32 架构上，除主线程外的所有线程，其栈的缺省大小均为 2MB。（在一些 64 位架构下，默认尺寸要大一些，例如，IA-64 有 32MB。）为了应对栈的增长（参考图 29-1），主线程栈的空间要大出许多。

偶尔，也需要改变线程栈的大小。在通过线程属性对象创建线程时，调用函数 `pthread_attr_setstacksize()` 所设置的线程属性（29.8 节）决定了线程栈的大小。而使用与之相关的另一函数 `pthread_attr_setstack()`，可以同时控制线程栈的大小和位置，不过设置栈的地址将降低程序的可移植性。手册页（manual page）提供了对这些函数的具体说明。

更大的线程栈可以容纳大型的自动变量或者深度的嵌套函数调用（也许是递归调用），这是改变每个线程栈大小的原因之一。而另一方面，应用程序可能希望减小每个线程栈，以便进程可以创建更多的线程。例如，在 x86-32 系统中，用户（模式）可访问的虚拟地址空间是 3GB，而 2MB 的缺省栈大小则意味着最多只能创建 1500 个线程。（更为准确的最大值还视乎文本段、数据段、共享函数库等对虚拟内存的消耗量。）特定架构的系统上，可采用的线程栈大小最小值可以通过调用 `sysconf(_SC_THREAD_STACK_MION)` 来确定。在 Linux/x86-32 上的 NPTL 实现中，该调用返回 16384。

在 NPTL 线程实现中，如果对线程栈尺寸资源限制（`RLIMIT_STACK`）的设置不同于 `unlimited`，那么创建线程时会以其作为默认值。对该限制的设置必须在运行程序之前，通常通过执行 shell 内建命令 `ulimit -s` 完成（在 C shell 下命令为 `limit stacksize`）。在主程序中调用 `setrlimit()` 来设置限制的办法可能行不通，因为 NPTL 在调用 `main()` 之前的运行时初始化期间就已经确定了默认的栈大小。

33.2 线程和信号

UNIX 信号模型是基于 UNIX 进程模型而设计的，问世比 Pthreads 要早几十年。自然而然，信号与线程模型之间存在一些明显的冲突。主要是因为，一方面，针对单线程进程要保持传统的信号语义（Pthreads 不应改变传统进程的信号语义），与此同时，又需要开发出适用于多线程进程环境的新信号模型。

信号与线程模型之间的差异意味着，将二者结合使用，将会非常复杂，应尽可能加以避免。尽管如此，有的时候还是必须在多线程程序中处理信号问题。本节将讨论信号与线程间的交互，并描述在多线程程序中处理信号的各种有效函数。

33.2.1 UNIX 信号模型如何映射到线程中

要了解 UNIX 信号如何映射到 Pthreads 模型，就需要了解，信号模型的哪些方面属于进程层面（由进程中的所有线程所共享），哪些方面是属于进程中的单个线程层面。如下是对其关键点的汇总。

- 信号动作属于进程层面。如果某进程的任一线程收到任何未经（特殊）处理的信号，且其缺省动作为 stop 或 terminate，那么将停止或者终止该进程的所有线程。
- 对信号的处置属于进程层面，进程中的所有线程共享对每个信号的处置设置。如果某一线程使用函数 `sigaction()` 为某类信号（比如，SIGINT）创建了处理函数，那么当收到 SIGINT 时，任何线程都会去调用该处理函数。与之类似，如果将对信号的处置设置为忽略（ignore），那么所有线程都会忽略该信号。
- 信号的发送既可针对整个进程，也可针对某个特定线程。满足如下三者之一的信号当属面向线程的。
 - 信号的产生源于线程上下文中对特定硬件指令的执行（即 22.4 节所描述的硬件异常：SIGBUS、SIGFPE、SIGILL 和 SIGSEGV）。
 - 当线程试图对已断开的（broken pipe）管道进行写操作时所产生的 SIGPIPE 信号。
 - 由函数 `pthread_kill()` 或 `pthread_sigqueue()` 所发出的信号，这些函数（由 33.2.3 节描述）允许线程向同一进程下的其他线程发送信号。
由其他机制产生的所有信号都是面向进程的。例如，其他进程通过调用 `kill()` 或者 `sigqueue()` 所发送的信号；用户键入特殊的终端字符所产生的信号，诸如 SIGINT 和 SIGTSTP；还有一些信号由软件事件产生，例如终端窗口大小的调整（SIGWINCH）或者定时器到期（例如，SIGALRM）。
- 当多线程程序收到一个信号，且该进程已然为此信号创建了信号处理程序时，内核会任选一条线程来接收这一信号，并在该线程中调用信号处理程序对其进行处理。这种行为与信号的原始语义保持了一致。让进程针对单个信号重复处理多次是没有意义的。
- 信号掩码（mask）是针对每个线程而言。（对于多线程程序来说，并不存在一个作用于整个进程范围的信号掩码，可以管理所有线程。）使用 Pthreads API 所定义的函数 `pthread_sigmask()`，各线程可独立阻止或放行各种信号。通过操作每个线程的信号掩码，应用程序可以控制哪些线程可以处理进程收到的信号。
- 针对为整个进程所挂起（pending）的信号，以及为每条线程所挂起的信号，内核都分别维护有记录。调用函数 `sigpending()` 会返回为整个进程和当前线程所挂起信号的并集。在再创

建的线程中，每线程的挂起信号集合初始时空。可将一个针对线程的信号仅向目标线程投送。如果该信号遭线程阻塞，那么它会一直保持挂起，直至线程将其放行（或者线程终止）。

- 如果信号处理程序中断了对 `pthread_mutex_lock()` 的调用，那么该调用总是会自动重新开始。如果一个信号处理函数中断了对 `pthread_cond_wait()` 的调用，则该调用要么自动重新开始（Linux 就是如此），要么返回 0，表示遭遇了假唤醒（如 30.2.3 节所述，此时，设计良好的应用程序会重新检查相应的判断条件并重新发起调用）。SUSv3 对这两个函数的行为要求与此处的描述一致。
- 备选信号栈是每线程特有的（参考 21.3 节对函数 `sigaltstack()` 的描述）。新创建的线程并不从创建者处继承备选信号栈。

更确切地说，SUSv3 规定每个内核调度实体（KSE）都有一个单独的备选信号栈。在按 1:1 比例实现线程的系统中，例如 Linux，每一个线程对应一个 KSE（见 33.4 节）。

33.2.2 操作线程信号掩码

刚创建的新线程会从其创建者处继承信号掩码的一份拷贝。线程可以使用 `pthread_sigmask()` 来改变或/并获取当前的信号掩码。

```
#include <signal.h>

int pthread_sigmask(int how, const sigset_t *set, sigset_t *oldset);

Returns 0 on success, or a positive error number on error
```

除了所操作的是线程信号掩码之外，`pthread_sigmask()` 与 `sigprocmask()` 的用法完全相同（见 20.10 节）。

SUSv3 特别指出，注明在多线程程序中使用函数 `sigprocmask()`，其结果是未定义的，也无法保证程序的可移植性。事实上，函数 `sigprocmask()` 和 `pthread_sigmask()` 在包括 Linux 在内的很多系统实现中是相同的。

33.2.3 向线程发送信号

函数 `pthread_kill()` 向同一进程下的另一线程发送信号 `sig`。目标线程由参数 `thread` 标识。

```
#include <signal.h>

int pthread_kill(pthread_t thread, int sig);

Returns 0 on success, or a positive error number on error
```

因为仅在同一进程中可保证线程 ID 的唯一性（参见 29.5 节），所以无法调用 `pthread_kill()` 向其他进程中的线程发送信号。

在实现函数 `pthread_kill()` 时，使用了 Linux 特有的 `tgkill(tgid, tid, sig)` 系统调用，将信号 `sig` 发送给由 `tid`（由 `gettid()` 所返回的内核线程 ID）标识的线程，该线程从属于由 `tgid` 标识的线程组中。更多细节，请参考 `tgkill(2)` 手册页。

Linux 特有的函数 `pthread_sigqueue()` 将 `pthread_kill()` 和 `sigqueue()` 的功能合二为一（见

22.8.1 节): 向同一进程中的另一线程发送携带数据的信号。

```
#define _GNU_SOURCE
#include <signal.h>

int pthread_sigqueue(pthread_t thread, int sig, const union signal value);

Returns 0 on success, or a positive error number on error
```

与函数 `pthread_kill()` 一样, `sig` 表示将要发送的信号, `thread` 标识目标线程。参数 `value` 则指定了伴随信号的数据, 其使用方式与函数 `sigqueue()` 中的对应参数相同。

函数 `pthread_sigqueue()` 从 2.11 版开始加入 `glibc` 函数库中, 同时需要内核的支持。始于 Linux 2.6.31, 内核通过系统调用 `rt_tgsigqueueinfo()` 来提供这一支持。

33.2.4 妥善处理异步信号

第 20 章至第 22 章所探讨的各种因素 (诸如, 可重入问题、重启遭中断的系统调用, 以及避免竞争条件), 当使用信号处理函数对异步产生的信号加以处理时, 这些都将导致情况变得复杂。另外, 没有任何 Pthreads API 属于异步信号安全 (async-signal-safe) 函数, 均无法在信号处理函数 (21.1.2 节) 中安全加以调用。因为这些原因, 所以当多线程应用程序必须处理异步产生的信号时, 通常不应该将信号处理函数作为接收信号到达的通知机制。相反, 推荐的方法如下。

- 所有线程都阻塞进程可能接收的所有异步信号。最简单的方法是, 在创建任何其他线程之前, 由主线程阻塞这些信号。后续创建的每个线程都会继承主线程信号掩码的一份拷贝。
- 再创建一个专用线程, 调用函数 `sigwaitinfo()`、`sigtimedwait()` 或 `sigwait()` 来接收收到的信号。22.10 节对 `sigwaitinfo()` 和 `sigtimedwait()` 做了说明。下面则对 `sigwait()` 有所描述。

这一方法的优势在于, 同步接收异步产生的信号。当接收到信号时, 专有线程可以安全地修改共享变量 (在互斥量的保护之下), 并可调用并非异步信号安全 (non-async-signal-safe) 的函数。也可以就条件变量发出信号, 并采用其他线程或进程的通讯及同步机制。

函数 `sigwait()` 会等待 `set` 所指信号集合中任一信号的到达, 接收该信号, 且在参数 `sig` 中将其返回。

```
#include <signal.h>

int sigwait(const sigset_t *set, int *sig);

Returns 0 on success, or a positive error number on error
```

除了以下不同以外, `sigwait()` 的操作与 `sigwaitinfo()` 相同。

- 函数 `sigwait()` 只返回信号编号, 而非返回一个描述信号信息的 `siginfo_t` 类型结构。
 - 并且返回值与其他线程相关函数保持一致 (而非传统 UNIX 系统调用返回的 0 或 -1)。
- 如有多个线程在调用 `sigwait()` 等待同一信号, 那么当信号到达时只有一个线程会实际接收到, 也无法确定收到信号的会是哪条线程。

33.3 线程和进程控制

与信号机制类似, `exec()`、`fork()` 和 `exit()` 的问世均早于 Pthreads API。接下来的段落将指出

在多线程程序中使用此类系统调用所应关注的细节。

线程和 exec()

只要有任一线程调用了 `exec()` 系列函数之一时，调用程序将被完全替换。除了调用 `exec()` 的线程之外，其他所有线程都将立即消失。没有任何线程会针对线程特有数据执行解构函数（`destructor`），也不会调用清理函数（`cleanup handler`）。该进程的所有互斥量（为进程私有）和属于进程的条件变量都会消失。调用 `exec()` 之后，调用线程的线程 ID 是不确定的。

线程和 fork()

当多线程进程调用 `fork()` 时，仅会将发起调用的线程复制到子进程中。（子进程中该线程的线程 ID 与父进程中发起 `fork()` 调用线程的线程 ID 相一致。）其他线程均在子进程中消失，也不会为这些线程调用清理函数以及针对线程特有数据的解构函数。这将导致如下一些问题。

- 虽然只将发起调用的线程复制到子进程中，但全局变量的状态以及所有的 `Pthreads` 对象（如互斥量、条件变量等）都会在子进程中得以保留。（因为在父进程中为这些 `Pthreads` 对象分配了内存，而子进程则获得了该内存的一份拷贝。）这会导致很棘手的问题。例如，假设在调用 `fork()` 时，另一线程已然锁定了某一互斥量，且对某一全局数据结构的更新也做到了一半。此时，子进程中的该线程无法解锁这一互斥量（因为其并非该互斥量的属主），如果试图获取这一互斥量，线程会遭阻塞。此外，子进程中的全局数据结构拷贝可能也处于不一致状态，因为对其进行更新的线程在执行到一半时消失了。
- 因为并未执行清理函数和针对线程特有数据的解构函数，多线程程序的 `fork()` 调用会导致子进程的内存泄漏。另外，子进程中的线程很可能无法访问（父进程中）由其他线程所创建的线程特有数据项，因为（子进程）没有相应的引用指针。

由于这些问题，推荐在多线程程序中调用 `fork()` 的唯一情况是：其后紧跟对 `exec()` 的调用。因为新程序会覆盖原有内存，`exec()` 将导致子进程的所有 `Pthreads` 对象消失。

对于那些必须执行 `fork()`，而其后又无 `exec()` 跟随的程序来说，`Pthreads` API 提供了一种机制：`fork` 处理函数（`handler`）。可以利用函数 `pthread_atfork()` 来创建 `fork` 处理函数，格式如下：
`pthread_atfork(prepare_func, parent_func, child_func);`

每一次 `pthread_atfork()` 调用都会将 `prepare_func` 添加到一个函数列表中，在调用 `fork()` 创建新的子进程之前，会（按与注册次序相反的顺序）自动执行该函数列表中的函数。与之类似，会将 `parent_func` 和 `child_func` 添加到一个函数列表中，在 `fork()` 返回前，将分别在父、子进程中（按注册顺序）自动运行。

在使用线程的函数库中，有时候 `fork` 处理函数很实用。如果没有这一机制，对于那些随意调用了此函数库和 `fork()`，又对函数库创建的其他线程一无所知的应用程序，函数库还真是无计可施。

调用 `fork()` 所产生的子进程从调用 `fork()` 的线程处继承 `fork` 处理函数。执行 `exec()` 期间，`fork` 处理函数将不再保留（因为处理函数的代码会在执行 `exec()` 的过程中遭到覆盖）。

关于 `fork` 处理函数及其使用的更多细节可以参考 [Butenhof, 1996]。

在 Linux 上，如果使用 NPTL 线程库的程序执行了 `vfork()`，那么将不再调用 `fork` 处理函数。不过，在使用 `LinuxThreads` 程序的同一情况下却有效。

线程与 exit()

如果任何线程调用了 `exit()`，或者主线程执行了 `return`，那么所有线程都将消失，也不会

执行线程特有数据的解构函数以及清理函数。

33.4 线程实现模型

本节将涉及一些理论知识，简要阐述实现线程 API 的 3 种不同模型，从而为 33.5 节中关于 Linux 线程实现的讨论提供必要的背景知识。这 3 种实现模型的差异主要集中在线程如何与内核调度实体（KSE, Kernel Scheduling Entity）相映射。KSE 是内核分配 CPU 以及其他系统资源的（对象）单位。（在早于线程而出现的传统 UNIX 中，KSE 等同于进程。）

多对一（M:1）实现（用户级线程）

在 M:1 线程实现中，关乎线程创建、调度以及同步（互斥量的锁定，条件变量的等待等）的所有细节全部由进程内用户空间（user-space）的线程库来处理。对于进程中存在的多个线程，内核一无所知。

M:1 实现的优势不多，其中最大的优点在于，许多线程操作（例如线程的创建和终止、线程上下文间的切换、互斥量以及条件变量操作）速度都很快，因为无需切换到内核模式。此外，由于线程库无需内核支持，所以 M:1 实现在系统间的移植相对要容易一些。

不过，M:1 实现也存在一些严重缺陷。

- 当一线程发起系统调用（如 `read()`）时，控制由用户空间的线程库转交给内核。这就意味着，如果 `read()` 调用遭到阻塞，那么所有的线程都会被阻塞。
- 内核无法调度进程中的这些线程。因为内核并不知晓进程中存在这些线程，也就无法在多处理器平台上将各线程调度给不同的处理器。另外，也不可能将一进程中某线程的优先级调整为高于其他进程中的线程，这是没有意义的，因为对线程的调度完全在进程中处理。

一对一（1:1）实现（内核级线程）

在 1:1 线程实现中，每一线程映射一个单独的 KSE。内核分别对每个线程做调度处理。线程同步操作通过内核系统调用实现。

1:1 实现消除了 M:1 实现的种种弊端。遭阻塞的系统调用不会导致进程的所有线程被阻塞，在多处理器硬件平台上，内核还可以将进程中的多个线程调度到不同的 CPU 上。

不过，因为需要切换到内核模式，所以诸如线程创建、上下文切换以及同步操作就要慢一些。另外，为每个线程分别维护一个 KSE 也需要开销，如果应用程序包含大量线程，则可能对内核调度器造成严重的负担，降低系统的整体性能。

尽管有这些缺点，1:1 实现通常更胜于 M:1 实现。LinuxThreads 和 NPTL 都采用 1:1 模型。

在 NPTL 的开发期间，为了使得包含数以千计线程的进程得以高效运行，投入了巨大的努力，对内核调度器进行了重写并设计了新的线程实现。后续的测试也显示了预期目标的达成。

多对多（M:N）实现（两级模型）

M:N 实现旨在结合 1:1 和 M:1 模型的优点，避免二者的缺点。

在 M:N 模型中，每个进程都可拥有多个与之相关的 KSE，并且也可以把多个线程映射到一个 KSE。这种设计允许内核将同一应用的线程调度到不同的 CPU 上运行，同时也解决了随线程数量而放大的性能问题。

M:N 模型的最大问题是过于复杂。线程调度任务由内核及用户空间的线程库共同承担，二者之

间势必要进行分工协作和信息交换。在 M:N 模型下，按照 SUSv3 标准要求来管理信号也极为复杂。

最初曾考虑采用 M:N 模型来实现 NPTL 线程库，但若要保证 Linux 调度器即使在处理大量 KSE 的情况下也能应对自如，则需要对内核所作的改动范围过大，可能也没有必要，故而否决了这一方案。

33.5 Linux POSIX 线程的实现

针对 Pthreads API：Linux 下有两种实现。

- **LinuxThreads**：这是最初的 Linux 线程实现，由 Xavier Leroy 开发。
- **NPTL** (Native POSIX Threads Library)：这是 Linux 线程实现的现代版，由 Ulrich Drepper 和 Ingo Molnar 开发，以取代 LinuxThreads。NPTL 的性能优于 LinuxThreads，也更符合 SUSv3 的 Pthreads 标准。对 NPTL 的支持需要修改内核，这始于 Linux 2.6。

一度，人们曾将由 IBM 开发的另一线程实现——NGPT (Next Generation POSIX Threads) 视为 LinuxThreads 的继任者。NGPT 采用 M:N 模型设计，性能明显优于 LinuxThreads。不过，NPTL 的开发者决意推出新的实现。NPTL 的设计方法有所调整，采用 1:1 模型，性能也优于 NGPT。随着 NPTL 的发布，对 NGPT 的开发也随之终止。

后续各节将讨论这两种实现的更多细节，并将二者对 SUSv3 Pthreads 标准的背离之处一一指出。

此处值得强调的是：LinuxThreads 实现已经过时，并且 glibc 从 2.4 版本开始也已不再支持它，所有新的线程库开发都基于 NPTL。

33.5.1 LinuxThreads

多年以来，LinuxThreads 曾一直是 Linux 上的主流线程实现，也能够满足各种线程化应用程序实现的需要。LinuxThreads 实现的要点如下。

- 线程的创建使用了 `clone()`，并指定有如下标志：
`CLONE_VM | CLONE_FILES | CLONE_FS | CLONE_SIGHAND`
这意味着，LinuxThreads 线程共享虚拟内存、文件描述符、文件系统相关信息 (umask, 根目录和当前工作目录) 以及信号处置。不过，线程间并不共享进程 ID 和父进程 ID。
- 除了由应用程序创建的线程以外，LinuxThreads 还会创建一个附加的管理线程，负责处理其他线程的创建和终止。
- LinuxThreads 利用信号来处理内部的操作。对于支持实时信号的内核来说 (Linux 2.2 及以后版本)，会使用头 3 个实时信号。对于老版本内核，则使用 SIGUSR1 和 SIGUSR2。应用程序不能使用这些信号。(对于各种线程同步操作而言，使用信号会导致较高延迟。)

LinuxThreads 对标准行为的背离之处

LinuxThreads 在很多方面与 SUSv3 Pthreads 标准并不一致。(LinuxThreads 实现受限于其开发时可用的内核特性，而在此范围内，则会尽量保持一致。) 以下列表对背离之处作了概述。

- 在同一进程的不同线程中调用 `getpid()` 会返回不同的值。调用 `getppid()` 则反应了如下事实：除了主线程以外的所有线程都由进程的管理线程创建 (`getppid()` 会返回管理线程

的进程 ID)。在主线程和其他线程中，对 `getppid()`调用的返回值相同。

- 如果某线程调用 `fork()`创建了一子进程，那么任何其他线程都可使用 `wait()`或类似技术来获取子进程的终止状态。不过，事实并非如此，只有创建子进程的线程才能使用 `wait()`。
- 如果某线程执行了 `exec()`，那么按照 SUSv3 要求，将终止所有其他线程。不过，只要调用 `exec()`的是主线程之外的线程，那么产生进程则与调用线程拥有相同的进程 ID，而与主线程的进程 ID 不同。而依照 SUSv3 标准，该进程 ID 应与主线程的进程 ID 保持一致。
- 线程之间不会共享凭证（用户 ID 与用户组 ID）。当一多线程进程执行一个 `set-user-ID` 程序时，将导致线程之间无法通过 `pthread_kill()`来发送信号，因为这两个线程的凭证已经发生了改变，发送线程不再有权发信号给目标线程（请参考图 20-2）。另外，由于 `LinuxThreads` 实现在内部使用了信号，一旦线程改变了自身的凭证，那么各种 `Pthreads` 操作有可能失败或者挂起（hang）。
- 还未能顾及 SUSv3 关于线程与信号间交互规范的各个方面。
 - 采用 `kill()`或者 `sigqueue()`向某进程发送的信号，应该由目标进程中不阻塞该信号的任意线程来接收和处理。不过，因为 `LinuxThreads` 线程的进程 ID 不同，所以只能将信号送给特定线程。要是该线程阻塞这一信号，即使其他线程并不阻塞此信号，它也会一直保持挂起（pending）。
 - `LinuxThreads` 并不支持信号为整个进程挂起的概念，只支持每个线程分别挂起信号。
 - 如果信号针对包含某个多线程应用的进程组，那么应用中的所有线程（即每个创建了信号处理函数的线程）都会处理该信号，而不是由（任意）某个线程去处理。例如，键入某个终端字符就会产生针对前台进程组的任务控制（`job-control`）信号。
 - 备选信号栈设置（由 `sigaltstack()`建立）是针对每个线程的。不过，如果新线程不慎从 `pthread_create()`调用者那里继承了备选信号栈设置，那么这两个线程将共享同一备选信号栈。SUSv3 规定新线程启动时不应定义备选信号栈。`LinuxThreads` 这一“违规”操作的后果是，如果两个线程恰巧在它们共享的备选信号栈中同时处理不同的信号，很可能会导致混乱（例如，程序崩溃）。这一问题可能非常难以重现，也很难调试，因为问题的出现依赖于在同一时点处理两个信号，这种情况极其罕见。

在使用 `LinuxThreads` 的程序中，新线程可以通过调用 `sigaltstack()`来确保使用与其创建者线程不同的备选信号栈（或许根本就没有栈）。不过，可移植程序（以及创建线程的库函数）并不知道这些，因为在其他实现上这并非必须。另外，即使采用这种技术，依然可能产生竞争条件：新线程在有机会调用 `sigaltstack()`之前依然有可能接收并处理备选栈上的信号。

- 线程不共享一般的任务号以及进程组号。不能利用 `setsid()`和 `setpgid()`系统调用去改变多线程程序中的会话号以及进程组号。
- 使用 `fcntl()`建立的记录锁也不能共享。重复地对同一类型的锁的请求是无法合并在一起执行的。
- 线程不共享资源的限制。SUSv3 定义资源限制是一种进程范围的属性。
- 函数 `times()`返回的 CPU 时间以及由 `getrusage()`返回的资源使用信息也都是针对每一个线程的。这些系统调用应该返回这个进程的总量。
- 一些版本的 `ps(1)`会显示进程中的所有线程（包括管理线程），作为单独的项，且它们的进程号也不相同。
- 线程之间并不共享由 `setpriority()`设置的 `nice` 值。

- 使用 `setitimer()` 创建的间隔定时器也无法在线程之间共享。
- 线程之间不能共享 System V 信号量还原 (`semadj`) 值。

LinuxThreads 的其他问题

除了以上与 SUSv3 标准的偏差外，LinuxThreads 实现还有如下问题。

- 如果管理线程被杀掉，那么余下的线程只能手工清理。
- 多线程程序的核心转储 (`core dump`) 可能并不包含所有的线程（甚至可能也不包含触发转储的线程）。
- 只有从主线程调用非标准的 `ioctl()` `TIOCNOTTY` 操作才能移除进程与控制终端的关联。

33.5.2 NPTL

设计 NPTL 是为了弥补 LinuxThreads 的大部分的缺陷。特别是如下部分。

- NPTL 更接近 SUSv3 Pthreads 标准。
- 使用 NPTL 的有大量线程的应用程序的性能要远优于 LinuxThreads。

NPTL 允许应用程序创建大量的线程。NPTL 实现的测试程序可以创建 10 万个线程。对于 LinuxThreads，实际线程数量的限制大约是一两千个。（应当承认，很少有程序需要创建这个多的线程。）

NPTL 实现的开发从 2002 年开始，大约在第 2 年完成。同时，Linux 内核为适应 NPTL 也做了各种调整。这些变动出现在 Linux 2.6 内核，并对 NPTL 的如下方面提供支持。

- 改进线程组的实现（28.2.1 节）。
- 增加 `futex` 作为一种同步机制（`futex` 作为一种通用机制，并不只是为 NPTL 而设计）。
- 增加新的系统调用（`get_thread_area()` 和 `set_thread_area()`）以便支持线程本地存储。
- 支持线程化的核心转储和对多线程程序的调试功能。
- 修改并支持与 Pthreads 模型一样的信号处理。
- 增加新的系统调用 `exit_group()`，可以终止进程中的所有线程（从 `glibc2.3` 开始，库函数 `exit()` 是 `exit_group()` 的包装函数，而函数 `pthread_exit()` 调用真正的内核系统调用 `_exit()`，仅终止调用的线程）。
- 重写内核调度程序以便能够有效地调度和处理大量（上千个）KSE 的情况。
- 提升内核进程终止的执行效率。
- 扩展系统调用 `clone()`（28.2 节）。

NPTL 实现最基本的部分如下：

- 线程使用函数 `clone()` 创建并指定如下标志。

```
CLONE_VM | CLONE_FILES | CLONE_FS | CLONE_SIGHAND |  
CLONE_THREAD | CLONE_SETTLS | CLONE_PARENT_SETTID |  
CLONE_CHILD_CLEARTID | CLONE_SYSVSEM
```

NPTL 比 LinuxThreads 可以共享更多的信息。标志 `CLONE_THREAD` 意思是新线程与创建者线程属于同一个线程组，并且共享同样的进程号以及父进程号。`CLONE_SYSVSEM` 表示新线程与创建者共享 System V 信号量还原值。

使用 `ps(1)` 列出一个运行在 NPTL 下的多线程程序时，只会输出一条记录。为了看到进程中的线程信息，可以使用 `ps-L` 选项。

- 实现的内部使用前两个实时信号。应用程序不能使用这些信号。

其中一个信号用来实现线程的取消功能。另一个信号用于确保进程中的所有线程拥有同样的用户号和用户组号。而在内核模式，线程有不同的用户和组凭证。所以 NPTL 实现对每一个改变用户号和用户组号的系统调用 (`setuid()`、`setresuid()` 等) 以及类似的组操作函数的包装函数都做了修改，以确保进程中的所有线程都做了相应的改变。

- 与 LinuxThreads 不同，NPTL 并不需要管理线程。

NPTL 标准一致性

这些改变意味着 NPTL 比 LinuxThreads 更接近 SUSv3 标准。在作者撰写本书的时候，遗留以下不一致的地方。

- 线程之间不共享 nice 值。
- 在早期的 2.6.x 内核中，还有一些额外的不一致的地方。
- 内核版本 2.6.16 之前，备选信号栈是针对每个线程的，但是新的线程从调用 `pthread_create()` 函数的线程那里错误地继承了备选信号栈设置（通过 `sigaltstack()` 产生），导致出现两个线程共享同一备选信号栈的问题。
 - 内核 2.6.16 之前，只有一个线程组的组长（即主线程）可以通过调用函数 `setsid()` 启动一个新的会话。
 - 内核 2.6.16 之前，只有一个线程组的组长可以使用函数 `setpgid()` 让宿主进程成为进程组主进程。
 - 早于 2.6.12 的内核版本，在同一进程的线程之间无法共享使用 `setitimer()` 创建的间隔定时器。
 - 早于 2.6.10 的内核版本，同一进程中的所有线程并不共享资源限制的设置。
 - 早于 2.6.9 的内核版本，函数 `times()` 返回的 CPU 时间以及函数 `getrusage()` 返回的资源使用信息都是针对每个线程的。

NPTL 设计与 LinuxThreads ABI 兼容。那些与提供 LinuxThreads 的 GNU C 库链接的程序换用 NPTL 时无需再重新编译。不过当程序运行在 NPTL 环境时某些行为可能会有些不同，主要是因为 NPTL 更接近于 SUSv3 Pthreads 标准。

33.5.3 哪一种线程实现

一些 Linux 发布版本附带包换 LinuxThreads 和 NPTL 的 GNU C 库，依据系统运行在何种内核上动态地确定链接哪一种 GNU C 库。（这些发布版本有其历史原因，自版本 2.4 后 glibc 不再提供 LinuxThreads。）所以有时候可能需要回答以下的问题。

- 特定的 Linux 发布版本中，哪一种线程实现是有效的？
- 在既提供 LinuxThreads 也提供 NPTL 的 Linux 发布版本中，缺省使用哪一种？如何明确地选择一个应用程序所使用的线程库？

找出线程实现

可以通过一些技术去找出某个特定系统使用的线程实现，也可以发现在提供两种线程实现的系统上运行的程序默认使用的实现版本。

在提供 glibc 2.3.2 或后续版本的系统上，可以使用如下命令找出系统提供的线程实现，如

果提供两种实现的，则显示默认的那个：

```
$ getconf GNU_LIBPTHREAD_VERSION
```

在只有 NPTL 或者将其作为默认实现的系统上，将会显示类似下面的信息：

```
$ ldd /bin/ls | grep libc.so
libc.so.6 => /lib/tls/libc.so.6 (0x40050000)
```

NPTL 2.3.4

自 glibc 2.3.2 以来，程序可以通过 `confstr(3)` 获得类似的信息，并取得 glibc 特定的配置变量 `_CS_GNU_LIBPTHREAD_VERSION` 的值。

使用老版本 Glibc 库的系统上，需要做一些额外的动作。首先，下面的命令可以被用来显示程序运行时使用的 Glibc 库的路径（这里使用标准程序 `ls` 作为例子，其位置为 `/bin/ls`）：

```
$ /lib/tls/libc.so.6 | egrep -i 'threads|nptl'
Native POSIX Threads Library by Ulrich Drepper et al
```

GNU C 库的路径显示在 `=>` 之后。如果作为命令执行这个路径的程序，那么 glibc 将会显示关于自身的一系列信息。可以通过 `grep` 选取与线程实现相关的信息。

在 `egrep` 正则表达式（regular expression）中包含 `nptl`，是因为某些包含 NPTL 的 glibc 发布显示如下的字符串信息：

```
NPTL 0.61 by Ulrich Drepper
```

因为 glibc 路径会随着不同的 Linux 发布而改变，可以使用 shell 的替换功能来产生一个显示 Linux 系统上使用的线程实现信息的命令行：

```
$ $(ldd /bin/ls | grep libc.so | awk '{print $3}') | egrep -i 'threads|nptl'
Native POSIX Threads Library by Ulrich Drepper et al
```

选择程序使用的线程实现

在即提供 NPTL 也提供 LinuxThreads 的 Linux 系统上，能够明确地控制具体使用的线程实现有时候非常地有用。最常见的例子是，当遇到一个旧有的依赖于某些 LinuxThreads 行为（可能非标准）的程序时，要能够强制程序使用指定的线程实现，而不是默认的 NPTL。

出于这个目的，可以使用一个动态链接器（dynamic linker）能够理解的特定的环境变量 `LD_ASSUME_KERNEL`。顾名思义，这个环境变量告诉动态链接器就好像运行在特定的 Linux 内核版本上一样。通过指定并不提供 NPTL 支持的内核版本（例如 2.2.5）可以确保 LinuxThreads 被使用到。所以可以使用如下命令运行一个基于 LinuxThreads 的多线程应用程序：

```
$ LD_ASSUME_KERNEL=2.2.5 ./prog
```

当环境变量设置与之前提及的显示使用的线程实现信息的命令行一起使用时，可以看到如下的一些信息：

```
$ export LD_ASSUME_KERNEL=2.2.5
$ $(ldd /bin/ls | grep libc.so | awk '{print $3}') | egrep -i 'threads|nptl'
linuxthreads-0.10 by Xavier Leroy
```

可以通过 `LD_ASSUME_KERNEL` 设置的内核版本号的范围受一些限制因素的制约。在一些提供 NPTL 和 LinuxThreads 的一般发布版本中，将版本号指定为 2.2.5 已经足够保证会使用 LinuxThreads。此环境变量更完整的描述请参考 <http://people.redhat.com/drepper/assumkernel.html>。

33.6 Pthread API 的高级特性

Pthreads API 还包括一些如下的高级特性。

- 实时调度 (Realtime scheduling): 可以对线程设置实时调度策略以及优先级。类似于 35.3 节中描述的进程的实时调度的系统调用。
- 进程共享互斥量和条件变量: SUSv3 规定进程之间共享互斥量和条件变量是可选的 (不只是针对进程中的线程而言)。这种情况, 条件变量或者互斥量必须在进程间的共享内存中分配。NPTL 支持这种特性。
- 高级线程同步原语: 这些功能包括障碍 (barrier)、读写锁 (read-write lock) 以及自旋锁 (spin lock)。

关于这些特性的更多的细节请参考[Butenhof, 1996]。

33.7 总结

不要将线程与信号混合使用, 只要可能多线程应用程序的设计应该避免使用信号。如果多线程应用必须处理异步信号的话, 通常最简洁的方法是所有的线程都阻塞信号, 创建一个专门的线程调用 `sigwait()` 函数 (或者类似的函数) 来接收收到的信号。这个线程就可以安全地执行像修改共享内存 (处于互斥量的保护之下) 和调用非异步信号安全的函数。

一般有两种有效的 Linux 线程实现: LinuxThreads 和 NPTL。LinuxThreads 多年以来一直为 Linux 使用, 但是很多方面并不遵循 SUSv3 的标准, 而且已经过时。全新的 NPTL 实现更接近 SUSv3 标准并且提供更优的性能, 现今的 Linux 发布也都提供这种实现。

更多的信息

请参考列于 29.10 节的更多的信息来源。

LinuxThreads 的作者编写了实现文档, 可以在如下地址找到: <http://pauillac.inria.fr/~xleroy/linuxthreads/>。NPTL 实现在如下的论文 (有些过时) 中有所描述: <http://people.redhat.com/drepper/nptl-design.pdf>。

33.8 练习

- 33-1. 编写程序以便证明: 作为函数 `sigpending()` 的返回值, 同一个进程中的的不同线程可以拥有不同的 pending 信号。可以使用函数 `pthread_kill()` 分别发送不同的信号给阻塞这些信号的两个不同的线程, 接着调用 `sigpending()` 方法并显示这些 pending 信号的信息。(可能会发现程序清单 20-4 中函数的作用。)
- 33-2. 假设一个线程使用 `fork()` 创建了一个子进程。当子进程终止时, 可以保证由此产生的 SIGCHLD 信号一定会发送给调用 `fork()` 的线程吗 (可以用进程中的其他线程做对比)?

Linux/UNIX

系统编程手册 (上册)

本书是 Linux 和 UNIX 系统编程接口方面的权威指南, 该接口几乎被 Linux 或 UNIX 系统上运行的所有应用所采用。

在这本权威巨著中, Linux 编程专家 Michael Kerrisk 针对掌握系统编程技艺所需的系统调用和库函数进行了巨细靡遗的描述, 并用清晰、完整的程序示例来补充完善其讲解。

本书囊括了 500 多个系统调用和库函数, 外加 200 多个程序实例、88 张表格和 115 幅图片。你将学到如何:

- 高效读写文件;
- 使用信号、时钟和定时器;
- 创建进程并执行程序;
- 编写安全的程序;
- 使用 POSIX 线程编写多线程程序;
- 构建和使用共享库;
- 使用管道、消息队列、共享内存以及信号量进行进程间通信;
- 使用套接字 API 编写网络应用程序。

本书在涵盖大量 Linux 专有特性 (比如, epoll、inotify、/proc 文件系统) 的同时, 对 UNIX 标准 (POSIX.1-2001/SUSv3 以及 POSIX.1-2008/SUSv4) 也极为重视, 这使得本书对于在其他 UNIX 平台下工作的程序员也同样极具价值。

本书是 Linux 和 UNIX 编程接口方面覆盖最为广泛全面的著作, 注定将成为新的经典。



Michael Kerrisk (<http://man7.org>) 具有 20 多年的 UNIX 系统使用和编程经验, 所开设的 UNIX 系统编程周训课程更是不计其数。自 2004 年起, 他开始维护手册页项目, 该项目旨在生成描述 Linux 内核以及 glibc 编程 API 的手册页。他已经撰写或与他人合著了 250 多篇手册页, 至今仍积极参与对 Linux 内核 / 用户空间接口的测试和设计评审工作。Michael 与家人居住在德国慕尼黑。

美术编辑: 王建国



读者可通过<http://www.man7.org/tlpi>下载本书中的所有源代码, 获知更多信息。



分类建议: 计算机 / 操作系统 / Linux

人民邮电出版社网址: www.ptpress.com.cn

异步社区会员 flyman150(2410757683@qq.com) 专享 尊重版权