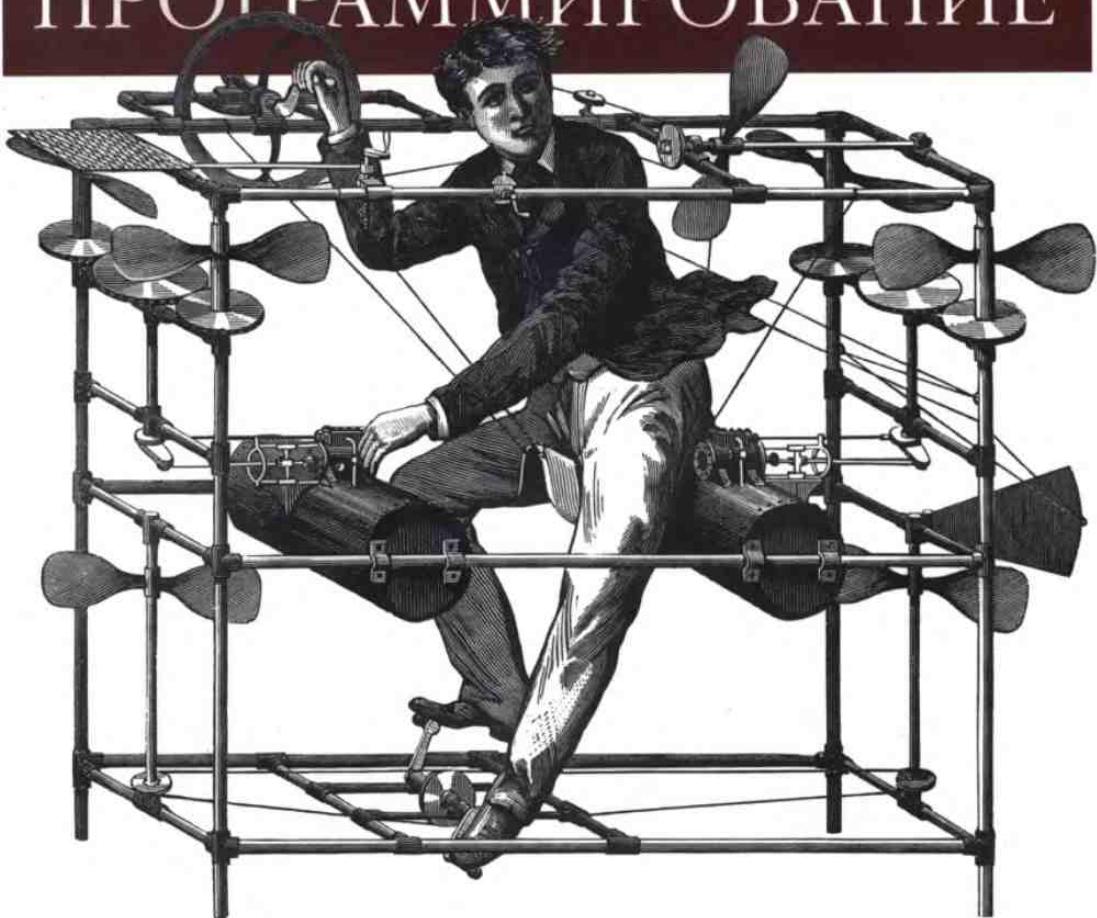


ВСЁ О СИСТЕМНЫХ ВЫЗОВАХ В LINUX

LINUX

СИСТЕМНОЕ ПРОГРАММИРОВАНИЕ



РОБЕРТ ЛАВ

O'REILLY®

 ПИТЕР®



Robert Love

LINUX

System Programming

O'REILLY®

Beijing • Cambridge • Farnham • Köln • Paris • Sebastopol • Taipei • Tokyo

РОБЕРТ ЛАВ

LINUX

СИСТЕМНОЕ
ПРОГРАММИРОВАНИЕ



Москва · Санкт-Петербург · Нижний Новгород · Воронеж
Ростов-на-Дону · Екатеринбург · Самара · Новосибирск
Киев · Харьков · Минск

2008

ББК 32.973.2-018.2

УДК 004.451

Л13

Лав Р.

Л13 Linux. Системное программирование. — СПб.: Питер, 2008. — 416 с.: ил. — (Серия «Бестселлеры O'Reilly»).

ISBN 978-5-388-00014-9

Эта книга о том, как создавать программное обеспечение под Linux, эффективно используя возможности системы — функции ядра и базовые библиотеки, включая оболочку, текстовый редактор, компилятор, отладчик и системные процедуры. Большая часть программного кода для Unix и Linux написана на системном уровне, поэтому в книге основное внимание сфокусировано на приложениях, находящихся вне ядра, таких как Apache, bash, cp, vim, Emacs, gcc, gdb, glibc, ls, mv и прочих. Книга написана специально для разработчиков и является необходимым инструментом любого программиста.

ББК 32.973.2-018.2

УДК 004.451

Права на издание получены по соглашению с O'Reilly.

Все права защищены. Никакая часть данной книги не может быть воспроизведена в какой бы то ни было форме без письменного разрешения владельцев авторских прав.

Информация, содержащаяся в данной книге, получена из источников, рассматриваемых издательством как надежные. Тем не менее, имея в виду возможные человеческие или технические ошибки, издательство не может гарантировать абсолютную точность и полноту приводимых сведений и не несет ответственности за возможные ошибки, связанные с использованием книги.

ISBN 978-0596009588 (англ.)

ISBN 978-5-388-00014-9

© Robert Love, 2007

© Перевод на русский язык ООО «Питер Пресс», 2008

© Издание на русском языке, оформление ООО «Питер Пресс», 2008

Краткое содержание

Об авторе	14
Предисловие	15
Глава 1. Введение и важнейшие концепции	17
Глава 2. Файловый ввод-вывод	44
Глава 3. Буферизованный ввод-вывод	88
Глава 4. Расширенный файловый ввод-вывод	112
Глава 5. Управление процессами	159
Глава 6. Расширенное управление процессами	200
Глава 7. Управление файлами и каталогами	239
Глава 8. Управление памятью	291
Глава 9. Сигналы	332
Глава 10. Время	364
Приложение. Расширения GCC языка C	397
Библиография	410

Содержание

Об авторе	14
Предисловие	15
Глава 1. Введение и важнейшие концепции	17
Системное программирование	18
Системные вызовы	19
Библиотека С	20
Компилятор С	20
Интерфейсы API и ABI	21
Интерфейсы API	21
Интерфейсы ABI	22
Стандарты	23
История POSIX и SUS	23
Стандарты языка С	24
Linux и стандарты	25
Эта книга и стандарты	25
Концепции программирования в Linux	26
Файлы и файловая система	26
Процессы	34
Пользователи и группы	36
Разрешения	37
Сигналы	38
Взаимодействие процессов	39
Заголовки	39
Обработка ошибок	39
Знакомство с системным программированием	43
Глава 2. Файловый ввод-вывод	44
Открытие файлов	45
Системный вызов open()	45
Владельцы новых файлов	48
Разрешения новых файлов	48
Функция creat()	50
Возвращаемые значения и коды ошибок	51
Чтение файла при помощи системного вызова read()	51
Возвращаемые значения	52
Считывание всех байтов	53
Считывание без блокировки	54
Прочие значения ошибки	54
Лимиты размера для вызова read()	55

Запись при помощи системного вызова write()	55
Частичная запись	56
Режим присоединения	57
Запись без блокировки	57
Прочие коды ошибок	57
Ограничения размера для вызова write()	58
Поведение вызова write()	58
Синхронизированный ввод-вывод	60
fsync() и fdatasync()	60
Флаг O_SYNC	62
Флаги O_DSYNC и O_RSYNC	63
Прямой ввод-вывод	64
Закрытие файлов	64
Значения ошибок	65
Поиск при помощи lseek()	66
Поиск за пределами конца файла	67
Значения ошибок	67
Ограничения	68
Позиционное чтение и запись	68
Значения ошибок	69
Усечение файлов	69
Мультиплексированный ввод-вывод	71
select()	72
Системный вызов poll()	78
Сравнение poll() и select()	82
Внутреннее устройство ядра	83
Виртуальная файловая система	83
Страницочный кэш	84
Страницчная отложенная запись	86
Заключение	87
Глава 3. Буферизованный ввод-вывод	88
Ввод-вывод с пользовательским буфером	88
Размер блока	89
Стандартный ввод-вывод	91
Указатели файла	91
Открытие файлов	92
Режимы	92
Открытие потока данных при помощи дескриптора файла	93
Закрытие потоков данных	94
Закрытие всех потоков данных	94
Чтение из потока данных	94
Чтение по одному символу за раз	94
Считывание строки целиком	95
Считывание двоичных данных	97
Запись в поток данных	98
Запись одного символа	98
Запись строки символов	99
Запись двоичных данных	99

Пример программы с использованием буферизованного ввода-вывода	100
Поиск в потоке данных	101
Получение текущей позиции в потоке данных	103
Сброс потока данных	103
Ошибки и конец файла	104
Получение связанного дескриптора файла	105
Управление буферизацией	105
Безопасность потоков выполнения	107
Блокировка файла вручную	108
Операции над потоками без блокировки	109
Недостатки стандартного ввода-вывода	110
Заключение	111
Глава 4. Расширенный файловый ввод-вывод	112
Разбросанный ввод-вывод	113
Системные вызовы readv() и writev()	113
Интерфейс event poll	118
Создание нового экземпляра epoll	118
Управление epoll	119
Ожидание событий при помощи epoll	122
События, запускаемые фронтом, и события, запускаемые уровнем	123
Отображение файлов в память	124
Системный вызов mmap()	124
Системный вызов munmap()	129
Пример отображения	130
Преимущества системного вызова mmap()	131
Недостатки системного вызова mmap()	132
Изменение размера отображения	133
Изменение защиты отображения	134
Синхронизация файла с отображением	135
Добавление советов относительно отображения	136
Советы относительно обычного файлового ввода-вывода	139
Системный вызов posix_fadvise()	139
Системный вызов readahead()	141
Советы стоят недорого	142
Синхронизированные, синхронные и асинхронные операции	142
Асинхронный ввод-вывод	144
Планировщики ввода-вывода и производительность ввода-вывода	145
Адресация диска	146
Жизнь планировщика ввода-вывода	147
Помощь при считывании	148
Выбор и настройка планировщика ввода-вывода	151
Оптимизация производительности ввода-вывода	152
Заключение	158
Глава 5. Управление процессами	159
Идентификатор процесса	159
Выделение идентификатора процесса	160

Иерархия процессов	161
Тип pid_t	161
Получение идентификатора процесса и идентификатора родительского процесса	161
Запуск процесса	162
Семейство вызовов exec	162
Системный вызов fork()	166
Завершение процесса	170
Прочие способы завершения	171
Библиотечный вызов atexit()	172
Функция on_exit()	173
SIGCHLD	173
Ожидание завершенных дочерних процессов	174
Ожидание определенного процесса	176
Еще больше гибкости при ожидании процессов	179
BSD говорит свое слово: wait3() и wait4()	181
Запуск и ожидание нового процесса	182
Зомби	184
Пользователи и группы	185
Реальные, действительные и сохраненные идентификаторы пользователя и группы	186
Изменение реального и сохраненного идентификатора пользователя или группы	187
Изменение действительного идентификатора пользователя или группы	188
Изменение идентификаторов пользователя и группы, стиль BSD	189
Изменение идентификаторов пользователя и группы, стиль HP-UX	189
Манипулирование предпочтительными идентификаторами пользователя/группы	190
Поддержка сохраненных идентификаторов пользователя	190
Получение идентификаторов пользователя и группы	191
Сеансы и группы процессов	191
Системные вызовы сеансов	193
Системные вызовы для группы процессов	195
Устаревшие функции для группы процессов	196
Демоны	197
Заключение	199
Глава 6. Расширенное управление процессами	200
Планирование процессов	200
Запись с О большим	202
Кванты времени	202
Процессы, ограниченные вводом-выводом, и процессы, ограниченные процессором	203
Вытесняющее планирование	204
Управление потоками	205
Уступка процессора	205
Правомерное использование	206
Уступка, прошлое и настоящее	207

Приоритеты процессов	208
Системный вызов nice()	209
Системные вызовы getpriority() и setpriority()	210
Приоритеты ввода-вывода	211
Привязка процессов к процессорам	212
Системные вызовы sched_getaffinity() и sched_setaffinity()	213
Системы реального времени	216
Жесткие и мягкие системы реального времени	216
Задержка, дрожание и предельные сроки	217
Поддержка процессов реального времени в Linux	218
Политики планирования и приоритеты в Linux	219
Установка параметров планирования	223
Интерфейс sched_rr_get_interval()	227
Предосторожности при использовании процессов реального времени	228
Детерминизм	229
Лимиты ресурсов	232
Лимиты	233
Установка и проверка лимитов	237
Глава 7. Управление файлами и каталогами	239
Файлы и их метаданные	239
Семейство stat	239
Разрешения	243
Владение	244
Расширенные атрибуты	247
Каталоги	256
Текущий рабочий каталог	257
Создание каталогов	263
Удаление каталогов	264
Чтение содержимого каталога	266
Ссылки	269
Жесткие ссылки	270
Символические ссылки	271
Удаление ссылки	273
Копирование и перемещение файлов	275
Копирование	275
Перемещение	275
Узлы устройств	278
Специальные узлы устройств	278
Генератор случайных чисел	279
Внеполосная коммуникация	280
Мониторинг событий файлов	281
Инициализация inotify	282
Стражи	283
События inotify	285
Расширенные параметры стражей	288
Удаление стража inotify	289
Проверка размера очереди событий	289
Разрушение экземпляра inotify	290

Глава 8. Управление памятью	291
Адресное пространство процесса	291
Страницы и подкачка страниц	292
Области памяти	293
Выделение динамической памяти	294
Выделение массивов	296
Изменение размера выделенной памяти	297
Освобождение динамической памяти	299
Выравнивание	301
Управление сегментом данных	305
Анонимные отображения в памяти	306
Создание анонимных отображений в памяти	307
Отображение /dev/zero	309
Расширенное выделение памяти	310
Тонкая настройка при помощи malloc_usable_size() и malloc_trim()	313
Отладка выделения памяти	314
Получение статистики	314
Выделение памяти на основе стека	315
Дублирование строк в стеке	317
Массивы переменной длины	317
Выбор механизма выделения памяти	319
Манипулирование памятью	320
Установка байтов	320
Сравнение байтов	321
Перемещение байтов	322
Поиск байтов	323
Прочее манипулирование байтами	324
Блокировка памяти	324
Блокировка части адресного пространства	325
Блокировка всего адресного пространства	326
Разблокировка памяти	327
Лимиты блокирования	328
Находится ли страница в физической памяти?	328
Уступающее выделение памяти	329
Принятие чрезмерных обязательств и ООМ	329
Глава 9. Сигналы	332
Концепции сигналов	333
Идентификаторы сигналов	334
Сигналы, поддерживаемые в Linux	334
Простейшее управление сигналами	340
Ожидание сигнала, любого сигнала	341
Примеры	342
Запуск и наследование	344
Сопоставление номеров сигналов и строк	345
Отправка сигнала	345
Разрешения	346
Примеры	347

Отправка сигнала себе	347
Отправка сигнала всей группе процессов	347
Повторный вход	348
Функции, гарантированно поддерживающие повторный вход	349
Наборы сигналов	351
Больше функций для наборов сигналов	351
Блокирование сигналов	352
Извлечение ожидающих сигналов	353
Ожидание набора сигналов	353
Расширенное управление сигналами	354
Структура siginfo_t	356
Великолепный мир si_code	358
Отправка сигнала с полезной нагрузкой	361
Пример	362
Заключение	363
Глава 10. Время	364
Структуры данных для работы со временем	367
Исходное представление	367
А теперь точность до микросекунд	367
Еще лучше: точность до наносекунд	368
Разбор времени	368
Тип для времени процесса	370
Часы POSIX	370
Разрешение источника времени	370
Получение текущего времени дня	372
Лучший интерфейс	372
Расширенный интерфейс	373
Получение времени процесса	374
Установка текущего времени дня	374
Установка времени с точностью	375
Расширенный интерфейс для установки времени	376
Различные действия со временем	376
Подстройка системных часов	378
Засыпание и ожидание	382
Засыпание с точностью до микросекунд	382
Засыпание с точностью до наносекунд	383
Расширенный подход к засыпанию	385
Переносимый способ засыпания	387
Превышение пределов	387
Альтернативы засыпанию	388
Таймеры	388
Простые сигнализации	388
Интервальные таймеры	389
Расширенные таймеры	391
Приложение. Расширения GCC языка С	397
GNU C	397
Подставляемые функции	398

Подавление подстановки	399
Чистые функции	399
Постоянные функции	400
Функции, не возвращающие результат	400
Функции, выделяющие память	400
Как заставить вызывающего проверять возвращаемое значение	401
Как пометить функцию как устаревшую	401
Как пометить функцию как используемую	401
Как пометить функции или параметры как неиспользующиеся	402
Упаковка структуры	402
Увеличение границы выравнивания переменной	402
Помещение глобальных переменных в регистр	403
Комментирование ветвей	404
Получение типа выражения	405
Получение границы выравнивания типа	405
Смещение члена структуры	406
Получение адреса возврата функции	407
Диапазоны оператора case	407
Арифметика указателей типа void и указателей на функции	408
Более переносимо и красиво	408
Библиография	410
Книги по языку программирования C	410
Книги по программированию в Linux	411
Книги, посвященные ядру Linux	412
Книги по дизайну операционной системы	413

Об авторе

Роберт Лав (Robert Love) — пользователь Linux и хакер с юных лет. Он питает страсть к сообществам, интересующимся ядром Linux и рабочими столами GNOME, и принимает активное участие в их жизни. Его последним вкладом в ядро Linux стала работа над событиями уровня ядра и программой inotify. А вклад в GNOME включает Beagle, GNOME Volume Manager, NetworkManager и Project Utopia. В настоящее время Роберт работает в отделе Open Source Program Office в Google.

Как автор Роберт несет ответственность за книгу «Linux Kernel Development» (издательство Novell Press), теперь уже второе издание. Также он соавтор пятого издания «Linux in a Nutshell» издательства O'Reilly. Будучи редактором журнала Linux Journal, Роберт написал множество статей и прочитал огромное количество лекций о Linux по всему миру.

Роберт окончил университет Флориды, получив степень бакалавра в области математики и вычислительной техники. Родом он из южной Флориды, а сегодня называет своим домом Бостон.

Предисловие

Когда разработчики ядра Linux в плохом настроении и им хочется побрюзгать, они недовольно бросают: «Пользовательское пространство — это всего лишь испытательная нагрузка для ядра».

Бормоча эти слова, разработчики пытаются снять с себя всю ответственность за любые ошибки, возникающие при выполнении пользовательского кода. Они считают, что разработчики из пользовательского пространства должны отойти в сторону и заниматься собственным кодом, потому что проблемы пользовательского кода — это, определенно, не проблемы ядра.

Для того чтобы доказать, что обычно ошибки действительно заключаются не в ядре, один из ведущих разработчиков ядра Linux уже три года читает лекцию «Почему пользовательское пространство — это сборище неудачников» в заполненных конференц-залах, приводя реальные примеры отвратительного пользовательского кода, на который каждый полагается в своей ежедневной работе. Прочие разработчики ядра создали инструменты, демонстрирующие, как сильно программы из пользовательского пространства обижают аппаратное обеспечение и истощают аккумуляторы ничего не подозревающих ноутбуков.

Но хотя код из пользовательского пространства может быть всего лишь «испытательной нагрузкой» для разработчиков ядра, над которым они насмехаются, выясняется, что все эти разработчики также зависят в своей каждодневной работе от пользовательского кода. Если бы пользовательских программ не было, ядро годилось бы только для того, чтобы выводить на экран шаблон из чередующихся букв АВАВАВ.

Сегодня Linux — это самая гибкая и мощная операционная система из всех когда-либо созданных, работающая везде, начиная от крошечных мобильных телефонов и встроенных устройств и заканчивая 70 % суперкомпьютеров из списка 500 лучших в мире. Никакая другая операционная система никогда не умела настолько хорошо масштабироваться и удовлетворять требованиям разнообразных типов аппаратного обеспечения и сред.

И вместе с ядром код, выполняющийся в пользовательском пространстве Linux, также функционирует на всех этих платформах, обеспечивая мир реальными приложениями и предоставляя людям надежные утилиты.

В этой книге Роберт Лав взялся за серьезную задачу рассказать читателю практически обо всех системных вызовах в системе Linux. Он произвел на свет сочинение, которое позволяет вам четко понять, как работает ядро Linux с точки зрения пользовательского пространства, и укротить мощь этой системы.

Из этой книги вы узнаете, как создавать код, который будет работать на всех существующих дистрибутивах Linux и на всех типах аппаратного обеспечения. Вы поймете, как функционирует Linux и как пользоваться преимуществами гибкости этой операционной системы.

В конце концов, эта книга учит вас писать код, который никто и никогда не сможет назвать отвратительным, а это лучшее, чего можно желать.

Greg Kroah-Hartman

1

Введение и важнейшие концепции

Эта книга о *системном программировании*, то есть об искусстве написания системного программного обеспечения. Системное программное обеспечение живет на низком уровне, общаясь напрямую с ядром и системными библиотеками. Системное программное обеспечение включает в себя оболочку и текстовый редактор, компилятор и отладчик, утилиты ядра и системные демоны. Все эти компоненты работают с кодом ядра и библиотеками С. Большая часть остального программного обеспечения (например, высокоуровневые приложения с графическим интерфейсом пользователя) живет на высоком уровне, ныряя на низкий уровень лишь изредка или вообще не касаясь его. Одни программисты каждый день все время проводят за написанием системного программного обеспечения; другие тратят на эту задачу лишь часть своего рабочего дня. Но нет такого программиста, которому во вред бы были знания о системном программировании. Будь это *raison d'être*¹ программиста или просто основание для самоутверждения, системное программирование находится в сердце всего программного обеспечения, которое мы пишем.

В частности, эта книга посвящена системному программированию в Linux. Linux – это современная Unix-подобная система, написанная Линусом Торвальдсом (Linus Torvalds) и сообществом хакеров со всего света. Хотя Linux разделяет цели и идеологию Unix, Linux – это не Unix. Linux следует по собственному пути, отклоняясь где удобно и присоединяясь, когда этого требуют вопросы практичности. В целом, суть системного программирования в Linux также, что и в любой другой Unix-системе. Но если выйти за границы основ, то Linux много делает для того, чтобы отличаться. По сравнению с традиционными Unix-системами Linux изобилует дополнительными системными вызовами, непохожим поведением и новыми возможностями.

¹ Raison d'être – фр.: смысл жизни (дословно: «оправдание для существования»). В англоязычной литературе также используется в качестве идиомы, означающей достаточно вескую причину для реализации какого-нибудь начинания.

Системное программирование

Говоря традиционно, все программирование в Unix – это программирование системного уровня. Исторически системы Unix включали в себя совсем немного высокуровневых абстракций. Даже программирование в среде для разработки, такой как X Window System, в полной мере обнажало корневой API (Application Programming Interface, интерфейс прикладного программирования) системы Unix. Следовательно, можно сказать, что эта книга рассказывает о программировании в Linux в целом. Но обратите внимание, что в этой книге не говорится о среде программирования Linux – вы не найдете руководства по использованию команды `make` на этих страницах. Речь пойдет об API системного программирования на современной машине под управлением Linux.

Системное программирование чаще всего противопоставляется прикладному программированию. Программирование на системном уровне и программирование на прикладном уровне различаются в некоторых аспектах, но в других схожи. Особенность системного программирования в том, что системные программисты должны обладать глубокими знаниями об аппаратном обеспечении и операционной системе, в которой они работают. Конечно же, различия есть и между используемыми библиотеками и совершаемыми вызовами. В зависимости от «уровня» стека, на котором пишется приложение, эти две области могут быть неравнозначными, но, в целом, перейти от прикладного программирования к системному программированию (или наоборот) совсем несложно. Даже когда приложение живет на самом верху стека, далеко от нижних уровней системы, знания, связанные с системным программированием, очень важны. Одни и те же хорошие практические приемы применяются во всех сферах программирования.

В последние годы наблюдается такая тенденция, что прикладное программирование отходит от программирования на системном уровне и превращается в высокоуровневую разработку с использованием либо веб-приложений (таких, как JavaScript и PHP), либо управляемого кода (например, на C# или Java). Такая разработка, однако, не предвещает смерть системного программирования. Действительно, кому-то же нужно писать интерпретатор JavaScript и среду исполнения C#, а это тоже программирование. Помимо этого, разработчики, пишущие на PHP или Java, также могут извлекать выгоду из знания системного программирования, так как понимание базовых элементов позволяет писать лучший код, неважно, на каком уровне он создается.

Несмотря на эту тенденцию в прикладном программировании, большая часть кода для Unix и Linux все так же пишется на системном уровне. Обычно это код на языке C, существующий в основном на базе интерфейсов, предоставляемых библиотекой C и ядром. Эти традиционные инструменты: Apache, bash, cp, Emacs, init, gcc, gdb, glibc, ls, mv, vim и X – совершенно не планируют умирать в скором времени.

Вершиной системного программирования зачастую называют разработку ядра или, по крайней мере, написание драйверов устройств. Но в этой книге, как и в большей части литературы по системному программированию, разра-

ботка ядра не рассматривается. Данная книга фокусируется на программировании системного уровня в пользовательском пространстве, то есть на всем, что лежит выше ядра (хотя знание внутреннего устройства ядра может стать полезным дополнением). Аналогичным образом, сетевое программирование — сокеты и подобные вещи — также не обсуждается в книге. Написание драйверов устройств и сетевое программирование — это огромные и постоянно расширяющиеся темы, которые рассматриваются в узкоспециализированных книгах.

Что такое интерфейс системного уровня, и как писать приложения системного уровня в Linux? Что именно предоставляют ядро и библиотека C? Как писать оптимальный код, и какие трюки скрываются в Linux? Какие искусные системные вызовы есть в Linux, которых нет в других вариантах Unix? Как все это работает? Вот на какие вопросы пытается ответить эта книга.

Существует три краеугольных камня, на которых держится системное программирование в Linux: *системные вызовы, библиотека C и компилятор C*. Каждый из них заслуживает официального представления.

Системные вызовы

Системное программирование начинается с системных вызовов. Системные вызовы — это обращения к функциям, которые делаются из пользовательского пространства — текстового редактора, вашей любимой игры и т. д. — в ядро (главный внутренний орган системы) для того, чтобы запросить у операционной системы какую-либо службу или ресурс. Системные вызовы могут быть как распространенными и часто используемыми, такими как `read()` и `write()`, так и экзотическими, как `get_thread_area()` или `set_tid_address()`.

В Linux реализовано намного меньше системных вызовов, чем в ядрах большинства других операционных систем. Например, число системных вызовов в архитектуре i386 приближается к 300, а в Microsoft Windows, по слухам, их тысячи. В ядре Linux для каждой из архитектур (таких, как Alpha, i386 или PowerPC) реализован собственный список доступных системных вызовов. Следовательно, системные вызовы, имеющиеся в одной архитектуре, могут отличаться от списка системных вызовов в другой архитектуре. Тем не менее очень большой поднабор системных вызовов — более 90 % — реализован одинаково во всех архитектурах. В данной книге я рассматриваю именно этот общий поднабор — наиболее распространенные интерфейсы.

Выполнение системных вызовов

Невозможно напрямую связать приложения пользовательского пространства с пространством ядра. По причинам, относящимся к безопасности и надежности, приложения пользовательского пространства не могут непосредственно выполнять код ядра или манипулировать данными ядра. Вместо этого предоставляется механизм, при помощи которого приложение пользовательского пространства может «сигнализировать» ядру, что ему необходимо выполнить системный вызов. Приложение может отправить прерывание ядру и исполнить только тот код, к которому ему разрешает обращаться ядро. Конкретная реализация механизма варьируется в зависимости от архитектуры. В архитектуре i386, на-

пример, приложение пользовательского пространства исполняет инструкцию программного прерывания `int 0x80`. Эта инструкция порождает переход в пространство ядра — защищенную область ядра, где ядро исполняет обработчик программного прерывания. А что такое обработчик для прерывания `0x80`? Не что иное, как обработчик системного вызова!

Приложение сообщает ядру, какой системный вызов нужно выполнить и с какими параметрами, используя *аппаратные регистры* (machine register). Системные вызовы обозначаются номерами начиная с нуля (0). В архитектуре i386, для того чтобы запросить системный вызов 5 (то есть вызов `open()`), приложение пользовательского пространства помещает значение 5 в регистр `eax` и только после этого исполняет инструкцию `int`.

Передача параметров обрабатывается схожим образом. В архитектуре i386, например, для каждого из возможных параметров используется отдельный регистр: регистры `ebx`, `ecx`, `edx`, `esi` и `edi` содержат, в указанном порядке, первые пять параметров. В редких случаях, когда системный вызов включает в себя более пяти параметров, используется один регистр, который указывает на буфер в пользовательском пространстве, где хранятся все параметры вызова. Конечно же, большинство системных вызовов выполняются всего лишь с парой параметров.

В других архитектурах выполнение системных вызовов обрабатывается по-другому, хотя общий дух сохраняется. Как системному программисту, обычно вам не нужно знать, каким способом ядро обрабатывает выполнение системных вызовов. Это знание закодировано в стандартных соглашениях о вызовах для каждой конкретной архитектуры и автоматически обрабатывается компилятором и библиотекой C.

Библиотека C

Библиотека C (`libc`) находится в самом сердце приложений Unix. Даже когда вы программируете на другом языке, библиотека C, вероятнее всего, все равно принимает участие в работе, обернутая более высокоуровневыми библиотеками, и предоставляет корневые службы, упрощая выполнение системных вызовов. В современных системах Linux библиотека C — это GNU `libc`, сокращенно `glibc`, что произносится как *джи-либ-си* (gee-lib-see) или, реже, *глиб-си* (glib-see).

В библиотеке GNU C можно найти намного больше, чем подразумевает ее имя. Помимо реализации стандартной библиотеки C, `glibc` предоставляет обертки для системных вызовов, поддержку многопоточной обработки и базовые возможности приложений.

Компилятор C

В Linux стандартный компилятор C предоставляется в форме коллекции компиляторов GNU (GNU Compiler Collection, `gcc`). Первоначально `gcc` был версией `cc`, компилятора C (C Compiler — `cc`), разработанной в проекте GNU. Таким образом, `gcc` расшифровывается как GNU C Compiler. Однако со временем

была добавлена поддержка многих других языков. Поэтому сегодня `gcc` используется в качестве общего имени для семейства компиляторов GNU. Однако `gcc` – это также двоичный файл, при помощи которого вызывается компилятор C. В этой книге, говоря о `gcc`, обычно я имею в виду программу `gcc`, если контекст не подразумевает иное.

Компилятор, используемый в системах Unix, в том числе в Linux, тесно связан с системным программированием, так как компилятор помогает реализовывать стандарт C (см. раздел «Стандарты языка C») и системный интерфейс ABI (см. раздел «Интерфейсы API и ABI»), о которых рассказывается далее в этой главе.

Интерфейсы API и ABI

Программисты, естественно, заинтересованы в том, чтобы их программы выполнялись во всех системах, поддержка которых декларируется, и не только сегодня, но и в будущем. Им необходима уверенность в том, что программы, которые они пишут для своих дистрибутивов Linux, будут работать и в других дистрибутивах, а также в прочих поддерживаемых архитектурах Linux и в более новых (или старых) версиях Linux.

На системном уровне существует два отдельных набора определений и описаний, влияющих на переносимость программ. Первый из них – это *интерфейс прикладного программирования* (Application Programming Interface, API), а второй – *двоичный интерфейс приложений* (Application Binary Interface, ABI). Оба определяют и описывают интерфейсы между различными частями компьютерного программного обеспечения.

Интерфейсы API

В интерфейсе API определяются способы, при помощи которых один фрагмент программного обеспечения общается с другим на уровне исходных текстов. Он предоставляет абстракцию в виде стандартного набора интерфейсов – обычно функций, которые какая-то часть программного обеспечения (обычно, хотя и не всегда, высокоуровневая) может вызывать из другой части программного обеспечения (обычно низкоуровневой). Например, API может абстрагировать концепцию вывода текста на экран при помощи семейства функций, обеспечивающих все необходимое для отображения текста. В API просто определяется интерфейс. Фрагмент программного обеспечения, который фактически представляет интерфейс API, называется *реализацией API* (*implementation*).

Очень часто API называют «контрактом». Это не совсем верно, по крайней мере в юридическом смысле термина, так как API – это не двухстороннее соглашение. Пользователь API (обычно высокоуровневое программное обеспечение) не может ничего внести в API и его реализацию. Он может использовать API в том виде, в котором интерфейс существует, или же вообще не использовать его – третьего не дано! API всего лишь гарантирует, что если обе части программного обеспечения будут удовлетворять требованиям API, то они будут

совместимы на уровне исходного кода (*source compatible*). Это означает, что приложение-пользователь API будет успешно компилироваться с данной реализацией API.

В качестве примера из реального мира можно привести API, определенный стандартом С и реализованный стандартной библиотекой С. Этот API определяет семейство базовых и обязательных функций, таких, как программы для манипулирования строками.

В этой книге мы полагаемся на существование разнообразных API, таких, как библиотека стандартного ввода-вывода, о которой речь пойдет в главе 3. Самые важные API для системного программирования в Linux обсуждаются в разделе «Стандарты» далее в этой главе.

Интерфейсы ABI

Тогда как API определяет исходный интерфейс, ABI определяет низкоуровневый двоичный интерфейс между двумя или несколькими частями программного обеспечения в конкретной архитектуре. Интерфейс ABI формулирует, как приложение взаимодействует с самим собой, как приложение взаимодействует с ядром и как приложение взаимодействует с библиотекой. ABI обеспечивает *совместимость на двоичном уровне* (*binary compatibility*), то есть гарантирует, что фрагмент конечной программы будет функционировать в любой системе, где имеется тот же ABI, не требуя перекомпиляции.

Интерфейсы ABI связаны с такими вещами, как соглашения о вызовах, порядок следования байтов, использование регистров, выполнение системных вызовов, компоновка, поведение библиотек и формат двоичных объектов. Например, соглашения о вызовах определяют способы вызова функций, способ передачи аргументов, какие регистры сохраняются, а какие искажаются и какзывающий код получает возвращаемое значение.

Хотя было предпринято несколько попыток определить единый ABI для одной архитектуры, охватывающей несколько операционных систем (в частности, для i386 в системах Unix), усилия не увенчались особым успехом. Наоборот, для всех операционных систем, в том числе Linux, обычно определяются собственные интерфейсы ABI. Любой ABI тесно привязан к архитектуре; большую часть ABI составляют машинно-ориентированные концепции, такие, как назначение определенных регистров или инструкции сборки. Таким образом, для каждой машинной архитектуры существует собственный интерфейс ABI для Linux. В действительности, мы обычно называем отдельные ABI по их машинным именам, например alpha или x86-64.

Системный программист должен быть знаком с ABI, но обычно необходимости запоминать подробности не возникает. ABI реализуется *цепочкой инструментов* (*toolchain*) — компилятором, компоновщиком и т. д. — и обычно не выходит на поверхность никакими другими способами. Знание ABI, однако, может помочь научиться создавать оптимальные программы и необходимо для написания кода сборки или для того, чтобы влезть в саму цепочку инструментов (что также относится к системному программированию).

Описание ABI для любой конкретной архитектуры в Linux можно найти в Интернете; этот интерфейс реализуется цепочкой инструментов соответствующей архитектуры и ядром.

Стандарты

Системное программирование в Unix – это древнее искусство. Основы программирования в Unix остаются неприкословенными вот уже десятилетия. Системы Unix, однако, весьма динамичны, и их поведение меняется с добавлением новых функций. Для того чтобы внести свой вклад в процесс превращения порядка в хаос, группы стандартов объединяются в официальные стандарты, включающие в себя варианты системных интерфейсов. Существует множество подобных стандартов, но, технически говоря, Linux официально не удовлетворяет ни одному из них. Однако система Linux *нацелена* на соответствие двум наиболее важным и распространенным стандартам: POSIX и Single UNIX Specification (SUS – единая спецификация Unix).

Стандарты POSIX и SUS документируют, среди прочего, API языка C для интерфейса Unix-подобной операционной системы. Фактически, они определяют системное программирование или, по крайней мере, самую общеупотребительную его часть для совместимых систем Unix.

История POSIX и SUS

В середине 80-х годов XX в. Институт инженеров по электротехнике и радиоэлектронике (Institute of Electrical and Electronics Engineers, IEEE) возглавил усилия по стандартизации интерфейсов системного уровня в системах Unix. Ричард Столлман (Richard Stallman), основатель движения Free Software (Свободное программное обеспечение), предложил для стандарта название POSIX (произносится как «пазикс» (pahz-icks)), которое сегодня расшифровывается как Portable Operating System Interface (Интерфейс переносимых операционных систем).

Первым результатом этих усилий, появившимся на свет в 1988 году, стал стандарт IEEE Std 1003.1-1988 (короче – POSIX 1988). В 1990 году IEEE не рассмотрел стандарт POSIX, выпустив стандарт IEEE Std 1003.1-1990 (POSIX 1990). Необязательная поддержка обработки в реальном времени и поточной обработки была документирована в стандартах IEEE Std 1003.1b-1993 (POSIX 1993 или POSIX.1b) и IEEE Std 1003.1c-1995 (POSIX 1995 или POSIX.1c) соответственно. В 2001 году необязательные стандарты были объединены с базовым стандартом POSIX 1990, благодаря чему родился единый стандарт IEEE Std 1003.1-2001 (POSIX 2001). Последняя версия, выпущенная в апреле 2004 года, носит название IEEE Std 1003.1-2004. Все ключевые стандарты POSIX обозначаются аббревиатурой POSIX.1, и последняя версия датирована 2004 годом.

В конце 80-х и начале 90-х годов производители систем Unix боролись друг с другом в так называемых «Юниксовых войнах», в которых каждый стремился

к тому, чтобы его вариант Unix был признан *настоящей* операционной системой Unix. Несколько основных производителей Unix объединились вокруг Open Group — промышленного консорциума, который был сформирован путем объединения Open Software Foundation (OSF) и X/Open. Консорциум Open Group обеспечивал сертификацию, официальные документы и проверку соответствия. В начале 90-х годов, когда «Юниксовы войны» были в самом разгаре, Open Group выпустила единую спецификацию Unix — Single UNIX Specification. SUS быстро завоевала популярность, большей частью благодаря своей цене (она бесплатная), в противоположность высокой стоимости стандарта POSIX. Сегодня SUS включает в себя новейший стандарт POSIX.

Первая версия SUS была опубликована в 1994 году. Системы, соответствующие спецификации SUSv1, маркируются как UNIX 95. Вторая версия SUS вышла в 1997 году, и совместимые системы маркируются как UNIX 98. Третья, и последняя, версия SUS, SUSv3, была опубликована в 2002 году. Совместимые системы маркируются как UNIX 03. В SUSv3 были выполнены пересмотр и объединение стандарта IEEE Std 1003.1–2001 и нескольких других стандартов. В книге я буду указывать, стандартизированы ли системные вызовы и другие интерфейсы по POSIX. Я упоминаю POSIX, а не SUS потому, что POSIX входит в состав SUS.

Стандарты языка C

Прославленная книга Дениса Ричи и Брайана Кернигана (Dennis Ritchie, Brian Kernighan) «The C Programming Language» (издательство Prentice Hall) многие годы после публикации в 1978 году служила неформальной спецификацией языка C. Эта версия C заслужила прозвище K&R C. Язык C уже в те времена быстро отвоевывал у BASIC и других языков звание *lingua franca*¹ микрокомпьютерного программирования. Поэтому, чтобы стандартизировать уже заслуживший популярность язык, в 1983 году Американский национальный институт стандартов (American National Standards Institute, ANSI) сформировал комитет, задачей которого было разработать официальную версию C, включающую возможности и усовершенствования, предложенные разнообразными производителями, а также новый язык C++. Процесс был длинным и трудоемким, но стандарт ANSI C был завершен в 1989 году. В 1990 году Международная организация по стандартизации (International Standardization Organization, ISO) ратифицировала стандарт ISO C90, основанный на ANSI C и включающий несколько изменений.

В 1995 году ISO выпустила обновленную (хотя и редко реализуемую) версию языка C, ISO C95. В 1999 году последовало крупное обновление языка — ISO C99, в ходе которого было добавлено множество новых возможностей, включая подстановку функций, новые типы данных, массивы переменной длины, комментарии в стиле C++ и новые библиотечные функции.

¹ Дословно «язык франков» — идиома, обозначающая язык, используемый как средство межэтнического общения в определенной сфере деятельности. — Примеч. ред.

Linux и стандарты

Как уже говорилось ранее, операционная система Linux нацелена на соответствие стандартам POSIX и SUS. Она предоставляет интерфейсы, документированные в SUSv3 и POSIX.1, включая необязательную поддержку обработки в реальном времени (POSIX.1b) и необязательную поддержку поточной обработки (POSIX.1c). Еще важнее то, что Linux пытается обеспечить функционирование в соответствии с требованиями POSIX и SUS. В целом, несоблюдение стандартов считается ошибкой. Считается, что Linux соответствует POSIX.1 и SUSv3, но официальная сертификация POSIX или SUS не проводилась (отдельно для каждой версии Linux), поэтому я не могу утверждать, что Linux официально отвечает требованиям POSIX или SUS.

Что касается стандартов языков, то в Linux все хорошо. Компилятор C gcc поддерживает стандарт ISO C99. Помимо этого, gcc предоставляет множество собственных расширений языка C. Эти расширения носят коллективное название GNU C и документированы в приложении к книге.

Linux не может похвастаться выдающейся историей совместимости¹, хотя сегодня в этом смысле дела в Linux обстоят лучше. Интерфейсы, документированные согласно стандартам, такие, как стандартная библиотека C, очевидно, всегда будут оставаться совместимыми на уровне источника. Совместимость на двоичном уровне обеспечивается, как минимум, для старшей версии glibc. И, так как язык C стандартизован, gcc всегда компилирует допустимый код на C правильно, хотя уникальные для gcc расширения могут быть в конечном итоге удалены из новых релизов gcc. Важнее всего, что ядро Linux гарантирует стабильность системных вызовов. Если системный вызов реализован в стабильной версии ядра Linux, то он может существовать в веках.

Среди разнообразных дистрибутивов Linux большую часть систем Linux стандартизирует Linux Standard Base (LSB). LSB – это совместный проект нескольких производителей Linux под покровительством консорциума Linux Foundation (ранее известного как Free Standards Group). Проект LSB расширяет спецификации POSIX и SUS и добавляет несколько собственных стандартов; он нацелен на стандартизацию двоичного кода, чтобы конечные программы могли выполняться на совместимых системах без необходимости вносить изменения. Большинство производителей Linux соблюдают требования LSB в той или иной степени.

Эта книга и стандарты

В этой книге я намеренно стараюсь избегать выражения преданности любым стандартам. Слишком часто встречаются книги о системном программировании в Unix, в которых в деталях рассказывается, как интерфейс ведет себя в одном стандарте по сравнению с другим, реализован ли определенный системный

¹ Опытные пользователи Linux могут помнить переход с a.out на ELF, переход с libc5 на glibc, изменения в gcc и тому подобное. К счастью, эти дни остались в прошлом.

вызов в этой системе и в той, в общем, страницы заполнены пустой болтовней. Эта книга посвящена системному программированию в современной системе Linux, которая обеспечивается новейшими версиями ядра Linux (2.6), компилятором C gcc (4.2) и библиотекой C (2.5).

Поскольку системные интерфейсы в целом определены раз и навсегда — например, разработчики ядра Linux проходят через огромные муки, чтобы никогда не нарушать интерфейсы системных вызовов — и обеспечивают определенный уровень совместимости на уровне источников и на двоичном уровне, этот подход позволяет мне глубже погрузиться в детали системного интерфейса Linux, не ограничивая себя вопросами совместимости со множеством других систем и стандартов Unix. Благодаря такому фокусу с Linux я также могу в этой книге предложить всестороннее рассмотрение самых современных уникальных для Linux интерфейсов, которые будут оставаться актуальными еще очень долго. Эта книга основывается на близком знакомстве с Linux и, в частности, знании реализации и поведения таких компонентов, как gcc и ядро, что дает мне возможность демонстрировать вам взгляд изнутри и делиться, как опытному ветерану, лучшими практическими советами по оптимизации.

Концепции программирования в Linux

Далее приводится краткий обзор служб, предоставляемых системой Linux. Все Unix-системы, включая Linux, предоставляют общий набор абстракций и интерфейсов. Именно эта общность и *определяет Unix*. Такие абстракции, как файл и процесс, интерфейсы для управления конвейерами и сокетами и так далее, представляют собой суть того, чем является Unix.

Я предполагаю, что вы знакомы со средой Linux. Вы должны уметь ориентироваться в оболочке, использовать базовые команды и компилировать простые программы на C. Это *не* обзор операционной системы Linux или ее среды программирования, а, скорее, перечисление того, что формирует основу системного программирования в Linux.

Файлы и файловая система

Файл — это самая простая и фундаментальная абстракция в Linux. Linux придерживается философии «все есть файл» (хотя и не так строго, как некоторые другие системы, такие, как Plan9¹). Следовательно, большая часть взаимодействия реализуется через считывание и запись файлов, даже когда интересующий вас объект никак нельзя назвать файлом в общепринятом смысле этого слова.

Для того чтобы обратиться к файлу, его сначала нужно открыть. Файл можно открыть для чтения, записи или для чтения и записи одновременно. Обра-

¹ Plan9 — операционная система от Bell Labs, часто называемая преемником Unix. Она базируется на нескольких инновационных идеях и стойко придерживается философии «все есть файл».

щение к открытому файлу осуществляется с использованием уникального дескриптора, обеспечивающего отображение метаданных, связанных с открытым файлом, на сам файл. Внутри ядра Linux этот дескриптор обрабатывается при помощи целочисленного значения (типа `int` в C), называемого *дескриптором файла* (*file descriptor*) или сокращенно `fd`. Файловые дескрипторы совместно используются пользовательским пространством и ядром, а пользовательские программы применяют их, чтобы напрямую обращаться к файлам. Большая доля системного программирования в Linux состоит из открытия, манипулирования, закрытия и использования дескрипторов файлов различными способами.

Обычные файлы

То, что большинство из нас называют «файлами», в Linux носит название *обычного файла* (*regular file*). Обычный файл содержит байты данных, организованные в линейный массив, который называется потоком байтов. В Linux для файла не определяется никакая дополнительная организация или формат. Байты могут иметь любые значения и следовать внутри файла в любом порядке. На системном уровне Linux не требует никакой структуризации от файлов, за исключением потока байтов. Некоторые операционные системы, например VMS, допускают только структурированные файлы, поддерживающие такие концепции, как *записи* (*record*). В Linux это не так.

Любой байт из файла можно считать и в любой байт в файле можно записать значение. Выполнение операций начинается на определенном байте, представляющем умозрительное «местоположение» внутри файла. Это местоположение называется *позицией в файле* (*file position*) или *смещением в файле* (*file offset*). Позиция в файле – это ценный фрагмент метаданных, которые ядро ассоциирует с каждым открытым файлом. Когда файл открывается впервые, значение позиции в файле равно нулю. Обычно, по мере того как файл побайтово считывается или записывается, позиция в файле, так или иначе, увеличивается. Значение позиции в файле также можно установить вручную, причем можно даже выбрать значение, находящееся дальше конца файла. Если записать байт в позицию, которая находится дальше конца файла, то промежуточные байты будут заполнены нулями. Хотя этот способ позволяет записывать байты дальше конца файла, невозможно таким же образом записывать байты в позициях, предшествующих началу файла. Это звучит бессмысленно, и, действительно, от этой практики было бы мало пользы. Значения позиции в файле начинаются с нуля, позиция не может быть отрицательной. При записи в байт в середине файла значение, которое ранее находилось по этому смещению, заменяется новым. Поэтому невозможно увеличить файл, записав что-то внутрь него. Обычно запись в файл начинают в конце файла. Максимальное значение позиции в файле ограничивается только размером типа данных C, который используется для ее хранения, и в современных системах Linux представляет 64-битное целое.

Размер файла измеряется количеством байтов и называется *длиной файла* (*file length*). Другими словами, длина – это просто число байтов в линейном

массиве, составляющем файл. Длину файла можно изменить при помощи операции, называемой *усечением* (truncation). Файл можно усечь до размера, меньшего его первоначальной длины, причем байты удаляются в конце файла. Хотя это звучит странно, учитывая название оператора, файл также можно «усечь» до размера, большего его первоначальной длины. В этом случае новые байты (добавляемые в конце файла) заполняются нулями. Файл может быть пустым (иметь длину, равную нулю) и вообще не содержать значимых байтов. Максимальная длина файла, как и максимальное значение позиции в файле, ограничивается только пределами типов C, которые ядро Linux применяет для управления файлами. Конкретные файловые системы, однако, могут накладывать собственные ограничения, не допуская создания файлов размером больше определенного значения.

Один и тот же файл может быть открыт более одного раза, причем другим или тем же самым процессом. Каждому открытому экземпляру файла присваивается уникальный дескриптор файла; несколько процессов могут совместно использовать один дескриптор. Ядро не накладывает никаких ограничений на одновременный доступ к файлу. Несколько процессов могут свободно считывать и записывать один и тот же файл в одно и то же время. В результате этого при одновременном доступе процессам приходится полагаться на упорядочение отдельных операций, что обычно непредсказуемо. Программам из пользовательского пространства нужно координировать деятельность между собой, чтобы гарантировать, что одновременные обращения к файлу будут в достаточной степени синхронизированы.

Хотя доступ к файлам обычно осуществляется при помощи *имен файлов* (filename), в действительности файлы напрямую не связываются со своими именами. Вместо этого обращение к файлу производится через структуру *inode* (первоначально information node, информационный узел), которой присваивается уникальное числовое значение. Это значение называется *номером inode* (inode number); и чаще всего его название сокращают до *i-number* или *ino*. В inode хранятся связанные с файлом метаданные, такие, как его временная метка, указывающая момент последней модификации, его владелец, тип, длина и местоположение данных файла — но не имя файла! Inode — это и физический объект, который можно обнаружить в Unix-подобных файловых системах, и умозрительная сущность, представляемая структурой данных в ядре Linux.

Каталоги и ссылки

Обращаться к файлу при помощи его номера inode затруднительно (и это потенциальная брешь в безопасности), поэтому из пользовательского пространства файлы всегда открываются по имени, а не по номеру inode. *Каталоги* (directory) используются для определения имен, по которым выполняется обращение к файлам. Каталог играет роль отображения между удобными для человека именами и номерами inode. Пара из имени и номера inode называется *ссылкой* (link). Физическая форма этого отображения на диске — простая таблица, хэш или что-либо еще — реализуется и управляет кодом ядра, который

поддерживает данную файловую систему. По существу, каталог — это обычный файл, но его отличие от обычного файла состоит в том, что он содержит только карту отображения имен в номера inode. Ядро напрямую использует эти отображения для разрешения имен в номера inode.

Когда приложение из пользовательского пространства запрашивает открытие определенного имени файла, ядро открывает каталог, в котором хранится данное имя файла, и ищет указанное имя. Исходя из имени файла, ядро получает номер inode. Inode содержит метаданные, связанные с файлом, включая местоположение данных файла на диске.

Первоначально на диске есть только один каталог, *корневой каталог* (root directory). Этот каталог обычно обозначается путем /. Но, как мы все знаем, чаще всего в системе хранится много каталогов. Так как же ядро узнает, в *какой* каталог нужно заглянуть, чтобы найти определенное имя файла?

Как я уже сказал выше, каталоги очень похожи на обычные файлы. И с ними также связаны структуры inode. Следовательно, ссылки внутри каталогов могут указывать на inode других каталогов. Это означает, что каталоги могут вкладываться друг в друга, формируя иерархию каталогов, что, в свою очередь, позволяет использовать *полные пути* (pathname), с которыми знакомы все пользователи Unix, например: /home/blackbeard/landscaping.txt.

Когда ядро получает запрос на открытие такого полного пути, оно проходит по всем *записям каталога* (directory entry; внутри ядра — *dentry*) в пути, каждый раз находя inode следующей записи. В предыдущем примере ядро начинает просмотр с каталога /, получает inode для каталога home, переходит туда, получает inode для blackbeard, переходит в этот каталог и, наконец, получает inode для landscaping.txt. Эта операция называется *разрешением каталога*, или *разрешением полного пути* (directory, или pathname resolution). Ядро Linux также использует кэш, называемый *кэшем dentry* (dentry cache), для хранения результатов разрешения каталогов, что обеспечивает более быстрый поиск в будущем, с учетом сосредоточенности во времени¹.

Полный путь, начинающийся с корневого каталога, является *полностью уточненным* (fully qualified) и называется *абсолютным полным путем*. Некоторые полные пути не полностью уточнены; они указываются по отношению к какому-то другому каталогу (например, todo/plunder). Такие пути называются *относительными полными путями*. Если ядро предоставляет относительный полный путь, то разрешение его начинается с *текущего рабочего каталога* (current working directory). В текущем рабочем каталоге ядро ищет каталог todo. Попав туда, ядро получает inode для каталога plunder.

Хотя каталоги обрабатываются как обычные файлы, ядро не позволяет открывать их и манипулировать ими как обычными файлами. Работать с каталогами можно, только используя специальный набор системных вызовов. Такие

¹ Сосредоточенность во времени (temporal locality) — это термин, обозначающий высокую вероятность того, что за доступом к определенному ресурсу вскоре последует еще один доступ к этому ресурсу. Сосредоточенность во времени распространяется на многие ресурсы в компьютере.

системные вызовы позволяют добавлять и удалять ссылки – все равно только эти две операции и имеют смысл. Если бы пользовательскому пространству было разрешено манипулировать каталогами без посредничества ядра, то файловая система могла бы быть полностью разрушена из-за единственной ошибки.

Жесткие ссылки

По существу, ничто из вышеописанного не запрещает разрешать несколько имен в один и тот же номер inode. И это действительно разрешено. Когда несколько ссылок отображают разные имена на одну и ту же структуру inode, они называются *жесткими ссылками* (hard link).

Жесткие ссылки позволяют создавать в файловой системе сложные структуры, в которых несколько полных имен указывают на одни и те же данные. Жесткие ссылки могут находиться в одном каталоге или в нескольких разных каталогах. В любом случае ядро просто разрешает полный путь в правильную inode. Например, две жесткие ссылки – `/home/bluebeard/map.txt` и `/home/black-beard/treasure.txt` – могут разрешать определенный номер inode в один и тот же фрагмент данных на диске.

При удалении файла выполняется отсоединение его от структуры каталогов, и делается это всего лишь путем удаления из каталога пары из имени и inode, соответствующей этому файлу. Так как Linux поддерживает жесткие ссылки, файловая система не может разрушать inode и связанные данные при каждой операции отсоединения. Что, если где-то в файловой системе существует другая жесткая ссылка? Для того чтобы гарантировать, что файл не будет разрушен до тех пор, пока *все* ссылки на него не будут удалены, каждая структура inode содержит *счетчик ссылок* (link count), позволяющий отслеживать в системе число ссылок, указывающих на нее. Когда полное имя отсоединяется, счетчик ссылок уменьшается на единицу; только когда он достигает нуля, inode и связанные данные фактически удаляются из файловой системы.

Символические ссылки

Жесткие ссылки не могут охватывать разные файловые системы, потому что номер inode не имеет никакого смысла за пределами собственной файловой системы. В системах Unix также реализованы *символические ссылки* (symbolic link, название которых часто сокращается до symlink), которые могут охватывать разные файловые системы, которые проще сами по себе и менее прозрачны.

Символические ссылки выглядят как обычные файлы. У символической ссылки есть собственная структура inode и какие-то данные, содержащие полный путь к файлу, на который эта ссылка указывает. Это означает, что символические ссылки могут указывать куда угодно, в том числе на файлы и каталоги в других файловых системах, и даже на файлы и каталоги, которых вообще не существует. Символическая ссылка, указывающая на несуществующий файл, называется *сломанной ссылкой* (broken link).

Символические ссылки оказывают большую нагрузку на систему, чем жесткие ссылки, потому что разрешение символической ссылки фактически включает

чает в себя разрешение двух файлов: самой символьической ссылки, а затем файла, на который она указывает. Жесткие ссылки не создают дополнительную нагрузку — нет никакой разницы между доступом к файлу, привязанному к файловой системе один раз, и к файлу, привязанному несколько раз. Нагрузка, вызываемая символьическими ссылками, минимальна, но она тем не менее считается негативной.

Символьические ссылки не так прозрачны, как жесткие ссылки. Использование жестких ссылок абсолютно прозрачно; в действительности, нужно только выяснить, что файл привязан более одного раза. Манипулирование символьическими ссылками, с другой стороны, требует особых системных вызовов. Такое отсутствие прозрачности зачастую считается позитивным фактором, а символьические ссылки выступают в роли скорее ярлыков (shortcut), нежели внутренних ссылок файловой системы.

Специальные файлы

Специальные файлы (special file) — это объекты ядра, которые представляются как файлы. За время своего существования системы Unix поддерживали несколько различных видов специальных файлов. Linux поддерживает четыре:

- файлы блочных устройств;
- файлы устройств посимвольного ввода-вывода;
- именованные конвейеры;
- и доменные сокеты Unix.

Специальные файлы позволяют определенным абстракциям встраиваться в файловую систему, внося вклад в парадигму «все есть файл». В Linux предусмотрен системный вызов для создания специального файла.

Доступ к устройствам в системах Unix осуществляется через файлы устройств, которые ведут себя и выглядят как обычные файлы, находящиеся в файловой системе. Файлы устройств можно открывать, считывать из них данные и записывать данные в них, и именно таким способом пользовательское пространство обращается с устройствами (физическими и виртуальными) в системе и манипулирует ими. Устройства Unix обычно подразделяют на две группы: **блочные устройства** (block device) и **устройства посимвольного ввода-вывода** (character device). У каждого типа устройств — собственный специальный файл устройства.

Доступ к устройству посимвольного ввода-вывода осуществляется как к линейной очереди байтов. Драйвер устройства помещает байты в одну очередь, один за другим, а пользовательское устройство считывает байты в том порядке, в котором они были в очередь поставлены. Пример устройства посимвольного ввода-вывода — клавиатура. Например, если пользователь вводит «рег», то приложение считает с устройства клавиатуры символ р, потом е и, наконец, г. Когда больше символов для считывания не остается, устройство возвращает код конца файла (end-of-file, EOF). Пропуск символов или считывание их

в любом другом порядке не имели бы никакого смысла. Для доступа к таким устройствам вывода используются *файлы устройств посимвольного ввода-вывода*.

К блочному устройству, в противоположность этому, обращаются как к массиву байтов. Драйвер устройства сопоставляет байты на диске с устройством с возможностью поиска, а пользовательское пространство может обращаться к любым значимым байтам в массиве в любом порядке: оно может считать байт 12, затем байт 7, а затем снова байт 12. Блочные устройства – это обычно устройства хранения данных. Жесткие диски, дисководы гибких дисков, приводы компакт-дисков и flash-память – все это примеры блочных устройств. Доступ к ним осуществляется через *файлы блочных устройств*.

Именованные конвейеры (named pipe, зачастую называемые FIFO – first in, first out – первым вошел, первым вышел) – это механизм *взаимодействия процессов* (interprocess communication, IPC), обеспечивающий коммуникационный канал через файловый дескриптор, доступ к которому выполняется через специальный файл. Обычные конвейеры используются для «перекачки» вывода одной программы на вход другой; они создаются в памяти при помощи системного вызова и не присутствуют ни в каких файловых системах. Именованные конвейеры функционируют как обычные конвейеры, но доступ к ним осуществляется через файл, называемый *специальным файлом FIFO* (FIFO special file). Не связанные между собой процессы могут обращаться к этому файлу и общаться между собой.

Последний тип специальных файлов – *сокеты* (socket). Сокеты представляют собой расширенную форму IPC, обеспечивающую коммуникацию между двумя различными процессами, причем не только на одной машине, но и на двух разных машинах. Фактически, сокеты формируют базис сетевого программирования и программирования для Интернета. Существует множество разновидностей сокетов, включая доменные сокеты Unix – эта форма сокетов используется для коммуникации в пределах локальной машины. Сокеты, обменивающиеся данными через Интернет, могут использовать для идентификации цели пару из имени компьютера и порта, а доменные сокеты Unix используют специальный файл, находящийся в файловой системе, обычно называемый просто файлом сокета.

Файловые системы и пространства имен

В Linux, как и во всех системах Unix, предусмотрено глобальное унифицированное *пространство имен* (namespace) файлов и каталогов. В некоторых операционных системах разные задачи и диски выносятся в разные пространства имен. Например, к файлу на диске можно обратиться при помощи полного пути A:\plank.jpg, тогда как путь к жесткому диску – это C:. В Unix к одному и тому же файлу на диске можно обратиться через полный путь /media/floppy/plank.jpg или даже через путь /home/captain/stuff/plank.jpg точно так же, как к файлам с других носителей. Таким образом, пространство имен в Unix унифицировано.

Файловая система (filesystem) – это набор файлов и каталогов в формальной и допустимой иерархии. Файловые системы можно индивидуально добавлять и удалять из глобального пространства имен файлов и каталогов. Эти операции называются *монтажированием* и *размонтированием* (mounting и unmounting). Каждая файловая система монтируется к определенному местоположению в пространстве имен, называемому точкой монтажа. В этой точке монтажа выполняется доступ к корневому каталогу файловой системы. Например, компакт-диск можно монтировать в каталог /media/cdrom, чтобы в этой точке монтажа обращаться к корню файловой системы компакт-диска. Первая подмонтированная файловая система находится в корне пространства имен, /, и называется корневой файловой системой. В системах Linux всегда есть корневая файловая система. Монтаж других файловых систем в других точках монтажа необязательно.

Файловые системы обычно существуют физически (то есть хранятся на диске), хотя Linux также поддерживает *виртуальные файловые системы* (virtual filesystem), существующие только в памяти, и *сетевые файловые системы* (network filesystem), находящиеся на других машинах в сети. Физические файловые системы содержатся на блочных устройствах хранения, таких как компакт-диски, диски, компактные карты памяти (flash-память) или жесткие диски. Некоторые из этих устройств поддерживают *разбиение на разделы*, то есть вы можете поделить устройство на несколько файловых систем, каждой из которых можно будет манипулировать индивидуально. Linux поддерживает широкий диапазон файловых систем и, определенно, любые файловые системы, с которыми может столкнуться средний пользователь, включая файловые системы, зависящие от носителя (например ISO9660), сетевые файловые системы (NFS), неуправляемые файловые системы (ext3), файловые системы из других систем Unix (XFS) и даже файловые системы из систем, не относящихся к семейству Unix (FAT).

Наименьший адресуемый модуль на блочном устройстве – это *сектор* (sector). Сектор относится к физическим свойствам устройства. Размер сектора может быть любой степенью двойки, и довольно распространены секторы объемом 512 байт. Блочное устройство не может передавать модули данных, размер которых меньше сектора, или обращаться к ним; весь ввод-вывод осуществляется в терминах одного или нескольких секторов.

Схожим образом, наименьший логически адресуемый модуль в файловой системе – это *блок* (block). Блок – это абстракция файловой системы, а не физического носителя, на котором хранится файловая система. Обычно размер блока равен размеру сектора, умноженному на степень двойки. Обычно блоки больше секторов, но они должны быть меньше *размера страницы*¹ (page size – наименьший блок, адресуемый *блоком управления памяти* (memory management unit), аппаратным компонентом). Распространены блоки объемом 512 байт, 1 и 4 Кбайт.

¹ Это искусственное, введенное в целях упрощения работы, ограничение ядра, от которого могут отказаться в будущем.

Исторически в системах Unix присутствует только одно общее пространство имен, которое могут просматривать все пользователи и все процессы в системе. В Linux реализован инновационный подход, и эта операционная система поддерживает *пространства имен отдельных процессов* (reg-process namespace), позволяя каждому процессу при необходимости работать с уникальным представлением иерархии файлов и каталогов системы¹. По умолчанию каждый процесс наследует пространство имен своего предка, но процесс может создать и собственное пространство имен с собственным набором точек монтирования и уникальным корневым каталогом.

Процессы

Если файлы – это самая фундаментальная абстракция в системе Unix, то сразу за ними следуют процессы. *Процессы* (process) – это выполняющийся код конечных программ: активные, живые, работающие программы. Но процессы – это больше чем просто конечные программы. Процессы состоят из данных, ресурсов, состояния и виртуализированного процессора.

Процессы начинают свою жизнь как исполняемые конечные программы, то есть машинный код в формате, который понимает ядро (наиболее распространенный в Linux формат – это ELF). Исполняемый формат содержит метаданные и несколько *разделов* (section) кода и данных. Разделы – это линейные фрагменты конечной программы, которые загружаются в линейные фрагменты памяти. Все байты в разделе обрабатываются одинаково, им даются одинаковые разрешения, и, в целом, они используются для схожих целей.

Самые важные и наиболее часто встречающиеся разделы – это *текстовый раздел* (text section), *раздел данных* (data section) и *раздел bss* (bss section). В текстовом разделе содержится исполняемый код и данные, доступные только для чтения, такие, как константы. Обычно этот раздел помечается как доступный только для чтения и исполняемый. В разделе данных хранятся инициализированные данные, такие, как переменные С с уже определенными значениями, и обычно он помечается как доступный для чтения и записи. В разделе bss содержатся неинициализированные глобальные данные. Так как стандарт С требует, чтобы переменные С по умолчанию инициализировались нулями, нет никакой необходимости хранить нули в конечной программе на диске. Вместо этого неинициализированные переменные можно просто перечислить в разделе bss, а при загрузке раздела в память ядро отобразит на него *нулевую страницу* (zero page, страница, содержащая только нули). Раздел bss был задуман исключительно в качестве оптимизационной техники. Название раздела – это пережиток прошлого, оно расшифровывается как *block started by symbol* (блок, начинающийся с символа) или *block storage segment* (сегмент хранения блока). Прочие разделы, которые обычно присутствуют в исполняемых программах ELF, – это *абсолютный раздел* (absolute section), в котором содержатся непере-

¹ Этот подход впервые был реализован в операционной системе Plan9 от Bell Labs.

мещаемые символы, и *неопределенный раздел* (undefined section), также называемый ловушкой или отстойником.

Процесс ассоциируется с разнообразными системными ресурсами, которые выделяются и управляются ядром. Обычно процессы запрашивают ресурсы и манипулируют ими только при помощи системных вызовов. Ресурсы включают в себя таймеры, ожидающие сигналы, открытые файлы, сетевые соединения, аппаратные ресурсы и механизмы IPC. Ресурсы процесса вместе с относящимися к процессу данными и статистикой хранятся внутри ядра в *дескрипторе процесса* (process descriptor).

Процесс – это абстракция виртуализации. Ядро Linux, поддерживающее и приоритетную многозадачность, и виртуальную память, предоставляет процессу и виртуализованный процессор, и виртуализованное представление памяти. С точки зрения процесса система выглядит так, как будто только он управляет ею. То есть несмотря на то, что данный процесс может выполняться по расписанию, включающему множество других процессов, он работает так, как если бы у него в руках был единоличный контроль над системой. Ядро незаметно и прозрачно меняет приоритеты и графики выполнения процессов, деля ресурсы системных процессоров на все выполняющиеся процессы. Процессы ничего об этом не знают. Аналогичным образом каждому процессу выделяется собственное линейное адресное пространство, как если бы он один управлял всей памятью в системе. Благодаря виртуальной памяти и подкачке страниц ядро позволяет множеству процессов совместно существовать в системе, и каждый из них работает в своем адресном пространстве. Ядро управляет такой виртуализацией при помощи аппаратной поддержки, предоставляемой современными процессорами, что дает операционной системе возможность одновременно управлять состоянием множества независимых процессов.

Потоки выполнения

Каждый процесс состоит из одного или нескольких *потоков выполнения* (thread of execution), часто называемых просто *потоками* (thread). Поток – это единица активности в процессе, абстракция, несущая ответственность за исполнение кода и поддержание состояния выполнения процесса.

Большинство процессов включают в себя только один поток; они называются *однопоточными* (single-threaded). Процессы, содержащие несколько потоков, называются *многопоточными* (multithreaded). Традиционно программы Unix относятся к однопоточным. Причиной этого являются историческая простота Unix, быстрое создание процессов и надежные механизмы IPC – все это снижает необходимость в дополнительных потоках.

Поток состоит из *стека* (stack), где хранятся его локальные переменные (так же, как в стеках процессов в системах, не поддерживающих потоки), состояния процессора и текущего местоположения в конечной программе (которое обычно хранится в *указателе команд* (instruction pointer) процессора). Большинство остальных составляющих процесса совместно используется всеми потоками.

В глубине системы ядро Linux реализует уникальное представление потоков: это просто обычные процессы, которые совместно используют некоторые ресурсы (например, адресное пространство). В пользовательском пространстве Linux реализует потоки согласно стандарту POSIX 1003.1c (обычно их называют *pthread*). Текущая реализация потоков Linux, являющаяся частью glibc, носит название Native POSIX Threading Library (NPTL).

Иерархия процессов

Каждый процесс идентифицируется уникальным положительным целым числом, называемым *идентификатором процесса* (process ID, pid). Pid первого процесса равен единице, и каждый последующий процесс получает новое уникальное значение pid.

В Linux процессы формируют жесткую иерархию, известную как дерево процессов. Корень дерева процессов – это первый процесс, называемый процессом инициализации или процессом init и обычно принадлежащий программе init(8). Новые процессы создаются при помощи системного вызова fork(). Этот системный вызов создает дубликат вызывающего процесса. Исходный процесс называется *предком* (parent); новый процесс называется *потомком* (child). У каждого процесса, за исключением первого, есть предок. Если предок – родительский процесс – завершается до того, как заканчивается выполнение его потомка, то ядро *переназначает* предка (parent) для потомка, делая его потомком процесса инициализации.

Когда процесс завершается, он не удаляется немедленно из системы. Наоборот, ядро сохраняет в памяти часть резидента процесса, чтобы предок процесса мог запросить состояние своего потомка после завершения. Это называется *обслуживанием* (waiting on) завершенного процесса. Как только родительский процесс обслужит завершенный дочерний процесс, дочерний процесс полностью удаляется. Процесс, который уже завершился, но еще не был обслужен, называется *зомби* (zombie). Как положено, процесс инициализации обслуживает все свои дочерние процессы, гарантируя, что процессы с переназначенным предком не превратятся навсегда в зомби.

Пользователи и группы

Авторизация в Linux обеспечивается на базе концепции *пользователей* (user) и *групп* (group). С каждым пользователем связывается уникальное положительное целое значение, называемое *идентификатором пользователя* (user ID, uid). Каждому процессу, в свою очередь, назначается в точности один uid, идентифицирующий пользователя, запустившего процесс, который называется *реальным uid* (real uid) процесса. Внутри ядра Linux идентификатор uid – это единственная концепция пользователя. Однако сами пользователизываются на себя и других пользователей при помощи *имен пользователя* (username), а не числовых значений. Имена пользователя и соответствующие значения uid хранятся в файле /etc/passwd, а библиотечные процедуры сопоставляют указываемые пользователями имена с соответствующими идентификаторами uid.

Во время входа в систему пользователь предоставляет имя и пароль программе `login(1)`. Если это существующее имя пользователя и правильный пароль, то программа `login(1)` запускает *оболочку входа* (`login shell`) пользователя, которая также указана в файле `/etc/passwd`, и назначает оболочке идентификатор `uid`, значение которого равно значению `uid` этого пользователя. Дочерние процессы наследуют значения `uid` своих предков.

Идентификатор `uid` со значением 0 связан с особым пользователем, известным как `root`. У пользователя `root` особые привилегии, и он может делать в системе практически все, что угодно. Например, только он может изменить `uid` процесса. Следовательно, программа `login(1)` выполняется под именем пользователя `root`.

Помимо реального `uid`, у каждого процесса также есть *действительный uid* (effective `uid`), *сохраненный uid* (saved `uid`) и *uid файловой системы* (filesystem `uid`). Значение реального `uid` всегда равно значению `uid` пользователя, запустившего процесс, а значение действительного `uid` может меняться согласно различным правилам, позволяя процессу выполнять с правами других пользователей. В сохраненном `uid` содержится исходное значение действительного `uid`; оно используется при принятии решения, какие другие действительные `uid` разрешены для пользователя. `Uid` файловой системы, значение которого обычно равно значению действительного `uid`, применяется для проверки доступа к файловой системе.

Каждый пользователь может принадлежать одной или нескольким группам, включая *первичную группу* (primary group) или *группу входа в систему* (`login group`), перечисленным в `/etc/passwd`, и, возможно, нескольким *дополнительным группам* (supplemental group), перечисленным в `/etc/group`. Таким образом, с каждым процессом связывается соответствующий *идентификатор группы* (group ID, `gid`), и для него создается *реальный gid*, *действительный gid*, *сохраненный gid* и *gid файловой системы*. Процессы обычно связываются с группой входа в систему пользователя, а не с какими-либо дополнительными группами.

Проверки безопасности позволяют процессам выполнять определенные операции, только если они отвечают каким-то существующим критериям. Исторически в Unix это решение было реализовано очень строго: у процессов со значением `uid`, равным 0, был доступ, а у других нет. Недавно в Linux эта система безопасности была заменена более общей системой *возможностей* (capabilities). Вместо простой бинарной проверки возможности позволяют ядру определять разрешения доступа на базе более мелкозернистых настроек.

Разрешения

Стандартно в Linux используется такой же механизм разрешений файлов и безопасности, какой исторически существовал в Unix.

С каждым файлом связывается владелец-пользователь, владелец-группа и набор битов разрешений. Значения битов определяют, может ли владелец-пользователь, владелец-группа и кто угодно читать, записывать и выполнять файл. Для каждого из этих трех классов используется по три бита, итого битов

разрешений девять. Сведения о владельцах и разрешениях хранятся в inode файла.

Для обычных файлов разрешения очевидны: они определяют возможность открывать файл для чтения, открывать файл для записи и выполнять файл. Разрешения на чтение и запись для специальных файлов такие же, как и для обычных файлов, хотя что именно считывается или записывается в файл — зависит от конкретного специального файла. Разрешения на выполнение для специальных файлов игнорируются. Для каталогов разрешение на чтение позволяет перечислять содержимое каталога, разрешение на запись — добавлять в каталог новые ссылки, а разрешение на выполнение — входить в каталог и использовать его в полных путях. В табл. 1.1 перечислены все девять битов разрешений, их восьмеричные значения (распространенный способ представления девяти битов), их текстовые значения (то, что вы можете увидеть в выводе команды `ls`), а также дается описание каждого разрешения.

Таблица 1.1. Биты разрешений и их значения

Бит	Восьмеричное значение	Текстовое значение	Соответствующее разрешение
8	400	r—	Владелец может читать
7	200	-w--	Владелец может писать
6	100	-x—	Владелец может выполнять
5	040	--r--	Группа может читать
4	020	--w—	Группа может писать
3	010	--x--	Группа может выполнять
2	004	—r-	Кто угодно может читать
1	002	—w-	Кто угодно может писать
0	001	—x	Кто угодно может выполнять

Помимо исторических разрешений Unix, Linux также поддерживает списки управления доступом (access control list, ACL). Списки ACL позволяют использовать более детальные и точные разрешения и управление безопасностью, но за счет увеличения сложности и места на диске.

Сигналы

Сигналы — это механизм односторонних асинхронных уведомлений. Сигнал может отправляться ядром процессу, процессом другому процессу или процессом самому себе. Сигналы обычно извещают процесс о каком-то событии, таком, как сбой сегментации, или таком, как нажатие пользователем клавишного сочетания `Ctrl+C`.

Ядро Linux реализует около 30 сигналов (точное их число зависит от архитектуры). Каждый сигнал представляется числовой константой и текстовым именем. Например, у сигнала SIGHUP, который уведомляет о терминальном зави- сании, значение в архитектуре i386 равно 1.

За исключением сигналов SIGKILL (который всегда завершает процесс) и SIGSTOP (который всегда останавливает процесс), остальные сигналы позволяют процессам решать, что должно происходить при получении сигналов. В зависимости от сигнала, процесс может выполнить действие по умолчанию, которым может быть завершение процесса, завершение процесса с созданием дампа, остановка процесса или отсутствие действия. Иначе, процесс может явно проигнорировать или обработать сигнал. Если сигнал проигнорирован, то он бесшумно удаляется. Обработка сигнала включает в себя выполнение определенной пользователем функции *обработчика сигналов* (signal handler). Программа переходит к этой функции, как только получает сигнал, а после возвращения обработчика сигнала выполнение программы продолжается с той инструкции, на которой было прервано.

Взаимодействие процессов

Позволять процессам обмениваться информацией и уведомлять друг друга о событиях – это одна из самых важных задач операционной системы. Ядро Linux реализует большинство из исторических механизмов взаимодействия процессов Unix, в том числе определенные и стандартизованные как System V, так и POSIX, а также включает один или два собственных механизма.

Механизмы IPC, поддерживаемые в Linux, – это конвейеры, именованные конвейеры, семафоры, очереди сообщений, общая память и фьютексы (futex).

Заголовки

Системное программирование в Linux вращается вокруг нескольких заголовков. И само ядро, и glibc предоставляют заголовки, которые используются в программировании на системном уровне. Эти заголовки включают в себя стандартные заголовки C (например, `<string.h>`) и обычные предложения Unix (скажем, `<unistd.h>`).

Обработка ошибок

Нет необходимости говорить, что проверка возникновения ошибок и их обработка имеют первостепенную важность. В системном программировании ошибка обозначается при помощи возвращаемого значения функции и описывается специальной переменной `errno`. Glibc обеспечивает прозрачную поддержку `errno` как для библиотечных, так и для системных вызовов. Большая часть интерфейсов, о которых пойдет речь в данной книге, использует этот механизм для того, чтобы сообщать об ошибках.

Функция уведомляет вызывающего об ошибке при помощи специального возвращаемого значения, которое обычно равно -1 (точное значение зависит от функции). Значение ошибки уведомляет вызывающего о том, что ошибка произошла, но не предоставляет никаких сведений, позволяющих понять, почему ошибка возникла. Переменная `errno` используется для поиска причины ошибки.

Эта переменная определена в заголовке `<errno.h>` следующим образом:

```
extern int errno;
```

Ее значение имеет смысл только сразу же после того, как функция, в которой устанавливается переменная `errno`, сообщает об ошибке (обычно возвращая значение -1). Во время успешного выполнения функции значение `errno` может меняться.

Переменную `errno` можно считывать и записывать напрямую — это модифицируемое именуемое выражение. Значение `errno` позволяет также получать текстовое описание конкретной ошибки. Помимо этого, препроцессор `#define` определяет числовое значение `errno`. Например, числовое значение 1 соответствует текстовому значению `EACCESS` и ошибке «отсутствует разрешение» (`permission denied`). В табл. 1.2 перечислены стандартные определения и соответствующие описания ошибок.

Таблица 1.2. Ошибки и их описания

Определение препроцессора	Описание
<code>E2BIG</code>	Слишком длинный список аргументов
<code>EACCES</code>	Отсутствует разрешение
<code>EAGAIN</code>	Повторите попытку
<code>EBADF</code>	Неправильный номер файла
<code>EBUSY</code>	Устройство или ресурс заняты
<code>ECHILD</code>	Дочерние процессы отсутствуют
<code>EDOM</code>	Математический аргумент за пределами домена функции
<code>EEXIST</code>	Файл уже существует
<code>EFAULT</code>	Неправильный адрес
<code>EFBIG</code>	Слишком большой файл
<code>EINTR</code>	Системный вызов был прерван
<code>EINVAL</code>	Недопустимый аргумент
<code>EIO</code>	Ошибка ввода-вывода
<code>EISDIR</code>	Это каталог
<code>EMFILE</code>	Слишком много открытых файлов
<code>EMLINK</code>	Слишком много ссылок
<code>ENFILE</code>	Переполнение таблицы файлов

Определение препроцессора	Описание
ENODEV	Устройство отсутствует
ENOENT	Файл или каталог отсутствует
ENOEXEC	Ошибка формата запуска
ENOMEM	Нехватка памяти
ENOSPC	На устройстве нет свободного пространства
ENOTDIR	Это не каталог
ENOTTY	Недопустимая операция управления ввода-вывода
ENXIO	Устройство или адрес отсутствует
EPERM	Операция не разрешена
EPIPE	Сломанный конвейер
ERANGE	Слишком большой результат
EROFS	Файловая система, доступная только для чтения
ESPIPE	Недопустимый поиск
ESRCH	Процесс отсутствует
ETXTBSY	Текстовый файл занят
EXDEV	Неправильная ссылка

Библиотека C предоставляет несколько функций, позволяющих преобразовывать значение errno в соответствующее текстовое представление. Это необходимо только для того, чтобы сообщать об ошибках, и для подобных ситуаций; проверка и обработка ошибок может выполняться с использованием определений препроцессора и переменной errno напрямую.

Первая из подобных функций – perror():

```
#include <stdio.h>
```

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

Эта функция печатает на устройстве stderr (standard error, стандартная ошибка) строковое представление текущей ошибки, взятое из errno, добавляя в качестве префикса строку, на которую указывает параметр str, за которой следует двоеточие. Чтобы получать удобный и полезный результат, в эту строку следует включать имя функции, возвратившей ошибку, например:

```
if (close (fd) == -1)
    perror ("close").
```

Библиотека C также предоставляет функции strerror() и strerror_r(), прототипы которых

```
#include <string.h>
```

```
char * strerror (int errnum);
```

и

```
#include <string.h>

int strerror_r (int errnum, char *buf, size_t len);
```

Первая функция возвращает указатель на строку, описывающую ошибку, которая передается при помощи `errnum`. Приложение не может менять эту строку, но это могут делать последующие вызовы `perror()` и `strerror()`. Таким образом, данная функция не поточно-ориентирована.

Функция `strerror_r()` поточно-ориентирована. Она заполняет буфер длиной `len`, на который указывает параметр `buf`. Вызов `strerror_r()` возвращает 0 при успешном завершении и -1 в случае сбоя. Интересно, что, когда возникает ошибка, функция устанавливает все ту же переменную `errno`.

Несколько функций могут возвращать любые значения из диапазона типа возврата. В этих случаях, перед тем как использовать переменную `errno`, ей необходимо присваивать значение 0, а затем после завершения функции проверять ее новое значение (эти функции в случае ошибки должны возвращать только значения, не равные нулю). Например:

```
errno = 0;
arg = strtoul (buf, NULL, 0);
if (errno)
    perror ("strtoul").
```

Очень часто при проверке переменной `errno` забывают, что любой библиотечный или системный вызов может модифицировать ее значение. Например, это неправильный код:

```
if (fsync (fd) == -1) {
    fprintf (stderr, "fsync failed!\n");
    if (errno == EIO)
        fprintf (stderr, "I/O error on %d!\n", fd).
}
```

Для того чтобы сохранять значение `errno` между вызовами функций, передавайте его другой переменной:

```
if (fsync (fd) == -1) {
    int err = errno;
    fprintf (stderr, "fsync failed: %s\n", strerror (err));
    if (err == EIO) {
        /* если ошибка связана с вводом-выводом.
         * дезертировать с корабля */
        fprintf (stderr, "I/O error on %d!\n", fd);
        exit (EXIT_FAILURE);
    }
}
```

В однопоточных программах `errno` — это глобальная переменная, как показано ранее в этом разделе. В многопоточных программах, однако, переменная `errno` сохраняется для каждого потока и, следовательно, поточно-ориентирована.

Знакомство с системным программированием

В этой главе мы познакомились с основами системного программирования Linux и получили общее представление о системе Linux, необходимое любому программисту. В следующей главе рассматривается базовый файловый ввод-вывод. Сюда входят, конечно же, чтение и запись файлов. Однако, так как в Linux множество интерфейсов реализованы как файлы, файловый ввод-вывод принципиально важен не только для, м-м-м, файлов, но и для многоного другого.

Теперь, когда предварительное обсуждение осталось позади, настало время окунуться в настоящее системное программирование. Вперед!

2 Файловый ввод-вывод

В этой главе мы будем говорить об основах чтения файлов и записи в файлы. Эти операции формируют основу системы Unix. В следующей главе рассматривается стандартный ввод-вывод из стандартной библиотеки C, а в четвертой главе обсуждение продолжается на уровне более сложных и специализированных интерфейсов ввода-вывода. Завершается эта тема в главе 7, которая посвящена манипулированию файлами и каталогами.

Перед тем как считать файл или записать что-нибудь в него, файл необходимо открыть. Ядро ведет списки открытых файлов для всех процессов. Эти списки называются *таблицами файлов* (file table), которые индексируются при помощи неотрицательных целых значений, известных как *дескрипторы файлов* (file descriptor, что часто сокращают до fd). Каждая запись в списке содержит информацию об открытом файле, включая указатель на находящуюся в памяти копию inode файла и связанные метаданные, такие, как позиция в файле и режимы доступа. И пользовательское пространство, и пространство ядра используют файловые дескрипторы в качестве уникальных для каждого процесса маркеров. При открытии файла возвращается файловый дескриптор, а последующие операции (чтывание, запись и т. п.) принимают файловый дескриптор как основной аргумент.

По умолчанию дочерний процесс получает копию таблицы файлов своего предка. Список открытых файлов и их режимов доступа, текущие позиции в файлах и т. д. остаются одними и теми же, но изменение в одном процессе – например, когда дочерний процесс закрывает файл, – не влияет на таблицу файлов другого процесса. Однако, как вы увидите в главе 5, дочерний и родительский процессы могут совместно использовать таблицу файлов родителя (как это делают потоки выполнения).

Файловые дескрипторы представляются типом C `int`. Тот факт, что не используется специальный тип – например, `fd_t`, – часто считается недоразумением, но так в Unix сложилось исторически. У каждого процесса в Linux есть максимальное число файлов, которые он может открыть. Нумерация файловых

дескрипторов начинается с нуля и продолжается до значения, на единицу меньшего максимального числа открытых файлов. По умолчанию максимальное значение равно 1024, но его можно увеличить до 1 048 576. Так как отрицательные значения не являются допустимыми файловыми дескрипторами, -1 часто используется в качестве возвращаемого значения ошибки в функциях, которые в противном случае возвращают допустимый файловый дескриптор.

Если только процесс явно не закроет их, у каждого процесса по определению есть по крайней мере три открытых файловых дескриптора: 0, 1 и 2. Файловый дескриптор 0 относится к *стандартному вводу* (standard in, stdin), файловый дескриптор 1 – это *стандартный вывод* (standard out, stdout), а файловый дескриптор 2 – *стандартная ошибка* (standard error, stderr).

Вместо того чтобы напрямую ссылаться на эти целые значения, библиотека C предоставляет определения препроцессора STDIN_FILENO, STDOUT_FILENO и STDERR_FILENO.

Обратите внимание, что файловые дескрипторы могут ссылаться не только на обычные файлы. Они используются для доступа к файлам устройств и конвейерам, каталогам и фьютексам, конвейерам FIFO и сокетам – согласно философии «все является файлом», практически ко всему, что можно прочитать или записать, обращение идет через файловый дескриптор.

Открытие файлов

Самый простой способ обратиться к файлу – использовать системный вызов `read()` или `write()`. Однако, перед тем как обращаться к файлу, его необходимо открыть при помощи системного вызова `open()` или `creat()`. После того как вы завершите работу с файлом, его нужно закрыть, используя системный вызов `close()`.

Системный вызов `open()`

При помощи системного вызова `open()` вы открываете файл и получаете дескриптор файла:

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

int open (const char *name, int flags);
int open (const char *name, int flags, mode_t mode);
```

Системный вызов `open()` сопоставляет переданный ему в аргументе `name` полный путь к файлу с дескриптором файла, который он и возвращает в случае успешного завершения операции. Позиция в файле устанавливается в ноль, а разрешения на доступ к файлу определяются исходя из флагов, переданных при помощи аргумента `flags`.

Флаги для системного вызова open()

Аргумент `flags` может принимать значения `O_RDONLY`, `O_WRONLY` и `O_RDWR`. Используя эти значения, вы, соответственно, запрашиваете открытие файла только для чтения, только для записи или для чтения и записи.

Например, следующий код открывает файл `/home/kidd/madagascar` для чтения:

```
int fd;  
  
fd = open ("/home/kidd/madagascar", O_RDONLY).  
if (fd == -1)  
    /* ошибка */
```

Невозможно считывать данные из файла, открытого только для записи, или записывать в файл, открытый только для чтения. Процесс, выполняющий системный вызов `open()`, должен обладать достаточными разрешениями, чтобы получить запрошенный уровень доступа.

Аргумент `flags` можно сочетать со следующими значениями, применяя к ним операцию побитового «ИЛИ», чтобы модифицировать поведение запроса на открытие:

O_APPEND

Файл открывается в *режиме присоединения* (append mode). Это означает, что перед каждой записью указатель позиции в файле устанавливается на конец файла. Указатель передвигается, даже если другой процесс что-то записал в файл после последней операции записи вызывающего процесса (см. раздел «Режим присоединения» далее в этой главе).

O_ASYNC

Когда указанный файл становится доступным для чтения или записи, генерируется определенный сигнал (по умолчанию `SIGIO`). Этот флаг доступен только для терминалов и сокетов, но не для обычных файлов.

O_CREAT

Если файл, указанный при помощи аргумента `name`, не существует, то ядро создает его. Если файл уже существует, то этот флаг не оказывает никакого эффекта, если только одновременно не используется флаг `O_EXCL`.

O_DIRECT

Файл открывается для прямого ввода-вывода (см. раздел «Прямой ввод-вывод» далее в этой главе).

O_DIRECTORY

Если файл, указанный при помощи аргумента `name`, не является каталогом, то вызов `call()` завершается ошибкой. Этот флаг используется библиотечным вызовом `opendir()` внутри системы.

O_EXCL

Когда этот флаг добавляется к флагу `O_CREAT`, то вызов `open()` завершается ошибкой, если файл, указанный при помощи аргумента `name`, уже существует.

вует. Используется для предотвращения условий состязания при создании файлов.

O_LARGEFILE

Указанный файл открывается с использованием 64-битных смещений, что позволяет открывать файлы, размер которых превышает два гигабайта. Предназначается для 64-разрядных архитектур.

O_NOCTTY

Если указанный при помощи аргумента `name` файл относится к терминальному устройству (например, `/dev/tty`), то он не становится управляющим терминалом процесса, даже если у процесса в данный момент нет управляющего терминала. Используется редко.

O_NOFOLLOW

Если аргумент `name` содержит в конце пути символьическую ссылку, то вызов `open()` завершается ошибкой. В обычной ситуации ссылка разрешается и открывается целевой файл. Если другие компоненты указанного пути являются ссылками, то вызов завершается успешно. Например, если в аргументе `name` содержится путь `/etc/ship/plank.txt`, то в случае, когда `plank.txt` – это символьская ссылка, вызов завершается ошибкой. Если же `etc` или `ship` – символьская ссылка, а `plank.txt` – нет, то вызов выполняется успешно.

O_NONBLOCK

Если возможно, файл открывается в режиме без блокировки. Ни вызов `open()`, ни любые другие операции не приводят к блокировке (засыпанию) процесса во время ввода-вывода. Это поведение может определяться только для конвейеров FIFO.

O_SYNC

Файл открывается для синхронного ввода-вывода. Ни одна операция записи не завершается, пока данные физически не записываются на диск; обычные операции считывания и так выполняются синхронно, поэтому данный флаг никак не влияет на чтение файлов. В POSIX дополнительно определяются флаги `O_DSYNC` и `O_RSYNC`; в Linux эти флаги – синонимы `O_SYNC` (см. раздел «Флаг `O_SYNC`» далее в этой главе).

O_TRUNC

Если файл существует и это обычный файл, а указанные флаги допускают запись, то файл усекается до нулевой длины. Флаг `O_TRUNC` для конвейера FIFO или терминального устройства игнорируется. Использование для файлов других типов не определено. Применение `O_TRUNC` с флагом `O_RDONLY` также не определено, так как, для того чтобы выполнить усечение файла, необходим доступ к нему на запись.

Например, следующий код открывает для записи файл `/home/teach/pearl`. Если файл уже существует, то он усекается до нулевой длины. Так как флаг `O_CREAT` не используется, если файл не существует, вызов возвращает ошибку:

```
int fd;
fd = open ("/home/teach/pear1", O_WRONLY | O_TRUNC);
if (fd == -1)
    /* ошибка */
```

Владельцы новых файлов

Определение пользователя, который должен стать владельцем нового файла, выполняется просто: uid владельца файла — это действительный uid процесса, создающего файл.

Определение владельца-группы немного сложнее. По умолчанию файл относится к той группе, которой принадлежит действительный gid процесса, создающего файл. Это поведение, заложенное в System V (поведенческая модель для Linux), и стандартный *modus operandi* Linux.

Чтобы усложнить ситуацию, BSD определяет собственное поведение: группу для файла выбирается на основе gid родительского каталога. Это поведение можно включить в Linux при помощи параметра времени монтирования (*mount-time option*)¹; также оно работает в Linux по умолчанию, если для родительского каталога файла установлен бит setgid (set group ID, установка идентификатора группы). Хотя большинство систем Linux работают по принципу System V (когда новые файлы получают gid создавшего процесса), возможность поведения в стиле BSD (когда новые файлы получают gid родительского каталога) подразумевает, что, когда это действительно необходимо, группу для файла нужно выбирать явно при помощи системного вызова *chown()* (см. главу 7).

К счастью, вопрос выбора владельца-группы файла встает не часто.

Разрешения новых файлов

Допустимы и могут использоваться обе из указанных ранее форм системного вызова *open()*. Аргумент mode игнорируется во всех случаях, когда файл не создается; он необходим, когда передается флаг O_CREAT. Если вы забудете добавить аргумент mode при использовании флага O_CREAT, то результат будет непредсказуемым и, весьма вероятно, довольно неприятным, поэтому старайтесь этого не забывать.

При создании файла вы передаете разрешения для него при помощи аргумента mode. Во время операции открытия файла режим не проверяется, вследствие чего можно выполнить противоречивое действие, скажем, открыть файл для записи, но присвоить ему разрешения только на чтение.

Аргумент mode представляет собой уже знакомый набор битов разрешений Unix. Например, восьмеричное значение 0644 обозначает, что владелец может считывать и записывать файл, а все остальные пользователи — только считывать. Технически говоря, стандарт POSIX допускает варьирование точных значений в зависимости от реализации, чтобы в разных системах Unix биты раз-

¹ Параметры монтирования bsdgroups и sysvgroups.

решений компоновались так, как это необходимо авторам. Для того чтобы компенсировать невозможность переноса позиций битов в обозначении режима в другие системы, в POSIX был представлен следующий набор констант, которые можно объединять в аргументе `mode`, применяя к ним операцию бинарного ИЛИ:

`S_IRWXU`

У владельца есть разрешения на чтение, запись и выполнение.

`S_IRUSR`

У владельца есть разрешение на чтение.

`S_IWUSR`

У владельца есть разрешение на запись.

`S_IXUSR`

У владельца есть разрешение на выполнение.

`S_IRWXG`

У группы есть разрешения на чтение, запись и выполнение.

`S_IRGRP`

У группы есть разрешение на чтение.

`S_IWGRP`

У группы есть разрешение на запись.

`S_IXGRP`

У группы есть разрешение на выполнение.

`S_IRWXO`

У любых пользователей есть разрешения на чтение, запись и выполнение.

`S_IROTH`

У любых пользователей есть разрешение на чтение.

`S_IWOTH`

У любых пользователей есть разрешение на запись.

`S_IXOTH`

У любых пользователей есть разрешение на выполнение.

Фактические биты разрешений, которые записываются для данного файла, определяются путем применения операции бинарного И между значением аргумента `mode` и пользовательской *маской создания файла* (*file creation mask*, `umask`) в дополнительном коде. Другими словами, биты `umask` *выключают* соответствующие биты в значении аргумента `mode`, передаваемом вызову `open()`. Таким образом, значение маски `umask`, равное 022, превращает аргумент `mode` 0666 в 0644 (0666 & ~022). Как системный программист, вы обычно не учитываете `umask` при определении разрешений — эта маска существует для того, чтобы пользователь мог ограничивать разрешения, которые его программа устанавливает для новых файлов.

В качестве примера следующий код открывает файл, указанный при помощи аргумента `file`, для записи. Если файл не существует, то с учетом `umask 022` он создается с разрешением 0644 (хотя аргумент `mode` равен 0664). Если файл существует, то он усекается до нулевой длины:

```
int fd;
fd = open (file, O_WRONLY | O_CREAT | O_TRUNC,
           S_IWUSR | S_IRUSR | S_IWGRP | S_IRGRP | S_IROTH);
if (fd == -1)
    /* ошибка */
```

Функция `creat()`

Комбинация `O_WRONLY | O_CREAT | O_TRUNC` настолько распространена, что существует отдельный системный вызов, обеспечивающий как раз это поведение:

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

int creat (const char *name, mode_t mode);
```

ПРИМЕЧАНИЕ

Да, все правильно, в имени функции отсутствует буква `e` (`creat` вместо `create`). Кен Томпсон, создатель Unix, однажды пошутил, что эта отсутствующая буква — самая большая оплошность в дизайне Unix.

Это типичный вызов `creat()`:

```
int fd;
fd = creat (file, 0644);
if (fd == -1)
    /* ошибка */
```

Он идентичен следующему вызову:

```
int fd;
fd = open (file, O_WRONLY | O_CREAT | O_TRUNC, 0644);
if (fd == -1)
    /* ошибка */
```

В большинстве архитектур Linux¹ `creat()` — это системный вызов, хотя его довольно просто реализовать в пользовательском пространстве:

```
int creat (const char *name, int mode)
{
    return open (name, O_WRONLY | O_CREAT | O_TRUNC, mode);
```

¹ Вспомните, что системные вызовы определяются индивидуально для каждой архитектуры. Например, хотя в i386 есть системный вызов `creat()`, в Alpha он отсутствует. Конечно же, вы можете использовать `creat()` в любой архитектуре, но это может быть библиотечная функция, а не функция, имеющая собственный системный вызов.

Такое дублирование – это наследие прошлого, когда у системного вызова open() было только два аргумента. Сегодня системный вызов creat() существует только для обеспечения совместимости. В новых архитектурах вызов creat() может реализовываться, как определено в glibc.

Возвращаемые значения и коды ошибок

И open(), и creat() в случае успешного завершения возвращают дескриптор файла. В случае ошибки оба вызова возвращают -1 и устанавливают для переменной errno подходящее значение ошибки (переменная errno рассматривается в главе 1; там же перечисляются возможные значения ошибок). Ошибка открытия файла обрабатывается довольно просто, так как обычно до открытия выполняется лишь несколько шагов, которые необходимо отменить, или вообще никакие действия не выполняются. Обычный ответ на ошибку – запросить у пользователя другое имя файла или просто прервать программу.

Чтение файла при помощи системного вызова read()

Теперь, разобравшись с открытием файла, давайте посмотрим, как же его можно прочитать. А в следующем разделе поговорим о записи.

Самый простой и распространенный механизм чтения – это системный вызов read(), определенный в POSIX.1 следующим образом:

```
#include <unistd.h>
```

```
ssize_t read (int fd, void *buf, size_t len);
```

Каждый вызов считывает до len байтов в буфер buf, начиная с текущего смещения в файле, который указывается при помощи аргумента fd. В случае успеха возвращается число записанных в buf байтов. В случае ошибки вызов возвращает значение -1 и устанавливает переменную errno. Позиция в файле увеличивается на число байтов, считанных из fd. Если объект, представляемый fd, не поддерживает поиск (например, это может быть файл устройства посимвольного ввода-вывода), то чтение всегда начинается с «текущей» позиции.

Простейшее использование этого системного вызова банально. В следующем примере выполняется считывание из файлового дескриптора fd в переменную word. Считывается число байтов, равное размеру типа unsigned long, который составляет 4 байта в 32-разрядных системах Linux и 8 байт в 64-разрядных системах. После того как вызов возвращает результат, переменная nr содержит число считанных байтов или значение -1 в случае ошибки.

```
unsigned long word;
ssize_t nr;
```

```
/* считать пару байт в 'word' из 'fd' */
```

```
nr = read (fd, &word, sizeof (unsigned long)).  
if (nr == -1)  
    /* ошибка */
```

Однако с этой наивной реализацией связана пара проблем: вызов может возвратить какое-то значение, не считав все `len` байтов, и вызвать определенные ошибки, которые не проверяются и не обрабатываются в этом коде. Такой код, к сожалению, очень распространен. Давайте посмотрим, как его можно улучшить.

Возвращаемые значения

Системный вызов `read()` может возвращать положительные не равные нулю значения, меньшие `len`. Это может происходить по нескольким причинам: количество доступных байтов меньше `len`, системный вызов прерван сигналом, конвейер сломан (если `fd` – это конвейер) и т. д.

Вторая проблема, связанная с использованием системного вызова `read()`, – это возможность возврата значения 0. `Read()` возвращает 0, указывая на достижение конца файла (`end-of-file`, EOF); соответственно, в этом случае никакие байты нечитываются. Достижение EOF не считается ошибкой (и поэтому вызов не возвращает значение `-1`). Возвращаемое значение 0 просто указывает, что позиция в файле вышла за пределы последнего допустимого смещения в файле, поэтому считывать больше нечего. Если же делается вызов на считывание `len` байтов, но байты для считывания недоступны, то вызов блокируется (засыпает) до тех пор, пока не появятся доступные байты (предполагая, что файловый дескриптор был открыт не в режиме без блокировки; см. раздел «Считывание без блокировки»). Обратите внимание, что эта ситуация отличается от ситуации с EOF, когда достигается конец файла. В случае блокировки вызов ожидает дополнительные данные – такое может происходить при считывании из сокета или файла устройства.

Некоторые ошибки легко устраняются. Например, если вызов `read()` прерывается сигналом до того, как он успевает считать байты, то возвращается значение `-1` (значение 0 можно перепутать с ситуацией EOF) и переменной `errno` присваивается значение `EINTR`. В таком случае можно просто повторить считывание.

Действительно, существует много возможных завершений вызова `read()`:

- вызов возвращает значение, равное `len`. Все `len` считанных байтов записаны в буфер `buf`. Результат допустим и ожидаем;
- вызов возвращает значение, меньшее `len`, но большее нуля. Считанные байты записаны в `buf`. Это может происходить, когда сигнал прерывает незавершенную операцию считывания, когда в середине считывания происходит ошибка, когда доступно больше нуля, но меньше `len` байтов, когда конец файла достигается раньше, чем удается считать `len` байтов. Повторение вызова (с соответствующим образом обновленными значениями `buf` и `len`) по-

зволит считать оставшиеся байты в оставшееся свободное пространство буфера или узнать причину ошибки;

- вызов возвращает значение 0. Это указывает на достижение EOF. Считывать нечего;
- вызов блокируется, потому что данные в данный момент времени недоступны. Этого не происходит в режиме без блокировки;
- вызов возвращает значение -1, а для переменной errno устанавливается значение EINTR. Это указывает на то, что сигнал был получен до того, как удалось считать хотя бы один байт. Вызов можно повторить;
- вызов возвращает значение -1, а для переменной errno устанавливается значение EAGAIN. Это указывает на то, что считывание было заблокировано из-за отсутствия данных, поэтому запрос нужно повторить позже. Происходит только в режиме без блокировки;
- вызов возвращает значение -1, а переменной errno присваивается значение, отличное от EINTR или EAGAIN. Это указывает на более серьезную ошибку.

Считывание всех байтов

Перечисленные возможные ситуации демонстрируют, что предыдущее тривиальное и упрощенное использование системного вызова read() не подходит, если вы хотите обрабатывать все ошибки и действительно считывать все len необходимых байтов (по крайней мере, до достижения EOF). Для реализации правильного вызова вам необходим цикл и несколько условных операторов:

```
ssize_t ret.

while (len != 0 && (ret = read (fd, buf, len)) != 0) {
    if (ret == 1) {
        if (errno == EINTR)
            continue;
        perror ("read");
        break;
    }

    len -= ret;
    buf += ret.
}
```

В этом фрагменте кода обрабатываются все пять условий. Цикл считывает len байтов, начиная с текущей позиции в файле fd, в буфер buf. Чтение файла продолжается до тех пор, пока не считаются все len байтов или пока не достигается условие EOF. Если считывается больше нуля, но меньше len байтов, то значение len уменьшается на число полученных байтов, buf увеличивается на то же число, и вызов повторяется. Если вызов возвращает -1 и переменная errno содержит EINTR, то вызов повторяется без обновления параметров. Если вызов возвращает -1 и переменная errno содержит любое другое значение, то вызывается perror(), чтобы вывести в стандартной ошибке описание проблемы, и цикл прерывается.

Частичное считывание не только допустимо, но и встречается довольно часто. Бессчетное количество ошибок возникает из-за того, что программисты не выполняют правильную проверку и обработку коротких запросов на считывание. Вы не должны пополнять этот список!

Считывание без блокировки

Иногда программистам нужно, чтобы вызов `read()` в случае, когда данные недоступны, не блокировался. Они предпочитают, чтобы вызов сразу же возвращал какое-то значение, указывая, что данных для считывания нет. Это называется *вводом-выводом без блокировки* (*nonblocking I/O*) и позволяет приложениям выполнять ввод-вывод, возможно, даже на нескольких файлах, вообще не блокируясь и, следовательно, не упуская доступные данные в других файлах.

Следовательно, стоит проверять еще одно значение переменной `errno`: `EAGAIN`. Как говорилось ранее, если определенный дескриптор файла открывается в режиме без блокирования (вызову `open()` передан флаг `O_NONBLOCK`; см. раздел «Флаги для системного вызова `open()`»), но данных для считывания нет, то вызов `read()` возвращает значение `-1` и устанавливает для переменной `errno` значение `EAGAIN`, а не блокируется. При выполнении считывания без блокировки необходимо всегда проверять `errno` на значение `EAGAIN`, чтобы не спутать серьезную ошибку с простым отсутствием данных. Например, можно использовать такой код:

```
char buf[BUFSIZ];
ssize_t nr.

start
nr = read(fd, buf, BUFSIZ).
if (nr == 1) {
    if (errno == EINTR)
        goto start; /* о. тише! */
    if (errno == EAGAIN)
        /* повторить позже */
    else
        /* ошибка */
}
```

ПРИМЕЧАНИЕ

Обработка значения ошибки `EAGAIN` в этом примере с использованием инструкции `goto start` в действительности несет мало смысла — можно с таким же успехом вообще не использовать ввод-вывод без блокировки. Использование режима без блокировки экономит время, но приводит к увеличению нагрузки из-за повторяющихся циклов.

Прочие значения ошибки

Прочие коды ошибки относятся к программным ошибкам или (как `EIO`) к проблемам на низком уровне. Возможные значения `errno` после сбоя вызова `read()` включают в себя:

EBADF

Данный файловый дескриптор недопустим или не открыт для чтения.

EFAULT

Указатель, предоставленный вызову при помощи buf, находится за пределами адресного пространства вызывающего процесса.

EINVAL

Файловый дескриптор соответствует объекту, не допускающему чтение.

EIO

Произошла низкоуровневая ошибка ввода-вывода.

Лимиты размера для вызова read()

Типы `size_t` и `ssize_t` определены в POSIX. Тип `size_t` используется для хранения значений, обозначающих размеры, в байтах. Тип `ssize_t` – это версия `size_t` со знаком (отрицательные значения обозначают ошибки). В 32-разрядных системах дополнительные типы C – это обычно `unsigned int` и `int` соответственно. Так как эти два типа часто используются вместе, потенциально меньший диапазон `ssize_t` ограничивает диапазон `size_t`.

Максимальное значение типа `size_t` – это `SIZE_MAX`; максимальное значение типа `ssize_t` – `SSIZE_MAX`. Если значение `len` больше `SSIZE_MAX`, то результаты вызова `read()` не определены. В большинстве систем Linux размер `SSIZE_MAX` равен `LONG_MAX`, то есть `0xffffffff` на 32-разрядной машине. Это относительно много для одной операции считывания, но все равно не стоит забывать об этом ограничении. Если вы используете показанный выше цикл считывания для обычного объемного считывания, то стоит добавлять подобные строки:

```
if (len > SSIZE_MAX)
    len = SSIZE_MAX.
```

Вызов `read()` с аргументом `len`, равным пузлу, приводит к тому, что немедленно возвращается значение, равное 0.

Запись при помощи системного вызова write()

Самый простой и распространенный системный вызов, используемый для записи, – это `write()`. Вызов `write()` представляет собой эквивалент `read()` и также определен в POSIX.1:

```
#include <unistd.h>
```

```
ssize_t write (int fd, const void *buf, size_t count);
```

Когда вы делаете вызов `write()`, начиная с текущей позиции в файле, указанным при помощи файлового дескриптора `fd`, в него записывается до `count` байтов из буфера `buf`. Файлы, представляющие объекты, которые не поддерживают

поиск (например, устройства посимвольного ввода-вывода), всегда записываются начиная с «головы».

В случае успеха возвращается количество записанных байтов, а позиция в файле соответствующим образом обновляется. В случае ошибки возвращается значение `-1` и соответствующим образом устанавливается переменная `errno`. Вызов `write()` может вернуть значение `0`, но оно всего лишь указывает, что было записано ноль байт.

Как и для вызова `read()`, самый простой вариант использования `write()` выглядит так:

```
const char *buf = "My ship is solid!";
ssize_t nr.
```

```
/* записать строку из 'buf' в 'fd' */
nr = write (fd, buf, strlen (buf));
if (nr == -1)
    /* ошибка */
```

Но снова, как и для `read()`, этот вариант не совсем верен. Вызывающий должен проверить возможность того, что запись была выполнена лишь частично:

```
unsigned long word = 1720,
size_t count,
ssize_t nr.
```

```
count = sizeof (word),
nr = write (fd, &word, count);
if (nr == -1)
    /* ошибка, проверить переменную errno */
else if (nr != count)
    /* возможная ошибка, но переменная errno не установлена */
```

Частичная запись

Системный вызов `write()` выполняет лишь частичную запись с меньшей вероятностью, чем системный вызов `read()` делает частичное считывание. Помимо этого, для системного вызова `write()` не существует условия EOF. Что касается обычных файлов, `write()` гарантированно выполняет запрошеннюю операцию записи полностью, если только не происходит какая-либо ошибка.

Следовательно, для обычных файлов нет необходимости реализовывать запись в цикле. Однако для других типов файлов, например сокетов, цикл может потребоваться, чтобы гарантировать, что все запрошенные байты *действительно* будут записаны. Еще одно преимущество использования цикла состоит в том, что второй вызов `write()` может вернуть ошибку, поясняющую, почему первый вызов выполнил запись только частично (хотя, повторюсь, такая ситуация встречается не часто). Вот пример реализации с циклом:

```
ssize_t ret, nr.

while (len != 0 && (ret = write (fd, buf, len)) != 0) {
    if (ret == -1) {
        if (errno == EINTR)
```

```

        continue;
perror ("write");
break;
}

len -= ret;
buf += ret;
}

```

Режим присоединения

Когда файл fd открывается в режиме присоединения (при помощи флага `O_APPEND`), то запись начинается не с текущей позиции в файле, а с конца файла.

Например, предположим, что два процесса записывают данные в один и тот же файл. В другом режиме, отличном от режима присоединения, если первый процесс пишет данные в конец файла и потом то же самое сделает второй процесс, то позиция в файле для первого процесса уже не будет указывать на конец файла; она будет указывать в точку, отстающую от реального конца файла на длину данных, только что записанных вторым процессом. Это означает, что несколько процессов не могут присоединять данные к одному файлу без явной синхронизации, так как при этом они будут попадать в условия состязания.

Режим присоединения позволяет избежать этой проблемы. Он гарантирует, что позиция в файле всегда устанавливается на конец файла, поэтому операции записи всегда присоединяют данные в конце файла, даже если их выполняют разные процессы. Можно представлять себе это так, как будто перед каждым запросом на запись позиция в файле обновляется и после обновления указывает на самый конец последних записанных данных. Это не имеет значения для следующего вызова `write()`, так как `write()` автоматически обновляет позицию в файле, но может иметь значение, если после этого вы вызовете `read()` по какой-то причине.

Режим присоединения очень удобно использовать для определенных задач, таких, как обновление файлов журнала, но для большинства других он не нужен.

Запись без блокировки

Когда файл fd открывается в режиме без блокировки (при помощи флага `O_NONBLOCK`) и выполняется запись, которая в обычных условиях заблокировала бы, системный вызов `write()` возвращает значение `-1` и присваивает переменной `errno` значение `EAGAIN`. В такой ситуации запрос нужно повторить еще раз. С обычными файлами такая ситуация случается крайне редко.

Прочие коды ошибок

Прочие полезные значения `errno` включают в себя:

`EBADF`

Указанный файловый дескриптор недопустим или не открыт для записи.

EFAULT

Указатель в аргументе `buf` указывает за пределы адресного пространства процесса.

EFBIG

Данная операция записи сделала бы файл больше максимального допустимого размера для процесса или больше внутренних пределов реализации.

EINVAL

Указанный файловый дескриптор соответствует объекту, не допускающему запись.

EIO

Произошла низкоуровневая ошибка.

ENOSPC

В файловой системе, откуда взят данный дескриптор файла, недостаточно свободного пространства.

EPIPE

Данный файловый дескриптор связан с конвейером или сокетом,читывающая сторона которого закрыта. Процесс также получает сигнал `SIGPIPE`. Действие по умолчанию для сигнала `SIGPIPE` – завершить получивший его процесс. Таким образом, процессы получают данное значение `errno`, только если они явно игнорируют, блокируют или обрабатывают этот сигнал.

Ограничения размера для вызова `write()`

Если значение `count` превышает значение `SSIZE_MAX`, то результат вызова `write()` не определен.

Если выполнить `write()` со значением `count`, равным нулю, то вызов мгновенно вернет значение 0.

Поведение вызова `write()`

Когда `write()` возвращает значение, это означает, что ядро скопировало данные из предоставленного буфера в буфер ядра, но нет никакой гарантии, что данные были записаны в указанное целевое местоположение. Действительно, этот вызов возвращает значение слишком быстро, чтобы так действительно происходило. Несоответствие производительности процессоров и жестких дисков делало бы такое поведение мучительно очевидным.

Вместо этого, когда приложение из пользовательского пространства выполняет системный вызов `write()`, ядро Linux делает несколько проверок, а затем просто копирует данные в буфер. Позже в фоновом режиме ядро собирает все «грязные» буфера, сортирует их оптимальным способом и записывает их содержимое на диск (этот процесс называется *отложенной записью (writeback)*). Благодаря этому создается впечатление, что вызовы `write()` работают чрезвы-

чайно быстро, возвращая значение практически мгновенно, а ядро может откладывать фактическую запись на диск на периоды бездействия и потом сразу же обрабатывать большие пакеты запросов.

Отложенная запись не меняет семантику POSIX. Например, если выполняется вызов на считывание только что записанных данных, еще находящихся в буфере, то этот запрос удовлетворяется из буфера, а не возвращает «устаревшие» данные с диска. Такое поведение действительно повышает производительность, так как данныечитываются из кэша в памяти без обращения к диску. Запросы на чтение и запись чередуются, как положено, и результаты всегда возвращаются ожидаемые — конечно же, если в системе не происходит аварийный сбой до того, как данные попадают на диск! Даже если приложение будет считать, что запись произошла успешно, в такой ситуации данные никогда не окажутся в нужном месте на диске.

Еще один вопрос, связанный с отложенной записью, — это невозможность принудительно реализовывать *упорядочение записи* (write ordering). Хотя приложение может пытаться упорядочивать запросы на запись, чтобы эти данные оказывались на диске в определенном порядке, ядро все равно меняет порядок записи так, как считает необходимым, — в основном учитывая соображения производительности. Обычно это превращается в проблему только в случае системных сбоев, так как в конечном итоге данные из всех буферов записываются на диск и все хорошо. Как бы то ни было, большая часть приложений никогда не заботится о порядке записи буферов на диск.

Последняя проблема, связанная с отложенной записью, относится к сообщениям об определенных ошибках ввода-вывода. Процессу, который выполнил запрос на запись, невозможно сообщить ни о какой ошибке ввода-вывода, которая потенциально произойдет во время отложенной записи, например, о физическом сбое жесткого диска. Действительно, буферы совершенно не связаны с процессами. В один буфер могут попадать данные от нескольких процессов, и процессы могут завершаться после записи данных в буфер, но до того, как эти данные в итоге сбрасываются на диск. Помимо этого, как сообщить процессу, что запись не удалась, *ex post facto*?

Ядро делает попытки минимизировать риски отложенной записи. Для того чтобы гарантировать, что данные будут записываться своевременно, ядро определяет *максимальный возраст буфера* (maximum buffer age) и записывает на диск все данные из грязных буферов до того, как их срок жизни превышает этот предел. Пользователь может настроить это значение в файле /proc/sys/vm/dirty_expire_centiseconds; оно определяется в сотых долях секунды.

Также можно принудительно заставить ядро записать на диск буфер определенного файла или даже сделать все операции записи синхронными. Эти темы обсуждаются в следующем разделе, «Синхронизированный ввод-вывод».

Далее в этой главе в разделе «Внутреннее устройство ядра» мы глубже изучим систему отложенной записи буферов ядра Linux.

Синхронизированный ввод-вывод

Хотя синхронизация ввода-вывода — это очень важная тема, не следует бояться вопросов, связанных с отложенной записью. Буферизация записи обеспечивает *огромный выигрыш* в производительности, и, следовательно, в любой операционной системе, находящейся хотя бы на полпути к получению звания «современная», реализуется отложенная запись с использованием буферов. Тем не менее иногда в приложениях возникает необходимость контролировать момент, когда данные оказываются на диске. Для этого ядро Linux предоставляет несколько параметров, позволяющих обменивать производительность на синхронность операций.

fsync() и fdatasync()

Самый простой способ гарантировать, что данные достигнут диска, — использовать системный вызов `fsync()`, определенный в `POSIX.1b` следующим образом:

```
#include <unistd.h>
```

```
int fsync (int fd);
```

Вызов `fsync()` гарантирует, что все грязные данные, соответствующие файлу, который указан при помощи дескриптора файла `fd`, записываются обратно на диск. Файловый дескриптор `fd` должен быть открыт для записи. Этот вызов обеспечивает запись на диск и данных, и метаданных, например временных меток создания и других атрибутов, содержащихся в `inode`. Он не возвращает значение до тех пор, пока жесткий диск не сообщает, что данные и метаданные были успешно записаны на диск.

Если запись кэшируется на диске, то `fsync()` не может знать, попадают ли данные на диск физически. Жесткий диск может сообщить, что данные были записаны, пока они в действительности еще остаются в кэше записи диска. К счастью, данные в кэше жесткого диска должны фиксироваться на диске немедленно.

В Linux также предусмотрен системный вызов `fdatasync()`:

```
#include <unistd.h>
```

```
int fdatasync (int fd);
```

Этот системный вызов делает то же самое, что и `fsync()`, но он ограничивается сбросом только данных. Вызов не гарантирует, что метаданные также синхронизируются с диском, и поэтому работает потенциально быстрее. Зачастую этого бывает достаточно.

Обе функции применяются одинаково и очень просто:

```
int ret;
```

```
ret = fsync (fd);
if (ret == -1)
    /* ошибка */
```

Ни одна из функций не гарантирует, что обновленные записи каталога, содержащего файл, синхронизируются с диском. Это подразумевает, что если ссылка на файл только что обновлена, то данные файла уже могут быть успешно записаны на диск, но еще не связаны с записью каталога, что делает файл недоступным. Чтобы гарантировать, что все обновления в записях каталога также фиксируются на диске, необходимо вызывать `fsync()` для дескриптора файла, открытого для самого каталога.

Возвращаемые значения и коды ошибок

В случае успеха оба вызова возвращают значение 0. В случае сбоя оба вызова возвращают -1 и присваивают переменной `errno` одно из следующих трех значений:

`EBADF`

Указанный файловый дескриптор является недопустимым или не открыт для записи.

`EINVAL`

Указанный файловый дескриптор соответствует объекту, не поддерживающему синхронизацию.

`EIO`

Во время синхронизации произошла низкоуровневая ошибка ввода-вывода. Это говорит о реальной ошибке ввода-вывода, и именно здесь зачастую отлавливаются такие ошибки.

Вызов `fsync()` может легко завершиться ошибкой из-за того, что `fsync()` не реализован в конкретной файловой системе, даже если `fdatasync()` реализован. В параноидальных приложениях в случае, когда `fsync()` возвращает код ошибки `EINVAL`, можно пробовать делать то же самое с использованием `fdatasync()`. Например:

```
if (fsync (fd) == -1) {
    /*
     * Мы предпочитаем fsync(), но давайте попробуем fdatasync( ) на
     * случай, если fsync( ) не сработает, просто для гарантии
     */
    if (errno == EINVAL) {
        if (fdatasync (fd) == -1)
            perror ("fdatasync");
    } else
        perror ("fsync");
}
```

Так как POSIX требует использования `fsync()` и считает `fdatasync()` необязательным вызовом, системный вызов `fsync()` всегда должен реализовываться для обычных файлов во всех распространенных файловых системах Linux. Однако для необычных типов файлов (тех, в которых нет метаданных, требующих синхронизации) или странных файловых систем можно реализовывать только `fdatasync()`.

sync()

Не такой оптимальный, но более масштабный консервативный системный вызов sync() обеспечивает синхронизацию *всех* буферов с диском:

```
#include <unistd.h>
```

```
void sync (void);
```

У этой функции нет параметров, и она не возвращает никакое значение. Она всегда завершается успешно, и возврат значения указывает на то, что все буфера — как с данными, так и с метаданными — были гарантированно записаны на диск¹.

Стандарты не требуют, чтобы вызов sync() ожидал, пока все буфера будут сброшены на диск, чтобы возвратить значение; они требуют только, чтобы вызов инициировал процесс фиксации всех буферов на диске. По этой причине часто рекомендуется выполнять синхронизацию несколько раз, чтобы гарантировать, что все данные успешно оказались на диске. Linux, однако, *ожидается* фиксации буферов. Таким образом, одного вызова sync() достаточно.

Единственное реальное применение sync() можно найти в реализации утилиты sync(8). Для фиксации на диске данных для конкретных дескрипторов файлов в приложениях следует использовать fsync() и fdatasync(). Обратите внимание, что в загруженной системе для завершения работы sync() может потребоваться несколько минут.

Флаг O_SYNC

Системному вызову open() можно передавать флаг O_SYNC, указывая, что весь ввод-вывод для данного файла должен быть *синхронизированным*:

```
int fd;
fd = open (file, O_WRONLY | O_SYNC);
if (fd == -1) {
    perror ("open");
    return -1;
}
```

Запросы на чтение всегда синхронизируются. Если бы это было не так, то достоверность считываемых данных в предоставленном буфере всегда находилась бы под вопросом. Однако, как обсуждалось чуть выше, вызовы write() обычно не синхронизируются. Нет никакой связи между возвращением вызовом какого-то значения и фиксацией данных на диске. Флаг O_SYNC принудительно создает это взаимоотношение между действиями, гарантируя, что все вызовы write() будут выполнять синхронизированный ввод-вывод.

Можно представлять себе это так: O_SYNC заставляет систему неявно выполнять вызов fsync() после каждой операции write() до того, как вызов на запись

¹ Что ж, здесь скрывается та же ловушка: жесткий диск может солгать и сказать ядру, что содержимое буферов уже на диске, тогда как оно все еще остается в дисковом кэше.

возвращает значение. В действительности семантика именно такова, хотя в ядре Linux работа `O_SYNC` реализуется немного более эффективно.

При использовании флага `O_SYNC` время пользователя (user time) и время ядра (kernel time) (время, проведенное в пользовательском пространстве и пространстве ядра соответственно) для операций записи немного увеличивается. Более того, в зависимости от размера записываемого файла `O_SYNC` может вызывать увеличение общего затрачиваемого времени на один-два порядка, так как к общему времени работы добавляется еще и время ожидания ввода-вывода (I/O wait time) для процесса (время, которое тратится на ожидание завершения операции ввода-вывода). Издержки увеличиваются очень сильно, поэтому к синхронизированному вводу-выводу следует прибегать только после того, как опробованы все возможные альтернативы.

Обычно, когда требуется гарантия, что операции записи будут успешно завершаться фиксацией данных на диске, в приложениях применяются вызовы `fsync()` и `fdatasync()`. При этом издержки увеличиваются меньше, чем при использовании флага `O_SYNC`, так как эти вызовы можно делать реже (например, только после того, как завершаются определенные критически важные операции).

Флаги `O_DSYNC` и `O_RSYNC`

В стандарте POSIX определяются еще два флага системного `open()`, связанные с синхронизированным вводом-выводом: `O_DSYNC` и `O_RSYNC`. В Linux эти флаги считаются синонимами `O_SYNC` и предоставляют в точности ту же функциональность.

Флаг `O_DSYNC` указывает, что после каждой операции записи должны синхронизироваться только обычные данные, но не метаданные. Это можно представлять себе как неявное выполнение после каждого запроса на запись вызова `fdatasync()`. Так как `O_SYNC` обеспечивает более надежные гарантии, отсутствие явной поддержки `O_DSYNC` не приводит ни к каким потерям; единственное, что возможно небольшое падение производительности из-за более серьезных требований, предъявляемых `O_SYNC`.

Флаг `O_RSYNC` включает синхронизацию не только запросов на запись, но и запросов на чтение. Его необходимо использовать вместе с одним из других флагов — `O_SYNC` или `O_DSYNC`. Как я уже сказал выше, операции считывания всегда синхронизируются, — в конце концов, они не могут возвращать значение до тех пор, пока у них не будет чего-то, что можно предъявить пользователю. Флаг `O_RSYNC` предусматривает, что любые побочные эффекты операции считывания также были синхронизированы. Это означает, что метаданные, обновляемые после операций считывания, должны записываться на диск до того, как вызов возвращает значение. В практических терминах это требование, вероятнее всего, означает, что, перед тем как `read()` вернет значение, должно обновиться время последнего доступа к файлу в копии inode на диске. Linux определяет, что флаг `O_RSYNC` должен работать так же, как `O_SYNC`, хотя в действительности в этом

нет смысла (эти два флага не связаны между собой так, как `O_SYNC` и `O_DSYNC`). Пока что в Linux не существует способа добиться функциональности, которую должен обеспечивать флаг `O_RSYNC`; максимум, что может делать разработчик, — это вызывать `fdatasync()` после каждого вызова `read()`. Но в таком поведении необходимость возникает крайне редко.

Прямой ввод-вывод

Ядро Linux, как и ядро любой другой современной операционной системы, реализует сложный пласт возможностей кэширования, буферизации и управления вводом-выводом между устройствами и приложением (см. раздел «Внутреннее устройство ядра» в конце этой главы). Когда речь идет о быстродействующих приложениях, часто возникает желание обойти этот уровень сложности и реализовать собственные функции управления вводом-выводом.

Однако развертывание собственной системы ввода-вывода обычно не стоит затрачиваемых усилий, и в действительности инструменты, доступные на уровне операционной системы, скорее всего, позволяют добиться намного более высокой производительности, чем средства уровня приложений. И все же системы управления базами данных зачастую используют собственные функции кэширования, чтобы минимизировать присутствие в них операционной системы, как только это возможно.

Когда вы добавляете флаг `O_DIRECT` в вызов `open()`, ядро минимизирует управление вводом-выводом и ввод-вывод осуществляется напрямую из буферов в пользовательском пространстве в устройство, обходя страничный кэш. Все операции ввода-вывода синхронизируются и не возвращают значение до тех пор, пока полностью не завершаются.

При использовании прямого ввода-вывода длина запроса, выравнивание буфера и смещения файлов — все должны быть целыми числами, кратными размеру сектора соответствующего устройства — обычно он составляет 512 байт. До появления ядра Linux 2.6 это требование было еще строже: в версии 2.4 все должно было быть выровнено по размеру логического блока файловой системы (обычно 4 Кбайт). Для обеспечения совместимости приложения должны использовать выравнивание по размеру самого большого логического блока (что чаще всего наименее удобно).

Закрытие файлов

После того как работа с дескриптором файла в программе завершена, можно удалить связь между дескриптором и соответствующим файлом, применив системный вызов `close()`:

```
#include <unistd.h>

int close (int fd);
```

Когда `close()` срабатывает, проецирование открытого дескриптора файла `fd` на сам файл прекращается и процесс отсоединяется от файла. Данный файловый дескриптор больше не является допустимым, и ядро может повторно использовать его в качестве возвращаемого значения для последующих вызовов `open()` и `creat()`. Вызов `close()` в случае успеха возвращает значение 0. В случае ошибки он возвращает -1 и соответствующим образом устанавливает переменную `errno`. Применять его просто:

```
if (close (fd) == -1)
    perror ("close").
```

Обратите внимание, что закрытие файла не имеет никакого отношения к тому моменту, когда файл сбрасывается на диск. Если в приложении нужно удостовериться, что файл будет зафиксирован на диске до того, как он будет закрыт, то для этого необходимо использовать один из вариантов синхронизации, которые были рассмотрены в разделе «Синхронизированный ввод-вывод».

У операции закрытия файла есть тем не менее некоторые побочные эффекты. Когда закрывается последний открытый дескриптор файла, ссылающийся на файл, структура данных, представляющая этот файл внутри ядра, освобождается. Когда структура данных освобождается, она отсоединяет находящуюся в памяти копию структуры `inode`, связанной с файлом. Если больше ничего не удерживает `inode` в памяти, эта структура также может быть освобождена (она может и остаться там, так как ядро кэширует структуры `inode` для улучшения производительности, но не обязана). Если файл отсоединен от диска в момент, когда он еще открыт, то он физически не удаляется до тех пор, пока не будет закрыт, а его `inode` не будет удалена из памяти. Следовательно, вызов `close()` может приводить к итоговому физическому удалению отсоединенного файла с диска.

Значения ошибок

Очень часто в программах не проверяют возвращаемое значение вызова `close()`. Если не выполнить такую проверку, то можно пропустить критическое условие ошибки, так как ошибки, связанные с отложенными операциями, могут не проявляться довольно долго, а при помощи `close()` их вполне можно отлавливать.

Существует несколько возможных значений переменной `errno`, возвращаемых в случае ошибки. Помимо `EBADF` (данный дескриптор файла недопустим), наиболее важное значение ошибки — это `EIO`, указывающее на низкоуровневую ошибку ввода-вывода, возможно, не связанную с фактической операцией закрытия. Независимо от того, какая ошибка возвращается, дескриптор файла, если он допустим, всегда закрывается, а связанные с ним структуры данных освобождаются.

Хотя POSIX допускает это, вызов `close()` никогда не возвращает значение `EINTR`. Разработчики ядра Linux лучше знают, что такая реализация не блещет смыслом.

Поиск при помощи lseek()

Обычно ввод-вывод выполняется в файле линейно, и все действия поиска, в которых возникает необходимость, относятся только к неявному обновлению позиции в файле после операций считывания и записи. Некоторым приложениям, однако, приходится перемещаться по файлу. Системный вызов lseek() предназначен для установки позиции в файле, соответствующем файловому дескриптору. Помимо обновления позиции в файле, он больше ничего не делает и вообще никакие действия ввода-вывода не инициирует:

```
#include <sys/types.h>
#include <unistd.h>

off_t lseek (int fd, off_t pos, int origin);
```

Поведение lseek() зависит от аргумента `origin`, который может принимать следующие значения:

`SEEK_CUR`

Текущая позиция в файле `fd` увеличивается на значение аргумента `pos`, который может содержать отрицательное, положительное и нулевое значение. Если `pos` равен нулю, то возвращается значение текущей позиции в файле.

`SEEK_END`

В качестве текущей позиции в файле `fd` устанавливается значение, равное текущей длине файла плюс значение `pos`. Аргумент `pos` может содержать отрицательное, положительное и нулевое значение. Если `pos` равен нулю, то смещение устанавливается на конец файла.

`SEEK_SET`

Текущая позиция в файле `fd` приравнивается значению аргумента `pos`. Если `pos` равен нулю, то смещение устанавливается на начало файла.

В случае успеха вызов возвращает новую позицию в файле. В случае ошибки он возвращает значение `-1` и соответствующим образом устанавливает переменную `errno`.

Например, чтобы установить позицию в файле `fd` на значение 1825, используйте следующий код:

```
off_t ret;
```

```
ret = lseek (fd, (off_t) 1825, SEEK_SET);
if (ret == (off_t) -1)
    /* ошибка */
```

Этот код устанавливает позицию в файле `fd` на конец файла:

```
off_t ret;
```

```
ret = lseek (fd, 0, SEEK_END);
if (ret == (off_t) -1)
    /* ошибка */
```

Так как вызов lseek() возвращает обновленное значение позиции в файле, его можно использовать для поиска текущей позиции в файле, используя значение pos, равное нулю, и значение offset, равное SEEK_CUR:

```
int pos;
pos = lseek(fd, 0, SEEK_CUR);
if (pos == (off_t) -1)
    /* ошибка */
else
    /* 'pos' -- это текущая позиция в файле fd */
```

Чаще всего lseek() используют для поиска начала файла, конца файла или определения текущей позиции в файле, связанном с данным дескриптором файла.

Поиск за пределами конца файла

Вызов lseek() можно применять также для перемещения указателя файла за пределы файла. Например, следующий код переносит указатель на 1688 байт дальше, чем находится конец файла, соответствующий дескриптору fd:

```
int ret;
ret = lseek(fd, (off_t) 1688, SEEK_END);
if (ret == (off_t) -1)
    /* ошибка */
```

Сам по себе поиск дальше конца файла ничего не дает -- запрос на чтение с новой позиции в файле вернет значение EOF. Если же делается запрос на запись с этой новой позиции, то между старым концом файла и новым фрагментом создается новое пространство, которое заполняется нулями.

Это заполненное нулями пространство называется *дырой* (hole). В файловых системах типа Unix дыры не занимают место на диске. Это означает, что общий размер всех файлов в файловой системе может превышать физический размер диска. Файлы с дырами называются *разреженными файлами* (sparse file). Разреженные файлы экономят пространство и повышают производительность, так как манипулирование дырами не требует никакого физического ввода-вывода.

Запрос на чтение, сделанный в тот фрагмент файла, который является дырой, возвращает соответствующее число двоичных нулей.

Значения ошибок

В случае ошибки вызов lseek() возвращает значение -1 и присваивает переменной errno одно из следующих четырех значений:

EBADF

Указанный дескриптор файла не открыт.

EINVAL

Значение аргумента offset не равно SEEK_SET, SEEK_CUR или SEEK_END или результирующая позиция в файле меньше нуля. Очень неудачно, что одна

ошибка `EINVAL` представляет оба этих случая. Первая из них практически точно является программной ошибкой времени компиляции, тогда как вторая может представлять собой более коварную логическую ошибку времени выполнения.

EOVERFLOW

Результатирующее смещение не может быть представлено переменной типа `off_t`. Такое происходит только в 32-разрядных архитектурах. На данный момент позиция в файле *обновляется*; эта ошибка указывает только на то, что вернуть позицию невозможно.

ESPIPE

Указанный дескриптор файла связан с объектом, не поддерживающим поиск, таким, как конвейер, FIFO или сокет.

Ограничения

Максимальное значение позиции в файле ограничивается размером типа `off_t`. Большинство машинных архитектур определяют его как тип C `long`, размер которого в Linux всегда равен одному слову (обычно это размер регистров общего назначения машины). Внутри, однако, ядро хранит смещения, используя тип C `long`. Это не создает никаких проблем на 64-разрядных машинах, но означает, что на 32-разрядных машинах при выполнении относительных операций поиска могут возникать ошибки `EOVERFLOW`.

Позиционное чтение и запись

Вместо вызова `lseek()` можно использовать два существующих в Linux варианта системных вызовов — `read()` и `write()`, принимающих в качестве одного из параметров позицию в файле, с которой должны начинаться чтение или запись. После завершения работы они *не* обновляют позицию в файле.

Вариант вызова `read()` называется `pread()`:

```
#define _XOPEN_SOURCE 500  
  
#include <unistd.h>  
  
ssize_t pread (int fd, void *buf, size_t count, off_t pos);
```

Этот вызов считывает из файла, переданного при помощи дескриптора файла `fd`, до `count` байтов в буфер `buf`, начиная с позиции в файле `pos`.

Вариант вызова `write()` называется `pwrite()`:

```
#define _XOPEN_SOURCE 500  
  
#include <unistd.h>  
  
ssize_t pwrite (int fd, const void *buf, size_t count, off_t pos);
```

Этот вызов записывает в файл, переданный при помощи дескриптора файла `fd`, до `count` байтов из буфера `buf`, начиная с позиции в файле `pos`.

Поведение этих вызовов практически идентично их собратьям без буквы `r` в начале названия, за исключением того, что они полностью игнорируют текущую позицию в файле; вместо того чтобы использовать ее, они начинают работу с позиции, содержащейся в аргументе `pos`. Помимо этого, по завершении они не обновляют позицию в файле. Другими словами, вызовы `read()` и `write()` могут потенциально вредить работе позиционных вызовов, если перемешивать эти два вида в коде программы.

Оба позиционных вызова можно использовать только с дескрипторами файлов, поддерживающими поиск. Они обеспечивают семантику, схожую с применением перед вызовом `read()` или `write()` вызова `lseek()`, но с тремя различиями. Во-первых, эти вызовы проще в использовании, особенно при выполнении таких хитрых операций, как перемещение по файлу назад или случайным образом. Во-вторых, они не обновляют указатель в файле после завершения. Наконец, самое главное, они позволяют избегать условий состязания, которые могут возникать при использовании `lseek()`. Если несколько потоков выполнения одновременно работают с одним дескриптором файла, то сразу же после того, как первый поток сделает вызов `lseek()`, второй поток в той же программе может успеть обновить позицию в файле, и тогда, когда первый поток выполнит операцию считывания или записи, это будет сделано совсем в другом месте, а не в том, которое планировалось изначально. Избежать таких ситуаций можно, применив системные вызовы `pread()` и `pwrite()`.

Значения ошибок

В случае успеха оба вызова возвращают число считанных или записанных байтов. Если `pread()` возвращает значение 0, то это указывает на условие EOF; если `pwrite()` возвращает значение 0, то этот вызов ничего не записал в файл. В случае ошибки оба вызова возвращают значение `-1` и соответствующим образом устанавливают переменную `errno`. Вызов `pread()` может устанавливать любые значения `errno`, допустимые для вызовов `read()` и `lseek()`. Вызов `pwrite()` может устанавливать любые значения, допустимые для `write()` и `lseek()`.

Усечение файлов

Linux предоставляет два системных вызова, предназначенные для усечения файлов. Оба они определяются в разнообразных стандартах POSIX и являются в них обязательными (до той или иной степени). Определения этих вызовов:

```
#include <unistd.h>
#include <sys/types.h>

int ftruncate (int fd, off_t len);
```

и

```
#include <unistd.h>
#include <sys/types.h>

int truncate (const char *path, off_t len);
```

Оба системных вызова усекают указанный файл до длины `len`. Системный вызов `ftruncate()` работает на дескрипторе файла, который передается при помощи аргумента `fd`, и этот дескриптор должен быть открыт для записи. Системный вызов `truncate()` работает на имени файла, который передается в аргументе `path`, и указанный файл должен быть доступным для записи. Оба вызова в случае успеха возвращают значение 0. Если происходит ошибка, то они возвращают значение -1 и соответствующим образом устанавливают переменную `errno`.

Чаще всего эти системные вызовы применяются для уменьшения длины файла. В случае успеха новая длина файла становится равной значению, которое содержится в аргументе `len`. Данные, находившиеся между новым и старым концом файла, удаляются, и к ним больше нельзя обратиться при помощи запроса на чтение.

Эти функции также можно использовать для «усечения» файлов до длины, превышающей их первоначальный размер, — примерно то же самое делает комбинация поиска и записи, описанная ранее в разделе «Поиск за пределами конца файла». Добавленные байты заполняются нулями.

Ни одна из этих операций не обновляет текущую позицию в файле.

Например, у нас есть файл `pirate.txt` длиной 74 байт со следующим содержимым:

```
Edward Teach was a notorious English pirate.  
He was nicknamed Blackbeard.
```

В том же каталоге, где находится этот файл, выполним следующую программу:

```
#include <unistd.h>
#include <stdio.h>

int main( )
{
    int ret;

    ret = truncate ("./pirate.txt", 45);
    if (ret == -1) {
        perror ("truncate");
        return -1;
    }

    return 0.
}
```

В результате мы получим файл длиной 45 байт с таким содержимым:

```
Edward Teach was a notorious English pirate.
```

Мультиплексированный ввод-вывод

В приложениях часто возникает необходимость фиксироваться на нескольких файловых дескрипторах, жонглируя вводом-выводом между клавиатурой (вход `stdin`), взаимодействием процессов и несколькими файлами. Современные основанные на событиях приложения с графическим интерфейсом пользователя (*graphical user interface, GUI*) могут сражаться в своих главных циклах¹ буквально с сотнями ожидающих обработки событий.

Без помощи потоков выполнения, которые фактически обслуживают каждый дескриптор файла отдельно, один процесс не мог бы разумно справляться с фиксацией более чем на одном файловом дескрипторе одновременно. Работать с несколькими дескрипторами просто, если они всегда готовы к операции чтения или записи. Но как только встречается дескриптор файла, который еще невозможно использовать, — предположим, выполняется системный вызов `read()`, а данных для считывания еще нет, — процесс блокируется и не может обслужить другие файловые дескрипторы. Блокировка может продолжаться несколько секунд, делая приложение малопроизводительным и раздражая пользователя. А если для данного дескриптора файла доступные данные не появляются, то процесс может заблокироваться навсегда. Так как ввод-вывод для различных дескрипторов файла зачастую взаимосвязан — вспомните конвейеры, — вполне возможно, что один файловый дескриптор будет приходить в состояние готовности только после обслуживания другого. В частности, это может превращаться в большую проблему для сетевых приложений, которые одновременно открывают множество сокетов.

Представьте себе блокировку дескриптора файла, связанного с взаимодействием процессов, пока данные ожидают на входе `stdin`. Приложение не узнает, что с клавиатуры уже поступили какие-то данные, пока заблокированный файловый дескриптор в итоге не вернет данные; но что делать, если заблокированная операция никогда не возвращает?

Ранее в этой главе мы уже рассмотрели ввод-вывод в режиме без блокировки как одно из решений этой проблемы. Применяя ввод-вывод без блокировки, приложения могут создавать запросы ввода-вывода, которые вместо блокировки будут возвращать особое условие ошибки. Однако это решение неэффективно по двум причинам. Во-первых, процессу приходится постоянно выполнять операции ввода-вывода в каком-то произвольном порядке, ожидая, когда его открытые дескрипторы файла будут готовы для ввода-вывода. Это плохая конструкция программы. Во-вторых, было бы намного эффективнее, если бы программа могла засыпать, освобождая процессор для других задач, и просыпаться, только когда один или несколько файловых дескрипторов становятся доступными для ввода-вывода.

¹ Термин «главный цикл» (*mainloop*) должен быть знаком любому, кому приходилось писать GUI-приложения. Например, приложения GNOME используют главный цикл, предоставляемый GLib — их базовой библиотекой. Главный цикл позволяет отслеживать и отвечать на несколько событий из одной точки блокировки.

Познакомьтесь с мультиплексированным вводом-выводом.

Мультиплексированный ввод-вывод (multiplexed I/O) позволяет приложению одновременно фиксироваться на нескольких файловых дескрипторах и получать уведомления, когда один из них становится доступным для чтения или записи без блокировки. Таким образом, мультиплексированный ввод-вывод становится центральной точкой для приложений, сконструированных на основе приблизительно такого каркаса:

- 1) мультиплексированный ввод-вывод: сообщи мне, когда любой из этих дескрипторов файла будет готов для ввода-вывода;
- 2) заснуть до тех пор, пока один или несколько файловых дескрипторов не будут готовы;
- 3) проснуться: какие дескрипторы готовы?
- 4) обработать все готовые к вводу-выводу дескрипторы файлов без блокировки;
- 5) вернуться к шагу 1 и начать сначала.

В Linux предусмотрено три решения для мультиплексированного ввода-вывода: интерфейсы `select`, `poll` и `epoll`. Здесь мы рассмотрим два первых, а последний, относящийся к более сложным и уникальным для Linux решениям, мы обсудим в главе 4.

select()

Системный вызов `select()` предоставляет механизм для реализации синхронного мультиплексированного ввода-вывода:

```
#include <sys/time.h>
#include <sys/types.h>
#include <unistd.h>

int select (int n,
            fd_set *readfds,
            fd_set *writefds,
            fd_set *exceptfds,
            struct timeval *timeout);

FD_CLR(int fd, fd_set *set);
FD_ISSET(int fd, fd_set *set);
FD_SET(int fd, fd_set *set);
FD_ZERO(fd_set *set);
```

Вызов `select()` блокируется до тех пор, пока указанный дескриптор файла не становится готовым к выполнению ввода-вывода, или же до тех пор, пока не истекает необязательный интервал тайм-аута.

Наблюдаемые файловые дескрипторы разбиваются на три набора, в каждом из которых ожидается свое событие. Дескрипторы файла в наборе `readfds` наблюдаются в ожидании момента, когда данные станут доступными для чтения (то есть можно будет выполнить операцию считывания без блокировки). Деск-

рипторы файла в наборе writefds наблюдаются в ожидании момента, когда можно будет завершить без блокировки операцию записи. Наконец, файловые дескрипторы в наборе exceptfds наблюдаются с тем, чтобы отлавливать создание исключений или появление внеполосных данных (эти состояния распространяются только на сокеты). Каждый набор может содержать значение NULL, и в этом случае вызов select() не отслеживает соответствующее событие.

В случае успешного возвращения вызовом результата каждый набор модифицируется, и в нем остаются только дескрипторы файла, готовые к операции ввода-вывода, относящейся к этому набору. Например, предположим, что два дескриптора файла со значениями 7 и 9 помещены в набор readfds. Когда вызов возвращает результат, если 7 все еще находится в наборе, это означает, что соответствующий дескриптор файла готов к чтению без блокировки. Если значение 9 удалено из набора, вероятно, считать соответствующий файл без блокировки невозможно. (Я говорю «вероятно», потому что данные вполне могут стать доступными уже после того, как вызов завершится. В этом случае повторный вызов select() вернет файловый дескриптор как готовый к считыванию¹.)

Первый параметр, n, на единицу превышает максимальное во всех наборах значение дескриптора файла. Следовательно, вызывающий функцию select() должен проверить, у какого дескриптора максимальное значение, прибавить к нему единицу и передать полученное в качестве первого параметра.

Параметр timeout — это указатель на структуру timeval, которая определяется следующим образом:

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

struct timeval {
    long tv_sec;      /* секунды */
    long tv_usec;     /* микросекунды */
};
```

Если значение этого параметра не равно NULL, то вызов select() возвращает значение через tv_sec секунд и tv_usec микросекунд, даже если ни один дескриптор файла не готов к вводу-выводу. После возврата состояние этой структуры в различных системах Unix не определено, поэтому ее необходимо повторно инициализировать (вместе с наборами дескрипторов файла) перед каждым вызовом. Действительно, текущие версии Linux автоматически модифицируют этот параметр, устанавливая оставшееся время. Таким образом, если первоначально интервал тайм-аута был равен 5 секундам и готовый дескриптор файла появился через 3 секунды, то после возврата вызова в поле tv.tv_sec будет содержаться значение 2.

¹ Причина этого заключается в том, что вызовы select() и poll() запускаются уровнем (level-triggered), а не фронтом (edge-triggered). Вызов epoll(), о котором мы будем говорить в четвертой главе, может работать в любом из этих режимов. Запускаемая фронтом операция проще, но может пропускать события ввода-вывода, если об этом специально не заботиться.

Если в обоих полях структуры `tmeval` находится значение 0, то вызов возвращает результат немедленно, сообщая обо всех ожидающих событиях, существующих на момент вызова, но не дожидаясь последующих событий.

Манипулирование набором дескрипторов файла выполняется не напрямую, а через макросы-помощники. Это позволяет реализовывать наборы в различных системах Unix так, как удобно разработчикам. Обычно, однако, они имеют форму простых битовых массивов. `FD_ZERO` удаляет все дескрипторы файла из указанного набора. Этот вызов необходимо делать перед каждым обращением к `select()`:

```
fd_set writefds;
```

```
FD_ZERO(&writefds);
```

`FD_SET` добавляет дескриптор файла в указанный набор, а `FD_CLR` удаляет дескриптор из указанного набора:

```
FD_SET(fd, &writefds); /* добавить 'fd' в набор */
FD_CLR(fd, &writefds); /* удалить 'fd' из набора */
```

В хорошо продуманном коде никогда не должна возникать необходимость в вызове `FD_CLR`, и он применяется крайне редко, если вообще когда-либо такое случается.

`FD_ISSET` проверяет, является ли дескриптор файла частью указанного набора. Этот вызов возвращает не равное нулю целое число, если дескриптор находится в наборе, и 0, если это не так. `FD_ISSET` применяется после того, как `select()` возвращает результат, для проверки, что определенный дескриптор файла готов к действию:

```
if (FD_ISSET(fd, &readfds))
    /* дескриптор 'fd' готов к чтению без блокировки! */
```

Так как наборы дескрипторов файла создаются статически, они накладывают ограничение на максимальное число дескрипторов и максимальное значение файлового дескриптора, который может быть помещен в набор. Оба этих значения равны константе `FD_SETSIZE`, значение которой в Linux составляет 1024. Мы взглянем на последствия такого ограничения далее в этой главе.

Возвращаемые значения и коды ошибок

В случае успеха вызов `select()` возвращает число дескрипторов файла во всех трех наборах, готовых к вводу-выводу. Если вызову передавалось значение тайм-аута, то может быть возвращено значение 0. В случае ошибки вызов возвращает -1 и присваивает переменной `errno` одно из следующих значений:

`EBADF`

В одном из наборов содержится недопустимый дескриптор файла.

`EINTR`

Во время ожидания был получен сигнал, и вызов можно повторить.

EINVAL

Параметр *p* имеет отрицательное значение, или указанное значение тайм-аута недопустимо.

ENOMEM

Недостаточно памяти для выполнения запроса.

Пример использования вызова select()

Рассмотрим пример простой, но полностью функциональной программы, иллюстрирующей использование вызова `select()`. В этом примере ожидание ввода на `stdin` фиксируется на срок до 5 секунд. Так как здесь наблюдается только один дескриптор файла, фактически это не мультиплексированный ввод-вывод, но пример понятно демонстрирует, как использовать этот системный вызов:

```
#include <stdio.h>
#include <sys/time.h>
#include <sys/types.h>
#include <unistd.h>

#define TIMEOUT 5      /* установка тайм аута в секундах */
#define BUF_LEN 1024   /* длина буфера считывания в байтах */

int main (void)
{
    struct timeval tv;
    fd_set readfds;
    int ret;

    /* Ожидать ввода на stdin. */
    FD_ZERO(&readfds);
    FD_SET(STDIN_FILENO, &readfds);

    /* Ждать до пяти секунд */
    tv.tv_sec = TIMEOUT;
    tv.tv_usec = 0;

    /* Все в порядке, можно фиксировать! */
    ret = select (STDIN_FILENO + 1,
                  &readfds,
                  NULL,
                  NULL,
                  &tv);

    if (ret == -1) {
        perror ("select");
        return 1;
    } else if (!ret) {
        printf ("%d seconds elapsed \n", TIMEOUT);
        return 0;
    }

    /*
     * Наш дескриптор файла готов к считыванию?
    */
```

```

    /* (Он должен быть готов, так как мы передали
     * вызову только один дескриптор и вызов
     * вернул не нулевое значение, но эта проверка
     * нужна для самоуспокоения.)
    */
    if (FD_ISSET(STDIN_FILENO, &readfds)) {
        char buf[BUF_LEN+1];
        int len;

        /* гарантированно без блокировки */
        len = read (STDIN_FILENO, buf, BUF_LEN);
        if (len == -1) {
            perror ("read");
            return 1;
        }

        if (len) {
            buf[len] = '\0';
            printf ("read %s\n", buf);
        }
    }

    return 0;
}

fprintf (stderr, "Это не должно происходить!\n");
return 1;
}

```

Краткое засыпание при помощи select()

Так как вызов `select()` исторически более охотно реализуется в разнообразных системах Unix, чем механизмы засыпания меньше чем на секунду, очень часто для решения задачи засыпания применяется данный вызов. Вы просто передаете ему не равное NULL значение тайм-аута и значения NULL для всех трех наборов:

`struct timeval tv;`

```

tv.tv_sec = 0.
tv.tv_usec = 500.

```

```

/* заснуть на 500 микросекунд */
select (0, NULL, NULL, NULL, &tv);

```

Конечно же, в Linux есть интерфейсы для обеспечения засыпания на любые интервалы времени. Мы поговорим о них в десятой главе.

Вызов pselect()

Системный вызов `select()`, впервые представленный в 4.2BSD, очень распространен, но POSIX определяет собственное решение, `pselect()`, в стандарте POSIX 1003.1g-2000 и позднее в POSIX 1003.1-2001:

```

#define _XOPEN_SOURCE 600
#include <sys/select.h>

```

```

int pselect (int n,
             fd_set *readfds,
             fd_set *writefds,
             fd_set *exceptfds,
             const struct timespec *timeout,
             const sigset_t *sigmask);

FD_CLR(int fd, fd_set *set);
FD_ISSET(int fd, fd_set *set);
FD_SET(int fd, fd_set *set);
FD_ZERO(fd_set *set);

```

Между вызовами `select()` и `pselect()` три различия:

- 1) `pselect()` для параметра `timeout` использует структуру `timespec`, а не структуру `timeval`. В структуре `timespec` вы указываете секунды и наносекунды, а не секунды и микросекунды, что теоретически должно давать лучший контроль над длительностью тайм-аута. На практике, однако, ни один из этих вызовов не обеспечивает точности даже на уровне микросекунд;
- 2) при вызове `pselect()` значение параметра `timeout` не меняется. Следовательно, его не нужно повторно инициализировать перед последующими вызовами;
- 3) у системного вызова `select()` нет параметра `sigmask`. По отношению к сигналам, когда значение этого параметра равно `NULL`, `pselect()` работает так же, как и `select()`.

Структура `timespec` определяется следующим образом:

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

```

struct timespec {
    long tv_sec: /* секунды */
    long tv_nsec: /* наносекунды */
}:

```

Основная причина, почему в комплект инструментов Unix добавлен вызов `pselect()`, скрывается в параметре `sigmask`, при помощи которого была сделана попытка решить вопрос возникновения условий состязания между ожиданием дескрипторов файла и сигналами (сигналы подробно рассматриваются в главе 9). Предположим, что обработчик сигнала устанавливает глобальный флаг (как делает большинство из них) и процесс проверяет этот флаг перед тем, как выполнять вызов `select()`. Теперь предположим, что сигнал прибывает после проверки, но перед самим вызовом. Приложение может заблокироваться на бесконечное время и никогда не отреагировать на установленный флаг. `pselect()` устраняет эту проблему, позволяя приложению вызывать `pselect()` с набором блокируемых сигналов. Заблокированные сигналы не обрабатываются до тех пор, пока они не будут разблокированы. После того как `pselect()` возвратит значение, ядро восстановит старую маску сигналов. Серьезно об этом говорится в главе 9.

До появления ядра 2.6.16 `pselect()` в Linux реализовывался не как системный вызов, а как простая обертка вокруг `select()`, предоставляемая `glibc`. Эта

обертка минимизировала — но полностью не устранила — риск возникновения условия состязания. Когда появился настоящий системный вызов, состязание исчезло.

Несмотря на (относительно небольшие) усовершенствования в pselect(), в большинстве приложений все так же используется вызов select(), возможно, по привычке, а возможно, потому, что этот вызов обладает лучшей переносимостью.

Системный вызов poll()

Системный вызов poll() — это решение для мультиплексированного ввода-вывода из System V. Оно устраняет несколько недостатков вызова select(), хотя select() до сих пор используется очень часто (и снова, вероятнее всего, по привычке, но может быть и во имя переносимости):

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

```
int poll (struct pollfd *fds, unsigned int nfds, int timeout);
```

В отличие от select() с его тремя неэффективными наборами дескрипторов файлов на основе битовых масок, вызов poll() использует один массив структур nfds pollfd, на который указывает аргумент fds. Эта структура определяется следующим образом:

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

```
struct pollfd {
    int fd;          /* дескриптор файла */
    short events;   /* запрошенные события для наблюдения */
    short revents;  /* возвращаемые случившиеся события */
};
```

Каждая структура pollfd относится к одному наблюдаемому дескриптору файла. Можно передавать вызову несколько структур, чтобы poll() наблюдал за несколькими дескрипторами файла. Поле events в структуре представляет собой битовую маску наблюдаемых событий для данного дескриптора файла. Это поле устанавливает пользователь. Поле revents — это битовая маска событий, которые встретились для данного дескриптора файла. Его устанавливает ядро, когда вызов возвращает значение. Все события, запрошенные в поле events, могут быть возвращены в поле revents. Допустимы следующие значения:

POLLIN

Есть данные для считывания.

POLLRDNORM

Есть обычные данные для считывания.

POLLRDBAND

Есть приоритетные данные для считывания.

POLLPRI

Есть срочные данные для считывания.

POLLOUT

Операция записи не заблокируется.

POLLWRNORM

Операция записи обычных данных не заблокируется.

POLLWRBAND

Операция записи приоритетных данных не заблокируется.

PLLMSG

Доступно сообщение SIGPOLL.

Помимо этого, в поле `revents` могут возвращаться следующие значения:

POLLER

Ошибка в указанном дескрипторе файла.

POLLHUP

Зависшее событие для данного дескриптора файла.

POLLNVAL

Указанный дескриптор файла недопустим.

Эти события не будут иметь никакого смысла, если вставить их в поле `events`, так как они могут только возвращаться. При использовании вызова `poll()`, в отличие от `select()`, вам не нужно явно запрашивать сообщения об исключениях.

Сочетание `POLLIN | POLLPRI` эквивалентно событию считывания для вызова `select()`, а `POLLOUT | POLLWRBAND` — событию записи для вызова `select()`. Значение `POLLIN` эквивалентно сочетанию `POLLRDNORM | POLLRDBAND`, а `POLLOUT` эквивалентно `POLLWRNORM`.

Например, чтобы наблюдать файловый дескриптор с целью отлавливания моментов, когда он становится доступен для чтения и записи, нужно установить для параметра `events` значение `POLLIN | POLLOUT`. После возврата вызовом результата нужно проверить эти флаги в поле `revents` структуры, соответствующей данному дескриптору файла. Если установлен флаг `POLLIN`, то дескриптор доступен для чтения без блокировки. Если установлен флаг `POLLOUT`, то дескриптор доступен для записи без блокировки. Эти флаги не взаимоисключающие: могут быть установлены оба, указывая, что и операции чтения, и операции записи успешно вернут результат, не блокируясь на этом дескрипторе файла.

При помощи параметра `timeout` определяется длина интервала ожидания в миллисекундах. Когда этот интервал истекает, вызов возвращает значение, даже если дескрипторы, готовые к вводу-выводу, не обнаружены. Отрицательное значение параметра указывает на бесконечный интервал тайм-аута. Значение 0 заставляет вызов вернуть результат немедленно, перечислив все дескрипторы файлов, ожидающие ввода-вывода, но не дожидаясь новых событий. Таким образом, вызов `poll()` оправдывает свое имя, выполняя опрос ровно один раз и немедленно возвращая результат.

Возвращаемые значения и коды ошибок

В случае успеха `poll()` возвращает число дескрипторов файла, в структурах которых поле `revents` не равно нулю. Вызов возвращает значение 0, если интервал тайм-аута истекает до того, как встречается какое-либо событие. В случае ошибки вызов возвращает значение -1 и присваивает переменной `errno` одно из следующих значений:

`EBADF`

В одной или нескольких структурах содержится недопустимый дескриптор файла.

`EFAULT`

Указатель на `fds` указывает за пределы адресного пространства процесса.

`EINTR`

Сигнал был получен до того, как произошло любое из запрошенных событий. Вызов можно выполнить повторно.

`EINVAL`

Значение параметра `nfds` превышает значение `RLIMIT_NOFILE`.

`ENOMEM`

Недостаточно памяти для выполнения запроса.

Пример использования вызова `poll()`

Давайте взглянем на пример программы, в которой вызов `poll()` применяется для одновременной проверки, не будет ли блокироваться считывание из `stdin` и запись в `stdout`:

```
#include <stdio.h>
#include <unistd.h>
#include <sys/poll.h>

#define TIMEOUT 5          /* тайм-аут опроса в секундах */

int main (void)
{
    struct pollfd fds[2];
    int ret.

    /* ожидать ввод на stdin */
    fds[0].fd = STDIN_FILENO;
    fds[0].events = POLLIN;

    /* наблюдать stdout – возможна ли запись
     * (практически всегда это так)
     */
    fds[1].fd = STDOUT_FILENO;
    fds[1].events = POLLOUT;

    /* Все установлено, можно фиксировать! */
    ret = poll (fds, 2, TIMEOUT * 1000);
```

```
if (ret == -1) {
    perror ("poll").
    return 1.
}

if (!ret) {
    printf ("%d seconds elapsed.\n", TIMEOUT);
    return 0;
}

if (fds[0].revents & POLLIN)
    printf ("stdin is readable\n").

if (fds[1].revents & POLLOUT)
    printf ("stdout is writable\n").

return 0.
}
```

Если выполнить эту программу, то, вполне понятно, результат получится следующим:

```
$ ./poll
stdout is writable
```

Если выполнить ее повторно, но в этот раз перенаправить на вход какой-то файл, то можно увидеть оба события:

```
$ ./poll < ode_to_my_parrot.txt
stdin is readable
stdout is writable
```

При использовании вызова `poll()` в реальном приложении не нужно заново создавать структуры `pollfd` при каждом обращении к нему. Одну и ту же структуру можно передавать вызову повторно; ядро обрабатывает обнуление поля `revents`, как необходимо.

Системный вызов `ppoll()`

В Linux есть еще один брат вызова `poll()` — `ppoll()`, который отличается от него так же, как `pselect()` от `select()`. В отличие от `pselect()`, однако, `ppoll()` представляет собой уникальный для Linux интерфейс:

```
#define _GNU_SOURCE
#include <sys/poll.h>

int ppoll (struct pollfd *fds,
           nfds_t nfds,
           const struct timespec *timeout,
           const sigset_t *sigmask);
```

Как и для `pselect()`, параметр `timeout` позволяет указать значение тайм-аута в секундах и наносекундах, а параметр `sigmask` содержит набор сигналов, которые вызов должен ожидать.

Сравнение poll() и select()

Хотя эти два вызова выполняют одну и ту же базовую функцию, системный вызов `poll()` превосходит по своим возможностям вызов `select()`, и вот почему:

- вызов `poll()` не требует, чтобы пользователь вычислял и передавал ему в качестве параметра увеличенное на единицу максимальное значение файлового дескриптора в наборе;
- `poll()` эффективнее для дескрипторов файла, имеющих большое значение. Представьте себе наблюдение при помощи `select()` за одним дескриптором файла со значением 900 – ядру придется проверять каждый бит переданного вызова набора, вплоть до девятисотого;
- размер наборов файловых дескрипторов для `select()` определяется статически, что заставляет решаться на компромисс: либо они слишком маленькие, что ограничивает максимальное значение дескриптора файла, за которым может наблюдать `select()`, либо они неэффективные. Операции на больших битовых масках неэффективны, особенно если неизвестно, разрежены ли они¹. Используя `poll()`, вы можете создавать массив в точности того размера, который необходим. Нужно наблюдать только за одним элементом? Просто передайте единственную структуру;
- при использовании вызова `select()`, когда он возвращает результат, наборы дескрипторов файлов меняются, поэтому перед каждым последующим вызовом их необходимо заново инициализировать. Системный вызов `poll()` отделяет вход (поле `events`) от выхода (поле `revents`), позволяя повторно использовать массив без изменений;
- параметр `timeout` вызова `select()` имеет неопределенное значение после того, как вызов возвращает результат. В переносимом коде его нужно заново инициализировать. Однако это не проблема для `pselect()`.

У системного вызова `select()`, однако, есть и пара преимуществ:

- вызов `select()` лучше переносится, так как некоторые системы Unix не поддерживают вызов `poll()`;
- `select()` обеспечивает большую точность тайм-аута: до микросекунд. И `poll()`, и `pselect()` теоретически должны обеспечивать точность до наносекунд, но на практике ни один из этих вызовов не достигает даже точности на уровне микросекунд.

Превосходит оба вызова, и `poll()`, и `select()`, интерфейс `epoll` – уникальное для Linux решение мультиплексирования ввода-вывода, с которым мы познакомимся в главе 4.

¹ Если битовая маска разрежена, то каждое слово, составляющее маску, можно проверить по отношению к нулю; только если эта операция вернет ложное значение, нужно будет проверить каждый бит. Эта операция, однако, непроизводительна, если битовая маска плотная.

Внутреннее устройство ядра

В этом разделе мы узнаем, как ядро Linux реализует ввод-вывод, сконцентрировавшись на трех основных подсистемах ядра: виртуальной файловой системе (virtual filesystem, VFS), страничном кэше (page cache) и страничной отложенной записи (page writeback). Работая вместе, эти три подсистемы помогают делать ввод-вывод бесшовным, эффективным и оптимальным.

ПРИМЕЧАНИЕ

В четвертой главе мы познакомимся с четвертой подсистемой, планировщиком ввода-вывода.

Виртуальная файловая система

Виртуальная файловая система, которую также иногда называют виртуальным файловым коммутатором (virtual file switch), – это механизм абстракции, позволяющий ядру Linux вызывать функции файловой системы и манипулировать данными файловой системы, не зная о том, какой конкретно тип файловой системы используется, и даже не заботясь об этом.

VFS реализует эту абстракцию благодаря общей модели файлов (common file model), представляющей собой основу всех файловых систем в Linux. Через указатели функций и разнообразные объектно-ориентированные практики¹ общая модель файлов предоставляет каркас, которому должны соответствовать файловые системы в ядре Linux. Это позволяет VFS отправлять в файловую систему генерализованные вызовы. У этого каркаса есть «якоря» для поддержки считывания, создания ссылок, синхронизации и других действий, и каждая файловая система, используя их, регистрирует собственные функции для обработки операций.

Данный подход обеспечивает определенную схожесть между файловыми системами. Например, VFS работает в терминах структур inode, суперблоков и записей каталога. Файловой системе, не относящейся к Unix и, вероятно, лишенной таких концепций Unix, как структура inode, просто приходится справляться с этим. И они справляются: Linux поддерживает, например, файловые системы FAT и NTFS безо всяких проблем.

У VFS огромное количество преимуществ. Единственный системный вызов умеет считывать данные из любых файловых систем и с любых носителей; единственная утилита может копировать данные из любой файловой системы в любую другую. Все файловые системы поддерживают одни и те же концепции, одни и те же интерфейсы и одни и те же вызовы. Все просто-напросто работает, и работает хорошо.

Когда приложение выполняет системный вызов `read()`, он проделывает интересное путешествие. Библиотека С предоставляет определения системного вызова, которые во время компиляции преобразуются в соответствующие операторы ловушки. Как только процесс из пользовательского пространства

¹ Да, в языке С.

отлавливается ядром, передается через обработчик системного вызова и попадает к системному вызову `read()`, ядро выясняет, какой объект *служит базой* (*back*) для указанного файлового дескриптора. После этого ядро вызывает функцию считывания, связанную именно с этим базовым объектом. Для файловых систем эта функция является частью кода файловой системы. После этого функция выполняет свою задачу, например физически считывает данные из файловой системы, и возвращает данные вызову `read()` из пользовательского пространства, который затем возвращает данные обработчику системного вызова, который, в свою очередь, копирует данные обратно в пользовательское пространство, где системный вызов `read()` возвращает результат и процесс продолжает исполнение.

VFS дает очень важные преимущества системным программистам. Им не приходится заботиться о том, в файловой системе какого типа или на каком носителе находится файл. Генерализованные системные вызовы — `read()`, `write()` и т. д. — умеют манипулировать файлами в любых поддерживаемых файловых системах и на любых поддерживаемых носителях.

Страницный кэш

Страницный кэш — это хранилище в памяти, где содержатся данные из дисковой файловой системы, к которым недавно происходило обращение. Доступ к диску выполняется ужасно медленно, особенно по сравнению с быстрым действием современных процессоров. Хранение запрошенных данных в памяти позволяет ядру удовлетворять последующие запросы, направленные на те же самые данные в памяти, избегая повторного обращения к диску.

Страницный кэш использует концепцию *сосредоточенности во времени* (*temporal locality*), один из типов локальности ссылок, которая утверждает, что если к ресурсу обратились в какой-то точке времени, то высока вероятность того, что обращение повторится в ближайшем будущем. Затраты памяти, которая расходуется на кэшированные данные при первом обращении, таким образом, окупают себя и предотвращают будущие дорогостоящие обращения к диску.

Страницный кэш — это первое место, где ядро начинает поиск данных из файловой системы. Ядро вызывает подсистему памяти для считывания данных с диска только в том случае, если их не удается обнаружить в кэше. Следовательно, когда элемент данных считывается впервые, он переносится с диска в страницный кэш и возвращается приложению из кэша. Если к тем же данным вскоре обращаются повторно, то они просто возвращаются из кэша. Все операции прозрачно выполняются через страницный кэш, гарантируя, что все данные всегда актуальны и верны.

В Linux реализуется страницный кэш динамического размера. По мере того как операции ввода-вывода передают в память все больше и больше данных, размер страницного кэша увеличивается, потребляя всю свободную память. Если страницный кэш в итоге разрастается на всю имеющуюся свободную память, то, после того как запрашивается операция, требующая дополнительной памяти, страницный кэш *обрезается* (*prune*) и наименее часто используемые

страницы освобождаются, чтобы дать место новым данным. Такое обрезание выполняется бесплатно и автоматически. Кэш динамического размера позволяет Linux действовать всю память в системе и кэшировать максимальный объем данных.

Однако часто бывают ситуации, когда лучше скидывать на диск редко используемые фрагменты данных, а не обрезать необходимые части страничного кэша, которые с большой вероятностью вернутся обратно в память буквально при следующем же запросе на считывание (*подкачка (swapping)*) позволяет ядру хранить данные на диске, чтобы поддерживать отпечаток памяти большего размера, чем объем доступной на машине оперативной памяти). Ядро Linux балансирует подкачку данных с обрезанием страничного кэша (и других резервов памяти) при помощи эвристического анализа. Путем такого анализа оно принимает решения, когда стоит не обрезать страничный кэш, а сбрасывать данные на диск, особенно если сбрасываемые данные в данный момент не используются.

Балансировка подкачки и кэширования настраивается в файле `/proc/sys/vm/swappiness`. В этом файле может содержаться значение от 0 до 100, а значение по умолчанию равно 60. Чем большее значение, тем выше предпочтение хранению страничного кэша в памяти и тем охотнее ядро выполняет подкачку. Меньшее значение позволяет ядру чаще обрезать страничный кэш, не подкачивая данные.

Еще одна форма локальности ссылок — это *сосредоточенность последовательности (sequential locality)*. Это означает, что к данным часто обращаются последовательно. Для того чтобы воспользоваться преимуществом данного принципа, ядро также реализует *опережающее чтение (readahead)* страничного кэша. Опережающее чтение — это акт считывания дополнительных данных с диска в страничный кэш, который выполняется после каждого запроса на чтение. Фактически ядро заранее считывает данные, которые могут потребоваться в будущем. Когда ядро читает фрагмент данных с диска, оно заодно считывает еще один-два фрагмента, следующие за первым. Операция считывания большей последовательности фрагментов данных за раз достаточно эффективна, так как обычно при этом не возникает необходимости в поиске на диске. Помимо этого, ядро может удовлетворять запрос на опережающее чтение, пока процесс обрабатывает первый фрагмент данных. Если процесс после этого передаст новый запрос на считывание следующего фрагмента, как это часто бывает, то ядро просто отдаст ему данные, считанные заранее, не запрашивая считывание их с диска.

Как и в случае со страничным кэшем, ядро управляет опережающим считыванием динамически. Если оно замечает, что процесс постоянно использует данные от предыдущего запроса на опережающее считывание, то оно увеличивает окно опережающего считывания, позволяя считывать заранее большие данных. Размер окна может составлять от 16 до 128 Кбайт. И наоборот, если ядро замечает, что опережающее считывание не дает попаданий, то есть приложение ищет данные в разных местах файла, а не читает их последовательно, то оно может полностью отключать эту возможность.

Работа страничного кэша должна быть прозрачной. В целом, системный программист не может оптимизировать свой код, чтобы лучше пользоваться преимуществами страничного кэша, за исключением, пожалуй, самостоятельной реализации подобного кэша в пользовательском пространстве. Эффективный код – это все, что требуется для того, чтобы наилучшим образом применять страничный кэш. Оптимизация опережающего считывания, однако, возможна. Последовательный файловый ввод-вывод всегда предпочтительнее, чем случайный доступ, хотя и не всегда возможен.

Страницчная отложенная запись

Как обсуждалось ранее в разделе «Поведение вызова write()», ядро откладывает фактическую запись данных на диск, помещая их в буферы. Когда процесс отправляет запрос на запись, данные копируются в буфер, который помечается как грязный, что указывает на то, что копия данных в памяти новее, чем копия их на диске. После этого запрос на запись просто возвращает результат. Если делается другой запрос на запись в тот же фрагмент файла, то буфер обновляется и в него помещаются новые данные. Запросы на запись в других местах того же файла приводят к созданию новых грязных буферов.

В конце концов, возникает необходимость зафиксировать содержимое грязных буферов на диске, чтобы синхронизировать файлы на диске с данными в памяти. Это называется отложенной записью. Синхронизация происходит в двух случаях:

- когда количество свободной памяти становится меньше настраиваемого порогового значения, содержимое грязных буферов записывается обратно на диск, чтобы очищенные буферы можно было удалить, освободив память;
- когда возраст грязного буфера превышает настраиваемое пороговое значение, его содержимое записывается обратно на диск. Это не позволяет данным перманентно оставаться грязными.

Отложенная запись выполняется только набором потоков ядра, которые называются потоками pdflush (вероятно, это расшифровывается как page dirty flush – сброс грязных страниц – но кто знает?). Когда выполняется одно из указанных выше условий, потоки pdflush просыпаются и начинают фиксировать данные из грязных буферов на диск, продолжая делать это до тех пор, пока оба условия не станут ложными.

Одновременно отложенную запись могут выполнять несколько потоков pdflush. Благодаря этому ядро пользуется преимуществами параллелизма, а также реализуется техника *избегания скученности* (congestion avoidance). Эта техника направлена на то, чтобы во время ожидания какого-то одного блочного устройства данные не скапливались. Если присутствуют грязные буферы для разных блочных устройств, то несколько потоков pdflush работают для того, чтобы полностью задействовать ресурсы каждого из этих устройств. Благодаря этому устраняется дефект более ранних версий ядра: предшественник потоков pdflush (единственный поток bdfflush) мог тратить все свое время на ожидание,

пока освободится одно блочное устройство, хотя в то же самое время другие блочные устройства могли бездействовать. На современной машине ядро Linux способно одновременно подпитывать очень большое число дисков.

Буферы в ядре представляются структурой данных `buffer_head`. При помощи этой структуры данных отслеживаются разнообразные метаданные, связанные с буфером, например, проверяется, находится буфер в чистом или грязном состоянии. Также она содержит указатель на фактические данные. Фактические данные хранятся в страничном кэше, что объединяет буферную подсистему со страничным кэшем.

В ранних версиях ядра Linux, до версии 2.3, буферная подсистема не была связана со страничным кэшем, то есть существовало два разных типа кэша: страничный и буферный. Таким образом, данные могли одновременно находиться и в буферном кэше (в грязном буфере), и в страничном кэше (в виде кэшированных данных). Конечно же, синхронизация двух отдельных кэшей требовала определенных усилий. Унифицированный страничный кэш, впервые представленный в ядре Linux 2.4, стал очень приятным усовершенствованием.

Отложенная запись и подсистема буферов в Linux обеспечивают возможность быстрой записи, но ценой утери данных в случае сбоя питания. Для того чтобы избежать такого риска, в критически важных приложениях (а также приложениях программистов, страдающих паранойей) можно применять синхронизированный ввод-вывод (о котором говорилось ранее в этой главе).

Заключение

В этой главе мы рассмотрели основу системного программирования Linux файловый ввод-вывод. В такой системе, как Linux, которая пытается представлять все, что только возможно, в виде файлов, очень важно знать, как открывать, читать, записывать и закрывать файлы. Все эти операции представляют собой классику Unix и определены во многих стандартах.

Следующая глава посвящена буферизованному вводу-выводу и интерфейсам ввода-вывода стандартной библиотеки C. Стандартная библиотека C существует не только для удобства: буферизация ввода-вывода в пользовательском пространстве обеспечивает огромный выигрыш в производительности.

3 Буферизованный ввод-вывод

Вспомните – в первой главе говорилось, что блок (абстракция файловой системы) является *lingua franca* ввода-вывода, потому что все дисковые операции выполняются в терминах блоков. Следовательно, производительность ввода-вывода оптимальна тогда, когда запросы делаются в пределах целого количества блоков.

Проблема снижения производительности обостряется при увеличении числа системных вызовов, которым требуется, скажем, читать по одному байту 1024 раза, вместо того чтобы за один раз считать все 1024 байта. Даже производительность последовательности операций над фрагментами данных, размер которых больше блока, может быть ниже оптимальной, если этот размер не кратен размеру блока. Например, если в данной файловой системе размер блока составляет 1 Кбайт, то операции над фрагментами размером 1130 байтов могут выполняться медленнее, чем операции над фрагментами размером 1024 байта.

Ввод-вывод с пользовательским буфером

Программы, которым приходится выполнять множество небольших запросов ввода-вывода к обычным файлам, зачастую используют ввод-вывод с пользовательским буфером. Это означает, что буферизация выполняется в пользовательском пространстве либо вручную самим приложением, либо прозрачно в библиотеке, в противоположность буферизации, которую осуществляет ядро. Как говорилось во второй главе, в целях улучшения производительности системы, ядро буферизует данные в глубине системы, откладывая операции записи, объединяя следующие друг за другом запросы ввода-вывода и выполняя упреждающее считывание. Пользовательская буферизация также нацелена на повышение производительности, хотя и реализована по-другому.

Рассмотрим пример с программой dd из пользовательского пространства:

```
dd bs=1 count=2097152 if=/dev/zero of=pirate
```

Так как здесь значение аргумента bs равно 1, эта команда копирует два мегабайта данных с устройства /dev/zero (виртуальное устройство, обеспечивающее

бесконечный поток нулей) в файл `pirate`, делая это с использованием 2 097 152 однобайтовых фрагментов. Таким образом, будет выполнено около двух миллионов операций чтения и записи — по одному байту за раз.

Теперь рассмотрим ту же команду копирования двух мегабайтов, но блоками по 1024 байта:

```
dd bs=1024 count=2048 if=/dev/zero of=pirate
```

Эта операция копирует те же два мегабайта нулей в тот же файл, но делает в 1024 раза меньше операций чтения и записи. Как можно видеть из табл. 3.1, производительность повышается во много раз. В этой таблице я зафиксировал время, затраченное четырьмя командами `dd` на выполнение одинаковой работы с разными размерами блока. Реальное время — это общее потраченное время, которое я засек по своим часам; пользовательское время — это время, в течение которого код программы выполнялся в пользовательском пространстве; системное время — это время, потраченное на выполнение процессом системных вызовов в пространстве ядра.

Таблица 3.1. Влияние размера блока на производительность

Размер блока, байты	Реальное время, секунды	Пользовательское время, секунды	Системное время, секунды
1	18,707	1,118	17,549
1024	0,025	0,002	0,023
1130	0,035	0,002	0,027

Использование фрагментов размером 1024 байта дает *огромный* прирост производительности по сравнению с однобайтовыми фрагментами. Однако таблица также демонстрирует, что использование фрагментов большего размера, приводящее к уменьшению количества системных вызовов, может понижать производительность, если размер фрагментов не кратен размеру дискового блока. Несмотря на меньшее число системных вызовов, операции с фрагментами по 1130 байт требуют создания невыровненных запросов и, таким образом, менее эффективны, чем операции на фрагментах по 1024 байта.

Для того чтобы пользоваться преимуществами зависимости производительности от размера фрагментов данных, необходимо заранее знать вероятный размер физического блока. Исходя из показанной выше таблицы, можно предположить, что размер блока, скорее всего, равен 1024 байт, величине, кратной 1024, или делителю 1024. В случае с `/dev/zero` размер блока в действительности равен 4096 байт.

Размер блока

На практике чаще всего используются блоки размером 512, 1024, 2048 и 4096 байт.

Как демонстрирует табл. 3.1, большого выигрыша в производительности можно добиться, выполняя операции над фрагментами данных, размер которых кратен размеру блока или является делителем этого значения. Причина заключается в том, что ядро и аппаратное обеспечение говорят в терминах блоков. Таким образом, использование точного размера блока или значения, аккуратно вписывающегося внутрь блока, гарантирует, что запросы ввода-вывода будут выравниваться по размеру блока, а внутри ядра не будет выполняться ненужная работа.

Для того чтобы понять, чему равен размер блока для конкретного устройства, нужно просто применить системный вызов `stat()` (подробнее о нем в главе 7) или команду `stat(1)`. Выясняется, однако, что обычно нет необходимости знать точный размер блока.

Основная цель при выборе размера фрагментов данных для операций ввода-вывода — избегать нетипичных значений, таких, как 1130. Никогда в истории Unix не применялись блоки размером 1130 байт, и если вы выберете это значение для своих операций, то после первого же запроса выравнивание ввода-вывода будет нарушено. Использование любого значения, кратного размеру блока, либо значения, которое является одним из делителей размера блока, однако, помогает избежать появления невыровненных запросов. Таким образом, пока выбранный вами размер не нарушает выравнивание по блокам, производительность остается на высоком уровне. Кратные размеру блока значения просто будут уменьшать число системных вызовов.

Другими словами, проще всего выполнять операции ввода-вывода, используя большой размер буфера, кратный типичным размерам блока. Отлично будут работать фрагменты размером 4096 или 8192 байта.

Конечно же, проблема заключается в том, что программы редко работают в терминах блоков. Программы работают с полями, строками и отдельными символами, а не с такими абстракциями, как блоки. Как упоминалось ранее, для исправления этой ситуации в программах может применяться ввод-вывод с пользовательской буферизацией. Когда данные записываются, они сохраняются в буфер внутри адресного пространства программы. Когда буфер достигает определенного размера — *размера буфера* (*buffer size*), все его содержимое записывается на диск одной операцией записи. Схожим образом данныечитываются фрагментами, размер которых равен размеру буфера, выровненными по размеру блока. Когда приложение отправляет запрос на чтение фрагментов, не совпадающих по размеру с буфером, они выдаются из буфера по частям. В конечном итоге, когда буфер опустошается, в него считывается очередной большой кусок данных размером с буфер. С правильно выбранным размером буфера можно достигать огромного выигрыша в производительности.

Можно реализовывать пользовательскую буферизацию в своих программах вручную. И это действительно делается во многих критически важных приложениях. Однако большая часть программ использует популярную библиотеку стандартного ввода-вывода (являющуюся частью стандартной библиотеки C), которая предоставляет надежное и гибкое решение для пользовательской буферизации.

Стандартный ввод-вывод

Стандартная библиотека С предоставляет библиотеку стандартного ввода-вывода (зачастую называемую просто `stdio`), которая, в свою очередь, предоставляет независимое от платформы решение пользовательской буферизации. Библиотеку стандартного ввода-вывода использовать легко, но это тем не менее очень мощный инструмент.

В отличие от таких языков программирования, как FORTRAN, в языке С отсутствует встроенная поддержка более сложной функциональности, чем управление потоком данных, арифметические действия и т. п., и, определенно, в этом языке нет внутренней поддержки ввода-вывода. По мере того как язык программирования С развивался, пользователи разрабатывали стандартные наборы процедур, обеспечивающих основную функциональность, например манипулирование строками, математические процедуры, функции времени и даты и операции ввода-вывода. Эта функциональность со временем совершенствовалась, а после утверждения в 1989 году стандарта ANSI C (C89) они в конечном итоге были формализованы в виде стандартной библиотеки С. Хотя в стандартах C95 и C99 были добавлены некоторые новые интерфейсы, стандартная библиотека С остается практически в неприкосновенном состоянии с момента появления на свет в 1989 году.

Далее в этой главе обсуждаются ввод-вывод с пользовательской буферизацией и его связь с файловым вводом-выводом, его реализация в стандартной библиотеке С, то есть открытие, закрытие, чтение и запись файлов с использованием стандартной библиотеки С. Будет ли приложение применять стандартный ввод-вывод, использовать собственное решение пользовательской буферизации или прямые системные вызовы — это решение разработчики должны принимать, тщательно взвешивая все требования и задачи каждого приложения.

Стандарты С всегда оставляют какие-то детали на усмотрение разработчиков, создающих конкретные реализации, и в этих реализациях зачастую добавляются какие-то новые функции. В этой главе, как и в остальных главах книги, рассматриваются интерфейсы и поведение системы в том виде, как оно реализовано в библиотеке `glibc` в современной системе Linux. Если Linux в чем-то отличается от базового стандарта, об это упоминается отдельно.

Указатели файла

Стандартные процедуры ввода-вывода не работают с дескрипторами файла, а используют собственный уникальный идентификатор, известный как *указатель файла* (*file pointer*). Внутри библиотеки С указатель файла отображается на дескриптор файла. Указатель файла представляется указателем на оператор описания типа `FILE`, определенный в `<stdio.h>`.

На языке стандартного ввода-вывода открытый файл называется *потоком данных* (*stream*). Потоки данных могут открываться для чтения (входные потоки), записи (выходные потоки) или для обеих операций (потоки ввода-вывода).

Открытие файлов

Файлы открываются для чтения или записи при помощи функции `fopen()`:

```
#include <stdio.h>
```

```
FILE * fopen (const char *path, const char *mode);
```

Эта функция открывает файл `path` согласно указанным режимам и связывает с ним новый поток данных.

Режимы

Аргумент `mode` описывает, каким образом нужно открыть указанный файл. Он может содержать в качестве значения одну из следующих строк:

`r`

Открыть файл для чтения. Указатель потока устанавливается на начало файла.

`r+`

Открыть файл для чтения и записи. Указатель потока устанавливается на начало файла.

`w`

Открыть файл для записи. Если файл существует, то он усекается до нулевой длины. Если файл не существует, он создается. Указатель потока устанавливается на начало файла.

`w+`

Открыть файл для чтения и записи. Если файл существует, то он усекается до нулевой длины. Если файл не существует, он создается. Указатель потока устанавливается на начало файла.

`a`

Открыть файл для записи в режиме присоединения. Если файл не существует, он создается. Указатель потока устанавливается на конец файла. Любые операции записи присоединяют данные к файлу.

`a+`

Открыть файл для чтения и записи в режиме присоединения. Если файл не существует, он создается. Указатель потока устанавливается на конец файла. Любые операции записи присоединяют данные к файлу.

ПРИМЕЧАНИЕ

Строка режима также может содержать символ `b`, хотя это значение всегда игнорируется в Linux. Некоторые операционные системы по-разному обрабатывают текстовые и двоичные файлы, и режим `b` обозначает открытие файла в двоичном режиме. Linux, как и все удовлетворяющие стандарту POSIX системы, обрабатывает текстовые и двоичные файлы одинаково.

В случае успеха `fopen()` возвращает действительный указатель FILE. В случае неудачи функция возвращает значение NULL и соответствующим образом устанавливает переменную errno.

Например, в следующем коде файл `etc/manifest` открывается для чтения и связывается с потоком `stream`:

```
FILE *stream;
```

```
stream = fopen ("/etc/manifest", "r").  
if (!stream)  
/* ошибка */
```

Открытие потока данных при помощи дескриптора файла

Функция `fdopen()` преобразовывает уже открытый дескриптор файла (fd) в поток данных:

```
#include <stdio.h>
```

```
FILE * fdopen (int fd, const char *mode);
```

Возможные режимы — те же, что и для функции `fopen()`, и они должны соответствовать режимам, с которыми первоначально открывался файл. Можно использовать режимы `w` и `w+`, но при этом файл усекаться не будет. Указатель потока устанавливается на позицию в файле, соответствующую данному дескриптору файла.

После того как дескриптор файла преобразуется в поток данных, ввод-вывод не следует выполнять напрямую на дескрипторе. Это, однако, допустимо. Обратите внимание, что дескриптор файла не дублируется, а просто связывается с новым потоком. Если закрыть поток, то будет закрыт и дескриптор файла.

В случае успеха `fdopen()` возвращает действительный указатель файла, в случае неудачи вызов возвращает значение NULL.

Например, в следующем коде файл `/home/kidd/map.txt` открывается системным вызовом `open()`, а затем при помощи связанного с этим файлом дескриптора создается поток:

```
FILE *stream.  
int fd.
```

```
fd = open (" /home/kidd/map.txt", O_RDONLY).  
if (fd == -1)  
/* ошибка */  
  
stream = fdopen (fd, "r").  
if (!stream)  
/* ошибка */
```

Закрытие потоков данных

Функция `fclose()` закрывает указанный поток данных:

```
#include <stdio.h>
```

```
int fclose (FILE *stream);
```

В первую очередь сбрасываются на диск все буферизованные, но еще не записанные данные. В случае успеха функция `fclose()` возвращает значение 0. В случае ошибки она возвращает EOF и соответствующим образом устанавливает переменную `errno`.

Закрытие всех потоков данных

Функция `fcloseall()` закрывает все потоки, связанные с текущим процессом, включая стандартный вход, стандартный выход и стандартный поток ошибок:

```
#define _GNU_SOURCE
```

```
#include <stdio.h>
```

```
int fcloseall (void);
```

Перед закрытием все потоки сбрасываются. Функция всегда возвращает значение 0 — она уникальна для Linux.

Чтение из потока данных

В стандартной библиотеке С реализовано несколько функций чтения из открытого потока данных, от самых распространенных до известных лишь посвященным. В этом разделе мы рассмотрим три наиболее популярных подхода к чтению: чтение по одному символу за раз, чтение всей строки за раз и чтение двоичных данных. Для того чтобы из потока можно было читать данные, его нужно открыть в качестве входного потока в подходящем режиме, то есть любом допустимом режиме, за исключением режимов `w` и `a`.

Чтение по одному символу за раз

Очень часто идеальным шаблоном ввода-вывода становится простое чтение по одному символу за раз. Функция `fgetc()` считывает один символ из потока:

```
#include <stdio.h>
```

```
int fgetc (FILE *stream);
```

Эта функция считывает очередной символ из потока `stream` и возвращает его в виде значения типа `unsigned char`, приведенного к типу `int`. Приведение выполняется для обеспечения достаточного диапазона, чтобы иметь возможность уведомлять об ошибке «конец файла»: в такой ситуации возвращается значение

EOF. Возвращаемое fgetc() значение нужно сохранять как тип int. Сохранение его в переменной типа char — это распространенная, но опасная ошибка.

В следующем примере один символ считывается из потока stream, проверяется на наличие ошибок, а затем результат выводится как значение типа char:

```
int c;
c = fgetc (stream);
if (c == EOF)
    /* ошибка */
else
    printf ("c=%c\n", (char) c);
```

Поток, на который указывает stream, должен быть открыт для чтения.

Вталкивание символа в поток данных

В реализации стандартного ввода-вывода также есть функция для вталкивания символа обратно в поток данных, что позволяет быстро «заглянуть» в поток и вернуть символ на место, если окажется, что он вам не подходит:

```
#include <stdio.h>

int ungetc (int c, FILE *stream);
```

Каждый вызов функции ungetc() помещает обратно в поток stream значение параметра c, приведенное к типу unsigned char. В случае успеха возвращается значение c, в случае неудачи возвращается значение EOF. Последующее считывание из stream возвращает c. Если вы вталкиваете обратно в поток несколько символов, то они возвращаются в обратном порядке, то есть последний помещенный в поток символ возвращается первым. В POSIX утверждается, что только для одного вталкивания успех гарантируется без вмешательства в запросы на считывание. Некоторые реализации, в свою очередь, допускают только одно вталкивание; Linux разрешает выполнять вталкивание бесконечное число раз, если для этого достаточно свободной памяти. Одно вталкивание, конечно же, всегда завершается успешно.

Если после вызова ungetc(), но перед запросом на считывание вы выполните вызов функции поиска (см. раздел «Поиск в потоке данных» далее в этой главе), то все втолкнутые в поток данных символы будут удалены. Это верно для всех потоков выполнения одного процесса, так как они совместно используют буфер.

Считывание строки целиком

Функция fgets() считывает строку из указанного потока данных:

```
#include <stdio.h>
```

```
char * fgets (char *str, int size, FILE *stream);
```

Эта функция читает из потока stream до size - 1 байт и сохраняет результат в строке str. Следом за считанными байтами в буфере также сохраняется символ

нуля (\0). Считывание останавливается после достижения конца файла или символа новой строки. Если считывается символ новой строки, то в str сохраняется символ \n.

В случае успеха возвращается строка str, в случае ошибки возвращается значение NULL.

Например:

```
char buf[LINE_MAX];
```

```
if (!fgets (buf, LINE_MAX, stream))
    /* ошибка */
```

В стандарте POSIX в файле заголовка <limits.h> определяется константа LINE_MAX, представляющая максимальный размер входной строки, который могут обрабатывать интерфейсы POSIX для манипулирования строками. В библиотеке C в Linux нет подобных ограничений и строки могут быть любой длины, но нет никакого способа соотнести это с определением LINE_MAX. В переносимых программах LINE_MAX можно использовать для обеспечения безопасности кода в других реализациях системы, эта константа имеет в Linux довольно большое значение. В программах, предназначенных только для Linux, не нужно беспокоиться об ограничениях на длину строк.

Чтение произвольных строк

Функция fgets(), читающая строки, часто оказывается полезной. Но практически так же часто она раздражает. Иногда разработчикам нужно использовать другие разделители, отличные от символа новой строки. Бывает, что разделитель вообще не нужен, и совсем редко у разработчиков возникает желание, чтобы разделитель хранился в буфере! Если бросить взгляд в историю, то решение сохранять символ новой строки в возвращаемом буфере не выглядит правильным.

Несложно написать замену функции fgets(), используя функцию fgetc(). Например, в следующем фрагменте кода считывается n - 1 байт из потока stream в строку str, а затем к строке присоединяется символ \0:

```
char *s,
int c.

s = str.
while (-n > 0 && (c = fgetc (stream)) != EOF)
    *s++ = c,
*s = '\0'.
```

Это решение можно расширить, чтобы чтение останавливалось на разделителе, указанном при помощи переменной d (которая в данном примере не может содержать символ нуля):

```
char *s,
int c = 0.

s = str.
while (-n > 0 && (c = fgetc (stream)) != EOF && (*s++ = c) != d)
```

```

    ;
if (c == d)
    *s = '\0';
else
    *s = '\0';

```

Если присвоить переменной `d` символ перевода строки, то этот фрагмент кода будет вести себя так же, как и `fgets()`, за исключением того, что он не сохраняет в буфере символ перевода строки.

В зависимости от реализации `fgets()` этот вариант, вероятно, будет работать медленнее, так как он много раз вызывает функцию `fgetc()`. Однако здесь речь идет не о той проблеме, которая была продемонстрирована выше в примере с командой `dd`. Хотя этот фрагмент создает дополнительную нагрузку повторным вызовом функции, он не нагружает систему системными вызовами и не создает сложностей невыровненного ввода-вывода, как команда `dd`, когда `bs = 1`. Вторая проблема намного серьезнее.

Считывание двоичных данных

В некоторых приложениях считывания отдельных символов или строк недостаточно. Иногда разработчикам приходится читать и записывать сложные двоичные данные, такие, как структуры С. Для этой цели библиотека стандартного ввода-вывода предоставляет функцию `fread()`:

```
#include <stdio.h>

size_t fread (void *buf, size_t size, size_t nr, FILE *stream);
```

Когда вы вызываете `fread()`, эта функция считывает из потока `stream` в буфер, на который указывает `buf`, до `nr` элементов данных, размер каждого из которых составляет `size` байт. Указатель файла передвигается вперед на число считанных байтов.

Функция возвращает число считанных элементов (а не число считанных байтов!). Возвращая значение `EOF` или значение, меньшее `nr`, функция указывает на сбой. К сожалению, невозможно узнать, какое именно из этих двух событий произошло, не применяя функции `ferror()` и `feof()` (см. раздел «Ошибки и конец файла» далее в этой главе).

Из-за различий в размере, выравнивании, заполнении и порядке байтов двоичные данные, записанные одним приложением, не всегда можно считать в другом приложении или даже в том же приложении на другой машине.

Простейший пример применения функции `fread()` — считывание одного элемента линейной последовательности байтов из указанного потока данных:

```
char buf[64];
size_t nr;

nr = fread (buf, sizeof(buf), 1, stream);
if (nr == 0)
    /* ошибка */
```

Мы рассмотрим более сложные примеры, когда будем изучать эквивалент `fread()` для записи — функцию `fwrite()`.

Запись в поток данных

Как и для чтения, стандартная библиотека С определяет множество функций для записи данных в открытый поток. В этом разделе мы изучим три наиболее популярных подхода к записи: запись одного символа, запись строки символов и запись двоичных данных. Такие разнообразные подходы к записи идеально подходят для буферизованного ввода-вывода. Чтобы записать данные в поток, его необходимо открыть как выходной поток в подходящем режиме, то есть любом режиме, за исключением `r`.

ПРОБЛЕМЫ ВЫРАВНИВАНИЯ

Во всех машинных архитектурах существуют требования к выравниванию данных (*data alignment*). Программисты часто представляют себе память просто как массив байтов. Но процессоры читают и записывают данные в память не однобайтовыми фрагментами. Они обращаются к памяти блоками, например, по 2, 4, 8 или 16 байт. Так как адресное пространство каждого процесса начинается по адресу 0, процессы должны инициировать доступ с адреса, являющегося кратным блоку.

Следовательно, переменные С необходимо хранить и обращаться к ним по выровненным адресам. В целом, переменные выравниваются естественным образом (*naturally aligned*) благодаря выравниванию соответствующих типов данных С. Например, 32-битное целое выровнено по 4-байтовой границе. Другими словами, переменная типа `int` хранится в памяти по адресу, который можно без остатка поделить на четыре.

С доступом к невыровненным данным связаны различные проблемы, зависящие от машинной архитектуры. Некоторые процессоры умеют обращаться к невыровненным данным, но за счет большого падения производительности. Другие процессоры вообще не умеют обращаться к невыровненным данным, и попытка сделать это вызывает аппаратное исключение. Что еще хуже, некоторые процессоры незаметно для пользователя удаляют биты младших разрядов, принудительно выравнивая адреса, что практически стопроцентно приводит к непредвиденным результатам.

Обычно компилятор естественным образом выравнивает все данные, и вопрос выравнивания совершенно не касается программиста. Работа со структурами, управление памятью вручную, сохранение двоичных данных на диск и коммуникация по сети могут приводить к необходимости решать какие-то вопросы выравнивания. Следовательно, системным программистам необходимо хорошо ориентироваться в этих вопросах!

Выравнивание подробнее обсуждается в восьмой главе.

Запись одного символа

Эквивалентом `fgetc()` для записи является функция `fputc()`:

```
#include <stdio.h>
```

```
int fputc (int c, FILE *stream);
```

Функция `fputc()` записывает байт, указанный при помощи аргумента `c` (это переменная, тип которой приведен к `unsigned char`), в поток, на который указывает `stream`. В случае успешного завершения функция возвращает значение `c`.

В противном случае она возвращает EOF и соответствующим образом устанавливает переменную errno.

Применять ее легко:

```
if (fputc ('p', stream) == EOF)
    /* ошибка */
```

В этом примере символ p записывается в поток stream, который должен быть открыт для записи.

Запись строки символов

Функция fputs() используется для записи целой строки символов в указанный поток данных:

```
#include <stdio.h>
```

```
int fputs (const char *str, FILE *stream);
```

Когда вы вызываете функцию fputs(), она записывает все содержимое отдельной нулём строки, на которую указывает str, в буфер, указанный при помощи stream. В случае успеха fputs() возвращает неотрицательное число. В случае неудачи она возвращает EOF.

В следующем примере файл открывается на запись в режиме присоединения, далее в связанный с файлом поток данных записывается указанная строка, после чего поток закрывается:

```
FILE *stream;
```

```
stream = fopen ("journal.txt", "a");
if (!stream)
    /* ошибка */

if (fputs ("The ship is made of wood.\n", stream) == EOF)
    /* ошибка */

if (fclose (stream) == EOF)
    /* ошибка */
```

Запись двоичных данных

Функций записи и чтения индивидуальных символов и строк недостаточно, когда в программе возникает необходимость записать сложные данные. Для прямого сохранения двоичных данных в библиотеке стандартного ввода-вывода предусмотрена функция fwrite():

```
#include <stdio.h>
```

```
size_t fwrite (void *buf,
              size_t size,
              size_t nr,
              FILE *stream);
```

Когда вы вызываете функцию `fwrite()`, то в поток `stream` записывается из буфера `buf` до `nr` элементов и длина каждого составляет `size` байт. Указатель файла перемещается вперед на общее число записанных байтов.

Функция возвращает число успешно записанных элементов (но не число байтов!). Если возвращенное значение меньше `nr`, то это указывает на ошибку.

Пример программы с использованием буферизованного ввода-вывода

Давайте теперь рассмотрим пример – на самом деле законченную программу, где используются многие из интерфейсов, которые мы уже рассмотрели в этой главе. Сначала в программе определяется структура `pirate`, а затем объявляются две переменные этого типа. Программа инициализирует одну из переменных и записывает ее на диск в файл `data` через выходной поток. Через другой поток программа считывает данные обратно из `data` напрямую в другой экземпляр `struct pirate`. Наконец, программа записывает содержимое структуры в стандартный вывод:

```
#include <stdio.h>

int main (void)
{
    FILE *in, *out;
    struct pirate {
        char           name[100]; /* настоящее имя */
        unsigned long  booty;     /* добыча в фунтах стерлингов */
        unsigned int   beard_len; /* длина бороды в дюймах */
    } p;
    blackbeard = { "Edward Teach", 950, 48 };

    out = fopen ("data", "w");
    if (!out) {
        perror ("fopen");
        return 1;
    }

    if (!fwrite (&blackbeard, sizeof (struct pirate), 1, out)) {
        perror ("fwrite");
        return 1;
    }

    if (fclose (out)) {
        perror ("fclose");
        return 1;
    }

    in = fopen ("data", "r");
    if (!in) {
        perror ("fopen");
        return 1;
    }
```

```

if (!fread (&p, sizeof (struct pirate), 1, in)) {
    perror ("fread");
    return 1;
}

if (fclose (in)) {
    perror ("fclose");
    return 1;
}

printf ("name=%s\" booty=%lu beard_len=%u\n".
p.name, p.booty, p.beard_len);

return 0;
}

```

На выходе мы получаем, конечно же, исходные значения:

```
name="Edward Teach" booty=950 beard_len=48
```

И снова важно помнить, что из-за различий в размерах переменных, выравнивании и т. п. двоичные данные, записанные одним приложением, не всегда можно считывать другими приложениями. Это означает, что другое приложение — или даже то же самое приложение на другой машине — может быть не в состоянии правильно считать обратно данные, записанные при помощи fwrite(). Что касается этого примера, представьте себе, какие последствия могли бы быть, если бы изменилась длина типа unsigned long или варьировался бы объем заполнения. Такие вещи гарантированно остаются постоянными только на определенном типе машины с определенным интерфейсом АВИ.

Поиск в потоке данных

Часто бывает полезно менять текущую позицию в потоке данных. Иногда в приложении, которое считывает сложный составленный из записей файл, приходится перемещаться туда и обратно по файлу. Также иногда приходится сбрасывать позицию в потоке на ноль. Независимо от того, для чего вам это нужно, библиотека стандартного ввода-вывода предлагает семейство интерфейсов, по функциональности эквивалентных системному вызову lseek() (о котором говорилось во второй главе). Функция fseek(), наиболее распространенная среди стандартных интерфейсов поиска, меняет позицию в файле stream согласно значениям параметров offset и whence:

```
#include <stdio.h>

int fseek (FILE *stream, long offset, int whence);
```

Если параметр whence равен SEEK_SET, то позиции в файле присваивается значение offset. Если whence равен SEEK_CUR, то позиция в файле увеличивается, начиная с текущего значения, на значение offset. Если whence равен SEEK_END, то указатель позиции в файле перемещается дальше конца файла на значение offset.

В случае успешного завершения функция `fseek()` возвращает 0, сбрасывает индикатор `EOF` и отменяет все результаты работы `ungetc()` (если эта функция применялась). В случае ошибки `fseek()` возвращает -1 и соответствующим образом устанавливает переменную `errno`. Наиболее распространенные ошибки — это недопустимый поток (`EBADF`) и недопустимый аргумент `whence` (`EINVAL`).

Также в библиотеке стандартного ввода-вывода есть функция `fsetpos()`:

```
#include <stdio.h>
```

```
int fsetpos (FILE *stream, fpos_t *pos);
```

Эта функция устанавливает для потока `stream` указатель позиции в потоке на значение `pos`. Она работает так же, как `fseek()`, когда параметр `whence` равен `SEEK_SET`. В случае успеха она возвращает 0, а в случае неудачи возвращает -1 и соответствующим образом устанавливает переменную `errno`. Эта функция (вместе со своим дополнением, функцией `fgetpos()`, о которой мы поговорим чуть далее) предоставляется исключительно для использования на других (не относящихся к Unix) платформах, в которых позицию в потоке данных представляют сложные типы. На этих plataформах данная функция — единственный способ установить позицию на произвольное значение, так как возможностей типа С `long`, по-видимому, недостаточно. В нацеленных исключительно на Linux приложениях этот интерфейс использовать нет необходимости, хотя и можно, если вы хотите обеспечить переносимость программы на все возможные платформы.

Также для стандартного ввода-вывода есть функция `rewind()`, позволяющая сокращать путь:

```
#include <stdio.h>
```

```
void rewind (FILE *stream);
```

Она вызывается так:

```
rewind (stream);
```

и сбрасывает позицию обратно на начало потока данных. Эта функция эквивалентна такой:

```
fseek (stream, 0, SEEK_SET);
```

за исключением того, что она также очищает индикатор ошибок.

Обратите внимание, что функция `rewind()` не возвращает никакое значение и поэтому она не может напрямую сообщать об ошибках. Если вызывающему необходимо удостовериться в существовании ошибки, то нужно очистить переменную `errno` перед вызовом и проверить ее значение после вызова, чтобы посмотреть, равна ли она все так же нулю. Например:

```
errno = 0;
rewind (stream);
if (errno)
    /* ошибка */
```

Получение текущей позиции в потоке данных

В отличие от системного вызова `lseek()`, функция `fseek()` не возвращает обновленную позицию в потоке данных. Для этого существует отдельный интерфейс. Функция `ftell()` сообщает текущую позицию в потоке `stream`:

```
#include <stdio.h>
```

```
long ftell (FILE *stream);
```

В случае ошибки она возвращает значение `-1` и соответствующим образом устанавливает переменную `errno`.

Также в библиотеке стандартного ввода-вывода есть функция `fgetpos()`:

```
#include <stdio.h>
```

```
int fgetpos (FILE *stream, fpos_t *pos);
```

В случае успешного завершения `fgetpos()` возвращает `0` и помещает значение текущей позиции в потоке `stream` в параметр `pos`. В случае ошибки она возвращает `-1` и устанавливает переменную `errno`. Как и `fsetpos()`, `fgetpos()` предоставляется исключительно для использования на не-Linux-платформах, где позицию в файле представляют сложные типы данных.

Сброс потока данных

В библиотеке стандартного ввода-вывода есть интерфейс для записи пользовательского буфера в ядро, позволяющий гарантировать, что все данные, записанные в поток, будут сброшены на диск вызовом `write()`. Этую функциональность обеспечивает функция `fflush()`:

```
#include <stdio.h>
```

```
int fflush (FILE *stream);
```

При вызове этой функции все незаписанные данные в потоке, на который указывает параметр `stream`, сбрасываются в ядро. Если параметр `stream` равен `NULL`, то сбрасываются *все* открытые входные потоки в процессе. В случае успеха `fflush()` возвращает значение `0`, а в случае ошибки она возвращает `EOF` и соответствующим образом устанавливает переменную `errno`.

Для того чтобы понять, как работает `fflush()`, необходимо знать, чем отличается буфер, поддерживаемый библиотекой С, от собственной системы буферизации ядра. Все вызовы в этой главе работают с буфером, который поддерживает библиотека С и который находится в пользовательском пространстве, а не в пространстве ядра. Именно за счет этого повышается производительность — вы остаетесь в пользовательском пространстве и, следовательно, исполняете пользовательский код, а не системные вызовы. Системный вызов делается только тогда, когда возникает необходимость обратиться к диску или какому-либо другому носителю.

Функция `fflush()` просто записывает данные из пользовательского буфера в буфер ядра. Результат получается такой же, как если бы пользовательская буферизация не применялась, а вы напрямую использовали бы системный вызов `write()`. Эта функция не гарантирует, что данные будут физически зафиксированы на любом носителе, — для этого необходимо применять что-либо, подобное системному вызову `fsync()` (см. раздел «Синхронизированный ввод-вывод» в главе 2). Вероятнее всего, вы будете вызывать `fflush()`, а затем сразу же `fsync()`, то есть сначала гарантировать, что пользовательский буфер будет записан в ядро, а затем, что данные из буфера ядра зафиксированы на диске.

Ошибки и конец файла

Некоторые из стандартных интерфейсов ввода-вывода, такие, как `fread()`, не умеют правильно сообщать о сбоях вызывающему и не предоставляют никакого механизма, который позволял бы отличать ситуацию ошибки от ситуации, когда встретился конец файла. Когда вы используете эти вызовы, а также в некоторых других случаях может быть удобно проверять статус интересующего вас потока данных, определяя, произошла в нем ошибка, или он достиг конца файла. В библиотеке стандартного ввода-вывода есть для этого два интерфейса. Функция `ferror()` проверяет, установлен ли для потока `stream` индикатор ошибки:

```
#include <stdio.h>
```

```
int ferror (FILE *stream);
```

Индикатор ошибки устанавливается другими стандартными интерфейсами ввода-вывода в ответ на условие ошибки. Эта функция возвращает ненулевое значение, если индикатор установлен, и 0 в противном случае.

Функция `feof()` проверяет, установлен ли индикатор EOF для потока `stream`:

```
#include <stdio.h>
```

```
int feof (FILE *stream);
```

Индикатор EOF устанавливается другими стандартными интерфейсами ввода-вывода, когда достигается конец файла. Эта функция возвращает ненулевое значение, если индикатор установлен, и 0 в противном случае.

Функция `clearerr()` очищает индикаторы ошибки и EOF для потока `stream`:

```
#include <stdio.h>
```

```
void clearerr (FILE *stream);
```

Она не возвращает никакое значение и не может завершаться сбоем (и способа проверить, что предоставлено допустимое значение потока данных, также нет). Вызов `clearerr()` следует делать только после проверки индикаторов ошибки и EOF, так как их значения безвозвратно удаляются. Например:

```
/* f – это допустимый поток */
```

```
if (ferror (f))
```

```
printf ("Error on f!\n").  
if (feof (f))  
    printf ("EOF on f!\n");  
clearerr (f);
```

Получение связанного дескриптора файла

Иногда бывает полезно знать, какой дескриптор файла соответствует данному потоку данных. Например, это может потребоваться, чтобы выполнить на потоке системный вызов через файловый дескриптор, когда аналогичной функции стандартного ввода-вывода не существует. Для получения дескриптора файла, соответствующего потоку, используйте функцию `fileno()`:

```
#include <stdio.h>  
  
int fileno (FILE *stream);
```

В случае успеха `fileno()` возвращает дескриптор файла, связанный с потоком `stream`. В случае сбоя она возвращает значение `-1`. Это может происходить только в том случае, когда функции передается недопустимый поток, и тогда `fileno()` присваивает переменной `errno` значение `EBADF`.

Обычно не рекомендуется перемешивать стандартные вызовы ввода-вывода с системными вызовами. Необходимо соблюдать осторожность при использовании `fileno()`, чтобы гарантировать правильное и предсказуемое поведение. В частности, вероятно, будет полезно сбрасывать поток данных перед тем, как делать что-либо с соответствующим ему файловым дескриптором. Фактические операции ввода-вывода не следует смешивать практически никогда.

Управление буферизацией

В библиотеке стандартного ввода-вывода реализовано три типа пользовательской буферизации, и она предоставляет разработчикам интерфейс, позволяющий управлять типом и размером буфера. Различные типы пользовательских буферов служат разным целям и подходят для разных ситуаций. Доступны следующие варианты:

Без буферизации

Пользовательская буферизация не выполняется. Данные передаются напрямую ядру. Так как это противоречит самой идее пользовательской буферизации, данный вариант используется крайне редко. Стандартный выход ошибок по умолчанию не буферизуется.

Строковая буферизация

Буферизация выполняется построчно. При получении очередного символа новой строки содержимое буфера передается ядру. Строковая буферизация имеет смысл для потоков данных, которые выводятся на экран. Следовательно,

это стандартный тип буферизации для терминалов (к стандартному выходу строковая буферизация применяется по умолчанию).

Блочная буферизация

Буферизация выполняется блоками. Этот тип буферизации мы обсуждали в начале главы, и он идеально подходит для файлов. По умолчанию все потоки данных, связанные с файлами, буферизуются блоками. Для стандартного ввода-вывода блочную буферизацию называют *полной буферизацией* (full buffering).

Обычно тип буферизации менять не нужно, так как тип, который используется по умолчанию, оптимален. Однако в библиотеке стандартного ввода-вывода есть интерфейс для управления типом буферизации:

```
#include <stdio.h>
```

```
int setvbuf (FILE *stream, char *buf, int mode, size_t size);
```

Функция `setvbuf()` устанавливает для потока `stream` тип буферизации, указанный при помощи параметра `mode`, и этот параметр может принимать следующие значения:

`_IONBF`

Без буферизации.

`_IOFBF`

Строковая буферизация.

`_IOLBF`

Блочная буферизация.

Если вы используете значение `_IONBF`, то параметры `buf` и `size` игнорируются. В остальных случаях `buf` может указывать на буфер размером `size` байт, который при операциях стандартного ввода-вывода используется в качестве буфера для указанного потока данных. Если же значение `buf` равно `NULL`, то буфер выделяется автоматически библиотекой `glibc`.

Функцию `setvbuf()` необходимо вызывать после открытия потока, но до того, как на нем выполняется любая другая операция. Она возвращает значение 0 после успешного завершения и ненулевое значение в противном случае.

Если вы указываете конкретный буфер, то он должен существовать на момент, когда поток закрывается. Очень распространенная ошибка — объявить буфер как автоматическую переменную в контексте, который удаляется до того, как закрывается поток. В частности, будьте осторожны и избегайте ситуации, когда буфер объявляется локально в функции `main()`, а явное закрытие потоков в той же функции отсутствует. Например, следующий код содержит ошибку:

```
#include <stdio.h>
```

```
int main (void)
```

```
{
```

```
    char buf[BUFSIZ];
```

```
/* включение для stdin блочной буферизации с буфером BUFSIZ */
setvbuf (stdout, buf, _IOFBF, BUFSIZ);

printf ("Arrr!\n");

return 0;
}
```

Эту ошибку можно исправить, явно закрыв поток данных перед тем, как программа выйдет из данного контекста, или же сделав buf глобальной переменной.

В целом, разработчикам обычно не приходится задумываться о буферизации потоков данных. За исключением стандартного вывода ошибок, для всех терминалов используется строковая буферизация, и это действительно имеет смысл. Для файлов применяется блочная буферизация, и это также разумно. Размер буфера по умолчанию для блочной буферизации равен значению BUFSIZ, которое определяется в `<stdio.h>`, и обычно это оптимальный выбор (это значение кратно типичному размеру блока).

Безопасность потоков выполнения

Потоки выполнения (*thread*) — это несколько нитей исполнения, составляющих каждый процесс. Их можно представлять себе как несколько процессов, совместно использующих одно адресное пространство. Потоки могут выполняться одновременно и менять общие данные, если только доступ к данным специально не синхронизируется или они не определяются как локальные для потока выполнения (*thread-local*). Операционные системы, поддерживающие потоки, предлагают также механизмы блокировки (программные конструкции, обеспечивающие взаимоисключение), позволяющие гарантировать, что потоки не будут наступать друг другу на пятки. Стандартный ввод-вывод выполняется с использованием этих механизмов. Однако они не всегда позволяют добиваться желаемого. Например, иногда возникает необходимость заблокировать группу вызовов, увеличив *критический регион* (*critical region* — фрагмент кода, который должен выполняться без вмешательства других потоков) с одной операции ввода-вывода до нескольких. В других ситуациях может понадобиться вообще убрать блокировку, чтобы повысить эффективность¹. В этом разделе мы посмотрим, как решить обе эти задачи.

Стандартные функции ввода-вывода *поточно-ориентированы* (*thread-safe*) по своей природе. Внутри системы они связывают блокировку, счетчик блокировок и поток-владелец (поток выполнения, *thread*) с каждым открытым потоком данных (*stream*). Перед тем как делать любые запросы ввода-вывода, любой

¹ Обычно удаление блокировки приводит к целому набору проблем. Однако иногда в программах работа потоков явно реализуется так, чтобы весь ввод-вывод выполнял какой-то один поток. В таких случаях можно избавляться от лишней нагрузки, создаваемой блокировкой.

поток должен сначала захватить блокировку и превратиться в поток-владелец. Два и более потока, работающие на одном и том же потоке данных, не могут смешивать свои операции стандартного ввода-вывода, поэтому операции стандартного ввода-вывода являются атомарными в пределах контекста вызова одной функции.

Конечно же, на практике многим приложениям требуется более высокая степень атомности, чем на уровне отдельных вызовов функций. Например, если несколько потоков выполнения делают запросы на запись, хотя отдельные операции записи смешиваются и создавать неразбериху в выводе не будут, приложению может быть необходимо, чтобы все запросы на запись завершились без перерыва. Для того чтобы обеспечить это, в библиотеке стандартного ввода-вывода есть семейство функций для манипулирования блокировкой потоков.

Блокировка файла вручную

Функция `flockfile()` дожидается момента, когда блокировка с потока данных `stream` снимается, и затем захватывает блокировку, увеличивает счетчик блокировок, делает свой поток потоком-владельцем данного потока данных и возвращает результат:

```
#include <stdio.h>
```

```
void flockfile (FILE *stream);
```

Функция `funlockfile()` уменьшает счетчик блокировок, связанный с потоком `stream`:

```
#include <stdio.h>
```

```
void funlockfile (FILE *stream);
```

Когда значение счетчика блокировок достигает нуля, текущий поток выполнения освобождает владение потоком данных и захватить блокировку после этого может другой поток выполнения.

Эти вызовы могут вкладываться один в другой. То есть один поток выполнения может сделать несколько вызовов `flockfile()`, и поток данных не разблокируется до тех пор, пока процесс не выполнит соответствующее число вызовов `funlockfile()`.

Функция `ftrylockfile()` – это версия `flockfile()`, которая захватывает или не захватывает блокировку, в зависимости от обстоятельств:

```
#include <stdio.h>
```

```
int ftrylockfile (FILE *stream);
```

Если поток `stream` в данный момент заблокирован, `ftrylockfile()` ничего не делает, но немедленно возвращает ненулевое значение. Если поток `stream` в данный момент не заблокирован, то функция захватывает блокировку, увеличивает счетчик блокировок, делает свой поток выполнения потоком-владельцем и возвращает значение 0.

Рассмотрим пример:

```
flockfile (stream);  
  
fputs ("List of treasure:\n", stream);  
fputs ("    (1) 500 gold coins\n", stream);  
fputs ("    (2) Wonderfully ornate dishware\n", stream);
```

```
unlockfile (stream);
```

Хотя отдельные операции fputs() никогда не состязаются, например, в середину строки «List of treasure» «чужие» данные никогда попасть не смогут, другая операция стандартного ввода-вывода, принадлежащая другому потоку выполнения, работающему на том же потоке данных, может успевать между двумя вызовами fputs(). В идеальном случае приложение должно быть разработано так, чтобы несколько потоков никогда не пытались выполнять операции ввода-вывода на одном и том же потоке данных. Если же в вашем приложении требуется такая функциональность и критические регионы должны охватывать более одной функции, то функция flockfile() и ее друзья придут вам на помощь.

Операции над потоками без блокировки

Есть и еще одна причина для блокировки на потоках данных вручную. Благодаря более тонкому и точному управлению блокировкой, которое способен обеспечить только программист, можно минимизировать нагрузку, создаваемую блокировкой, соответственно повышая производительность. Для этого в Linux предусмотрено семейство функций, родственных обычным интерфейсам стандартного ввода-вывода, которые вообще не выполняют блокировку. Фактически это эквиваленты операций стандартного ввода-вывода, но без блокировки:

```
#define _GNU_SOURCE  
  
#include <stdio.h>  
  
int fgetc_unlocked (FILE *stream);  
char *fgets_unlocked (char *str, int size, FILE *stream);  
size_t fread_unlocked (void *buf, size_t size, size_t nr,  
                      FILE *stream);  
int fputc_unlocked (int c, FILE *stream);  
int fputs_unlocked (const char *str, FILE *stream);  
size_t fwrite_unlocked (void *buf, size_t size, size_t nr,  
                      FILE *stream);  
int fflush_unlocked (FILE *stream);  
int feof_unlocked (FILE *stream);  
int ferror_unlocked (FILE *stream);  
int fileno_unlocked (FILE *stream);  
void clearerr_unlocked (FILE *stream);
```

Эти функции все работают точно так же, как их эквиваленты с блокировкой, за исключением того, что они не проверяют наличие блокировки и не захватывают блокировку для указанного потока — stream. Если блокировка необходима, то программист обязан гарантировать, что блокировка вручную устанавливается и снимается.

Хотя в POSIX определены некоторые варианты функций стандартного ввода-вывода без блокировки, ни одна из перечисленных выше функций в POSIX не входит. Они все уникальны для Linux, но и некоторые другие системы Unix поддерживают часть этого набора.

Недостатки стандартного ввода-вывода

Несмотря на то что операции стандартного вывода применяются очень широко, некоторые эксперты находят в них определенные недостатки. Часть функций, например `fgets()`, иногда ведет себя непредсказуемо. Другие функции, такие, как `gets()`, настолько ужасающи, что они практически были изгнаны из стандартов.

Наибольшую критику вызывает отрицательное воздействие на производительность, к которому приводит двойное копирование в операциях стандартного ввода-вывода. При считывании данных функция стандартного ввода-вывода отправляет ядру системный вызов `read()`, копируя данные из ядра в буфер стандартного ввода-вывода. Когда приложение после этого делает запрос на чтение при помощи функции стандартного ввода-вывода, например `fgetc()`, данные снова копируются — на этот раз из буфера стандартного ввода-вывода в указанный в функции буфер. Запросы на запись работают наоборот: сначала данные копируются из предоставленного буфера в буфер стандартного ввода-вывода, а затем — из этого буфера в ядро при помощи системного вызова `write()`.

В альтернативной реализации можно было бы избежать двойного копирования, заставив все запросы на чтение возвращать указатель на буфер стандартного ввода-вывода. Затем данные можно было бы считывать напрямую из буфера стандартного ввода-вывода, не прибегая к дополнительному копированию. В случае если приложению понадобится переместить данные в собственный локальный буфер, например, чтобы что-то дописать к ним, оно всегда могло бы копировать их вручную. Такая реализация предоставила бы «свободный» интерфейс, позволяющий приложениям сигнализировать, когда они заканчивают обработку каждого фрагмента считанных данных.

Реализовать запись сложнее, но все же можно было бы избежать двойного копирования. Когда делается запрос на запись, в реализации необходимо запоминать указатель. В конечном итоге, когда данные будут готовы к сбросу в ядро, нужно будет пройти по списку сохраненных указателей и отправить все данные в ядро. Это можно делать при помощи разбросанного ввода-вывода через `writev()`, применяя, таким образом, только один системный вызов (мы поговорим о разбросанном вводе-выводе в следующей главе).

Существуют высоко оптимизированные библиотеки пользовательской буферизации, решающие проблему двойного копирования за счет реализации, похожей на описанную выше. Также некоторые разработчики реализуют собственные решения пользовательской буферизации. Однако, несмотря на существующие альтернативы, популярность стандартного ввода-вывода все так же остается на высоком уровне.

Заключение

Библиотека стандартного ввода-вывода — это библиотека, реализующая пользовательскую буферизацию, которая предоставляется как часть стандартной библиотеки С. Если не обращать внимания на некоторые недостатки, то это мощное и широко используемое решение. На самом деле многие программисты на С не знают других решений, кроме стандартного ввода-вывода. Действительно, для терминального ввода-вывода, для которого идеально подходит строковая буферизация, стандартный ввод-вывод — это единственный вариант. Кто стал бы применять `write()` напрямую, чтобы записать данные в стандартный вывод?

Стандартный ввод-вывод — и пользовательская буферизация в целом, если уж на то пошло, — применяется, когда выполняется любое из следующих условий:

- предполагается, что вы будете выполнять много системных вызовов, и поэтому вам хотелось бы минимизировать нагрузку, объединив их и сократив число системных вызовов до минимума;
- вопрос производительности имеет первостепенную важность, и вы хотите убедиться, что все операции ввода-вывода будут выполняться фрагментами, кратными блоку, с границами, выровненными по блокам;
- вы используете символьный или строковый шаблон доступа, и вам нужны интерфейсы, упрощающие подобный доступ, чтобы не прибегать к излишним системным вызовам;
- вы предпочитаете высокоуровневый интерфейс низкоуровневым системным вызовам Linux.

Наибольшей гибкости, однако, можно добиться, работая напрямую с системными вызовами Linux. В следующей главе мы изучим более сложные формы ввода-вывода и связанные с ними системные вызовы.

4 Расширенный файловый ввод-вывод

В главе 2 мы изучили основные системные вызовы ввода-вывода в Linux. Все эти вызовы формируют базис не только файлового ввода-вывода, но и практически всех видов коммуникаций в Linux. В третьей главе мы узнали, как производится буферизация поверх базовых системных вызовов ввода-вывода в пользовательском пространстве, а также познакомились с конкретным решением буферизации – библиотекой стандартного ввода-вывода C. Темой этой главы будут более сложные системные вызовы ввода-вывода, предоставляемые Linux.

Разбросанный ввод-вывод

Позволяет одному вызову читать или записывать данные в несколько буферов одновременно; полезен для объединения полей разных структур данных и формирования единой транзакции ввода-вывода.

Интерфейс epoll

Эта возможность расширяет системные вызовы poll() и select(), о которых рассказывалось в главе 2; она полезна, когда в одной программе необходимо опросить сотни дескрипторов файла.

Ввод-вывод с проецированием в память

Файл отображается в память, позволяя осуществлять файловый ввод-вывод при помощи простых манипуляций с памятью. Этот тип ввода-вывода полезен для определенных шаблонов ввода-вывода.

Советы по использованию файлов

Эта техника позволяет процессу передавать ядру подсказки относительно сценариев использования файлов; может повышать производительность ввода-вывода.

Асинхронный ввод-вывод

Позволяет процессу отправлять запросы ввода-вывода, не дожидаясь их завершения; полезен для управления большой нагрузкой ввода-вывода без использования потоков.

Завершается глава обсуждением вопросов производительности и подсистем ввода-вывода ядра.

Разбросанный ввод-вывод

Разбросанный ввод-вывод (scatter/gather I/O) — это способ ввода и вывода данных, когда один системный вызов записывает данные из одного потока данных в вектор буферов или, иначе, считывает данные из одного потока данных в вектор буферов. Этот тип ввода-вывода носит такое название потому, что данные разбрасываются (scatter) по указанному вектору буферов или собираются (gather) из этого вектора. Альтернативное название данного подхода к вводу-выводу — *векторный ввод-вывод* (vectored I/O). В противоположность ему, стандартные системные вызовы чтения и записи, о которых мы говорили в главе 2, обеспечивают *линейный ввод-вывод* (linear I/O).

Разбросанный ввод-вывод дает несколько преимуществ по сравнению с методами линейного ввода-вывода:

Более естественная обработка

Если данные естественным образом сегментированы, — например, это поля предопределенного файла заголовка, то векторный ввод-вывод обеспечивает интуитивную манипуляцию данными.

Эффективность

Одна операция векторного ввода-вывода может заменить несколько операций линейного ввода-вывода.

Производительность

Помимо сокращения числа выполняемых системных вызовов, векторный ввод-вывод также обеспечивает более высокую производительность по сравнению с линейным за счет внутренней оптимизации.

Атомность

В отличие от нескольких операций линейного ввода-вывода, процесс может выполнять одну операцию векторного ввода-вывода, не беспокоясь о пересечении с операцией, исполняемой другим процессом.

Добиться более естественной реализации ввода-вывода и атомности операций чтения и записи можно и без применения механизма разбросанного ввода-вывода. Перед записью процесс может объединять различные векторы в единый буфер и разбивать возвращаемый буфер на несколько векторов после завершения считывания, то есть приложение из пользовательского пространства может выполнять разбрасывание и сбор данных вручную. Однако такое решение неэффективно и его неинтересно реализовывать.

Системные вызовы `readv()` и `writev()`

В стандарте POSIX 1003.1-2001 определяется, а в Linux реализуется пара системных вызовов, обеспечивающих разбросанный ввод-вывод. Реализация Linux удовлетворяет всем целям, перечисленным в предыдущем разделе.

Функция `readv()` считывает `count` сегментов из файла, указанного при помощи дескриптора файла `fd`, в буфера, на которые указывает `iov`:

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

ssize_t readv (int fd,
               const struct iovec *iov,
               int count);
```

Функция `writev()` записывает максимум `count` сегментов из буферов, на которые указывает аргумент `iov`, в файл, указанный при помощи дескриптора файла `fd`:

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

ssize_t writev (int fd,
                const struct iovec *iov,
                int count);
```

Функции `readv()` и `writev()` работают так же, как `read()` и `write()`, соответственно, за исключением того, что для считывания и записи используется несколько буферов.

Каждая структура `iovec` описывает независимый автономный буфер, который называется *сегментом* (*segment*):

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

struct iovec {
    void *iov_base; /* указатель на начало буфера */
    size_t iov_len; /* размер буфера в байтах */
};
```

Набор сегментов называется *вектором* (*vector*). Каждый сегмент в векторе описывает адрес и длину буфера в памяти, в который или из которого данные должны записываться или считываться. Функция `readv()`, перед тем как переходить к очередному буферу, полностью заполняет каждый буфер длиной `iov_len` байтов. Функция `writev()` всегда записывает в файл все `iov_len` байтов из одного буфера и только после этого переходит к следующему. Обе функции всегда обрабатывают сегменты по порядку, начиная с `iov[0]`, затем `iov[1]` и т. д., пока не достигают сегмента `iov[count - 1]`.

Возвращаемые значения

В случае успеха `readv()` и `writev()` возвращают число считанных или записанных байтов соответственно. Это значение должно равняться сумме всех `count` значений `iov_len`. В случае ошибки системные вызовы возвращают значение `-1` и соответствующим образом устанавливают переменную `errno`. Эти системные вызовы могут завершаться теми же ошибками, что и системные вызовы `read()` и `write()`, и коды ошибок для переменной `errno` также не отличаются. Помимо этого, в стандартах определяются две дополнительные ошибочные ситуации.

Во-первых, так как эти системные вызовы возвращают значение типа `ssize_t`, если сумма всех `count` значений `iov_len` превышает значение `SSIZE_MAX`, никакие данные не перемещаются, возвращается значение `-1`, а переменной `errno` присваивается значение `EINVAL`.

Во-вторых, стандарт POSIX диктует, что значение count должно быть больше нуля, но меньше или равно значению константы IOV_MAX, которая определяется в заголовке `<limits.h>`. В Linux значение IOV_MAX в настоящее время равно 1024. Если аргумент count равен нулю, системные вызовы возвращают 0¹. Если аргумент count больше IOV_MAX, то никакие данные не перемещаются, вызовы возвращают значение -1, а переменной errno присваивается значение EINVAL.

ОПТИМИЗАЦИЯ СЧЕТЧИКА

Во время выполнения операции векторного ввода-вывода ядро Linux должно выделять внутренние структуры данных для представления каждого сегмента. Обычно такое выделение происходит динамически на основе значения аргумента count. Однако для оптимизации ядро Linux создает в стеке небольшой массив сегментов, который оно использует, если значение count достаточно мало, устраняя необходимость динамически выделять сегменты и, следовательно, давая небольшой выигрыш в производительности. В настоящий момент пороговое значение равно восьми, поэтому если значение count меньше или равно восьми, то операции векторного ввода-вывода выполняются из стека ядра процесса для более эффективного использования памяти.

Скорее всего, у вас не будет возможности решать, сколько сегментов одновременно передавать в каждой операции векторного ввода-вывода. Однако если это возможно и вы раздумываете, насколько мал должен быть аргумент count, то значение, меньшее или равное восьми, определенно повысит эффективность.

Пример использования системного вызова writev()

Давайте рассмотрим простой пример, в котором в файл записывается вектор из трех сегментов, содержащих строки разных размеров. Эта законченная программа достаточно полна для того, чтобы продемонстрировать работу системного вызова `writev()`, но в то же время достаточно проста, чтобы служить полезным отрывком кода:

```
#include <stdio.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <fcntl.h>
#include <string.h>
#include <sys/uio.h>

int main ( )
{
    struct iovec iov[3];
    ssize_t nr;
    int fd, i;

    char *buf[] = {
        "The term buccaneer comes from the word boucan.\n",
        "A boucan is a wooden frame used for cooking meat.\n",
        "Buccaneer is the West Indies name for a pirate.\n" };

```

¹ Обратите внимание, что в других системах Unix переменная errno может получать значение EINVAL, если аргумент count равен нулю. Это явно допускается стандартами, утверждающими, что если это значение равно 0, то либо может создаваться ошибка EINVAL, либо система может обрабатывать нулевой случай каким-то другим способом без выбрасывания ошибки.

```

fd = open ("buccaneer.txt", O_WRONLY | O_CREAT | O_TRUNC);
if (fd == -1) {
    perror ("open");
    return 1;
}

/* заполнение трех структур iovec */
for (i = 0; i < 3; i++) {
    iov[i].iov_base = buf[i];
    iov[i].iov_len = strlen (buf[i]);
}

/* один вызов записывает все три структуры в файл */
nr = writev (fd, iov, 3);
if (nr == -1) {
    perror ("writev");
    return 1;
}
printf ("wrote %d bytes\n", nr);

if (close (fd)) {
    perror ("close");
    return 1;
}

return 0;
}

```

При выполнении программы мы получаем желаемый результат:

```
$ ./writev
wrote 148 bytes
```

Точно так же мы можем успешно считать файл:

```
$ cat buccaneer.txt
The term buccaneer comes from the word boucan.
A boucan is a wooden frame used for cooking meat.
Buccaneer is the West Indies name for a pirate.
```

Пример использования системного вызова readv()

Теперь рассмотрим пример программы, в котором системный вызов readv() используется для чтения данных из ранее сгенерированного текстового файла с применением векторного ввода-вывода. Это также простая, но полная программа:

```

#include <stdio.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <fcntl.h>
#include <sys/uio.h>

int main ( )
{
    char foo[48], bar[51], baz[49];
    struct iovec iov[3];

```

```

ssize_t nr;
int fd, i;

fd = open ("buccaneer.txt", O_RDONLY);
if (fd == -1) {
    perror ("open");
    return 1;
}

/* установка значений структур iovc */
iov[0].iov_base = foo;
iov[0].iov_len = sizeof (foo);
iov[1].iov_base = bar;
iov[1].iov_len = sizeof (bar);
iov[2].iov_base = baz;
iov[2].iov_len = sizeof (baz);

/* считывание данных в структуры одним вызовом */
nr = readv (fd, iov, 3);
if (nr == -1) {
    perror ("readv");
    return 1;
}

for (i = 0; i < 3; i++)
    printf ("%d: %s", i, (char *) iov[i].iov_base);

if (close (fd)) {
    perror ("close");
    return 1;
}

return 0;
}

```

Если выполнить эту программу после предыдущей, то результат будет таким:

```

$ ./readv
0: The term buccaneer comes from the word boucan.
1: A boucan is a wooden frame used for cooking meat.
2: Buccaneer is the West Indies name for a pirate

```

Реализация

Простейшую реализацию системных вызовов `readv()` и `writev()` в пользовательском пространстве можно создать, используя цикл, например, так:

```

#include <unistd.h>
#include <sys/uio.h>

ssize_t naive_writev (int fd, const struct iovec *iov, int count)
{
    ssize_t ret = 0;
    int i;

    for (i = 0; i < count; i++) {

```

```

ssize_t nr;

nr = write(fd, iov[i].iov_base, iov[i].iov_len);
if (nr == -1) {
    ret = -1;
    break;
}
ret += nr;
}

return ret;
}

```

К счастью, в Linux это *не так*: `readv()` и `writev()` реализованы как системные вызовы и внутри системы действительно выполняется разбросанный ввод-вывод. В действительности весь ввод-вывод внутри ядра Linux векторизован; системные вызовы `read()` и `write()` используют векторный ввод-вывод, но для них вектор содержит только один сегмент.

Интерфейс event poll

В ответ на ограничения вызовов `poll()` и `select()` в версии ядра Linux 2.6¹ была представлена возможность `event poll` (`epoll` — опрос событий). Несмотря на более высокую сложность по сравнению с предыдущими интерфейсами, `epoll` решает фундаментальную проблему производительности, характерную для них обоих, и добавляет несколько новых функций.

Обоим системным вызовам — и `poll()`, и `select()` (подробнее о них говорится в главе 2) — необходим полный список дескрипторов файла, за которыми они должны наблюдать при каждом вызове. Ядру приходится проходить по списку, проверяя каждый дескриптор файла. Когда список становится слишком большим — он может содержать сотни и даже тысячи файловых дескрипторов, — прохождение списка при каждом вызове `poll()` или `select()` превращается в «узкое место», мешающее масштабированию приложений.

`Epoll` обходит эту проблему, разделяя регистрацию монитора и фактический мониторинг. Один системный вызов инициализирует контекст `epoll`, второй добавляет в контекст или удаляет из него дескрипторы файла, мониторинг которых необходим, а третий вызов несет ответственность за фактическое ожидание событий.

Создание нового экземпляра epoll

Контекст `epoll` создается при помощи системного вызова `epoll_create()`:

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

```
int epoll_create (int size)
```

¹ Возможность `epoll` впервые была представлена в разрабатываемом ядре 2.5.44, а завершен интерфейс был в версии ядра 2.5.66.

При успешном выполнении вызов `epoll_create()` создает новый экземпляр `epoll` и возвращает дескриптор файла, связанный с этим экземпляром. Данный файловый дескриптор никак не связан с реальным файлом; это всего лишь дескриптор, который будет использоваться в последующих вызовах, относящихся к возможностям `epoll`. Параметр `size` – это подсказка ядру относительно числа дескрипторов файла, за которыми вы собираетесь наблюдать; это не максимальное количество. Чем более точное значение вы передадите, тем выше будет производительность, но точное значение не требуется. В случае ошибки вызов возвращает значение `-1` и присваивает переменной `errno` одно из следующих значений:

`EINVAL`

Параметр `size` содержит не положительное число.

`ENFILE`

В системе открыто максимально допустимое число файлов.

`ENOMEM`

Для завершения операции недостаточно памяти.

Типичный вызов:

```
int epfd;
```

```
epfd = epoll_create (100); /* plan to watch ~100 fds */  
if (epfd < 0)  
    perror ("epoll_create").
```

Дескриптор файла, который возвращает `epoll_create()`, необходимо разрушать при помощи системного вызова `close()` после того, как опрашивание завершается.

Управление `epoll`

Системный вызов `epoll_ctl()` можно использовать для добавления дескрипторов файла или удаления дескрипторов из указанного контекста `epoll`:

```
#include <sys/epoll.h>  
  
int epoll_ctl (int epfd,  
               int op,  
               int fd,  
               struct epoll_event *event).
```

В заголовке `<sys/epoll.h>` структура `epoll_event` определяется так:

```
struct epoll_event {  
    __u32 events; /* events */  
    union {  
        void *ptr;  
        int fd;  
        __u32 u32;  
        __u64 u64;  
    } data;  
};
```

Успешный вызов `epoll_ctl()` позволяет изменить определенную информацию в экземпляре `epoll`, связанном с дескриптором файла `epfd`. При помощи параметра `op` вы указываете операцию, которая должна быть выполнена для файла, которому соответствует дескриптор `fd`. Параметр `event` содержит дополнительное описание поведения операции.

Далее перечислены допустимые значения для параметра `op`:

EPOLL_CTL_ADD

Добавляет в экземпляр `epoll`, связанный с дескриптором `epfd`, задачу мониторинга файла, который указан при помощи дескриптора файла `fd`, относительно событий, определенных в `event`.

EPOLL_CTL_DEL

Удаляет задачу мониторинга файла, связанного с дескриптором файла `fd`, из экземпляра `epoll`, на который указывает `epfd`.

EPOLL_CTL_MOD

Модифицирует существующую задачу мониторинга `fd`, заменяя старые события новыми, которые указаны при помощи параметра `event`.

В поле `events` структуры `epoll_event` вы перечисляете события, которые должны отслеживаться для данного дескриптора файла. Для того чтобы указать несколько событий, примените к ним операцию побитового ИЛИ. Допустимые значения перечислены далее:

EPOLLERR

Условие ошибки для файла. Мониторинг этого события выполняется всегда, даже если данное значение не указано.

EPOLLET

Это значение переводит мониторинг файла в режим запуска фронтом (см. далее раздел «События, запускаемые фронтом, и события, запускаемые уровнем»). По умолчанию используется режим запуска уровнем.

EPOLLHUP

Зависание файла. Мониторинг этого события выполняется всегда, даже если данное значение не указано.

EPOLLIN

Файл доступен для чтения без блокировки.

EPOLLONESHOT

После того как событие генерируется и считывается, мониторинг файла автоматически прекращается. Для того чтобы возобновить наблюдение, новую маску событий необходимо указать при помощи `EPOLL_CTL_MOD`.

EPOLLOUT

Файл доступен для записи без блокировки.

EPOLLPRI

Доступны для чтения срочные внеполосные данные.

Поле `data` в структуре `event_poll` предназначено для частного использования пользователем. Содержимое возвращается пользователю после получения запрошенного события. Распространенная практика — присваивать полю `event.data.fd` значение `fd`, чтобы с легкостью узнавать, к какому дескриптору файла относится отловленное событие.

В случае успеха `epoll_ctl()` возвращает значение 0. В случае ошибки вызов возвращает значение -1 и присваивает переменной `errno` одно из следующих значений:

EBADF

`Epfid` указывает на недопустимый экземпляр `epoll`, или `fd` содержит недопустимый дескриптор файла.

EEXIST

Значение `op` равно `EPOLL_CTL_ADD`, но дескриптор файла `fd` уже связан с дескриптором `epfd`.

EINVAL

`Epfid` указывает не на экземпляр `epoll`, значение `epfd` равно значению `fd` или передано недопустимое значение `op`.

ENOENT

Значение `op` равно `EPOLL_CTL_MOD` или `EPOLL_CTL_DEL`, но дескриптор файла `fd` не связан с дескриптором `epfd`.

ENOMEM

Недостаточно памяти для обработки запроса.

EPERM

`Fd` не поддерживает `epoll`.

В качестве примера, чтобы добавить новую задачу наблюдения за файлом, связанным с `fd`, в экземпляр `epoll` `epfd`, нужно использовать такой код:

```
struct epoll_event event;
int ret;

event.data.fd = fd; /* return the fd to us later */
event.events = EPOLLIN | EPOLLOUT;

ret = epoll_ctl(epfd, EPOLL_CTL_ADD, fd, &event);
if (ret)
    perror("epoll_ctl");
```

Чтобы изменить существующее наблюдаемое событие для файла, связанного с `fd`, в экземпляре `epoll` `epfd`, используйте такой код:

```
struct epoll_event event;
int ret;

event.data.fd = fd; /* return the fd to us later */
event.events = EPOLLIN;

ret = epoll_ctl(epfd, EPOLL_CTL_MOD, fd, &event);
if (ret)
    perror("epoll_ctl");
```

Для того чтобы удалить существующее событие для файла, связанного с fd, из экземпляра epoll_event event, нужно использовать следующий код:

```
struct epoll_event event;
int ret;

ret = epoll_ctl(epfd, EPOLL_CTL_DEL, fd, &event);
if (ret)
    perror("epoll_ctl");
```

Обратите внимание, что, когда значение op равно EPOLL_CTL_DEL, параметр event может содержать значение NULL, так как маску событий передавать системному вызову не нужно. Версии ядра до 2.6.9, однако, ошибочно проверяют данный параметр на значение NULL, считая, что он никогда не должен быть равен NULL. Для обеспечения совместимости с более старыми версиями ядра всегда передавайте допустимый не равный NULL указатель; он не будет изменен. В ядре 2.6.9 данная ошибка была исправлена.

Ожидание событий при помощи epoll

Системный вызов epoll_wait() ожидает события для файловых дескрипторов, связанных с указанным экземпляром epoll:

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

int epoll_wait (int epfd,
                struct epoll_event *events,
                int maxevents,
                int timeout);
```

Этот системный вызов ожидает события для файлов, связанных с экземпляром epoll, указанным при помощи параметра epfd, до timeout миллисекунд. В случае успеха, когда вызов возвращает результат, events указывает на область в памяти, которая содержит структуры epoll_event, описывающие каждое событие; всего до maxevents событий. Возвращаемое значение — это число событий или -1 в случае ошибки. Также, если происходит ошибка, переменной errno присваивается одно из следующих значений:

EBADF

Epfd не является допустимым дескриптором файла.

EFAULT

У процесса нет доступа на запись к области памяти, на которую указывает events.

EINTR

Системный вызов был прерван сигналом до того, как он завершил свою работу.

EINVAL

Epfd не указывает на допустимый экземпляр epoll или значение maxevents равно либо меньше нуля.

Если параметр `timeout` равен 0, то вызов возвращает результат немедленно, даже если никакие события недоступны, и тогда возвращаемое значение равно нулю. Если параметр `timeout` равен -1, то вызов не возвращает результат до тех пор, пока не появится какое-либо событие.

После того как вызов возвратил результат, в поле `events` структуры `epoll_event` содержится описание произошедших событий. В поле `data` содержится то значение, которое пользователь присвоил ему до вызова `epoll_ctl()`.

Полный пример использования `epoll_wait()` выглядит так:

```
#define MAX_EVENTS 64

struct epoll_event *events;
int nr_events, i, epfd;

events = malloc (sizeof (struct epoll_event) * MAX_EVENTS);
if (!events) {
    perror ("malloc");
    return 1;
}

nr_events = epoll_wait (epfd, events, MAX_EVENTS, -1);
if (nr_events < 0) {
    perror ("epoll_wait");
    free (events);
    return 1;
}

for (i = 0; i < nr_events; i++) {
    printf ("event=%d on fd=%d\n",
           events[i].events,
           events[i].data.fd);

    /*
     * Теперь для каждого поля events[i].events мы можем работать
     * с дескриптором в поле events[i].data.fd без блокировки.
     */
}

free (events);
```

Подробнее о функциях `malloc()` и `free()` мы поговорим в главе 8.

События, запускаемые фронтом, и события, запускаемые уровнем

Если у параметра `event`, переданного `epoll_ctl()`, поле `events` содержит значение `EPOLLET`, то наблюдение за дескриптором файла `fd` запускается фронтом, в противоположность запуску уровнем.

Рассмотрим следующие события, происходящие между производителем и потребителем, которые общаются по конвейеру Unix:

- Производитель записывает 1 Кбайт данных в конвейер.
- Потребитель делает вызов `epoll_wait()` для этого конвейера, ожидая, когда в конвейере появятся данные и, таким образом, когда конвейер станет доступным для считывания.

В случае наблюдения, запускаемого уровнем, вызов `epoll_wait()` на шаге 2 возвращает результат немедленно, показывая, что конвейер готов к считыванию. Если же наблюдение запускается фронтом, то вызов не возвращает результат до тех пор, пока не выполнится шаг 1. То есть даже если на момент, когда вызывается `epoll_wait()`, конвейер доступен для считывания, этот вызов не возвращает результат, пока в конвейер не будут записаны данные.

Запуск уровнем — это поведение вызова по умолчанию. Так работают `poll()` и `select()`, и этого поведения ожидает большинство разработчиков. Запуск фронтом требует другого подхода к программированию, частого использования ввода-вывода без блокировки и тщательной проверки на ошибку `EAGAIN`.

ПРИМЕЧАНИЕ

Эта терминология пришла из электротехники. Прерывание по уровню возникает, когда сигнал в линии достигает определенного значения. Прерывание по фронту возникает только во время нарастающего или спадающего фронта сигнала. Прерывания по уровню полезны, когда вас интересует состояние события (уровень сигнала, установленного в линии). Прерывания по фронту полезны, когда вас интересует само событие (момент установки сигнала в линии).

Отображение файлов в память

Альтернативой стандартному файловому вводу-выводу является предоставляемый ядром интерфейс, который позволяет приложению отображать файл в память, то есть создавать соответствие «один к одному» между адресом в памяти и словом в файле. Благодаря этому программист может обращаться к файлу напрямую через память точно так же, как к любым другим находящимся в памяти данным. Можно даже делать так, чтобы запись в область памяти прозрачно отображалась обратно в файл на диске.

В POSIX.1 стандартизируется, а в Linux реализуется системный вызов `mmap()`, предназначенный для отображения объектов в память. В этом разделе мы обсудим использование `mmap()` для отображения файлов в память, чтобы при помощи этой техники выполнять ввод-вывод; в главе 8 мы изучим другие приложения `mmap()`.

Системный вызов `mmap()`

Когда вы вызываете `mmap()`, этот вызов просит у ядра отобразить в память `len` байтов объекта, представляемого дескриптором файла `fd`, начиная со смещения в файле `offset` байтов. При помощи параметра `addr` можно указать предпочтительный начальный адрес в памяти. Разрешения на доступ описываются

параметром `prot`, а дополнительное поведение можно определить флагами в параметре `flags`:

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

void * mmap (void *addr,
             size_t len,
             int prot,
             int flags,
             int fd,
             off_t offset);
```

Параметр `addr` помогает ядру выбрать место в памяти, куда лучше отображать файл. Это только подсказка; большинство пользователей передают значение 0. Вызов возвращает фактический адрес в памяти, где начинается отображение.

При помощи параметра `prot` вы описываете желаемую защиту памяти в отображении. Он может содержать значение `PROT_NONE` (в этом случае доступ к страницам в отображении запрещается, что несет мало смысла) или один или несколько из следующих флагов, объединенных операцией побитового ИЛИ:

`PROT_READ`

Страницы можно считывать.

`PROT_WRITE`

Страницы можно записывать.

`PROT_EXEC`

Страницы можно выполнять.

Желаемый вариант защиты памяти не должен конфликтовать с режимом открытия файла. Например, если в программе файл открывается только для чтения, нельзя использовать в параметре `prot` флаг `PROT_WRITE`.

ФЛАГИ ЗАЩИТЫ, АРХИТЕКТУРЫ И БЕЗОПАСНОСТЬ

Хотя в POSIX определяются четыре бита защиты (чтение, запись, выполнение и «держись подальше»), некоторые архитектуры поддерживают не все четыре. Очень часто процессор не отличает друг от друга действия чтения и выполнения. В этом случае у него может быть только один флаг «считывания». В таких системах `PROT_READ` подразумевает также `PROT_EXEC`. До недавних пор подобным поведением отличалась архитектура x86.

Конечно же, нельзя полагаться на такое поведение, если вам нужна хорошая переносимость. В переносимых программах всегда нужно устанавливать флаг `PROT_EXEC`, если вы собираетесь исполнять код, соответствующий отображению.

Обратная ситуация представляет одну из причин широкого распространения атак на основе переполнения буфера: даже если в конкретном отображении не используется флаг, разрешающий выполнение, процессор тем не менее это допускает.

В новейших процессорах x86 существует бит `NX` (no-execute, не выполнять), позволяющий создавать отображения, разрешающие считывание, но запрещающие выполнение. В таких новых системах `PROT_READ` больше не подразумевает `PROT_EXEC`.

Аргумент `flags` позволяет описать тип отображения, а также некоторые элементы его поведения. Его значение включает следующие флаги, объединенные операцией побитового ИЛИ:

MAP_FIXED

Заставляет системный вызов `mmap()` считать значение параметра `addr` обязательным, а не ориентировочным. Если ядро не в состоянии поместить отображение по указанному адресу, то вызов завершается ошибкой. Если параметры адреса и длины определяют отображение так, что оно накладывается на какое-то существующее отображение, то существующие страницы удаляются, а на их место записываются новые данные отображения. Так как данный параметр требует глубокого знания адресного пространства процесса, он делает программу непереносимой и использовать его не рекомендуется.

MAP_PRIVATE

Указывает, что отображение не может использоваться совместно. Файл отображается копированием при записи, и никакие изменения, которые данный процесс вносит в память, не повторяются в фактическом файле, а также в отображениях других процессов.

MAP_SHARED

Позволяет всем процессам, отображающим тот же файл, использовать данное отображение. Запись в отображение эквивалентна записи в файл. При считывании из отображения учитываются данные, записанные в него другими процессами.

Необходимо указывать флаг либо `MAP_SHARED`, либо `MAP_PRIVATE`, но не оба. Другие более сложные флаги рассматриваются в главе 8.

Когда вы отображаете дескриптор файла, значение счетчика ссылок файла увеличивается. Таким образом, вы можете создать отображение файла, закрыть дескриптор этого файла, но у вашего процесса все так же останется возможность доступа к файлу. Соответствующее уменьшение значения счетчика ссылок выполняется после того, как отображение закрывается, либо когда процесс прекращается.

В качестве примера рассмотрим следующий фрагмент кода, который отображает `len` байтов файла, указанного при помощи дескриптора `fd`, начиная с первого байта, в область памяти, доступную только для чтения:

```
void *p;  
  
p = mmap (0, len, PROT_READ, MAP_SHARED, fd, 0);  
if (p == MAP_FAILED)  
    perror ("mmap");
```

На рис. 4.1 показано, как при создании отображения между файлом и адресным пространством процесса работают параметры `offset` и `len`, передаваемые системному вызову `mmap()`.

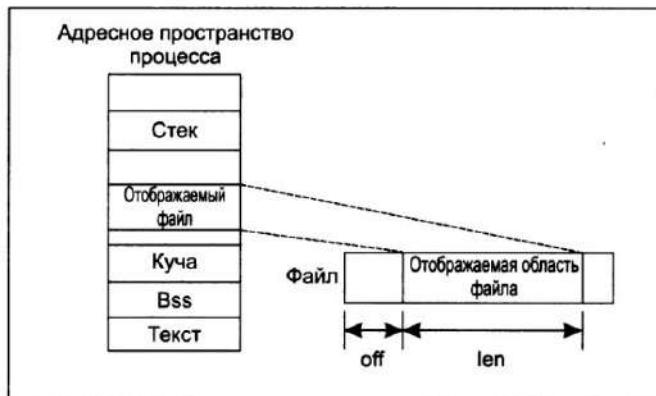


Рис. 4.1. Отображение файла в адресное пространство процесса

Размер страницы

Страница (page) – это наименьшая единица памяти, которая может обладать собственными разрешениями и собственным поведением. Следовательно, страницы представляют собой строительные кирпичики отображений в памяти, которые, в свою очередь, составляют адресное пространство процесса.

Системный вызов `mmap()` работает со страницами. Оба параметра, `addr` и `offset`, должны быть выровнены по размеру страницы. Это означает, что их значения должны быть кратными размеру страницы.

Таким образом, любое отображение включает целое число страниц. Если значение параметра `len`, предоставленное вызывающим, не выровнено по размеру страницы, например потому, что размер соответствующего файла включает не целое число страниц, то размер отображения округляется вверх до целого числа страниц. Байты в этом добавленном фрагменте памяти, между последним значащим байтом и концом отображения, заполняются нулями. Любые операции считывания из этой области возвращают нули. Никакие операции записи в эту память не влияют на содержимое отраженного файла, даже если он отображен с флагом `MAP_SHARED`. Только исходные `len` байтов записываются обратно в файл.

Метод `sysconf()`

Стандартный метод POSIX, позволяющий узнать размер страницы, – это `sysconf()`, умеющий извлекать разнообразную системную информацию:

```
#include <unistd.h>
```

```
long sysconf (int name);
```

Вызов `sysconf()` возвращает значение конфигурационного элемента `name` или значение `-1`, если `name` не является допустимым элементом. В случае ошибки вызов присваивает переменной `errno` значение `EINVAL`. Так как `-1` может быть допустимым значением некоторых элементов, рекомендуется сбрасывать значение

errno перед вызовом sysconf() и проверять ее новое значение после того, как вызов завершится.

В POSIX определяется конфигурационный элемент _SC_PAGESIZE (и его синоним _SC_PAGE_SIZE), представляющий размер страницы в байтах. Таким образом, узнать размер страницы легко:

```
long page_size = sysconf (_SC_PAGESIZE);
```

Функция getpagesize()

В Linux также есть функция получения размера страницы getpagesize():

```
#include <unistd.h>
```

```
int getpagesize (void);
```

Вызов getpagesize() возвращает размер страницы в байтах. Использовать эту функцию еще проще, чем sysconf():

```
int page_size = getpagesize ( );
```

Не все системы Unix поддерживают эту функцию; она была удалена из редакции стандарта POSIX 1003.1–2001. Здесь я привожу ее для полноты изложения.

Макрос PAGE_SIZE

Размер страницы также статически хранится в макросе PAGE_SIZE, который определяется в заголовочном файле `<asm/page.h>`. Таким образом, третий из возможных способов проверки размера страницы – это:

```
int page_size = PAGE_SIZE;
```

В отличие от первых двух вариантов, однако, в данном подходе размер системной страницы извлекается во время компиляции, а не во время выполнения. Некоторые архитектуры поддерживают несколько типов машин с разными размерами страницы, а некоторые типы машин даже поддерживают несколько размеров страницы! Каждый двоичный файл должен выполняться на всех типах машин в данной архитектуре, то есть вы однажды собираете его и он исполняется на любых других компьютерах. Жесткое кодирование размера страницы не допустило бы этого. Следовательно, определять размер страницы нужно во время исполнения. Так как значения параметров `addr` и `offset` обычно равны нулю, это требование выполнить несложно.

Более того, в будущих версиях ядра этот макрос, вероятно, не будет экспортироваться в пользовательское пространство. Я рассматриваю его в этой главе из-за того, что он часто встречается в коде Unix, но вам в ваших программах использовать его не следует. Лучший выбор среди перечисленного – sysconf().

Возвращаемые значения и коды ошибок

В случае успеха вызовов `mmap()` возвращает местоположение отображения. В случае ошибки он возвращает значение `MAP_FAILED` и соответствующим образом устанавливает переменную `errno`. Вызов `mmap()` никогда не возвращает 0.

Возможные значения переменной `errno` включают:

EACCES

Файл, указанный при помощи файлового дескриптора `fd`, не является обычным файлом или же режим, в котором файл открыт, конфликтует со значением параметра `prot` или `flags`.

EAGAIN

Файл заблокирован.

EBADF

Значение параметра `fd` не является допустимым дескриптором файла.

EINVAL

Один или несколько параметров вызова (`addr`, `len`, `offset`) имеют недопустимое значение.

ENFILE

В системе достигнуто максимальное количество открытых файлов.

ENODEV

Файловая система, в которой хранится файл, предназначенный для отображения, не поддерживает отображение в память.

ENOMEM

У процесса недостаточно памяти.

EOVERFLOW

Результат сложения значений параметров `addr` и `len` выходит за пределы адресного пространства.

EPERM

Для вызова был указан флаг `PROT_EXEC`, но файловая система подмонтирована с параметром `noexec`.

Связанные сигналы

С областями отображения связаны два сигнала:

SIGBUS

Этот сигнал генерируется, когда процесс пытается обратиться к области отображения, ставшей недействительной, например, вследствие того, что после отображения файл был усечен.

SIGSEGV

Этот сигнал генерируется, когда процесс пытается выполнить запись в область отображения, доступную только для чтения.

Системный вызов `munmap()`

В Linux существует системный вызов `munmap()`, позволяющий удалять отображения, создаваемые при помощи `mmap()`:

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

int munmap (void *addr, size_t len);
```

Вызов `munmap()` удаляет все отображения, которые содержат страницы, расположенные в области адресного пространства длиной `len` байт, начиная с адреса `addr`. Значения параметров `addr` и `len` должны быть выровнены по размеру страницы. После того как отображение удаляется, связанная с ней ранее область памяти становится недействительной и все попытки доступа приводят к созданию сигнала `SIGSEGV`.

Обычно системному вызову `munmap()` передают возвращенное предыдущим вызовом `mmap()` значение и его же параметр `len`.

В случае успеха `munmap()` возвращает 0; в случае ошибки он возвращает -1 и соответствующим образом устанавливает переменную `errno`. Единственное стандартное значение `errno` — это `EINVAL`, указывающее, что один или оба параметра вызова имеют недопустимые значения.

В качестве примера в следующем фрагменте кода удаляются все отображения, страницы которых находятся в интервале памяти `[addr,addr + len]`:

```
if (munmap (addr, len) == -1)
    perror ("munmap");
```

Пример отображения

Рассмотрим простую программу, в которой системный вызов `mmap()` используется для вывода на стандартный выход файла, указанного пользователем:

```
#include <stdio.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <fcntl.h>
#include <unistd.h>
#include <sys/mman.h>

int main (int argc, char *argv[])
{
    struct stat sb;
    off_t len;
    char *p;
    int fd;

    if (argc < 2) {
        fprintf (stderr, "usage: %s <file>\n", argv[0]);
        return 1;
    }

    fd = open (argv[1], O_RDONLY);
    if (fd == -1) {
        perror ("open");
        return 1;
    }

    if (fstat (fd, &sb) == -1) {
```

```
    perror ("fstat");
    return 1;
}

if (!S_ISREG (sb.st_mode)) {
    fprintf (stderr, "%s is not a file\n", argv[1]);
    return 1;
}

p = mmap (0, sb.st_size, PROT_READ, MAP_SHARED, fd, 0);
if (p == MAP_FAILED) {
    perror ("mmap");
    return 1;
}

if (close (fd) == -1) {
    perror ("close");
    return 1;
}

for (len = 0, len < sb.st_size; len++)
    putchar (p[len]);

if (munmap (p, sb.st_size) == -1) {
    perror ("munmap");
    return 1;
}

return 0;
}
```

В этом примере вам, возможно, незнаком только системный вызов `fstat()`, о котором мы поговорим в главе 7. Все, что вам пока нужно о нем знать, — это то, что `fstat()` возвращает информацию об указанном файле. Макрос `S_ISREG()` проверяет часть этой информации, позволяя нам перед созданием отображения гарантировать, что указанный файл является обычным файлом (а не файлом устройства или каталогом). Поведение остальных файлов при отображении зависит от соответствующих устройств. Некоторые файлы устройств поддерживают системный вызов `mmap()`, другие нет, и при попытке отобразить их возникает ошибка `EACCES`.

Все остальное в примере должно быть понятно. В качестве аргумента программе передается имя файла. Программа открывает файл, проверяет, что это обычный файл, отображает его, закрывает файл и побайтово печатает содержимое файла в стандартный вывод, после чего удаляет отображение из памяти.

Преимущества системного вызова `mmap()`

У манипулирования файлами при помощи системного вызова `mmap()` есть несколько преимуществ по сравнению со стандартными системными вызовами `read()` и `write()`. Вот некоторые из них:

- при считывании и записи в файл, отображенный в памяти, вы избегаете лишнего копирования, которое выполняется при использовании системных вызовов `read()` и `write()` (данные записываются в буфер в пользовательском пространстве и читаются из буфера);
- помимо возможных сбоев страниц памяти, операции считывания и записи в отображенный в памяти файл не влекут за собой никакой дополнительной нагрузки, которую обычно вызывают системные вызовы и переключение контекста. Это — то же самое, что простое обращение к памяти;
- когда несколько процессов отображают в память один и тот же объект, все эти процессы совместно используют данные. Отображения, доступные только для чтения и доступные для записи и совместного использования, открываются полностью; в закрытых доступных для записи отображениях открываются страницы, еще не скопированные при записи (`not-yet-COW (copy-on-write)`);
- поиск в отображениях — это обычное манипулирование указателем. Необходимости применять системный вызов `lseek()` нет.

По этим причинам системный вызов `mmap()` представляет собой разумный выбор для множества приложений.

Недостатки системного вызова `mmap()`

При использовании системного вызова `mmap()` также следует помнить о нескольких его недостатках:

- отображения в памяти всегда включают целое число страниц. Таким образом, разница между размером самого файла и целым числом страниц — это потерянное место. Когда речь идет о небольших файлах, то значительная доля отображения может пустовать. Например, если размер страницы равен 4 Кбайт, то для отображения размером 7 байт просто так расходует 4 089 байт;
- отображение в памяти не должно выходить за пределы адресного пространства процесса. В 32-битном адресном пространстве наличие очень большого количества отображений разного размера может приводить к фрагментации адресного пространства, затрудняя поиск больших свободных смежных регионов. Конечно же, эта проблема намного реже наблюдается в 64-битном адресном пространстве;
- внутри ядра создание и манипулирование отображениями в памяти и соответствующими структурами данных создает определенную нагрузку. Обычно от нее избавляются за счет устранения двойного копирования, о котором упоминалось в предыдущем разделе, особенно в отношении больших и часто используемых файлов.

По этим причинам преимущества системного вызова `mmap()` наиболее очевидны при отображении больших файлов (когда потерянное место составляет лишь небольшую долю всех отображений) или же когда общий размер отобра-

жаемого файла делится без остатка на размер страницы (и проблема потерянного места вообще не возникает).

Изменение размера отображения

В Linux предусмотрен системный вызов `mremap()`, предназначенный для увеличения или уменьшения размера указанного отображения. Эта функция уникальна для Linux:

```
#define _GNU_SOURCE

#include <unistd.h>
#include <sys/mman.h>

void * mremap (void *addr, size_t old_size,
               size_t new_size, unsigned long flags);
```

Вызов `mremap()` расширяет или сжимает отображение в области памяти `[addr,addr+old_size]` до нового размера `new_size`. В зависимости от наличия места в адресном пространстве процесса и значения параметра `flags`, ядро может одновременно переместить отображение.

ПРИМЕЧАНИЕ

Открывающаяся квадратная скобка в выражении `[addr,addr+old_size]` указывает, что область начинается с меньшего адреса и включает его, а закрывающаяся круглая скобка указывает, что область заканчивается прямо перед большим адресом и не включает его. Это условное обозначение называется обозначением интервала (interval notation).

Параметр `flags` может содержать либо значение 0, либо значение `REMAP_MAYMOVE`, указывающее, что при необходимости, для того чтобы выполнить запрошенное изменение размера, ядро может переместить отображение. Изменение размера в большом масштабе с большей вероятностью завершится успехом, если позволить ядру перемещать отображение.

Возвращаемые значения и коды ошибок

В случае успеха системный вызов `mremap()` возвращает указатель на новое отображение памяти с измененным размером. В случае ошибки он возвращает значение `MAP_FAILED` и соответствующим образом устанавливает переменную `errno`, присваивая ей одно из следующих значений:

EAGAIN

Область памяти заблокирована, и изменить ее размер невозможно.

EFAULT

Некоторые страницы в данном диапазоне не являются допустимыми страницами в адресном пространстве процесса, или же возникла проблема при повторном отображении указанных страниц.

EINVAL

Недопустимый аргумент.

ENOMEM

Указанный диапазон невозможно расширить без перемещения (флаг `MREMAP_MAYMOVE` не был передан системному вызову) или же в адресном пространстве процесса недостаточно свободного места.

В таких библиотеках, как `glibc`, системный вызов `mremap()` часто применяется для эффективной реализации `realloc()` — интерфейса, предназначенного для изменения размера блока памяти, первоначально полученного при помощи `malloc()`. Например:

```
void * realloc (void *addr, size_t len)
{
    size_t old_size = look_up_mapping_size (addr);
    void *p;

    p = mremap (addr, old_size, len, MREMAP_MAYMOVE);
    if (p == MAP_FAILED)
        return NULL;
    return p;
}
```

Этот код работает только в том случае, если память, выделенная `malloc()`, представляет собой уникальные анонимные отображения; тем не менее он служит полезным примером возможного выигрыша в производительности. В данном примере подразумевается, что программист написал также функцию `look_up_mapping_size()`.

В библиотеке GNU C системный вызов `mmap()` и его семейство применяются для выполнения некоторых задач распределения памяти. Мы глубже изучим эту тему в главе 8.

Изменение защиты отображения

В POSIX определяется интерфейс `mprotect()`, позволяющий программам менять разрешения существующих областей памяти:

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

int mprotect (const void *addr,
              size_t len,
              int prot);
```

Вызов `mprotect()` меняет режим защиты страниц памяти в интервале `[addr, addr+len]`, где значение `addr` должно быть выровнено по размеру страницы. Параметр `prot` может принимать те же значения, что и для вызова `mmap()`: `PROT_NONE`, `PROT_READ`, `PROT_WRITE` и `PROT_EXEC`. Это не аддитивные значения: если область памяти доступна для считывания, а значение параметра `prot` равно только `PROT_WRITE`, то данный вызов делает область доступной только для записи.

В некоторых системах `mprotect()` работает только с отображениями в памяти, ранее созданными при помощи `mmap()`. В Linux `mprotect()` может обрабатывать любые области в памяти.

Возвращаемые значения и коды ошибок

В случае успеха системный вызов `mprotect()` возвращает значение 0. В случае ошибки он возвращает -1 и присваивает переменной `errno` одно из следующих значений:

EACCES

Невозможно присвоить области памяти разрешения, запрошенные при помощи параметра `prot`. Причиной этого может быть, например, попытка сделать доступной для записи область памяти, в которую отображен файл, открытый только для чтения.

EINVAL

Недопустимое значение параметра `addr` или оно не выровнено по размеру страницы.

ENOMEM

Недостаточно памяти ядра для выполнения запроса либо одна или несколько страниц в данной области памяти не являются допустимыми составляющими адресного пространства процесса.

Синхронизация файла с отображением

В POSIX определен эквивалент системного вызова `fsync()`, о котором мы говорили в главе 2, но предназначенный для обработки отображений в памяти:

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

```
int msync(void *addr, size_t len, int flags);
```

Вызов `msync()` сбрасывает на диск все изменения, внесенные в файл, отображенный при помощи `mmap()`, синхронизируя отображенный файл с его отображением. А именно на диск записывается файл или фрагмент файла, связанный с отображением по адресу в памяти `addr`, занимающим `len` байтов. Значение аргумента `addr` должно быть выровнено по размеру страницы; обычно этому аргументу присваивают значение, возвращенное предыдущим вызовом `mmap()`.

Без использования `msync()` невозможно гарантировать, что грязное отображение будет записано на диск до того, как связь между файлом и отображением будет удалена. Это отличается от поведения системного вызова `write()`, когда буфер загрязняется в ходе процесса записи и ставится в очередь для сброса содержимого на диск. Когда вы записываете данные в отображение в памяти, процесс напрямую модифицирует страницы файла в страничном кэше ядра, не вовлекая в обработку само ядро. Нет никакой гарантии, что ядро сразу же синхронизирует страничный кэш и диск.

Параметр `flags` управляет поведением операции синхронизации. Он может включать следующие значения, объединенные операцией побитового ИЛИ:

MS_ASYNC

Указывает, что синхронизация должна выполняться асинхронно. Обновление вносится в расписание, но вызов `msync()` возвращает результат немедленно, не ожидая, пока данные будут фактически записаны на диск.

MS_INVALIDATE

Указывает, что все остальные кэшированные копии отображения можно аннулировать. Все будущие операции доступа к любым отображениям данного файла будут отражать только что синхронизированное содержимое файла на диске.

MS_SYNC

Указывает, что синхронизация должна выполняться синхронно. Вызов `msync()` не возвращает результат до тех пор, пока все страницы не записываются на диск.

Необходимо указывать флаг либо `MS_ASYNC`, либо `MS_SYNC`, но не оба.

Использовать этот вызов просто:

```
if (msync(addr, len, MS_ASYNC) == -1)
    perror("msync");
```

В этом примере файл, отображенный в область `[addr..addr+len]`, асинхронно синхронизируется (быстро произнесите это словосочетание десять раз подряд) с содержимым на диске.

Возвращаемые значения и коды ошибок

В случае успеха системный вызов `msync()` возвращает значение 0. В случае ошибки он возвращает -1 и присваивает переменной `errno` одно из следующих допустимых значений:

EINVAL

Параметр `flags` включает оба значения, `MS_SYNC` и `MS_ASYNC`; установлен бит, отличный от трех допустимых флагов, или значение `addr` не выровнено по размеру страницы.

ENOMEM

Указанная область памяти (или ее часть) не входит в отображение. Обратите внимание, что Linux вернет значение `ENOMEM`, как диктует POSIX, если попытаться синхронизировать область, отображения в которой удалены лишь частично, но при этом все равно синхронизирует действующие отображения в этой области.

До версии ядра Linux 2.4.19 системный вызов `msync()` возвращал значение `EFAULT` вместо `ENOMEM`.

Добавление советов относительно отображения

В Linux существует системный вызов с именем `madvise()`, позволяющий процессам давать ядру советы и подсказки относительно того, как они собираются использовать отображение. Благодаря этой информации ядро может оптимизировать свои действия, основываясь на предназначении отображения. Хотя ядро Linux динамически подстраивает свое поведение и обычно обеспечивает опти-

мальную производительность, не нуждаясь в явных советах, добавление подобных советов помогает гарантировать желаемое поведение кэширования и упреждающего считывания для некоторых типов рабочей нагрузки.

При помощи вызова `madvise()` вы сообщаете ядру, как вести себя по отношению к страницам, отображение которых в памяти начинается с адреса `addr` и занимает `len` байтов:

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

int madvise (void *addr,
             size_t len,
             int advice);
```

Если значение аргумента `len` равно 0, то ядро распространяет данный совет на все отображение, начиная с адреса `addr`. Параметр `advice` содержит сам совет и может принимать следующие значения:

MADV_NORMAL

Приложение не дает никаких явных советов для данного диапазона памяти. Память должна обрабатываться, как обычно.

MADV_RANDOM

Приложение будет обращаться к страницам в указанном диапазоне случайнym образом (не последовательно).

MADV_SEQUENTIAL

Приложение будет обращаться к страницам в указанном диапазоне последовательно, начиная с низших и заканчивая высшими адресами.

MADV_WILLNEED

Приложение обратится к страницам в указанном диапазоне в ближайшем будущем.

MADV_DONTNEED

Приложение не планирует обращаться к страницам в указанном диапазоне в ближайшем будущем.

Фактические варианты изменения поведения, которые ядро выбирает в ответ на получаемые советы, зависят от реализации: в POSIX определяются только значения советов, а не возможные последствия. Текущее ядро 2.6 реагирует на значения `advice` следующим образом:

MADV_NORMAL

Ядро работает как обычно, объем упреждающего считывания умеренный.

MADV_RANDOM

Ядро отключает упреждающее считывание и читает только минимальный объем данных во время каждой операции физического считывания.

MADV_SEQUENTIAL

Ядро выполняет агрессивное упреждающее считывание.

MADV_WILLNEED

Ядро инициирует упреждающее считывание, читая указанные страницы в память.

MADV_DONTNEED

Ядро освобождает все ресурсы, связанные с указанными страницами, и удаляет все грязные и еще не синхронизированные страницы. При последующем доступе к отображенным данным эти данныечитываются из файла на диске в страницы памяти.

Типичный способ применения системного вызова:

```
int ret:  
  
ret = madvise (addr, len, MADV_SEQUENTIAL);  
if (ret < 0)  
    perror ("madvise").
```

Этот вызов сообщает ядру, что процесс планирует обращаться к области памяти [addr,addr+len) последовательно.

УПРЕЖДАЮЩЕЕ СЧИТЫВАНИЕ

Когда ядро Linux считывает файлы с диска, оно выполняет оптимизацию, известную под названием упреждающего считывания (readahead). То есть, когда делается запрос на считывание определенного фрагмента файла, ядро также считывает следующий фрагмент. Если после этого делается запрос на считывание второго фрагмента — а так происходит, когда файл читается последовательно, — то ядро возвращает необходимые данные немедленно. Так как у дисков есть собственные буферы отслеживания (проще говоря, жесткие диски также выполняют собственное упреждающее считывание) и так как файлы обычно записываются на диск последовательно, стоимость такой оптимизации очень невелика.

Небольшой объем упреждающего считывания чаще всего оказывается очень полезным, но оптимальный результат зависит от того, как много данных считывается заранее. Если доступ к файлу осуществляется последовательно, то лучше использовать большое окно упреждающего считывания, а при случайному доступе к файлу нагрузка, создаваемая упреждающим считыванием, только мешает.

Как говорилось в разделе «Внутреннее устройство ядра» в главе 2, ядро динамически меняет размер окна упреждающего считывания в ответ на коэффициент попадания в это окно. Чем больше успешных попаданий, тем благоприятнее увеличение размера окна; если попаданий мало, это означает, что размер окна нужно уменьшить. Системный вызов `madvise()` позволяет приложениям без промедления регулировать размер окна.

Возвращаемые значения и коды ошибок

В случае успеха системный вызов `madvise()` возвращает значение 0. В случае ошибки он возвращает -1 и соответствующим образом устанавливает переменную `errno`. Возможны следующие значения ошибки:

EAGAIN

Недоступен внутренний ресурс ядра (вероятно, память). Процесс может повторить попытку.

EBADF

Область существует, но в нее не отображен файл.

EINVAL

Параметр `len` имеет отрицательное значение; значение параметра `addr` не выровнено по размеру страницы; недопустимый параметр `advice`, или страницы заблокированы либо открыты для совместного использования при помощи `MADV_DONTNEED`.

EIO

Произошла внутренняя ошибка, связанная с флагом `MADV_WILLNEED`.

ENOMEM

Указанная область не является допустимым отображением в адресном пространстве данного процесса, либо был передан флаг `MADV_WILLNEED`, но памяти оказалось недостаточно для того, чтобы поместить данные в страницы в указанной области.

Советы относительно обычного файлового ввода-вывода

В предыдущем подразделе мы рассмотрели, как предоставлять ядру советы относительно использования отображений в памяти. В этом разделе мы изучим, как сообщать ядру дополнительную информацию, касающуюся обычного файлового ввода-вывода. В Linux есть два интерфейса для подобных советов: `posix_fadvise()` и `readahead()`.

Системный вызов `posix_fadvise()`

Первый интерфейс, как свидетельствует его название, стандартизован в POSIX 1003.1-2003:

```
#include <fcntl.h>
```

```
int posix_fadvise (int fd,
                   off_t offset,
                   off_t len,
                   int advice);
```

Вызов `posix_fadvise()` дает ядру подсказку относительно дескриптора файла `fd` в интервале `[offset, offset+len]`. Сама подсказка содержится в параметре `advice`. Если значение `len` равно 0, то подсказка применяется к диапазону `[offset, длина файла]`. Чаще всего подсказку дают для всего файла, передавая для параметров `len` и `offset` значение 0.

Возможные значения параметра `advice` схожи с теми, которые используются для системного вызова `madvise()`. Этот параметр должен иметь в точности одно из следующих значений:

`POSIX_FADV_NORMAL`

Приложение не дает никакого конкретного совета для данного диапазона файла. Файл должен обрабатываться, как обычно.

POSIX_FADV_RANDOM

Приложение собирается обращаться к данным в указанном диапазоне случайным образом (не последовательно).

POSIX_FADV_SEQUENTIAL

Приложение собирается обращаться к данным в указанном диапазоне последовательно, начиная с низших и до высших адресов.

POSIX_FADV_WILLNEED

Приложение собирается обращаться к данным в указанном диапазоне в ближайшем будущем.

POSIX_FADV_NOREUSE

Приложение обратится к данным в указанном диапазоне в ближайшем будущем, но только один раз.

POSIX_FADV_DONTNEED

Приложение не собирается обращаться к страницам в указанном диапазоне в ближайшем будущем.

Как и в случае `madvice()`, фактическая реакция на каждый совет зависит от реализации — даже разные версии ядра Linux могут выполнять разные действия. В настоящее время реализованы следующие ответные действия:

POSIX_FADV_NORMAL

Ядро работает как обычно, выполняя средний объем упреждающего считывания.

POSIX_FADV_RANDOM

Ядро отключает упреждающее считывание, читая только минимальный объем данных при каждой операции физического считывания.

POSIX_FADV_SEQUENTIAL

Ядро выполняет агрессивное упреждающее считывание, удваивая размер окна упреждающего считывания.

POSIX_FADV_WILLNEED

Ядро инициирует упреждающее считывание, чтобы начать чтение указанных страниц в память.

POSIX_FADV_NOREUSE

В настоящее время поведение такое же, как для `POSIX_FADV_WILLNEED`; будущие версии ядра могут выполнять дополнительную оптимизацию, извлекая преимущества из однократного обращения. У этой подсказки нет аналога для системного вызова `madvice()`.

POSIX_FADV_DONTNEED

Ядро удаляет из страничного кэша все кэшированные данные, попадающие в указанный диапазон. Обратите внимание, что реакция ядра на эту подсказку отличается от реакции на аналогичную подсказку для `madvice()`.

В качестве примера следующий фрагмент кода сообщает ядру, что ко всему файлу, представляемому дескриптором `fd`, доступ будет осуществляться случайным образом:

```
int ret;  
  
ret = posix_fadvise (fd, 0, 0, POSIX_FADV_RANDOM);  
if (ret == -1)  
    perror ("posix_fadvise").
```

Возвращаемые значения и коды ошибок

В случае успеха системный вызов `posix_fadvise()` возвращает значение 0. В случае ошибки он возвращает `-1` и присваивает переменной `errno` одно из следующих значений:

`EBADF`

Указан недопустимый дескриптор файла.

`EINVAL`

Недопустимое значение параметра `advice`; указанный дескриптор файла ссылается на конвейер, или указанный совет невозможно применить к данному файлу.

Системный вызов `readahead()`

Системный вызов `posix_fadvise()` впервые появился в версии ядра Linux 2.6. До этого поведение, аналогичное подсказке `POSIX_FADV_WILLNEED`, обеспечивало системный вызов `readahead()`. В отличие от `posix_fadvise()`, `readahead()` — это уникальный для Linux интерфейс:

```
#include <fcntl.h>  
  
ssize_t readahead (int fd,  
                  off64_t offset,  
                  size_t count);
```

Вызов `readahead()` помещает в страничный кэш область `[offset..offset+count)` для дескриптора файла `fd`.

Возвращаемые значения и коды ошибок

В случае успеха системный вызов `readahead()` возвращает значение 0. В случае ошибки он возвращает `-1` и присваивает переменной `errno` одно из следующих значений:

`EBADF`

Указан недопустимый дескриптор файла.

`EINVAL`

Указанному дескриптору файла соответствует файл, не поддерживающий упреждающее считывание.

Советы стоят недорого

В нескольких распространенных ситуациях рабочей нагрузки можно получить большой выигрыш в производительности, если передать ядру небольшой хорошо продуманный совет. Благодаря такому совету можно добиться хороших результатов в уменьшении нагрузки, создаваемой вводом-выводом. Учитывая, что жесткие диски работают медленно, а современные процессоры — чрезвычайно быстро, любая помощь оказывается кстати, а удачный совет может иметь огромное значение.

Перед тем как читать фрагмент файла, процесс может передать ядру подсказку `POSIX_FADV_WILLNEED`, приказав ему считать файл в страничный кэш. Ввод-вывод будет выполняться асинхронно, в фоновом режиме. Когда приложение в конечном итоге обратится к файлу, операция завершится без создания блокировки ввода-вывода.

И наоборот, после считывания или записи большого фрагмента данных, например, во время потоковой записи видео на диск, процесс может передать подсказку `POSIX_FADV_DONTNEED`, приказывая ядру удалить данный фрагмент файла из дискового кэша. Большая потоковая операция может беспрерывно заполнять страничный кэш. Если приложение никогда больше не собирается обращаться к данным, это означает, что кэш наполняется ненужными данными, вероятно, за счет потенциально более полезных данных. Таким образом, приложение, работающее с потоковым видео, должно периодически запрашивать удаление потоковых данных из кэша.

Процесс, намеревающийся считать файл целиком, может передать подсказку `POSIX_FADV_SEQUENTIAL`, приказывая ядру выполнять агрессивное упреждающее считывание. И наоборот, процесс, который знает, что доступ к файлу будет осуществляться случайным образом, с переходами вперед и назад, может передать подсказку `POSIX_FADV_RANDOM`, чтобы ядро не выполняло упреждающее считывание, потому что это будет всего лишь ненужной дополнительной нагрузкой.

Синхронизированные, синхронные и асинхронные операции

В системах Unix термины «синхронизированный», «не синхронизированный», «синхронный» и «асинхронный» широко применяются, несмотря на тот факт, что это может приводить в замешательство: в английском языке, например, слова «синхронный» и «синхронизированный» практически не различаются!

Синхронная (*synchronous*) операция записи не возвращает результат до тех пор, пока записанные данные, по крайней мере, не сохраняются в буферном кэше ядра. Синхронная операция чтения не возвращает результат до тех пор, пока считанные данные не сохраняются в буфере в пользовательском пространстве, предоставленном приложением. С другой стороны, *асинхронная* (*asynchronous*) операция записи может возвращать результат еще до того, как данные покидают пользовательское пространство, а асинхронная операция чтения

может возвращать результат до того, как считанные данные становятся доступными. То есть операция фактически может ставиться в очередь и выполняться когда-то в будущем. Конечно же, в таком случае должен существовать какой-то механизм для определения, завершилась ли операция и с каким успехом.

Синхронизированная (*synchronized*) операция более ограничивающая и безопасная, чем просто синхронная операция. Синхронизированная операция записи сбрасывает данные на диск, гарантируя, что данные, находящиеся на диске, всегда синхронизированы по отношению к соответствующим буферам ядра. Синхронизированная операция чтения всегда возвращает самую новую копию данных, предположительно прямо с диска.

Короче говоря, термины «синхронный» и «асинхронный» описывают, ожидает ли операция ввода-вывода какого-то события (например, сохранения данных) перед тем, как вернуть результат. Термины «синхронизированный» и «не синхронизированный», однако, указывают, *какое именно* событие должно произойти (скажем, запись данных на диск).

Чаще всего операции записи в Unix синхронные и не синхронизированные, операции чтения синхронные и синхронизированные¹. Для большинства операций возможны любые комбинации этих характеристик, как показано в табл. 4.1.

Таблица 4.1. Синхронность операций записи

	Синхронизированные	Не синхронизированные
Синхронные	Операции записи не возвращают результат до тех пор, пока данные не сбрасываются на диск. Именно так они работают, если файл открывается с использованием флага O_SYNC	Операции записи не возвращают результат до тех пор, пока данные не сохраняются в буферах ядра. Это обычное поведение
Асинхронные	Операции записи возвращают результат, как только запрос ставится в очередь. После того как операция записи в конечном итоге завершается, данные гарантированно оказываются на диске	Операции записи возвращают результат сразу же, как только запрос ставится в очередь. После того как операция записи в конечном итоге завершается, данные гарантированно оказываются, по крайней мере, в буферах ядра

Операции чтения всегда синхронизированы, так как чтение устаревших данных не имеет смысла. Однако они могут быть либо синхронными, либо асинхронными, как показано в табл. 4.2.

¹ Технически операции чтения не являются синхронизированными, так же как и операции записи, но ядро гарантирует, что страничный кэш всегда содержит новейшие данные. Это означает, что данные в страничном кэше всегда идентичны или даже новее данных на диске. Таким образом, *на практике* поведение операций всегда синхронизировано. Действительно, сложно найти причины для иной реализации.

Таблица 4.2. Синхронность операций чтения

Синхронизированные	
Синхронные	Операции чтения не возвращают результат до тех пор, пока данные (новейшие) не сохраняются в предоставленном буфере (это обычное поведение)
Асинхронные	Операции чтения возвращают результат сразу же после того, как запрос ставится в очередь, но когда операция чтения в конечном итоге завершается, возвращаются самые новые данные

В главе 2 мы узнали, как сделать операцию записи синхронизированной (при помощи флага `O_SYNC`) и как гарантировать, что в каждый конкретный момент все операции ввода-вывода будут синхронизированными (при помощи системного вызова `fsync()` и ему подобных). Давайте теперь посмотрим, как можно реализовывать асинхронные операции чтения и записи.

Асинхронный ввод-вывод

Асинхронный ввод-вывод требует поддержки ядра на низших уровнях. В POSIX 1003.1-2003 определяются два интерфейса `aio`, которые, к счастью, реализованы в Linux. Библиотека `aio` предоставляет семейство функций для выполнения асинхронного ввода-вывода и получения уведомлений о завершении операций:

```
#include <aio.h>

/* Блок управления асинхронным вводом-выводом */
struct aiocb {
    int aio_filedes;           /* дескриптор файла */
    int aio_lio_opcode;        /* выполняемая операция */
    int aio_reqprio;           /* смещение приоритета запроса */
    volatile void *aio_buf;     /* указатель на буфер */
    size_t aio_nbytes;          /* длина операции */
    struct sigevent aio_sigevent; /* номер и значение сигнала */

    /* Далее идут внутренние закрытые члены... */
};

int aio_read (struct aiocb *aiocbp);
int aio_write (struct aiocb *aiocbp);
int aio_error (const struct aiocb *aiocbp);
int aio_return (struct aiocb *aiocbp);
int aio_cancel (int fd, struct aiocb *aiocbp);
int aio_fsync (int op, struct aiocb *aiocbp);
int aio_suspend (const struct aiocb * const cblist[],
                 int n,
                 const struct timespec *timeout);
```

Асинхронный ввод-вывод на основе потоков

Linux поддерживает интерфейс `aio` только для файлов, открытых с флагом `O_DIRECT`. Для того чтобы осуществлять асинхронный ввод-вывод на обычных

файлах, открытых без флага `O_DIRECT`, приходится применять собственные решения. Без поддержки ядра можно надеяться лишь на некое приближение к асинхронному вводу-выводу, которое даст результат, похожий на настоящий.

Преде всего, нужно понять, когда разработчику приложения может требоваться асинхронный ввод-вывод:

- для ввода-вывода без блокировки;
- для разделения актов постановки запроса на ввод-вывод в очередь, передачи запроса ядру и получения уведомления о завершении операции.

Первый пункт относится к вопросу производительности. Если операции ввода-вывода никогда не блокируются, то нагрузка, создаваемая вводом-выводом, стремится к нулю и не накладывает никаких ограничений на процесс. Второй пункт связан с вопросом процедуры, это просто другой способ обработки ввода-вывода.

Наиболее распространенный путь достижения этих целей основан на использовании потоков выполнения (вопросы составления расписаний и выполнения по расписанию обсуждаются в главах 5 и 6). Данный подход включает следующие задачи программирования:

1. Создайте пул «рабочих потоков выполнения» для обработки всего ввода-вывода.
2. Реализуйте набор интерфейсов для помещения операций ввода-вывода в рабочую очередь.
3. Сделайте так, чтобы каждый интерфейс возвращал дескриптор ввода-вывода, уникальным образом идентифицирующий связанную с ним операцию ввода-вывода. Каждый рабочий поток должен захватывать запрос ввода-вывода из головы очереди и передавать его на выполнение, а затем ожидать его завершения.
4. После завершения помещайте результаты операций (возвращаемые значения, коды ошибок и считанные данные) в очередь результатов.
5. Реализуйте набор интерфейсов для извлечения из очереди результатов информации о состоянии, используя дескрипторы ввода-вывода для идентификации каждой операции.

Такая реализация обеспечивает поведение, аналогичное интерфейсам `POSIX aio`, но за счет увеличения нагрузки, связанной с управлением потоками.

Планировщики ввода-вывода и производительность ввода-вывода

В современной системе относительный разрыв в производительности между дисками и остальными составляющими системы довольно велик — и продолжает увеличиваться. Худший компонент производительности диска — это процесс

перемещения головки чтения/записи с одной части диска на другую; эта операция называется подводом головки. В мире, где множество операций совершается всего лишь за несколько циклов процессора (каждый из которых может длиться треть наносекунды), один подвод головки может требовать в среднем более восьми миллисекунд. Конечно же, это небольшое значение, но оно в 25 миллионов раз больше, чем длительность одного цикла процессора!

Учитывая несоответствие производительности жесткого диска и остальных составляющих системы, было бы совершенно непродуманно и неэффективно отправлять диску запросы ввода-вывода в том порядке, в котором они отправляются приложениями. Поэтому в современных операционных системах применяются планировщики ввода-вывода, работой которых является минимизация числа и масштаба подводов головки. Они делают это, настраивая порядок обслуживания запросов ввода-вывода, и также время их реализации. Планировщики ввода-вывода проделывают огромную работу, уменьшая потери производительности, связанные с доступом к диску.

Адресация диска

Для того чтобы понять роль планировщика ввода-вывода, необходимо ознакомиться с некоторой общей информацией. Жесткие диски адресуют данные, используя знакомую геометрическую схему: *цилиндры* (cylinder), *головки* (head) и *секторы* (sector), называемую адресацией CHS (CHS addressing). Жесткий диск составлен из нескольких «тарелок» (platter), каждая из которых состоит из одного диска, шпинделя и головки чтения/записи. Можно представлять себе каждую тарелку как компакт-диск, а набор тарелок — как стопку компакт-дисков. Каждая тарелка делится на круговые *дорожки* (track), такие же, как на компакт-диске. Каждая дорожка затем подразделяется на целое число *секторов* (sector).

Для того чтобы найти на диске определенный блок данных, логике диска необходимо знать три значения: цилиндр, головку и сектор. Значение цилиндра указывает на дорожку, на которой записаны данные. Если положить тарелки одна на другую, то дорожки с одинаковым номером сформируют цилиндр, проходящий через все тарелки. Другими словами, цилиндр составляют дорожки, находящиеся на одинаковом расстоянии от центра каждого диска. Значение головки идентифицирует определенную головку чтения/записи (и, таким образом, конкретную тарелку). С учетом этих данных поиск сужается до одной дорожки на одной тарелке. После этого диск при помощи значения сектора идентифицирует конкретный сектор на дорожке. Таким образом, поиск завершается: жесткому диску известно, на какой тарелке, дорожке и в каком секторе искать данные. Он может установить головку чтения/записи на правильную тарелку и правильную дорожку и прочитать информацию или записать ее в указанный сектор.

К счастью, современные жесткие диски не заставляют компьютеры общаться с ними в терминах цилиндров, головок и секторов. Вместо этого теперешние жесткие диски присваивают каждой тройке, состоящей из цилиндра, головки

и сектора, уникальный *номер блока* (block number, также называемый *физическим блоком* (physical block) или *блоком устройства* (device block)). Фактически, блок отображается на определенный сектор. Современные операционные системы адресуют жесткие диски с использованием номеров блока — этот процесс называется *логической адресацией блоков* (logical block addressing, LBA), а жесткий диск сам переводит номер блока в правильный адрес CHS¹. Хотя это не гарантируется, отображение блоков в адреса CHS обычно оказывается последовательным: физический блок n чаще всего располагается на диске рядом с физическим блоком $n + 1$. Такое последовательное отображение очень важно, как вы увидите далее.

Файловые системы, однако же, существуют только в программном обеспечении. Они работают в своих единицах измерения, называемых *логическими блоками* (logical block; также иногда используется название *блок файловой системы* (filesystem block) или, чтобы запутать ситуацию, просто *блок* (block)). Размер логического блока должен быть кратен размеру физического блока. Другими словами, логические блоки файловой системы отображаются на один или несколько физических блоков диска.

ЖИЗНЬ ПЛАНИРОВЩИКА ВВОДА-ВЫВОДА

Планировщики ввода-вывода выполняют две основные операции: объединение и сортировку. *Объединение* (merging) — это процесс, в ходе которого два или более последовательных запросов ввода-вывода объединяются в один запрос. Рассмотрим два запроса: один — на чтение с дискового блока 5, а другой — на чтение с дисковых блоков с 6-го по 7-й. Эти запросы можно объединить в один запрос на чтение с дисковых блоков с 5-го по 7-й. Общий объем ввода-вывода может остаться тем же самым, но число операций ввода-вывода сократится наполовину.

Сортировка (sorting), наиболее важная из двух операций, — это процесс упорядочения ожидающих выполнения запросов ввода-вывода по возрастанию номеров блока. Например, если есть операции ввода-вывода на блоках 52, 109 и 7, планировщик ввода-вывода отсортирует эти запросы в следующем порядке: 7, 52, 109. Если после этого будет сделан запрос к блоку 81, то он будет вставлен между запросами к блокам 52 и 109. В итоге планировщик ввода-вывода отправляет запросы на диск в том порядке, в котором они находятся в очереди: 7, затем 52, затем 81 и, наконец, 109.

Такое поведение позволяет минимизировать перемещения головки диска. Вместо потенциально бессистемного перемещения — туда и обратно по всему диску — головка передвигается плавно и линейно. Так как подвод головки представляет собой самую дорогую составляющую дискового ввода-вывода, за счет этого значительно повышается производительность ввода-вывода.

¹ Ограничения на абсолютное значение этого номера блока несут основную ответственность за разнообразные ограничения на общий размер диска, существовавшие в течение многих лет.

Помощь при считывании

Все запросы на чтение должны возвращать новейшие данные. Это означает, что если запрошенные данные находятся не в страничном кэше, то процесс чтения должен блокироваться до момента, когда данные можно будет считывать с диска, — потенциально длительная операция. Такое падение производительности называется *задержкой чтения* (*read latency*).

Обычное приложение может инициировать несколько запросов на чтение за короткий период времени. Так как каждый запрос синхронизируется индивидуально, более поздние запросы зависят от завершения более ранних. Рассмотрим пример со считыванием всех файлов в каталоге. Приложение открывает первый файл, читает фрагмент этого файла, ожидает данные, затем читает еще один фрагмент и так далее до тех пор, пока не считает весь файл. После этого приложение повторяет то же самое на следующем файле. Запросы становятся сериализованными: каждый очередной запрос невозможно выполнить до тех пор, пока не завершится текущий.

Это поведение абсолютно противоположно поведению запросов на запись, которым (в своем обычном, не синхронизированном состоянии) не приходится инициировать дисковый ввод-вывод, так как он выполняется в какой-то момент времени в будущем. Таким образом, с точки зрения приложения из пользовательского пространства запросы на запись выполняются в потоковом режиме, не обремененные проблемами производительности диска. Такое потоковое поведение только усугубляет проблему считывания: потоковая запись может требовать больше внимания со стороны ядра и диска. Этот феномен зовется проблемой *записи, истощающей чтение* (*writes-starving-reads*).

Если бы планировщик ввода-вывода всегда сортировал новые запросы, вставляя их в очередь в порядке возрастания номеров блока, то запросы к блокам, находящимся далеко на диске, задерживались бы на неопределенное время. Вернемся к нашему предыдущему примеру. Если бы новые запросы постоянно обращались бы к блокам, скажем, с 50 по 60, то запрос к блоку 109 никогда не был бы обслужен. Так как задержка чтения критична для производительности, такое поведение значительно снижало бы эффективность системы. Таким образом, планировщики ввода-вывода применяют механизм предотвращения зависания запросов.

Простой подход, например такой, который был реализован в планировщике ввода-вывода ядра Linux 2.4, Linus Elevator¹, заключается в том, чтобы просто прекращать сортировку вставкой в момент, когда в очереди оказывается достаточно много старых запросов. За счет небольшого снижения общей производительности достигается честность обработки запросов и, в случае запросов на чтение, уменьшается задержка. Однако проблема этого варианта в излишней

¹ Да, планировщик ввода-вывода назвали в честь человека. Планировщики ввода-вывода иногда называются алгоритмами лифта (*elevator algorithm*), так как они решают ту же проблему, что и в задаче плавного перемещения лифта.

упрощенности эвристики. Вследствие этого в ядре Linux 2.6 на место Linus Elevator пришли несколько новых планировщиков ввода-вывода.

Планировщик ввода-вывода Deadline

Планировщик ввода-вывода Deadline был придуман для того, чтобы решить проблемы планировщика ввода-вывода версии 2.4 и традиционных алгоритмов лифта в целом. Linus Elevator ведет отсортированный список ожидающих выполнения запросов ввода-вывода. Обслуживаются те запросы ввода-вывода, которые находятся в голове очереди. Планировщик ввода-вывода Deadline также поддерживает такую очередь, но с небольшим усовершенствованием в виде двух дополнительных очередей: *очереди FIFO на чтение* и *очереди FIFO на запись*. Элементы в каждой из этих очередей сортируются по времени поступления (то есть первый прибывший обслуживается первым). Очередь FIFO на чтение, как подразумевает ее название, содержит только запросы чтения. Очередь FIFO на запись, аналогично, содержит только запросы записи. Каждому запросу в очередях FIFO присваивается значение срока действия. В очереди FIFO на чтение срок действия равен 500 миллисекундам, а в очереди FIFO на запись составляет пять секунд.

Когда подается новый запрос ввода-вывода, он вставляется в подходящее место в стандартную очередь и помещается в хвост соответствующей очереди FIFO (на чтение или на запись). Обычно жесткому диску отправляются запросы ввода-вывода из стандартной отсортированной очереди. Это максимизирует глобальную пропускную способность за счет минимизации количества подводов головки, так как нормальная очередь сортируется по номеру блока (как и для планировщика Linus Elevator).

Когда возраст элемента в голове одной из очередей FIFO превышает срок действия запросов в этой очереди, планировщик ввода-вывода прекращает отправку запросов ввода-вывода из стандартной очереди и начинает обслуживать запросы в этой очереди — обрабатывается запрос в голове очереди FIFO, а также еще парочка для ровного счета. Планировщику ввода-вывода нужно проверять и обрабатывать только запросы в голове очереди, так как они самые старые.

Таким образом, планировщик ввода-вывода Deadline может определять не-жесткий крайний срок существования запросов. Хотя он не гарантирует, что запрос ввода-вывода будет обслужен до окончания его срока действия, планировщик ввода-вывода чаще всего обрабатывает все запросы приблизительно около этого момента. Таким образом, планировщик ввода-вывода Deadline обеспечивает хорошую общую пропускную способность, не заставляя никакие запросы ожидать непозволительно долгое время. Так как срок действия запросов на чтение меньше, проблема записи, истощающей чтение, максимально смягчается.

Планировщик ввода-вывода Anticipatory

Планировщик ввода-вывода Deadline работает хорошо, но не идеально. Вспомните обсуждение зависимости запросов чтения. В случае с планировщиком Deadline первый запрос на чтение в последовательности таких запросов

обслуживается быстро, до истечения его срока действия или приблизительно в этот момент, после чего планировщик ввода-вывода возвращается к обслуживанию запросов ввода-вывода из отсортированной очереди — пока все хорошо. Но предположим, что приложение отправляет еще один запрос на чтение. В конце концов, его срок действия также истечет и планировщик ввода-вывода выполнит операцию считывания с диска: головка будет подведена к нужному месту для того, чтобы обработать запрос на чтение, а затем обратно туда, где обрабатываются запросы из отсортированной очереди. Подвод головки туда и обратно может продолжаться в течение некоторого времени, так как подобное поведение демонстрирует множество приложений. Хотя задержка остается на минимальном уровне, глобальная пропускная способность не очень высокая, поскольку запросы на чтение продолжают поступать и головку приходится подводить к разным местам диска, чтобы обрабатывать их. Производительность можно было бы повысить, если бы диск мог взять перерыв и подождать еще одного считывания, не возвращаясь к обслуживанию отсортированной очереди. Но, к сожалению, к тому моменту, когда приложение отправляет очередной зависимый запрос на чтение, планировщик ввода-вывода уже переключает передачу.

Проблема опять упирается в зависимые запросы на чтение — каждый очередной запрос на чтение отправляется только после того, как предыдущий возвращает результат, но к тому времени, как приложение получает считанные данные, ему разрешается продолжать выполнение и оно делает следующий запрос на чтение, планировщик ввода-вывода уже убегает далеко вперед и начинает обслуживать другие запросы. Это приводит к лишней паре подводов головки для каждого считывания: диск подводит головку к нужным данным, обслуживает запрос, а затем возвращает головку. Если бы планировщик ввода-вывода мог как-то знать — *предвидеть*, что в ту же часть диска скоро будет отправлен еще один запрос на чтение, то, вместо того чтобы переводить головку туда и обратно, он мог бы просто подождать следующего запроса. Экономия этих ужасных операций подвода головки определенно стоила бы нескольких миллисекунд ожидания.

Именно так работает планировщик ввода-вывода *Anticipatory*. Он начал свое существование как планировщик ввода-вывода *Deadline*, но ему подарили дополнительный механизм предвидения. Когда приложение отправляет запрос на чтение, планировщик ввода-вывода *Anticipatory* обслуживает его до завершения срока действия, как обычно. Но, в отличие от планировщика ввода-вывода *Deadline*, планировщик *Anticipatory* после этого до шести миллисекунд ничего не предпринимает и просто ожидает. Велика вероятность того, что в течение этих шести миллисекунд приложение отправит еще один запрос на чтение в ту же часть файловой системы. В таком случае этот запрос обслуживается мгновенно и планировщик ввода-вывода *Anticipatory* ждет еще какое-то время. Если шесть миллисекунд проходят, а новый запрос на чтение не приходит, то планировщик ввода-вывода *Anticipatory* полагает, что ждать больше нечего, и возвращается к тому, чем он занимался до этого (то есть к обслуживанию стандартной отсортированной очереди). Даже если правильно предугадывается

среднее число запросов, то экономится огромное количество времени — два дорогих подвода головки на каждый запрос. Так как большинство запросов на чтение зависимы, предугадывание большую часть времени окапает себя.

Планировщик ввода-вывода CFQ

Планировщик ввода-вывода CFQ (Complete Fair Queuing — абсолютно справедливое обслуживание очередей) достигает тех же целей, но применяя другой подход¹. С использованием CFQ каждому процессу назначается собственная очередь и каждой очереди присваивается квант времени. Планировщик ввода-вывода посещает все очереди по кругу, обслуживая запросы из одной очереди до тех пор, пока не истечет ее квант времени или пока запросы не закончатся. В последнем случае планировщик ввода-вывода CFQ бездействует некоторое время — по умолчанию 10 миллисекунд, ожидая новый запрос в этой очереди. Если предугадывание оправдывает себя, то планировщику ввода-вывода удается избежать подвода головки. Если нет, то время ожидания тратится впустую и планировщик переходит к очереди следующего процесса.

В очереди каждого процесса синхронизированные запросы (например, запросы на чтение) имеют приоритет по отношению к не синхронизированным запросам. Таким образом, CFQ отдает предпочтение считыванию и предотвращает проблему записи, источающей чтение. Так как для каждого процесса существует своя очередь, планировщик ввода-вывода CFQ справедливо обслуживает все процессы, в то же время обеспечивая хорошую общую производительность.

Планировщик ввода-вывода CFQ подходит для большинства типов рабочей нагрузки, что делает его превосходным первым выбором.

Планировщик ввода-вывода Noop

Планировщик ввода-вывода Noop — самый простой из существующих планировщиков. Он вообще не выполняет никакую сортировку — только базовое объединение. Этот планировщик применяется для специализированных устройств, которым не требуется собственная сортировка запросов (или на которых она выполняется отдельно).

Выбор и настройка планировщика ввода-вывода

Планировщик ввода-вывода по умолчанию можно выбирать во время загрузки при помощи параметра командной строки ядра `iosched`. Допустимые значения: `as`, `cfq`, `deadline` и `noop`. Планировщик ввода-вывода также можно выбирать во время выполнения для каждого конкретного устройства в файле `/sys/block/устройство/queue/scheduler`, где устройство — это конкретное блочное устройство. Если прочитать этот файл, то будет возвращен выбранный в данный момент

¹ Далее планировщик ввода-вывода CFQ обсуждается в том виде, как он реализован в настоящий момент. В предыдущих воплощениях планировщика не применялись кванты времени или эвристика предугадывания, но работали они схожим образом.

планировщик ввода-вывода; для выбора планировщика нужно записать в него одно из допустимых значений. Например, для того чтобы для устройства `hda` выбрать планировщик ввода-вывода `CFQ`, нужно выполнить следующую команду:

```
# echo cfq > /sys/block/hda/queue/scheduler
```

В каталоге `/sys/block/устройство/queue/iosched` находятся файлы, позволяющие администратору узнавать и устанавливать настраиваемые значения, относящиеся к планировщику ввода-вывода. Точный набор параметров зависит от текущего планировщика. Для изменения любых параметров необходимы привилегии пользователя `root`.

Хороший программист пишет программы, не зависящие от особенностей обеспечивающей их подсистемы ввода-вывода. Тем не менее знание этой подсистемы определенно помогает, если вы хотите создавать оптимальный код.

Оптимизация производительности ввода-вывода

Так как дисковый ввод-вывод выполняется медленно по отношению к производительности других компонентов системы, но одновременно является чрезвычайно важным аспектом современных вычислительных систем, максимизация производительности ввода-вывода представляет собой критическую задачу.

Минимизация операций ввода-вывода (путем объединения множества небольших операций в малое число более крупных), выравнивание операций ввода-вывода по размеру блока и применение пользовательской буферизации (см. главу 3), а также использование расширенных техник ввода-вывода, таких, как векторный ввод-вывод, позиционный ввод-вывод (см. главу 2) и асинхронный ввод-вывод, — это важные шаги, которые всегда необходимо предпринимать при написании системных программ.

В самых важных — требовательных и нагруженных операциями ввода-вывода программах, однако, можно применять дополнительные трюки максимизации производительности. Хотя ядро Linux, как говорилось выше, задействует сложные планировщики ввода-вывода для минимизации дорогих операций подвода головки, приложения из пользовательского пространства также могут вносить свой вклад, дополнительно повышая производительность.

Планирование ввода-вывода в пользовательском пространстве

Нагруженные операциями ввода-вывода приложения, отправляющие большое число запросов ввода-вывода и нуждающиеся в каждом грамме доступной производительности, могут сортировать и объединять свои ожидающие исполнения запросы ввода-вывода, выполняя ту же работу, что и планировщик ввода-вывода Linux¹.

¹ Описанные здесь техники следует применять только в критически важных приложениях с интенсивным вводом-выводом. Сортировка запросов ввода-вывода в предположении, что запросы для сортировки действительно есть, в приложении, отправляющем не так уж много подобных запросов, — это глупо и не нужно.

Но зачем же дважды выполнять одну и ту же работу, если вы знаете, что планировщик ввода-вывода отсортирует запросы в зависимости от номеров блока, минимизирует число подводов головки и сделает так, что головка будет двигаться линейно и плавно? Представьте себе приложение, отправляющее большое количество неотсортированных запросов ввода-вывода. Эти запросы прибывают в очередь планировщика ввода-вывода в случайном порядке. Планировщик делает свою работу, сортируя и объединяя запросы перед отправкой их на диск, но запросы начинают попадать на диск уже в то время, пока приложение продолжает генерировать и отправлять новые запросы. Планировщик в состоянии одновременно отсортировать только небольшой набор запросов, например несколько запросов этого приложения и еще какие-то ожидающие обработки запросы. Каждый пакет запросов приложения аккуратно сортируется, но, как только очередь заполняется, все будущие запросы исключаются из уравнения.

Таким образом, если приложение генерирует много запросов, в частности если оно собирается считывать данные, разбросанные по всему диску, то от сортировки запросов перед их отправкой оно только выигрывает, так как запросы попадают к планировщику ввода-вывода в правильном порядке.

Приложение из пользовательского пространства, однако, не обладает информацией, доступной ядру. На низших уровнях в глубине планировщика ввода-вывода запросы формулируются в терминах физических дисковых блоков. Сортировать их очень просто. Но в пользовательском пространстве запросы определяются в терминах файлов и смещений. Приложениям в пользовательском пространстве нужно зондировать почву, чтобы делать обоснованные предположения о размещении данных в файловой системе.

С учетом того, что стоит цель определения наиболее благоприятного в смысле подводов головки порядка запросов ввода-вывода к определенным файлам, у приложений из пользовательского пространства есть пара вариантов. Они могут выполнять сортировку на основе:

- полного пути;
- номера inode;
- физического дискового блока файла.

У каждого из этих вариантов свои преимущества и недостатки. Давайте вкратце рассмотрим их.

Сортировка по пути

Сортировка по полному пути — это самый простой, но и самый неэффективный способ аппроксимации сортировки с учетом блоков. Алгоритмы размещения, которые применяются в большинстве файловых систем, приводят к тому, что файлы в каждом каталоге и, следовательно, каталоги, находящиеся в одном родительском каталоге, оказываются на диске рядом друг с другом. Вероятность того, что файлы в одном каталоге были созданы приблизительно в одно время, только усиливает эту характеристику.

Таким образом, сортировка по пути грубо аппроксимирует физическое местоположение файлов на диске. Определенно, у двух файлов в одном каталоге больше шансов оказаться недалеко друг от друга, чем у двух файлов в совершенно разных частях файловой системы. Недостатком этого подхода является то, что он не учитывает фрагментацию: чем сильнее фрагментирована файловая система, тем бесполезнее оказывается сортировка по пути. Даже если игнорировать фрагментацию, сортировка по пути лишь приблизительно воспроизводит фактический порядок следования блоков. Но, с другой стороны, сортировку по пути можно применять практически во всех файловых системах. Вне зависимости от подхода к размещению файлов, сосредоточенность во времени диктует, что сортировка по пути должна обладать, как минимум, средней точностью. Помимо этого, такую сортировку легко реализовать.

Сортировка по inode

Inode – это конструкция в Unix, содержащая метаданные, связанные с каждым конкретным файлом. Хотя данные файла могут занимать несколько физических дисковых блоков, у каждого файла есть только одна структура inode, включающая такую информацию, как размер файла, разрешения, владелец и т. д. Мы подробнее поговорим о структурах inode в главе 7. Сейчас же вам нужно знать только два факта: с каждым файлом связана структура inode и у структур inode есть уникальные номера.

Сортировка по inode лучше, чем сортировка по пути, если предполагать, что следующее неравенство:

номер inode файла i < номера inode файла j

подразумевает, в целом, следующее:

номера физических блоков файла i < номеров физических блоков файла j

Определенно, это верно для Unix-подобных файловых систем, таких, как ext2 и ext3. В файловых системах, в которых не применяются фактические структуры inode, возможно все, что угодно, но номер inode (на что бы он ни отображался) все равно остается очень хорошей аппроксимацией первого порядка.

Далее показана простая программа, которая выводит номер inode указанного файла:

```
#include <stdio.h>
#include <stdlib.h>
#include <fcntl.h>
#include <sys/types.h>
#include <sys/stat.h>

/*
 * get_inode – возвращает inode файла, связанного с указанным
 * дескриптором файла, или значение -1 в случае ошибки
 */
int get_inode (int fd)
{
    struct stat buf;
    int ret;
```

```
ret = fstat (fd, &buf);
if (ret < 0) {
    perror ("fstat");
    return -1;
}

return buf.st_ino;
}
;

int main (int argc, char *argv[])
{
    int fd, inode;

    if (argc < 2) {
        fprintf (stderr, "usage: %s <file>\n", argv[0]);
        return 1;
    }

    fd = open (argv[1], O_RDONLY);
    if (fd < 0) {
        perror ("open");
        return 1;
    }

    inode = get_inode (fd);
    printf ("%d\n", inode);

    return 0;
}
```

Функцию `get_inode()` легко адаптировать для использования в собственных программах.

У сортировки по номеру `inode` есть несколько преимуществ: номер `inode` легко получить, по нему просто сортировать, и он представляет собой хорошую аппроксимацию физического размещения файлов. Основные его недостатки состоят в том, что фрагментация ухудшает аппроксимацию, что аппроксимация является всего лишь предположением и что она менее точна в не-Unix файловых системах. Тем не менее это наиболее часто используемый метод планирования запросов ввода-вывода в пользовательском пространстве.

Сортировка по физическим блокам

Лучший подход к разработке собственного алгоритма лифта, конечно же, включает сортировку по физическим дисковым блокам. Как говорилось ранее, каждый файл разбивается на логические блоки, представляющие собой наименьшие единицы размещения в файловой системе. Размер логического блока зависит от файловой системы: каждый логический блок соответствует одному физическому блоку. Таким образом, мы можем найти число логических блоков в файле, определить, каким физическим блокам они соответствуют, и выполнить сортировку на основе этих данных.

Ядро предоставляет метод, позволяющий получить физический дисковый блок, если известен номер логического блока файла. Это делается при помощи системного вызова `ioctl()`, о котором речь пойдет в главе 7, и команды `FIBMAP`:

```
ret = ioctl (fd, FIBMAP, &block);
if (ret < 0)
    perror ("ioctl");
```

Здесь `fd` представляет дескриптор файла, а `block` – логический блок, для которого мы хотим найти соответствующий физический блок. В случае успешного завершения значение `block` заменяется фактическим номером физического блока. Передаваемые вызову логические блоки индексируются с нуля и указываются по отношению к размеру файла. Это означает, что если файл составлен из восьми логических блоков, то допустимы значения от 0 до 7.

Поиск соответствия между логическими и физическими блоками, таким образом, представляет собой двухэтапный процесс. Во-первых, нужно определить число блоков в указанном файле. Это делается при помощи системного вызова `stat()`. Во-вторых, для каждого логического блока нужно выполнить запрос `ioctl()`, чтобы найти соответствующие физические блоки.

Далее показан пример программы, которая делает это для файла, передаваемого в командной строке:

```
#include <stdio.h>
#include <stdlib.h>
#include <fcntl.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <sys/ioctl.h>
#include <linux/fs.h>

/*
 * get_block – для файла, связанного с указанным дескриптором fd, возвращает
 * физический блок, соответствующий логическому блоку logical_block
 */
int get_block (int fd, int logical_block)
{
    int ret;

    ret = ioctl (fd, FIBMAP, &logical_block);
    if (ret < 0) {
        perror ("ioctl");
        return -1;
    }

    return logical_block;
}

/*
 * get_nr_blocks – возвращает число логических блоков, которые
 * занимает файл, связанный с дескриптором fd
 */
int get_nr_blocks (int fd)
{
    struct stat buf;
    int ret;

    ret = fstat (fd, &buf);
    if (ret < 0) {
```

```
    perror ("fstat");
    return -1;
}

return buf.st_blocks;
}

/*
 * print_blocks – для каждого логического блока, который занимает
 * файл, связанный с дескриптором fd, выводит на стандартный выход пару
 * значений "(логический блок, физический блок)"
 */
void print_blocks (int fd)
{
    int nr_blocks, i;

    nr_blocks = get_nr_blocks (fd);
    if (nr_blocks < 0) {
        fprintf (stderr, "get_nr_blocks failed!\n");
        return;
    }

    if (nr_blocks == 0) {
        printf ("no allocated blocks\n");
        return;
    } else if (nr_blocks == 1)
        printf ("1 block\n\n");
    else
        printf ("%d blocks\n\n", nr_blocks);

    for (i = 0; i < nr_blocks; i++) {
        int phys_block;

        phys_block = get_block (fd, i);
        if (phys_block < 0) {
            fprintf (stderr, "get_block failed!\n");
            return;
        }
        if (!phys_block)
            continue;

        printf ("(%u, %u) ", i, phys_block);
    }
    putchar ('\n');
}

int main (int argc, char *argv[])
{
    int fd;

    if (argc < 2) {
        fprintf (stderr, "usage: %s <file>\n", argv[0]);
        return 1;
    }

    fd = open (argv[1], O_RDONLY);
```

```
    if (fd < 0) {
        perror ("open");
        return 1;
    }

    print_blocks (fd);

    return 0;
}
```

Так как файлы обычно бывают непрерывными и было бы сложно (мягко говоря) сортировать все запросы ввода-вывода по логическим блокам, имеет смысл выполнять сортировку только по местоположению первого логического блока файла. Следовательно, необходимость в `get_nr_blocks()` отпадает и сортировать в приложениях можно только возвращаемые значения функции

```
get_block (fd, 0);
```

Недостатком `FIBMAP` является то, что эта команда требует характеристики `CAP_SYS_RAWIO` — фактически, привилегий пользователя `root`. Следовательно, в приложениях, не имеющих доступа к этим привилегиям, данный подход реализовать невозможно. Помимо этого, хотя команда `FIBMAP` стандартизована, ее фактическая реализация зависит от файловых систем. Наиболее распространенные системы, такие, как `ext2` и `ext3`, поддерживают ее, но какие-нибудь экзотические создания могут и не поддерживать. В такой ситуации вызов `ioctl()` возвращает ошибку `EINVAL`.

Среди преимуществ данного подхода, однако, находится то, что он возвращает фактический физический дисковый блок, в котором находится файл, что вам и необходимо для сортировки. Даже если сортировать все запросы ввода-вывода к одному файлу на основе местоположения только одного блока (планировщик ввода-вывода ядра сортирует все запросы на базе блоков), этот подход позволяет очень близко подходить к оптимальному результату. Однако требование, связанное с привилегиями пользователя `root`, конечно же, делает его невозможным для многих приложений.

Заключение

В последних трех главах мы коснулись всех аспектов файлового ввода-вывода в Linux. В главе 2 мы познакомились с основами файлового ввода-вывода в Linux — фактически с основами программирования в Unix и такими системными вызовами, как `read()`, `write()`, `open()` и `close()`. В главе 3 мы обсудили буферизацию в пользовательском пространстве и реализацию стандартной библиотеки C. В этой главе мы рассмотрели различные аспекты расширенного ввода-вывода, начиная от более мощных, но более сложных системных вызовов ввода-вывода и заканчивая техниками оптимизации и губительным для производительности подводом головки.

В следующих двух главах мы будем говорить об организации процессов: создание, разрушения и управления. Вперед!

5 Управление процессами

Как я сказал в первой главе, процессы — это вторая, после файлов, фундаментальная абстракция в системе Unix. Представляющие конечную программу, активную, живую и работающую, в ходе ее выполнения процессы — это нечто большее, чем просто язык ассемблера; они состоят из данных, ресурсов, состояния и виртуализированного процессора.

В этой главе мы рассмотрим базовые понятия, связанные с процессами: от их создания до завершения. Они пребывают в неизменном виде с первых дней существования систем Unix. Именно здесь, в тематике управления процессами, раскрываются во всей красе долговечность и прогрессивная концепция первоначального дизайна Unix. Система Unix прошла долгий путь, который редко кому удается осилить, и в ней были разделены акты создания нового процесса и загрузки нового двоичного образа. Хотя большую часть времени эти две задачи выполняются совместно, такое разделение обеспечило огромное поле для экспериментов и эволюционирования каждой из них. Эта редко применяемая концепция дожила до сегодняшнего дня, и, хотя в большинстве операционных систем предлагается один системный вызов для запуска новой программы, в Unix требуются два: fork и exec. Но перед тем как мы рассмотрим эти системные вызовы, давайте поближе познакомимся с самим процессом.

Идентификатор процесса

Каждый процесс представляется уникальным идентификатором процесса (process ID, название которого обычно сокращается просто до pid). Гарантируется, что в любой *конкретный момент времени* значение pid уникально. Это означает, что, хотя в момент времени t_0 может существовать только один процесс с pid 770 (или ни одного), нет никакой гарантии, что в момент времени t_1 не будет существовать совершенно другой процесс с pid 770. На практике обычно предполагается, что ядро не спешит повторно использовать и переназначать идентификаторы процессов, но, как вы совсем скоро узнаете, такое предположение небезопасно.

У *процесса бездействия* (*idle process*), который ядро «выполняет», когда нет никаких других доступных для выполнения процессов, идентификатор *pid* всегда равен 0. Первый процесс, который ядро выполняет после загрузки системы, называется *процессом init* — процессом инициализации, и значение *pid* для него равно 1. Обычно процесс *init* в Linux принадлежит программе *init*. Мы будем использовать термин «*init*» как для первоначального процесса, запускаемого ядром, так и для конкретной программы, которая для этого используется.

Если только пользователь явно не сообщает ядру, какой процесс нужно запустить (при помощи параметра командной строки ядра *init*), то ядру приходится самостоятельно идентифицировать подходящий процесс *init* — редкий случай, когда ядро диктует политику. Ядро Linux пробует воспользоваться четырьмя исполняемыми файлами в следующем порядке:

1. */sbin/init*: предпочтительное и наиболее вероятное местоположение процесса *init*;
2. */etc/init*: еще одно вероятное местоположение процесса *init*;
3. */bin/init*: возможное местоположение процесса *init*;
4. */bin/sh*: местоположение оболочки Bourne, которую ядро пытается запустить, если ему не удастся найти процесс *init*.

В качестве процесса инициализации запускается первый найденный из этих процессов. Если не удается запустить ни один из четырех, то ядро Linux останавливает систему с состоянием «паника».

После того как ядро запускает его, процесс инициализации обрабатывает оставшуюся часть процесса загрузки. Обычно она включает инициализацию системы, запуск различных служб и запуск программы входа в систему.

Выделение идентификатора процесса

По умолчанию ядро накладывает ограничение на максимальное значение идентификатора процесса, равное 32 768. Это необходимо для совместимости с более старыми системами Unix, в которых для идентификаторов процесса применялись меньшие 16-разрядные типы. Системный администратор может увеличить предел, установив его в файле */proc/sys/kernel/pid_max* и получив большее пространство идентификаторов *pid* в обмен на ухудшение совместимости.

Ядро выделяет процессам идентификаторы строго линейно. Если в определенный момент времени максимальное значение выделенного идентификатора *pid* равно 17, то следующим будет выделен идентификатор *pid*, равный 18, даже если на момент запуска нового процесса процесс с последним назначенным идентификатором *pid* 17 уже не будет выполняться. Ядро не начинает повторно использовать значения идентификаторов процесса до тех пор, пока не возвращается к началу — то есть меньшие значения не выделяются заново до тех пор, пока не достигается максимальное значение из */proc/sys/kernel/pid_max*. Таким образом, хотя Linux не гарантирует уникальность идентификаторов процесса в длинные периоды времени, поведение механизма выделения идентификаторов

ров обеспечивает, по крайней мере, краткосрочную уверенность в стабильности и уникальности значений pid.

Иерархия процессов

Процесс, запускающий новый процесс, называется *родительским процессом* или *предком* (parent); новый процесс называется *дочерним процессом* или *потомком* (child). Каждый процесс запускается каким-то другим процессом (за исключением, конечно же, процесса init). Таким образом, у каждого потомка есть предок. Это взаимоотношение фиксируется при помощи идентификатора предка каждого процесса (parent process ID, ppid), значение которого для каждого потомка равно pid родительского процесса.

Каждым процессом владеют *пользователь* (user) и *группа* (group). Информация о владении используется для управления правами доступа к ресурсам. Для ядра пользователи и группы — это просто целые значения, которые внутри файлов /etc/passwd и /etc/group сопоставляются с удобными для человека именами, знакомыми всем пользователям Unix, такими, как пользователь root или группа wheel (в целом ядро Linux совершенно не интересуют эти предназначенные для человека строки, и оно предпочитает идентифицировать объекты целыми числами). Каждый дочерний процесс наследует владельцев предка, то есть им владеют тот же пользователь и группа, что и его предком.

Каждый процесс также является частью *группы процессов* (process group), что просто выражает его взаимоотношения с другими процессами, — не путайте эту группу с вышеупомянутой концепцией пользователей и групп. Потомки обычно принадлежат той же группе процессов, что и их родители. Помимо этого, когда оболочка запускает конвейер (то есть когда пользователь выполняет команду ls | less), все команды в конвейере становятся членами одной группы процессов. Понятие группы процессов упрощает отправку сигналов и получение информации обо всем конвейере, а также обо всех потомках процессов в конвейере. С точки зрения пользователя, группа процессов тесно связана с понятием *задания* (job).

Тип pid_t

Программно идентификатор процесса представляется типом pid_t, который определен в файле заголовка <sys/types.h>. Точный тип C, поддерживающий идентификатор, зависит от архитектуры и не определяется ни в одном из стандартов C. Однако в Linux тип pid_t обычно основывается на типе C int.

Получение идентификатора процесса и идентификатора родительского процесса

Системный вызов getpid() возвращает идентификатор вызывающего процесса:

```
#include <sys/types.h>
#include <unistd.h>

pid_t getpid (void);
```

Системный вызов `getppid()` возвращает идентификатор процесса, являющегося предком вызывающего процесса:

```
#include <sys/types.h>
#include <unistd.h>
```

```
pid_t getppid (void);
```

Ни один из этих вызовов не может вернуть ошибку, поэтому использовать их просто:

```
printf ("Мой pid=%d\n", getpid ());
printf ("pid предка=%d\n", getppid ());
```

Откуда нам знать, что `pid_t` – это целое число со знаком? Хороший вопрос! Ответ прост: мы не знаем. Несмотря на то, что можно безопасно предполагать, что в Linux `pid_t` – это то же самое, что и `int`, такая догадка не вполне соответствует предназначению абстрактного типа и нарушает переносимость. К сожалению, как и для всех описаний типов в C, не существует простого способа вывести значение типа `pid_t` – это является частью абстракции, и технически нужно было бы использовать функцию `pid_to_int()`, которой не существует. Тем не менее чаще всего эти значения обрабатываются как целые, по крайней мере с функцией `printf()`.

Запуск процесса

В Unix акт загрузки в память и исполнения образа программы отделен от акта создания нового процесса. Один системный вызов (в действительности, один вызов из семейства вызовов) загружает двоичную программу в память, заменяя ею предыдущее содержимое адресного пространства, и начинает выполнение новой программы. Это называется *исполнением* (*executing*) новой программы, и функциональность обеспечивается семейством вызовов `exec`.

Другой системный вызов используется для создания нового процесса, который первоначально представляет собой практически копию своего родительского процесса. Часто новый процесс немедленно начинает выполнять новую программу. Акт создания нового процесса называется *ветвлением* (*forking*), и эту функциональность предоставляет системный вызов `fork()`. Таким образом, чтобы выполнить новый образ программы в новом процессе, необходимы два акта – сначала ветвление, для создания нового процесса, а затем исполнение, для загрузки нового образа в этот процесс. Сначала мы рассмотрим вызов `exec`, а затем `fork()`.

Семейство вызовов `exec`

Нет единой функции `exec`, вместо этого существует семейство функций `exec`, построенных на одном системном вызове. Давайте сначала изучим простейший из этих вызовов, `exec()`:

```
#include <unistd.h>

int exec1 (const char *path,
           const char *arg,
           ...);
```

Вызов `exec1()` заменяет текущий образ процесса новым путем загрузки в память программы, на которую указывает аргумент `path`. Параметр `arg` – это первый аргумент запускаемой программы. Многоточие указывает на переменное число аргументов – функция `exec1()` является функцией с переменным числом аргументов, то есть дополнительные аргументы можно указывать в скобках один за другим. Список аргументов необходимо завершать значением `NULL`.

Например, следующий код заменяет программу, которая исполняется в данный момент, программой `/bin/vi`:

```
int ret;

ret = exec1 ("/bin/vi", "vi", NULL);
if (ret == -1)
    perror ("exec1");
```

Обратите внимание, что я, следуя соглашениям Unix, передаю в качестве первого аргумента программы значение `vi`. Оболочка помещает последний компонент пути, `vi`, в первый аргумент, когда разветвляет или запускает процессы, поэтому программа может проверять свой первый аргумент, `argv[0]`, чтобы узнавать имя своего двоичного образа. Очень часто несколько системных утилит, которые пользователь знает под разными именами, в действительности представляют одну и ту же программу, но с жесткими ссылками для нескольких имен. Программа использует свой первый аргумент для определения нужного поведения.

Еще один пример: для того чтобы отредактировать файл `/home/kidd/hooks.txt`, можно выполнить следующий код:

```
int ret;

ret = exec1 ("/bin/vi", "vi", "/home/kidd/hooks.txt", NULL);
if (ret == -1)
    perror ("exec1");
```

Обычно системный вызов `exec1()` не возвращает значение. Успешный вызов завершается переходом к точке входа новой программы, а только что выполненный код пропадает из адресного пространства процесса. В случае ошибки, однако, `exec1()` возвращает значение `-1` и устанавливает переменную `errno`, указывая на ошибку. Я перечислю возможные значения `errno` далее в этом разделе.

Успешный вызов `exec1()` меняет не только адресное пространство и образ процесса, но и другие атрибуты процесса:

- любые ожидающие сигналы теряются;
- все сигналы, которые процесс отлавливает (см. главу 9), возвращаются к своему поведению по умолчанию, а обработчики сигналов удаляются из адресного пространства процесса;

- все блокировки памяти (см. главу 8) удаляются;
- для большинства атрибутов потока восстанавливаются значения по умолчанию;
- большая часть статистики процесса сбрасывается;
- все, связанное с памятью процесса, включая отображенные в память файлы, удаляется;
- все, что существует исключительно в пользовательском пространстве, включая особенности библиотеки C, такие, как поведение вызова `atexit()`, удаляется.

Многие свойства процесса, однако, не меняются. Например, идентификатор `pid`, приоритет и пользователь-владелец и группа-владелец остаются теми же самыми.

Обычно при работе системных вызовов из семейства `exec` открытые файлы наследуются. Это означает, что у запускаемых программ сохраняется полный доступ ко всем файлам, открытым в исходном процессе, если, конечно, им известны значения дескрипторов файлов. Однако часто такое поведение нежелательно. Обычная практика — закрывать файлы, перед тем как применять вызовы `exec`, хотя можно также заставить ядро делать это автоматически при помощи `fcntl()`.

Остальные члены семейства

Помимо `exec1()`, есть еще пять членов семейства `exec`:

```
#include <unistd.h>

int execvp (const char *file,
            const char *arg,
            ...);

int execle (const char *path,
            const char *arg,
            ...
            char * const envp[]);

int execv (const char *path, char *const argv[]);

int execvp (const char *file, char *const argv[]);

int execve (const char *filename,
            char *const argv[],
            char *const envp[]);
```

Мнемоника проста: `l` и `v` указывают, передаются аргументы в виде списка (*list*) или массива (*vector*); `p` обозначает, что поиск указанного файла осуществляется по полному пользовательскому пути (*path*). В командах, которые используют варианты с `r`, можно указывать только имя файла, если этот файл находится в пределах пользовательского пути. Наконец, `e` указывает, что для нового процесса также передается новая среда. Интересно, что, хотя техниче-

скогого обоснования для этого не существует, в семействе exec отсутствует вызов, который умел бы и искать файлы в пути, и принимать в качестве одного из аргументов новую среду. Вероятно, причиной послужило то, что варианты с *р* были реализованы для использования оболочками, а процессы, выполняемые в оболочках, обычно наследуют свои среды от оболочек.

Члены семейства exec, принимающие в качестве аргумента массив, работают практически так же, за исключением того, что нужно конструировать массив и передавать его вместо списка. Использование массива позволяет определять аргументы во время выполнения. Как и список аргументов переменной длины, массив должен заканчиваться значением NULL.

В следующем фрагменте кода вызов execvp() применяется для вызова программы *vi*, как мы уже делали ранее:

```
const char *args[] = { "vi", "/home/kidd/hooks.txt", NULL };
int ret;

ret = execvp ("vi", args);
if (ret == -1)
    perror ("execvp");
```

Если предполагать, что каталог */bin* находится в пользовательском пути, этот пример работает точно так же, как предыдущий.

В Linux только один член семейства exec является системным вызовом. Остальные — это обертки вокруг системного вызова в библиотеке C. Так как системные вызовы с переменной длиной списка аргументов в лучшем случае трудно реализовывать и так как концепция пользовательского пути существует исключительно в пользовательском пространстве, единственный кандидат на пост системного вызова — это execve(). Прототип системного вызова идентичен пользовательскому вызову.

Значения ошибок

В случае успеха системные вызовы exec не возвращают никакое значение. В случае ошибки вызовы возвращают значение *-1* и присваивают переменной *errno* одно из следующих значений:

E2BIG

Общее число байтов в предоставленном списке аргументов (*arg*) или среде (*environ*) слишком велико.

EACCES

У процесса отсутствует разрешение на поиск для компонента пути, указанного при помощи аргумента *path*; *path* не является обычным файлом; целевой файл не помечен как исполняемый; или файловая система, где находится *path* или *file*, подмонтирована с флагом *noexec*.

EFAULT

Указатель недопустим.

EIO

Произошла низкоуровневая ошибка ввода-вывода (это плохо).

EISDIR

Последний компонент в пути `path`, или интерпретатор, является каталогом.

ELOOP

При разрешении пути `path` встретилось слишком много символьических ссылок.

EMFILE

Для вызывающего процесса был достигнут предел числа открытых файлов.

ENFILE

Был достигнут системный предел числа открытых файлов.

ENOENT

Цель `path` или `file` не существует, либо не существует необходимая общая библиотека.

ENOEXEC

Цель `path` или `file` является недопустимым двоичным файлом, или же она предназначена для другой машинной архитектуры.

ENOMEM

Недостаточно памяти ядра для выполнения новой программы.

ENOTDIR

Один из компонентов в пути `path` не является каталогом.

EPERM

Файловая система, которой принадлежит `path` или `file`, подмонтирована с флагом `nosuid`, пользователь не является пользователем `root`, либо для `path` или `file` установлен бит `suid` или `sgid`.

ETXTBSY

Цель `path` или `file` открыта для записи другим процессом.

Системный вызов `fork()`

Новый процесс, выполняющий тот же образ, что и текущий, можно создать при помощи системного вызова `fork()`:

```
#include <sys/types.h>
#include <unistd.h>
```

```
pid_t fork (void);
```

Успешный вызов `fork()` создает новый процесс, практически во всех аспектах идентичный вызывающему процессу. Оба процесса продолжают выполняться, возвращаясь из `fork()` так, как будто бы ничего особенного не произошло.

Новый процесс называется «потомком» исходного процесса, который, в свою очередь, называется «предком». В потомке успешный вызов `fork()` возвращает значение 0. В предке `fork()` возвращает идентификатор `pid` потомка. Дочерний

и родительский процессы идентичны практически во всех аспектах, за исключением нескольких важных различий:

- для потомка, конечно же, выделяется новый идентификатор pid, значение которого отличается от значения pid его предка;
- для потомка в качестве pid предка устанавливается pid его родительского процесса;
- статистика ресурсов для потомка сбрасывается до нуля;
- любые ожидающие сигналы удаляются и не наследуются потомком (см. главу 9);
- никакие захваченные блокировки файлов не наследуются потомком.

В случае ошибки дочерний процесс не создается, системный вызов fork() возвращает значение -1, а переменной errno присваивается соответствующее значение. Существует два возможных значения errno с тремя возможными вариантами расшифровки:

EAGAIN

Ядру не удалось выделить определенные ресурсы, такие, как новый идентификатор pid, или же был достигнут лимит ресурсов RLIMIT_NPROC (rlimit) (см. главу 6).

ENOMEM

Недостаточно памяти ядра для завершения запроса.

Использовать этот системный вызов просто:

```
pid_t pid;  
  
pid = fork ( );  
if (pid > 0)  
    printf ("Я предок с pid=%d!\n", pid);  
else if (!pid)  
    printf ("Я младенец!\n");  
else if (pid == -1)  
    perror ("fork");
```

Чаще всего системный вызов fork() применяется для создания нового процесса, в который затем загружается новый двоичный образ, — представьте себе оболочку, в которой запускается новая программа для пользователя, или процесс, запускающий программу справки. Сначала процесс «ответствляет» новый процесс, а затем получившийся потомок начинает выполнять новый двоичный образ. Эта комбинация «fork плюс exec» часто встречается и очень проста. В следующем примере ответствуется новый процесс, выполняющий двоичный файл /bin/windlass:

```
pid_t pid;  
  
pid = fork ( );  
if (pid == -1)  
    perror ("fork");
```

```
/* потомок... */

if (!pid) {
    const char *args[] = { "windlass", NULL };
    int ret;

    ret = execv ("/bin/windlass", args);
    if (ret == -1) {
        perror ("execv");
        exit (EXIT_FAILURE);
    }
}
```

Родительский процесс продолжает выполняться без изменений, за исключением того, что теперь у него есть новый потомок. Вызов `execv()` изменяет потомок, заставляя его выполнять программу `/bin/windlass`.

Копирование при записи

В ранних системах Unix ветвление было простым, если не сказать примитивным. После вызова ядро создавало копии внутренних структур данных, дублировало записи таблицы страниц процесса, а затем выполняло постраничное копирование родительского адресного пространства в новое адресное пространство потомка. Но такое постраничное копирование было, по крайней мере, с точки зрения ядра слишком продолжительным процессом.

Поведение современных систем Unix отличается большей оптимальностью. Вместо оптового копирования родительского адресного пространства в современных системах Unix, таких, как Linux, применяются страницы *копирования при записи* (*copy-on-write*, COW).

Копирование при записи — это ленивая стратегия оптимизации, разработанная для того, чтобы смягчать нагрузку, создаваемую дублированием ресурсов. Предпосылка проста: если несколько потребителей запрашивают доступ на чтение собственных копий одного ресурса, никакого смысла в дублировании ресурса нет. Вместо этого каждому потребителю лучше просто предоставить указатель на один и тот же ресурс. До тех пор пока ни один из потребителей не пытается модифицировать свою «копию» ресурса, иллюзия эксклюзивного доступа к ресурсу сохраняется, что дает возможность избегать лишней нагрузки из-за копирования. Если потребитель все же пытается модифицировать свою копию ресурса, то в этот момент ресурс прозрачно дублируется и потребителю для модификации предоставляется копия. Потребитель, ничего не зная о происходящем за сценой, может спокойно редактировать свою копию ресурса, пока другие потребители продолжают работать с исходной, неизмененной версией. Отсюда и имя: *копирование* происходит только *при записи*.

Основное преимущество заключается в том, что если потребитель никогда не модифицирует свою копию ресурса, то потребность в копировании никогда не возникает. Общее преимущество ленивых алгоритмов — то, что они откладывают дорогостоящие действия до последнего момента, — здесь также работает.

В конкретном примере виртуальной памяти копирование при записи реализуется на базе страниц. Таким образом, если процесс не модифицирует все свое адресное пространство, целиком адресное пространство никогда не копируется. Когда ветвление завершается, предок и потомок уверены, что у каждого из них есть собственное уникальное адресное пространство, хотя в действительности они совместно используют оригинальные страницы предка, которые, в свою очередь, также могут быть предоставлены другим родительским или дочерним процессам, и так далее!

Реализация в ядре проста. В структурах данных ядра, связанных со страницами, страницы помечаются как доступные только для чтения и как доступные для копирования при записи. Когда любой из процессов пытается модифицировать страницу, происходит страничное прерывание или *ошибка из-за отсутствия страницы* (page fault). Ядро обрабатывает страничное прерывание, прозрачно создавая копию страницы; в этот момент атрибут копирования при записи для страницы очищается и она становится недоступна для совместного использования.

Так как современные машинные архитектуры обеспечивают поддержку копирования при записи на аппаратном уровне в *блоках управления памятью* (memory management unit, MMU), решение просто само по себе и легко реализуется.

У копирования при записи еще больше преимуществ в случае ветвления. Так как очень часто за ветвлением следует выполнение (exec), копирование адресного пространства предка в адресное пространство потомка обычно представляет собой пустую трату времени: если потомок сразу же приступает к выполнению нового двоичного образа, его предыдущее адресное пространство просто стирается. Копирование при записи позволяет оптимизировать этот случай.

Системный вызов vfork()

До появления страниц, реализующих копирование при записи, разработчикам Unix приходилось реализовывать трудоемкое и расточительное копирование адресного пространства во время ветвления, за которым сразу же следовало выполнение. Вследствие этого разработчики BSD представили в 3.0BSD системный вызов vfork():

```
#include <sys/types.h>
#include <unistd.h>

pid_t vfork (void);
```

При успешном выполнении вызов vfork() работает так же, как и fork(), за исключением того, что дочерний процесс должен немедленно сделать успешный вызов одной из функций exec либо завершиться, вызвав _exit() (об этом речь пойдет в следующем разделе). Системный вызов vfork() избегает копирования адресного пространства и таблиц страниц, приостанавливая родительский процесс до тех пор, пока дочерний не завершится или не начнет

выполнять новый двоичный образ. В промежутке предок и потомок совместно используют — без семантики копирования при записи — свое адресное пространство и записи в таблице страниц. Фактически, единственная работа, которая выполняется после запуска `vfork()`, — это копирование внутренних структур данных ядра. Следовательно, дочерний процесс не должен модифицировать никакую память в адресном пространстве.

Системный вызов `vfork()` — это исконное, которое никогда не должно было быть реализовано в Linux, хотя необходимо отметить, что даже с учетом копирования при записи `vfork()` работает быстрее `fork()`, так как отпадает необходимость копирования записей таблицы страниц¹. Тем не менее пришествие страниц, поддерживающих копирование при записи, ослабляет любые аргументы в пользу альтернатив `fork()`. Действительно, до версии ядра Linux 2.2.0 системный вызов `vfork()` был всего лишь оберткой вокруг `fork()`. Так как требования `vfork()` слабее, чем требования `fork()`, подобная реализация `vfork()` осуществима.

Строго говоря, ни одна реализация `vfork()` не защищена от ошибок: представьте себе ситуацию, когда вызов `exec` завершается сбоем! Предок будет находиться в замороженном состоянии неопределенное количество времени, пока потомок не сообразит, что делать дальше, или не завершится.

Завершение процесса

В POSIX и C89 определяется стандартная функция для завершения текущего процесса:

```
#include <stdlib.h>
```

```
void exit (int status);
```

Вызов `exit()` выполняет некоторые базовые шаги по завершению, а затем заставляет ядро прекратить процесс. У этой функции нет возможности возвращать ошибки — в действительности она вообще никогда не возвращает результат. Таким образом, следом за вызовом `exit()` нет смысла использовать никакие другие инструкции.

Параметр `status` используется для обозначения статуса выхода процесса. Другие программы, а также пользователь, работающий в оболочке, могут при необходимости проверять это значение. В частности, статус `status & 0377` возвращается предку. Мы узнаем, как извлекать возвращаемое значение, далее в этой главе.

Значения `EXIT_SUCCESS` и `EXIT_FAILURE` определяются в качестве способов представления успеха и неудачи, поддерживающих перенос на другие платформы.

¹Хотя в настоящее время оно не является частью ядра Linux 2.6, исправление, реализующее общие записи таблиц страниц с копированием при записи, было выпущено в списке рассылки Linux Kernel Mailing List (lkml). Если бы его не добавили, у системного вызова `vfork()` не осталось бы вообще никаких преимуществ.

мы. В Linux 0 обычно представляет успех; ненулевое значение, например 1 или -1, соответствует неудаче.

Следовательно, успешный выход определяется одной простой строкой:

```
exit (EXIT_SUCCESS);
```

Перед завершением процесса библиотека C выполняет следующие шаги в указанном порядке:

- 1) вызывает любые функции, зарегистрированные с atexit() или on_exit(), в порядке, обратном порядку регистрации (мы поговорим об этих функциях далее в этой главе);
- 2) сбрасывает все открытые стандартные потоки ввода-вывода (см. главу 3);
- 3) удаляет все временные файлы, созданные при помощи функции tmpfile().

Эти шаги завершают всю работу, которую процесс должен проделывать в пользовательском пространстве, и после этого вызов exit() делает системный вызов _exit(), давая ядру знать, что оно теперь может выполнить оставшуюся часть процесса завершения:

```
#include <unistd.h>

void _exit (int status);
```

Когда процесс завершается, ядро очищает все ресурсы, которые оно создало от имени процесса и которые более не используются. Сюда входят выделенная память, открытые файлы, семафоры System V и другие ресурсы. После очистки ядро разрушает процесс и уведомляет предка о кончине его потомка.

Приложения могут напрямую вызывать _exit(), но это редко имеет смысл: большинству приложений необходимо производить какую-то уборку, которую обеспечивает обычный выход, например сбрасывать поток stdout. Обратите внимание, однако, что после ветвления пользователи vfork() должны вызвать _exit(), а не exit().

ПРИМЕЧАНИЕ

В блестящей попытке обеспечения избыточности в стандарт ISO C99 была добавлена функция _Exit(), обладающая поведением, полностью идентичным системному вызову _exit():

```
#include <stdlib.h>

void _Exit (int status);
```

Прочие способы завершения

Классический способ завершить программу — не использовать явный системный вызов, а просто «сойти с пути». В случае языка C это происходит, когда функция main() возвращает результат. Этот подход, однако, все равно включает системный вызов: компилятор просто вставляет неявный вызов _exit() после собственного кода завершения. Хорошая практика кодирования — явно возвращать статус выхода либо через системный вызов exit(), либо через возвращаемое значение функции main(). Оболочка использует значение выхода для оценки успешности или для определения ошибок команд. Обратите внимание, что

успешное возвращение описывается при помощи вызова `exit(0)`; функция `main()` в случае успеха также возвращает значение 0.

Процесс также может завершаться, если ему отправляется сигнал, действие по умолчанию которого заключается в том, чтобы прекращать процесс. Среди таких сигналов можно назвать `SIGTERM` и `SIGKILL` (см. главу 9).

Последний способ завершить выполнение программы — вызвать гнев ядра. Ядро может убивать процессы за исполнение недопустимых инструкций, за нарушение сегментации, когда у него заканчивается память, и т. д.

Библиотечный вызов `atexit()`

В стандарте POSIX 1003.1-2001 определяется, а в Linux реализуется библиотечный вызов `atexit()`, используемый для регистрации функций, которые должны вызываться при завершении процесса:

```
#include <stdlib.h>
```

```
int atexit (void (*function)(void));
```

Успешный вызов `atexit()` регистрирует указанную функцию, чтобы она была запущена во время обычного завершения процесса, то есть когда процесс завершится при помощи системного вызова `exit()` или возврата из функции `main()`. Если процесс запускает функцию `exec`, список зарегистрированных функций очищается (так как функции более не существуют в адресном пространстве нового процесса). Если процесс завершается при помощи сигнала, то зарегистрированные функции не вызываются.

Указанная функция не принимает никакие параметры и не возвращает никакое значение. Прототип выглядит так:

```
void my_function (void);
```

Функции вызываются в порядке, обратном регистрации. Это означает, что функции хранятся в стеке и последняя попавшая в него выходит первой (*last in first out, LIFO*). Зарегистрированные функции не должны вызывать `exit()`, чтобы не начать бесконечную рекурсию. Если функция должна раньше запланированного прервать процесс завершения, то это необходимо делать с использованием системного вызова `_exit()`. Подобное поведение, однако, нежелательно, так как потенциально важные функции потом могут не запускаться.

Стандарт POSIX требует, чтобы `atexit()` поддерживала, по крайней мере, `ATEXIT_MAX` зарегистрированных функций, и это значение должно быть не менее 32. Точное максимальное значение можно получить при помощи `sysconf()` и значения `_SC_ATEXIT_MAX`:

```
long atexit_max;
```

```
atexit_max = sysconf (_SC_ATEXIT_MAX);
printf ("atexit_max=%ld\n", atexit_max);
```

В случае успеха `atexit()` возвращает значение 0. В случае ошибки она возвращает значение -1.

Простой пример:

```
#include <stdio.h>
#include <stdlib.h>

void out (void)
{
    printf ("atexit( ) succeeded!\n");
}

int main (void)
{
    if (atexit (out))
        fprintf(stderr, "atexit( ) failed!\n");

    return 0.
}
```

Функция on_exit()

В SunOS 4 определяется собственный эквивалент atexit(), и библиотека glibc в Linux поддерживает его:

```
#include <stdlib.h>

int on_exit (void (*function)(int , void *), void *arg);
```

Эта функция работает так же, как и atexit(), но прототип зарегистрированной функции отличается:

```
void my_function (int status, void *arg).
```

Аргумент `status` — это значение, передаваемое `exit()` или возвращаемое `main()`. Аргумент `arg` — это второй параметр, передаваемый `on_exit()`. Необходимо с большой внимательностью относиться к использованию этой функции, чтобы гарантировать, что на момент, когда функция в конечном итоге вызывается, в памяти, на которую указывает `arg`, содержатся допустимые данные.

Последняя версия Solaris уже не поддерживает эту функцию. Вместо нее следует использовать совместимую со стандартами функцию `atexit()`.

SIGCHLD

Когда процесс завершается, ядро отправляет сигнал SIGCHLD ее предку. По умолчанию этот сигнал игнорируется, и предок не предпринимает никакие действия. Однако процессы при необходимости могут обрабатывать этот сигнал при помощи системных вызовов `signal()` и `sigaction()`. Эти вызовы, а также остальная часть великолепного мира сигналов рассматриваются в главе 9.

Сигнал SIGCHLD может генерироваться и отправляться в любое время, так как завершение дочернего процесса выполняется асинхронно по отношению к его предку. Однако часто предок дожидается информации о завершении своего потомка или даже явно ожидает, пока не произойдет событие. Реализовать такое поведение позволяют системные вызовы, о которых речь пойдет далее.

Ожидание завершенных дочерних процессов

Получать уведомления при помощи сигналов удобно, но часто процессам-предкам требуется больше информации о завершении их потомков, например им необходимо получать возвращаемое значение.

Если бы дочерний процесс после завершения полностью исчезал, как можно было бы ожидать, то не оставалось бы никаких следов, которые мог бы изучить его родительский процесс. Следовательно, первые разработчики Unix решили, что, когда процесс умирает раньше своего предка, ядро должно переводить потомок в особое состояние. Процесс в таком состоянии называется *зомби* (zombie). В этом состоянии существует только минимальный скелет того, что раньше было процессом, — некоторые базовые структуры данных ядра, содержащие потенциально полезные данные. Процесс в этом состоянии ожидает того, что родитель запросит его статус (это называется обслуживанием (waiting on) процесса-зомби). Только после того, как предок получает сохраненную информацию о завершенном потомке, процесс формально удаляется и прекращает существовать даже в виде зомби.

Ядро Linux предоставляет несколько интерфейсов для получения информации о завершенных дочерних процессах. Простейший из этих интерфейсов, определенный в стандарте POSIX, — это `wait()`:

```
#include <sys/types.h>
#include <sys/wait.h>

pid_t wait (int *status);
```

Вызов `wait()` возвращает идентификатор `pid` завершенного дочернего процесса или значение `-1`, если произошла ошибка. Если ни один из дочерних процессов пока что не завершился, вызов блокируется до тех пор, пока это событие не произойдет. Если дочерний процесс уже завершился, то вызов возвращает результат немедленно. Следовательно, если вызвать `wait()` в ответ на сообщение о кончине потомка, — скажем, при получении сигнала `SIGCHLD`, то он вернет результат без блокировки.

В случае ошибки переменной `errno` присваивается одно из двух значений:

`ECHILD`

У вызывающего процесса нет дочерних процессов.

`EINTR`

Во время ожидания был получен сигнал, и вызов вернул результат раньше времени.

Если его значение не равно `NULL`, указатель `status` содержит дополнительную информацию о потомке. Так как POSIX допускает в реализациях определение битов `status` так, как это удобнее разработчикам, в стандарте предусматривает семейство макросов для интерпретации данного параметра:

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

int WIFEXITED (status);
int WIFSIGNALED (status);
int WIFSTOPPED (status);
int WIFCONTINUED (status);

int WEXITSTATUS (status);
int WTERMSIG (status);
int WSTOPSIG (status);
int WCOREDUMP (status);
```

Любой из первых двух макросов может возвращать значение `true` (не-нулевое значение) в зависимости от того, каким образом завершается процесс. Первый макрос, `WIFEXITED`, возвращает `true`, когда процесс завершается обычным образом, то есть когда он вызывает `_exit()`. В этом случае макрос `WEXITSTATUS` предоставляет восемь бит младших разрядов, которые были переданы `_exit()`.

Макрос `WIFSIGNALED` возвращает `true`, если завершение процесса было вызвано сигналом (подробное обсуждение сигналов вы найдете в главе 9). В этом случае `WTERMSIG` возвращает номер сигнала, приведшего к завершению, а `WCOREDUMP` возвращает `true`, если процесс сделал дамп ядра в ответ на получение сигнала. Макрос `WCOREDUMP` не определяется в стандарте POSIX, хотя многие системы Unix, включая Linux, поддерживают его.

Макросы `WIFSTOPPED` и `WIFCONTINUED` возвращают `true`, если процесс был остановлен или продолжен соответственно и в настоящий момент отслеживается при помощи системного вызова `ptrace()`. Эти условия обычно применимы только во время реализации отладчика, хотя совместно с `waitpid()` (см. следующий раздел) они также используются для реализации управления заданиями. Обычно системный вызов `wait()` применяется только для передачи информации о завершении процесса. Когда `WIFSTOPPED` возвращает значение `true`, `WSTOPSIG` позволяет получить номер сигнала, остановившего процесс. Макрос `WIFCONTINUED` не определяется стандартом POSIX, хотя будущие стандарты определяют его для `waitpid()`. В версии ядра Linux 2.6.10 данный макрос предоставляется также для `wait()`.

Давайте рассмотрим пример программы, в которой системный вызов `wait()` применяется для того, чтобы выяснить, что произошло с потомком программы:

```
#include <unistd.h>
#include <stdio.h>
#include <sys/types.h>
#include <sys/wait.h>

int main (void)
{
    int status;
    pid_t pid;

    if (!fork ( ))
        return 1;
```

```

pid = wait (&status).
if (pid == -1)
    perror ("wait");

printf ("pid=%d\n", pid);

if (WIFEXITED (status))
    printf ("Нормальное завершение со статусом выхода=%d\n",
        WEXITSTATUS (status));

if (WIFSIGNALED (status))
    printf ("Убит сигналом=%d%s\n",
        WTERMSIG (status),
        WCOREDUMP (status) ? " (дамп ядра)" : "");

if (WIFSTOPPED (status))
    printf ("Остановлен сигналом=%d\n",
        WSTOPSIG (status));

if (WIFCONTINUED (status))
    printf ("Продолжен\n");

return 0.
}

```

Эта программа отвечает дочерний процесс, который немедленно завершается. После этого родительский процесс выполняет системный вызов `wait()`, чтобы определить статус потомка. Процесс выводит `pid` потомка и сведения о его завершении. Так как в данном случае дочерний процесс завершается путем возврата из функции `main()`, мы знаем, что на выходе получим приблизительно такой результат:

```

$ ./wait
pid=8529
Нормальное завершение со статусом выхода=1

```

Если вместо обычного возврата дочернего процесса реализовать использование вызова `abort()`¹, который отправляет своему процессу сигнал `SIGABRT`, то результат будет примерно таким:

```

$ ./wait
pid=8678
Убит сигналом=6

```

Ожидание определенного процесса

Очень важно наблюдать за поведением дочерних процессов. Часто, однако, у процессов бывает несколько потомков и процессу нужно ожидать завершения не всех них, а какого-то конкретного дочернего процесса. Можно было бы несколько раз сделать системный вызов `wait()`, каждый раз проверяя возвращаемое значение. Однако это было бы крайне неудобно. Что, если бы позже вам по-

¹ Определен в заголовке `<stdlib.h>`.

надобилось проверить статус другого завершенного процесса? Родительскому процессу пришлось бы сохранять все результаты всех вызовов `wait()` на случай, если они понадобятся позже.

Если известно значение идентификатора `pid` интересующего вас процесса, то можно воспользоваться системным вызовом `waitpid()`:

```
#include <sys/types.h>
#include <sys/wait.h>

pid_t waitpid (pid_t pid, int *status, int options);
```

Вызов `waitpid()` — это более мощная версия системного вызова `wait()`. Его дополнительные параметры позволяют осуществлять более тонкую настройку.

При помощи параметра `pid` вы указываете, какой процесс или какие процессы нужно ожидать. Он может принимать значения из четырех категорий:

< -1

Ожидать любые дочерние процессы, значение идентификатора группы процессов для которых равно абсолютному значению данного параметра. Например, для того чтобы настроить ожидание любых процессов из группы процессов 500, нужно передать значение `-500`.

-1

Ожидать любые дочерние процессы. Поведение, аналогичное поведению системного вызова `wait()`.

0

Ожидать любые дочерние процессы, принадлежащие той же группе процессов, что и вызывающий процесс.

> 0

Ожидать дочерний процесс, значение идентификатора `pid` которого в точности равно значению данного параметра. Например, чтобы настроить ожидание дочернего процесса с `pid` 500, нужно передать значение 500.

Параметр `status` работает идентично единственному параметру системного вызова `wait()`, и его можно использовать с макросами, о которых говорилось выше.

В качестве значения параметра `options` можно передавать следующие значения, объединенные операцией двоичного ИЛИ, или же пустое значение:

WNOHANG

Вызов не должен блокироваться. Если уже завершенного (или остановленного, или продолженного) подходящего дочернего процесса нет, вызов должен возвращать результат немедленно.

UNTRACED

Если это значение передается системному вызову, то флаг `WIFSTOPPED` устанавливается, даже если вызывающий процесс не отслеживает дочерний процесс. Этот флаг помогает реализовывать более общее управление заданиями, как, например, в оболочке.

WCONTINUED

Если это значение передается системному вызову, то флаг WIFCONTINUED устанавливается, даже если вызывающий процесс не отслеживает дочерний процесс. Как и WUNTRACED, этот флаг полезен для реализации оболочки.

В случае успеха системный вызов `waitpid()` возвращает идентификатор pid процесса, состояние которого изменилось. Если был передан флаг `WNOHANG`, а указанный потомок или потомки еще не изменили состояние, то вызов `waitpid()` возвращает значение 0. В случае ошибки вызов возвращает -1 и присваивает переменной `errno` одно из трех значений:

ECHILD

Процесс или процессы, указанные при помощи аргумента `pid`, не существуют либо они не являются потомками вызывающего процесса.

EINTR

Параметр `WNOHANG` не был указан, и во время ожидания был получен сигнал.

EINVAL

Аргумент `options` имеет недопустимое значение.

В качестве примера предположим, что в вашей программе необходимо получать возвращаемое значение определенного дочернего процесса с pid 1742, но при этом вызов должен немедленно возвращать результат, если потомок еще не завершился. Можно реализовать это при помощи такого кода:

```
int status;
pid_t pid;

pid = waitpid(1742, &status, WNOHANG);
if (pid == -1)
    perror("waitpid");
else {
    printf("pid=%d\n", pid);

    if (WIFEXITED(status))
        printf("Обычное завершение со статусом выхода=%d\n",
               WEXITSTATUS(status));

    if (WIFSIGNALED(status))
        printf("Убит сигналом=%d%s\n",
               WTERMSIG(status),
               WCOREDUMP(status) ? "(дамп ядра)" : "");
}
```

В качестве последнего примера хочу обратить ваше внимание на то, что следующий вариант использования системного вызова `wait()`:

```
wait(&status);
```

идентичен такому варианту системного вызова `waitpid()`:

```
waitpid(-1, &status, 0);
```

Еще больше гибкости при ожидании процессов

Для приложений, в которых требуется еще более высокая гибкость функциональности ожидания дочерних процессов, расширение XSI стандарта POSIX определяет, а Linux реализует системный вызов `waitid()`:

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

```
int waitid (idtype_t idtype,
            id_t id,
            siginfo_t *infop,
            int options);
```

Как и `wait()` и `waitpid()`, системный вызов `waitid()` используется для реализации ожидания и получения информации о статусе изменения (завершение, остановка, продолжение) дочернего процесса. Он предоставляет еще больше параметров, но за счет большей сложности.

Аналогично `waitpid()`, `waitid()` позволяет разработчику указывать, какой процесс нужно ожидать. Однако при использовании `waitid()` для этого необходимо использовать не один, а два параметра. Аргументы `idtype` и `id` определяют, какие дочерние процессы ожидать, точно так же, как один аргумент `pid` системного вызова `waitpid()`. Параметр `idtype` может принимать одно из следующих значений:

P_PID

Ожидать дочерний процесс, идентификатор `pid` которого равен значению аргумента `id`.

P_GID

Ожидать дочерний процесс, идентификатор группы процессов которого равен значению аргумента `id`.

P_ALL

Ожидать все дочерние процессы; аргумент `id` игнорируется.

Аргумент `id` принадлежит редко используемому типу `id_t` — типу, представляющему общий идентификационный номер. Он используется на тот случай, если в будущих реализациях будет добавлено новое значение `idtype`, и, предположительно, дает гарантию того, что в предопределенном типе можно будет хранить новый идентификатор, так как он достаточно велик, чтобы вместить любое значение `pid_t`. В Linux разработчики могут применять его точно так же, как тип `pid_t`, например напрямую передавать значения `pid_t` или числовые константы. Педантичные программисты, конечно же, могут беспрепятственно добавлять приведение типов.

Параметр `options` включает одно или несколько следующих значений, объединенных операцией двоичного ИЛИ:

WEXITED

Вызов будет ожидать завершившиеся дочерние процессы (определенные при помощи аргументов `id` и `idtype`).

WSTOPPED

Вызов будет ожидать дочерние процессы, которые были приостановлены в ответ на получение сигнала.

WCONTINUED

Вызов будет ожидать дочерние процессы, которые продолжили выполнение в ответ на получение сигнала.

WNOHANG

Вызов никогда не будет блокироваться, и в случае если подходящие завершенные (или остановленные, или продолженные) дочерние процессы не обнаружены, он будет возвращать результат немедленно.

WNOWAIT

Вызов не будет освобождать найденный процесс из состояния зомби. Этот процесс можно будет обслужить в будущем.

В случае успеха системный вызов `waitid()` заполняет параметр `infop`, который должен указывать на допустимый тип `siginfo_t`. Точная конструкция структуры `siginfo_t` зависит от реализации¹, но после вызова `waitid()` заполненными остаются лишь несколько полей. Это означает, что при успешном вызове гарантируется, что в следующих полях будут содержаться допустимые значения:

si_pid

Идентификатор `pid` дочернего процесса.

si_uid

Идентификатор `uid` дочернего процесса.

si_code

Принимает значение `CLD_EXITED`, `CLD_KILLED`, `CLD_STOPPED` или `CLD_CONTINUED` в ответ на завершение дочернего процесса, смерть вследствие получения сигнала, остановку вследствие получения сигнала или продолжение вследствие получения сигнала соответственно.

si_signo

Принимает значение `SIGCHLD`.

si_status

Если значение `si_code` равно `CLD_EXITED`, то в этом поле содержится код выхода дочернего процесса. В противном случае в этом поле находится номер доставленного дочернему процессу сигнала, который вызвал изменение состояния.

В случае успеха системный вызов `waitid()` возвращает значение 0. В случае ошибки `waitid()` возвращает значение -1 и присваивает переменной `errno` одно из следующих значений:

¹ Действительно, в Linux структура `siginfo_t` очень сложна. Ее определение можно найти в файле `/usr/include/bits/siginfo.h`. Мы подробнее изучим эту структуру в главе 9.

ECHLD

Процесс или процессы, указанные при помощи аргументов `id` и `idtype`, не существуют.

EINTR

Значение параметра `options` не включает флаг `WNOHANG`, и выполнение было прервано сигналом.

EINVAL

Аргумент `options` или комбинация аргументов `id` и `idtype` недопустимы.

Функция `waitid()` предоставляет дополнительную полезную семантику, отсутствующую в системных вызовах `wait()` и `waitpid()`. В частности, информация, которую можно получить из структуры `siginfo_t`, часто оказывается довольно ценной. Если же необходимости в подобной информации нет, то имеет смысл использовать более простые функции, поддерживаемые в широком диапазоне систем и, таким образом, переносимые на большее число не относящихся к Linux систем.

BSD говорит свое слово: `wait3()` и `wait4()`

Тогда как системный вызов `waitpid()` — это наследник версии System V Release 4 от AT&T, BSD идет собственным путем и предоставляет две другие функции, которые используются для ожидания изменения состояния дочернего процесса:

```
#include <sys/types.h>
#include <sys/time.h>
#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);
```

Цифры 3 и 4 указывают, что эти две функции — это версии системного вызова `wait()` с тремя и четырьмя параметрами соответственно.

Функции работают аналогично системному вызову `waitpid()`, за исключением того, что у них есть еще один аргумент — `rusage`. Следующий вызов `wait3()`:

```
pid = wait3 (status, options, NULL);
```

эквивалентен такому вызову `waitpid()`:

```
pid = waitpid (-1, status, options);
```

А такой вызов `wait4()`:

```
pid = wait4 (pid, status, options, NULL);
```

эквивалентен следующему вызову `waitpid()`:

```
pid = waitpid (pid, status, options);
```

Таким образом, `wait3()` ожидает изменения состояния любого потомка, а `wait4()` ожидает изменения состояния определенного потомка, указанного при помощи параметра `pid`. Аргумент `options` работает так же, как в системном вызове `waitpid()`.

Как упоминалось ранее, большое различие между этими вызовами и вызовом `waitpid()` заключается в аргументе `rusage`. Если его значение не равно `NULL`, функция заполняет память, на которую указывает `rusage`, информацией о дочернем процессе. Следующая структура предоставляет сведения об использовании ресурсов потомком:

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

struct rusage {
    struct timeval ru_utime; /* затраченное пользовательское время */
    struct timeval ru_stime; /* затраченное системное время */
    long ru_maxrss; /* максимальный размер резидентной части */
    long ru_ixrss; /* размер общей памяти */
    long ru_idrss; /* размер собственных данных */
    long ru_isrss; /* размер собственного стека */
    long ru_minflt; /* восстановления страниц */
    long ru_majflt; /* страничные прерывания */
    long ru_nswap; /* операции подкачки */
    long ru_inblock; /* блочные операции ввода */
    long ru_oublock; /* блочные операции вывода */
    long ru_msgsnd; /* отправленные сообщения */
    long ru_msgrcv; /* полученные сообщения */
    long ru_nsigvals; /* полученные сигналы */
    long ru_nvcsw; /* добровольные переключения контекста */
    long ru_nivcsw; /* вынужденные переключения контекста */
};
```

Мы поговорим об использовании ресурсов далее в этой главе.

В случае успеха эти функции возвращают идентификатор `pid` процесса, который поменял состояние. В случае ошибки они возвращают значение `-1` и присваивают переменной `errno` одно из значений ошибки, которые может возвращать системный вызов `waitpid()`.

Так как `wait3()` и `wait4()` не определены в стандарте POSIX¹, рекомендуется применять их только тогда, когда информация об использовании ресурсов критически важна. Несмотря на отсутствие стандартизации POSIX, однако, практически все системы Unix поддерживают эти два вызова.

Запуск и ожидание нового процесса

В обоих стандартах – ANSI C и POSIX – определяется интерфейс, объединяющий запуск нового процесса и ожидание его завершения; можно представлять

¹ Вызов `wait3()` входил в первоначальную спецификацию Single UNIX Specification, но был исключен из нее.

себе это как синхронное создание процессов. Если процесс ответвляет дочерний процесс только для того, чтобы сразу начать ожидать его завершение, имеет смысл использовать следующий интерфейс:

```
#define _XOPEN_SOURCE /* если необходим WEXITSTATUS и т.д. */
#include <stdlib.h>
```

```
int system (const char *command);
```

Функция `system()` носит такое имя, потому что синхронный запуск процессов называют *выходом в систему* (*shelling out to the system*). Очень часто `system()` применяют для выполнения простых утилит или сценариев оболочки, обычно с явной целью просто получить их возвращаемое значение.

Когда вы вызываете функцию `system()`, запускается команда, переданная при помощи параметра `command`, включая любые дополнительные аргументы. Параметр `command` добавляется к аргументам `/bin/sh -c`. В этом смысле параметр полностью передается оболочке.

В случае успеха возвращаемым значением является статус возвращения команды в том виде, как он предоставляется вызовом `wait()`. Следовательно, код выхода выполненной команды захватывается при помощи макроса `WEXITSTATUS`. Если вызов `/bin/sh` завершается ошибкой, то значение, передаваемое `WEXITSTATUS`, совпадает со значением, возвращенным `exit(127)`. Так как вызываемая команда также может возвращать код 127, безошибочного метода проверки того, кто вернул ошибку — команда или оболочка, — не существует. В случае ошибки вызов возвращает значение `-1`.

Если значение параметра `command` равно `NULL`, `system()` возвращает ненулевое значение в случае, если оболочка `/bin/sh` доступна, и 0 в противном случае.

Во время выполнения команды сигнал `SIGCHLD` блокируется, а сигналы `SIGINT` и `SIGQUIT` игнорируются. У игнорирования сигналов `SIGINT` и `SIGQUIT` есть несколько последствий, в частности когда `system()` вызывается внутри цикла. В этом случае необходимо гарантировать, что программа правильно проверяет статус выхода дочернего процесса, например:

```
do {
    int ret.

    ret = system ("pidof rudder");
    if (WIFSIGNALED (ret) &&
        (WTERMSIG (ret) == SIGINT ||
         WTERMSIG (ret) == SIGQUIT))
        break; /* или другой вариант обработки */
} while (1).
```

Реализация функции `system()` с использованием `fork()`, функции из семейства `exec`, а также системного вызова `waitpid()` — это полезное упражнение. Вы должны попытаться выполнить его самостоятельно, так как оно объединяет множество концепций из этой главы. Но для завершенности я приведу пример простой реализации:

```

/*
 * my_system – синхронно ответвляет дочерний процесс и ожидает команды
 * "/bin/sh -c <cmd>".
 *
 * Возвращает -1 в случае любой ошибки или код выхода запущенного процесса.
 * Не блокирует и не игнорирует сигналы.
 */
int my_system (const char *cmd)
{
    int status;
    pid_t pid;

    pid = fork ();
    if (pid == -1)
        return -1;
    else if (pid == 0) {
        const char *argv[4];

        argv[0] = "sh";
        argv[1] = "-c";
        argv[2] = cmd;
        argv[3] = NULL;
        execv ("/bin/sh", argv);

        exit (-1);
    }

    if (waitpid (pid, &status, 0) == -1)
        return -1;
    else if (WIFEXITED (status))
        return WEXITSTATUS (status);

    return -1;
}

```

Обратите внимание, что, в отличие от официальной функции `system()`, в этом примере не блокируются и не отключаются никакие сигналы. Это может быть к лучшему или к худшему, в зависимости от конкретной программы, но, по крайней мере, не блокировать SIGING зачастую бывает удобно, так как это позволяет прерывать запущенную программу так, как пользователь ожидает в обычных ситуациях. В более сложной реализации можно было бы добавить дополнительные указатели в качестве параметров, которые, когда их значение отлично от `NULL`, указывали бы на ошибки, в настоящее время неразличимые. Например, это могли бы быть `fork_failed` и `shell_failed`.

Зомби

Как уже говорилось выше, завершенный, но еще не обслуженный своим предком процесс называется зомби. Процессы-зомби продолжают потреблять системные ресурсы, хотя и в минимальном количестве — только для поддержания скелета того, чем они когда-то были. Эти ресурсы остаются для того, чтобы процессы-предки, желающие проверить статус своих потомков, могли получить

информацию, относящуюся к жизни и завершению этих процессов. Как только предок получает эти сведения, ядро окончательно удаляет процесс и зомби прекращает существовать.

Однако любой, кто работает с Unix, время от времени замечает процессы-зомби. Эти процессы, часто называемые *привидениями* (*ghost*), — дети безответственных родителей. Если ваше приложение ответвляется дочерний процесс, то ответственностью именно вашего приложения (если только это не кратковременное приложение, о чем мы поговорим через секунду) является обслуживание потомка, даже если оно заключается только в удалении собранной информации. В противном случае все потомки ваших процессов будут превращаться в привидения и продолжать существовать, заполняя списки процессов системы и генерируя отвращение к нелепой реализации вашего приложения.

Что же происходит, если родительский процесс умирает раньше своего потомка или же если он умирает раньше, чем получает возможность обслужить своих потомков-зомби? Когда процесс завершается, ядро Linux проходит по списку его потомков и переназначает их предка, делая их дочерними процессами процесса инициализации (*init process* — процесс со значением *pid*, равным 1). Это гарантирует, что ни один процесс не остается сиротой без предка. Процесс инициализации, в свою очередь, периодически обслуживает все свои потомки, гарантуя, что ни один из них не останется зомби слишком долго, — и никаких привидений! Таким образом, если предок умирает раньше своих потомков или не обслуживает свои потомки до того, как завершится, дочерние процессы, в конечном итоге, становятся потомками процесса инициализации, который их и обслуживает, позволяя им полностью завершиться. Хотя обслуживание всех дочерних процессов до сих пор считается хорошей практикой реализации процессов, вышеописанная предосторожность означает, что краткосрочным процессам не нужно излишне беспокоиться об этом.

Пользователи и группы

Как упоминалось ранее в этой главе и как обсуждалось в главе 1, процессы связываются с пользователями и группами. Идентификаторы пользователя и группы — это числовые значения, предоставляемые типами *C uid_t* и *gid_t* соответственно. Отображение числовых значений в удобные для человека имена — например, у пользователя *root* значение *uid* равно 0 — выполняется в пользовательском пространстве с помощью файлов */etc/passwd* и */etc/group*. Ядро работает только с числами.

В системе Linux идентификаторы пользователя и группы для процесса определяют, какие операции позволяет выполнять процессу. Таким образом, процессы должны выполняться от имени подходящих пользователей и групп. Многие процессы работают от имени пользователя *root*. Однако лучшие практики разработки программного обеспечения поддерживают доктрину прав с наименьшими привилегиями, что означает, что процесс должен работать с минимальным из возможных уровнем прав. Это требование динамично: если процессу

требуются привилегии пользователя `root`, для того чтобы выполнить какую-то операцию в начале своей жизни, но после этого необходимость в расширенных правах отпадает, он должен избавляться от привилегий `root` как можно скорее. С этой целью многие процессы — в частности, те, которым для выполнения определенных операций необходимы привилегии пользователя `root`, — зачастую подделывают свои идентификаторы пользователя или группы.

Перед тем как мы посмотрим, как такое возможно, нужно разобраться со сложностями идентификаторов пользователя и группы.

Реальные, действительные и сохраненные идентификаторы пользователя и группы

ПРИМЕЧАНИЕ

В следующем обсуждении мы будем фокусироваться на идентификаторах пользователя, но для идентификаторов групп ситуация абсолютно та же.

В действительности с процессом связан не один, а четыре идентификатора пользователя: реальный, действительный, сохраненный и идентификатор файловой системы. *Реальный идентификатор пользователя* (real user ID) — это uid пользователя, который первоначально запустил процесс. Он равен реальному идентификатору пользователя для родительского процесса данного процесса и не меняется во время вызова `exec`. Обычно процесс входа в систему устанавливает для оболочки входа в систему каждого пользователя идентификатор, равный реальному идентификатору этого пользователя, и все процессы пользователя продолжают нести этот идентификатор. Суперпользователь (`root`) может менять реальный идентификатор пользователя на любые значения, но никакому другому пользователю это неподвластно.

Действительный идентификатор пользователя (effective user ID) — это идентификатор пользователя, которым в настоящий момент владеет процесс. Проверки разрешений обычно основываются на этом значении. Первоначально данный идентификатор равен реальному идентификатору пользователя, потому что, когда процесс ветвится, действительный идентификатор пользователя родительского процесса наследуется его потомком. Помимо этого, когда процесс делает вызов `exec`, действительный идентификатор пользователя обычно не меняется. Но именно во время вызова `exec` возникает ключевое различие между реальным и действительным идентификаторами: выполняя двоичный файл `setuid` (`suid`), процесс может менять действительный идентификатор пользователя. Точнее, в качестве действительного идентификатора пользователя устанавливается идентификатор пользователя владельца файла программы. Например, так как файл `/usr/bin/passwd` является файлом `setuid` и его владелец — пользователь `root`, когда оболочка обычного пользователя ответвляет процесс, выполняющий этот файл, этот процесс получает действительный идентификатор пользователя `root`, независимо от того, что за пользователь в действительности его выполняет.

Непrivилегированные пользователи могут устанавливать в качестве действительного идентификатора пользователя реальный или сохраненный идентификатор пользователя, как вы через мгновение увидите. Суперпользователь может устанавливать в качестве действительного идентификатора пользователя любое значение.

Сохраненный идентификатор пользователя (*saved user ID*) — это исходное значение действительного идентификатора пользователя для процесса. Когда процесс ветвится, потомок наследует сохраненный идентификатор пользователя своего предка. Во время вызова `exec`, однако, ядро устанавливает в качестве сохраненного идентификатора пользователя действительный идентификатор пользователя, таким образом сохраняя информацию о действительном идентификаторе пользователя на момент вызова `exec`. Непrivилегированные пользователи не могут менять сохраненный идентификатор пользователя; суперпользователь может присвоить ему значение реального идентификатора пользователя.

В чем смысл всех этих значений? Действительный идентификатор пользователя — это самое главное значение; это идентификатор пользователя, который проверяется во время подтверждения достоверности реквизитов процесса. Реальный и сохраненный идентификаторы пользователя играют роль заместителей или потенциальных значений идентификатора пользователя, которые позволено использовать процессам, не принадлежащим пользователю `root`. Реальный идентификатор пользователя — это действительный идентификатор пользователя, принадлежащий пользователю, фактически выполняющему программу, а сохраненный идентификатор пользователя — это действительный идентификатор пользователя, каким он был до того, как двоичный файл `suid` привел к изменению его значения во время вызова `exec`.

Изменение реального и сохраненного идентификатора пользователя или группы

Идентификаторы пользователя и группы устанавливаются при помощи двух системных вызовов:

```
#include <sys/types.h>
#include <unistd.h>
```

```
int setuid (uid_t uid);
int setgid (gid_t gid);
```

Вызов `setuid()` устанавливает действительный идентификатор пользователя текущего процесса. Если текущий действительный идентификатор пользователя для процесса равен 0 (пользователь `root`), то также устанавливаются реальный и сохраненный идентификаторы пользователя. Пользователь `root` имеет возможность передавать в параметре `uid` любые значения, таким образом устанавливая в качестве всех трех идентификаторов пользователя значение `uid`. Остальные пользователи, кроме `root`, могут передавать в качестве параметра `uid` только реальный или сохраненный идентификатор пользователя. Другими

словами, пользователь-не root может присваивать действительному идентификатору пользователя только одно из этих двух значений.

В случае успеха системный вызов setuid() возвращает значение 0. В случае ошибки он возвращает -1 и присваивает переменной errno одно из следующих значений:

EAGAIN

Значение uid отличается от реального идентификатора пользователя, и если присвоить реальному идентификатору пользователя значение uid, то он выйдет за предел NPROC rlimit (указывающий максимальное число процессов, которые могут принадлежать пользователю).

EPERM

Пользователь не является пользователем root, а значение uid не равно значению действительного или сохраненного идентификатора пользователя.

Все вышесказанное также распространяется на группы — просто замените setuid() на setgid(), а uid на gid.

Изменение действительного идентификатора пользователя или группы

В Linux есть две одобренные стандартом POSIX функции для установки действительных идентификаторов пользователя и группы для выполняющегося в данный момент процесса:

```
#include <sys/types.h>
#include <unistd.h>

int seteuid (uid_t euid);
int setegid (gid_t egid);
```

Вызов seteuid() устанавливает в качестве действительного идентификатора пользователя значение параметра euid. Пользователь root может в качестве euid передавать любые значения. Остальные пользователи, кроме пользователя root, могут устанавливать в качестве действительного идентификатора пользователя только реальный или сохраненный идентификатор пользователя. В случае успеха seteuid() возвращает значение 0. В случае ошибки этот вызов возвращает -1 и присваивает переменной errno значение EPERM, указывающее, что владельцем текущего процесса является не root и что значение euid не равно реальному или сохраненному идентификатору пользователя.

Обратите внимание, что для пользователей, не являющихся пользователем root, функции seteuid() и setuid() работают одинаково. Таким образом, стандартная практика и хорошая привычка всегда использовать функцию seteuid(), если только ваш процесс не будет выполняться от имени пользователя root, в такой ситуации setuid() имеют большие смысла.

Все вышесказанное также распространяется на группы — просто замените seteuid() на setegid(), а euid на egid.

Изменение идентификаторов пользователя и группы, стиль BSD

Разработчики BSD остановились на собственных интерфейсах для установки идентификаторов пользователя и группы. В Linux эти интерфейсы предоставляются для обеспечения совместимости:

```
#include <sys/types.h>
#include <unistd.h>

int setreuid(uid_t ruid, uid_t euid);
int setregid(gid_t rgid, gid_t egid);
```

Вызов `setreuid()` устанавливает в качестве реального и действительного идентификаторов пользователя для процесса значения `ruid` и `euid` соответственно. Если передать `-1` для любого из параметров, то соответствующий идентификатор пользователя изменен не будет. Процессы, не принадлежащие пользователю `root`, могут устанавливать только действительный идентификатор пользователя, присваивая ему реальный или сохраненный идентификатор пользователя, а также устанавливать в качестве реального идентификатора пользователя действительный идентификатор пользователя. Если реальный идентификатор пользователя меняется или если действительному идентификатору пользователя присваивается значение, не равное предыдущему значению реального идентификатора пользователя, то также меняется и сохраненный идентификатор пользователя, которому присваивается значение нового действительного идентификатора пользователя. По крайней мере, так Linux и большинство других систем Unix реагируют на подобные изменения; данное поведение не определяется стандартом POSIX.

В случае успеха `setreuid()` возвращает значение `0`. В случае сбоя функция возвращает `-1` и присваивает переменной `errno` значение `EPERM`, указывающее, что владельцем текущего процесса является пользователь-не `root` и что значение `euid` не равно реальному или сохраненному идентификатору пользователя или же что значение `ruid` не равно действительному идентификатору пользователя.

Все высказывание также распространяется на группы — просто замените `setreuid()` на `setregid()`, `ruid` на `rgid`, а `euid` на `egid`.

Изменение идентификаторов пользователя и группы, стиль HP-UX

Вам может казаться, что ситуация становится комичной, но в HP-UX, системе Unix от Hewlett-Packard, были представлены собственные механизмы установки идентификаторов пользователя и группы для процесса. Linux не отстает от них благодаря следующим интерфейсам:

```
#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);
```

Вызов `setresuid()` устанавливает в качестве реального, действительного и сохраненного идентификаторов пользователя значения параметров `ruid`, `euid` и `suid` соответственно. Если указать для любого из параметров значение `-1`, то соответствующий идентификатор изменен не будет.

Пользователь `root` может устанавливать для любых идентификаторов пользователя любые значения. Остальные пользователи могут присваивать любым идентификаторам пользователя текущие значения реального, действительного и сохраненного идентификаторов пользователя. В случае успеха `setresuid()` возвращает значение `0`. В случае ошибки функция возвращает `-1` и присваивает переменной `errno` одно из следующих значений:

EAGAIN

Значение `uid` не соответствует реальному идентификатору пользователя, и если присвоить реальному идентификатору пользователя значение `uid`, то он выйдет за предел `NPROC rlimit` (указывающий максимальное число процессов, которые могут принадлежать пользователю).

EPERM

Пользователь не является пользователем `root`, и он попытался установить новое значение для реального, действительного или сохраненного идентификатора пользователя, не соответствующее одному из текущих значений реального, действительного или сохраненного идентификаторов пользователя.

Все высказанное также распространяется на группы — просто замените `setresuid()` на `setresgid()`, `ruid` на `rgid`, `euid` на `egid`, а `suid` на `sgid`.

Манипулирование предпочтительными идентификаторами пользователя/группы

Процессы, не принадлежащие пользователю `root`, должны использовать для изменения своих действительных идентификаторов пользователя функцию `seteuid()`. Процессы, принадлежащие пользователю `root`, должны использовать `setuid()`, если нужно изменить все три идентификатора пользователя, и `seteuid()`, если необходимо временно модифицировать только действительный идентификатор пользователя. Эти функции просты и работают согласно стандарту POSIX, правильно учитывая сохраненные идентификаторы пользователя.

Несмотря на предоставление дополнительной функциональности, функции в стиле BSD и HP-UX не позволяют вносить полезные изменения, возможность использовать которые дают `setuid()` и `seteuid()`.

Поддержка сохраненных идентификаторов пользователя

Существование сохраненных идентификаторов пользователя и группы провозглашается в стандарте IEEE Std 1003.1–2001 (POSIX 2001), и Linux поддержи-

вают эти идентификаторы с момента появления ядра 1.1.38. В программах, написанных только для Linux, можно не беспокоиться о существовании сохраненных идентификаторов пользователя — они всегда будут на месте. В программах, написанных для более старых систем Unix, нужно проверять наличие макроса `_POSIX_SAVED_IDS`, прежде чем ссылаться на сохраненный идентификатор пользователя или группы.

В случае отсутствия сохраненного идентификатора пользователя и группы все предыдущие рассуждения остаются в силе — просто игнорируйте те части правил, в которых упоминаются сохраненные идентификаторы.

Получение идентификаторов пользователя и группы

Эти два системных вызова возвращают реальные идентификаторы пользователя и группы соответственно:

```
#include <unistd.h>
#include <sys/types.h>
```

```
uid_t getuid (void);
gid_t getgid (void);
```

Они не могут завершаться ошибкой. Аналогично, следующие два системных вызова возвращают действительные идентификаторы пользователя и группы соответственно:

```
#include <unistd.h>
#include <sys/types.h>
```

```
uid_t geteuid (void);
gid_t getegid (void);
```

Они также никогда не завершаются ошибкой.

Сеансы и группы процессов

Каждый процесс является членом *группы процессов* (process group) — набора из одного или нескольких процессов, обычно связанных друг с другом с целью *управления заданиями* (job control). Основное отличие группы процессов заключается в том, что сигналы можно отправлять одновременно всем процессам в группе: единственное действие может завершать, останавливать или продолжать выполнение всех процессов из одной группы.

Каждая группа процессов определяется при помощи *идентификатора группы процессов* (process group ID), и в ней есть *лидер группы процессов* (process group leader). Идентификатор группы процессов равен значению `pid` лидера группы процессов. Группы процессов существуют до тех пор, пока в них остается хотя бы один член. Даже если лидер группы процессов завершается, группа продолжает существовать.

Когда новый пользователь впервые входит в систему на машине, процесс входа в систему создает новый *сесанс* (session), состоящий из одного процесса — *оболочки входа в систему* (login shell) пользователя. Оболочка входа в систему играет роль *лидера сеанса* (session leader). Идентификатор pid лидера сеанса используется в качестве идентификатора сеанса. Сеанс — это набор из одной или нескольких групп процессов. Сеансы позволяют упорядочивать действия зарегистрированного в системе пользователя и связывают этого пользователя с *управляющим терминалом* (controlling terminal) — это особое устройство tty, управляющее терминальным вводом-выводом для данного пользователя. Следовательно, сеансы большей частью являются ответственностью оболочек. В действительности ничто больше не управляет ими.

Тогда как группы процессов предоставляют механизм для адресации сигналов одновременно всем членам группы, упрощая управление заданиями и прочие функции оболочки, сеансы существуют только для того, чтобы объединять входы в систему вокруг управляющих терминалов. Группы процессов в сеансе делятся на одну *приоритетную группу процессов* (foreground process group) и ноль или несколько *фоновых групп процессов* (background process group). Когда пользователь выходит из терминала, сигнал SIGQUIT отправляется всем процессам в приоритетной группе процессов. Когда терминал распознает отключение сети, всем процессам в приоритетной группе процессов отправляется сигнал STIGHUP. Когда пользователь вводит ключ прерывания (обычно клавишное сочетание Ctrl+C), всем процессам в приоритетной группе процессов отправляется сигнал SIGINT. Таким образом, сеансы упрощают для оболочек управление терминалами и входами в систему.

В качестве примера представим, что пользователь входит в систему и у его оболочки входа в систему, bash, идентификатор pid равен 1700. Экземпляр bash пользователя теперь является единственным членом и лидером новой группы процессов с идентификатором группы процессов 1700. Эта группа процессов принадлежит новому сеансу с идентификатором сеанса 1700, и bash — единственный член и лидер этого сеанса. Новые команды, которые пользователь запускает в оболочке, выполняются в новых группах процессов в пределах сеанса 1700. Одна из этих групп процессов — соединенная напрямую с пользователем и находящаяся под управлением терминала — это *приоритетная группа процессов*. Все прочие группы процессов являются *фоновыми группами процессов*.

В каждой конкретной системе существует множество сеансов: один сеанс входа в систему для каждого пользователя и прочие сеансы для процессов, не привязанных к сеансам входа пользователя в систему, таким, как демоны. Демоны обычно создают собственные сеансы, чтобы избежать проблем привязки к другим сеансам, которые могут существовать в системе.

Каждый из этих сеансов включает одну или несколько групп процессов, и каждая группа процессов содержит, по крайней мере, один процесс. Группы процессов, включающие более одного процесса, в общем виде реализуют управление заданиями.

Команда оболочки, подобная этой:

```
$ cat ship-inventory.txt | grep booty | sort
```

возвращает группу процессов, включающую три процесса. Таким образом, оболочка может отправлять сигналы всем трем процессам одновременно. Так как пользователь ввел эту команду в консоли без ведущего амперсанда, можно сказать, что данная группа процессов будет приоритетной. На рис. 5.1 иллюстрируются взаимоотношения между сессиями, группами процессов, процессами и управляющими терминалами.

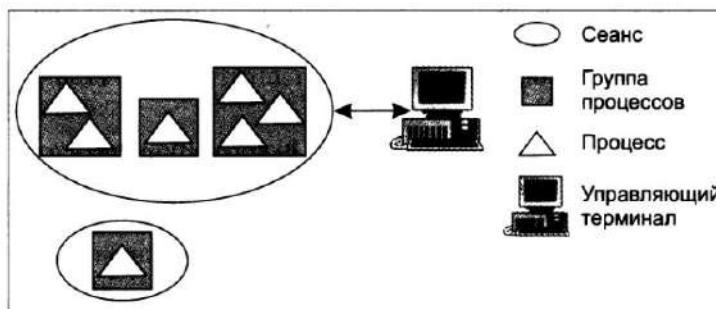


Рис. 5.1. Взаимоотношения между сессиями, группами процессов, процессами и управляющими терминалами

Linux предоставляет несколько интерфейсов для определения сессий и групп процессов, связанных с конкретным процессом, и извлечения информации о них. Они в основном используются оболочками, но могут также быть полезны процессам, подобным демонам, желающим вообще не связываться с детьми сессий и групп процессов.

Системные вызовы сессий

Оболочки создают новые сессии при входе в систему. Они делают это при помощи специального системного вызова, который упрощает создание нового сеанса до тривиальности:

```
#include <unistd.h>
```

```
pid_t setsid (void);
```

Вызов `setsid()` создает новый сеанс, предполагая, что процесс еще не является лидером группы процессов. Вызывающий процесс делается лидером сеанса и единственным членом нового сеанса, у которого нет управляющего устройства `tty`. Вызов также создает в сеансе новую группу процессов и делаетзывающий процесс лидером этой группы и единственным ее членом. В качестве идентификаторов нового сеанса и группы процессов устанавливается `pid` вызывающего процесса.

Другими словами, `setsid()` создает новую группу процессов внутри нового сеанса и делает вызывающий процесс лидером обоих. Это полезно для демонов,

которые не хотят быть членами существующих сеансов или иметь управляющие терминалы, и для оболочек, которые создают новые сеансы для каждого пользователя во время входа в систему.

В случае успеха setsid() возвращает идентификатор только что созданного сеанса. В случае ошибки вызов возвращает значение -1. Единственный возможный код errno для вызова — это EPERM, указывающий, что процесс в настоящее время является лидером группы процессов. Самый простой способ гарантировать, что любой конкретный процесс не будет лидером группы процессов, — это выполнить ветвление, завершить родительский процесс и заставить дочерний вызвать setsid(). Например:

```
pid_t pid;

pid = fork ( );
if (pid == -1) {
    perror ("fork");
    return -1;
} else if (pid != 0)
    exit (EXIT_SUCCESS);

if (setsid ( ) == -1) {
    perror ("setsid");
    return -1;
}
```

Получить идентификатор текущего сеанса тоже можно, хотя это и не часто бывает нужно:

```
#define _XOPEN_SOURCE 500
#include <unistd.h>

pid_t getsid (pid_t pid);
```

Вызов getsid() возвращает идентификатор сеанса для процесса, указанного при помощи параметра pid. Если значение аргумента pid равно 0, то вызов возвращает идентификатор сеанса для вызывающего процесса. В случае ошибки вызов возвращает значение -1. Единственное возможное значение переменной errno — это ESRCH, указывающее, что значение pid не соответствует никакому допустимому процессу. Обратите внимание, что в других системах Unix переменная errno также может принимать значение EPERM, указывающее, что pid и вызывающий процесс принадлежат разным сеансам; в Linux эта ошибка не возвращается, а вместо этого возвращается идентификатор сеанса любого процесса.

Применяется этот вызов редко и обычно в диагностических целях:

```
pid_t sid;

sid = getsid (0);
if (sid == -1)
    perror ("getsid"); /* не должно происходить */
else
    printf ("Мой идентификатор сеанса=%d\n", sid);
```

Системные вызовы для группы процессов

Вызов `setpgid()` присваивает значение `pgid` идентификатору группы процессов для процесса, указанного при помощи параметра `pid`:

```
#define _XOPEN_SOURCE 500
#include <unistd.h>

int setpgid (pid_t pid, pid_t pgid);
```

Если значение аргумента `pid` равно 0, то используется текущий процесс. Если значение `pgid` равно 0, то в качестве идентификатора группы процессов используется идентификатор процесса, указанного при помощи `pid`.

В случае успеха вызов `setpgid()` возвращает значение 0. Успех зависит от нескольких условий:

- процесс, идентифицированный при помощи `pid`, должен быть вызывающим процессом или потомком вызывающего процесса, который еще не выполнил вызов `exec` и принадлежит тому же сеансу, что и вызывающий процесс;
- процесс, идентифицированный при помощи `pid`, не должен быть лидером сеанса;
- если группа процессов с идентификатором `pgid` уже существует, она должна принадлежать тому же сеансу, что и вызывающий процесс;
- значение `pgid` должно быть неотрицательным.

В случае ошибки вызов возвращает значение -1 и присваивает переменной `errno` один из следующих кодов ошибки:

EACCES

Процесс, идентифицированный параметром `pid`, является потомком вызывающего процесса, который уже сделал вызов `exec`.

EINVAL

Значение `pgid` меньше нуля.

EPERM

Процесс, идентифицированный параметром `pid`, является лидером сеанса или принадлежит другому сеансу, отличному от сеанса вызывающего процесса. Или же была сделана попытка переместить процесс в группу процессов, принадлежащую другому сеансу.

ESRCH

Идентификатор `pid` принадлежит не текущему процессу и не его потомку, и его значение не равно 0.

Как и для сеансов, получить идентификатор группы процессов для процесса можно, но обычно не нужно:

```
#define _XOPEN_SOURCE 500
#include <unistd.h>

pid_t getpgid (pid_t pid);
```

Вызов `getpgid()` возвращает идентификатор группы процессов для процесса, указанного при помощи параметра `pid`. Если `pid` равен 0, то используется идентификатор группы процессов для текущего процесса. В случае ошибки вызов возвращает значение `-1` и присваивает переменной `errno` единственно возможное значение `ESRCH`, указывающее, что `pid` — это недопустимый идентификатор процесса.

Точно так же, как `getsid()`, `getpgid()` используется в основном в диагностических целях:

```
pid_t pgid;

pgid = getpgid (0);
if (pgid == -1)
    perror ("getpgid"); /* не должно происходить */
else
    printf ("Мой идентификатор группы процессов=%d\n", pgid);
```

Устаревшие функции для группы процессов

В Linux поддерживаются два старых интерфейса из BSD для манипулирования или получения значения идентификатора группы процессов. Так как они не приносят особой пользы, в отличие от ранее перечисленных системных вызовов, использовать их в новых программах рекомендуется только тогда, когда стоит жесткое требование о переносимости. `setpgrp()` используется для установки идентификатора группы процессов:

```
#include <unistd.h>

int setpgrp (void);
```

Такой вызов:

```
if (setpgrp ( ) == -1)
    perror ("setpgrp");
```

идентичен следующему варианту вызова:

```
if (setpgid (0,0) == -1)
    perror ("setpgid").
```

Оба пытаются назначить текущий процесс группе процессов с тем же номером, что и `pid` текущего процесса, возвращая 0 в случае успеха и `-1` в случае ошибки. Все значения переменной `errno` для вызова `setpgid()` распространяются и на `setpgrp()`, за исключением `ESRCH`.

Схожим образом вызов `getpgrp()` можно использовать для получения идентификатора группы процессов:

```
#include <unistd.h>
```

```
pid_t getpgrp (void);
```

Следующий вариант вызова:

```
pid_t pgid = getpgrp ( ).
```

идентичен такому:

```
pid_t pgid = getpgid (0);
```

Оба возвращают идентификатор группы процессов для вызывающего процесса. Функция `getpgid()` не может вернуть ошибку.

Демоны

Демон (*daemon*) – это процесс, выполняющийся в фоновом режиме, не подключающийся ни к одному из управляющих терминалов. Демоны обычно запускаются во время загрузки, выполняются от имени пользователя `root` или другого специального пользователя (такого, как `apache` или `postfix`) и обрабатывают задачи системного уровня. В качестве общепринятого соглашения имена демонов обычно заканчиваются на букву `d` (как в `crond` и `sshd`), но это не обязательное требование и даже не универсальная особенность.

Это название берет начало в демоне Максвелла — мысленном эксперименте, проведенном в 1867 году физиком Джеймсом Максвеллом. Демоны — это сверхъестественные существа в греческой мифологии, существующие где-то между людьми и богами, наделенные силой и пророческим знанием. В отличие от демонов в иудейско-христианской культуре, греческий демон не обязательно злой. Действительно, демоны в мифологии обычно были помощниками богов, выполнившими задания, которые не хотелось самостоятельно делать жителям горы Олимп, — точно так же, как демоны в Unix справляются с задачами, которых привилегированные пользователи стараются избегать.

К демону предъявляются два общих требования: он должен выполняться как потомок процесса инициализации и не должен быть соединен с терминалом.

В целом, для того чтобы стать демоном, программа выполняет следующие шаги:

1. Выполняется вызов `fork()`. Это создает новый процесс, который затем будет превращен в демон.
2. В родительском процессе выполняется вызов `exit()`. Это гарантирует, что первоначальный предок (дедушка демона) удовлетворен завершением своего потомка, что предок демона более не существует и что демон не является лидером группы процессов. Последнее требование обязательно для успеха следующего шага.
3. Выполняется вызов `setsid()`, предоставляя демону новую группу процессов и сеанс, в обоих из которых он становится лидером. Это также гарантирует, что с данным процессом не связан управляющий терминал (так как процесс только что создал новый сеанс и не будет назначать терминал).
4. Рабочий каталог меняется на корневой каталог при помощи `chdir()`. Это делается потому, что унаследованный рабочий каталог может находиться в любом месте системы. Демоны обычно продолжают выполняться в течение всего времени работы системы, и вам не нужно, чтобы какой-то случайный каталог оставался открытым. Этим шагом вы предотвращаете размонтирование администратором файловой системы, содержащей данный каталог.

5. Закрываются все дескрипторы файлов. Нет никакой необходимости наследовать открытые файловые дескрипторы и, не подозревая об этом, оставлять их открытыми.
6. Открываются дескрипторы файлов 0, 1 и 2 (стандартный ввод, стандартный вывод и стандартный вывод для ошибок), и все они перенаправляются на `/dev/null`.

Следуя этим правилам, можно написать такую программу, которая будет превращать в демона сама себя:

```
#include <sys/types.h>
#include <sys/stat.h>
#include <stdlib.h>
#include <stdio.h>
#include <fcntl.h>
#include <unistd.h>
#include <linux/fs.h>

int main (void)
{
    pid_t pid;
    int i;

    /* создание нового процесса */
    pid = fork ( );
    if (pid == -1)
        return -1;
    else if (pid != 0)
        exit (EXIT_SUCCESS);

    /* создание нового сеанса и группы процессов */
    if (setsid ( ) == -1)
        return -1;

    /* установка в качестве рабочего каталога корневого каталога */
    if (chdir ("") == -1)
        return -1;

    /* закрытие всех открытых файлов */
    /* NR_OPEN – это слишком, но это работает */
    for (i = 0; i < NR_OPEN, i++)
        close (i);

    /* перенаправление дескрипторов файла 0,1,2 в /dev/null */
    open ("/dev/null", O_RDWR);      /* stdin */
    dup (0);                      /* stdout */
    dup (0);                      /* stderr */

    /* всякие действия демона... */

    return 0;
}
```

В большинстве систем Unix в составе библиотеки С есть функция `daemon()`, автоматизирующая эти шаги и превращающая все вышеизложенное в простой вызов:

```
#include <unistd.h>

int daemon (int nochdir, int noclose);
```

Если значение nochdir не равно нулю, то демон не будет менять текущий рабочий каталог на корневой. Если значение noclose не равно нулю, то демон не будет закрывать все открытые дескрипторы файлов. Данные параметры полезны, если родительский процесс уже позаботился об этих аспектах процедуры создания демона. Обычно, однако, для обоих параметров передается значение 0.

В случае успеха вызов возвращает 0. В случае ошибки вызов возвращает -1 и присваивает переменной errno один из кодов ошибки системных вызовов fork() и setsid().

Заключение

В этой главе мы познакомились с основами управления процессами в Unix, от создания процесса до его завершения. В следующей главе мы более подробно изучим сложные интерфейсы управления процессами, а также интерфейсы для настройки планирования процессов.

6

Расширенное управление процессами

В главе 5 вы познакомились с абстракцией процесса, и мы обсудили интерфейсы ядра, применяемые для создания, управления и разрушения процессов. Эта глава основана на идеях предыдущей: мы начнем с обсуждения планировщика процессов Linux и его алгоритма планирования, а затем перейдем к сложным интерфейсам управления процессами. Эти системные вызовы предназначены для манипулирования планированием и семантикой процессов, позволяя влиять на поведение планировщика для достижения целей, диктуемых приложением или пользователем.

Планирование процессов

Планировщик процессов (process scheduler) – это компонент ядра, выбирающий, какой процесс будет выполняться следующим. Другими словами, планировщик процессов, или просто планировщик, представляет собой подсистему ядра, которая делит ограниченный ресурс времени процессора среди всех процессов в системе. Принимая решения о том, какие процессы выполнять и когда, планировщик несет ответственность за максимизацию загрузки процессора, обеспечивая впечатление, что множество процессов выполняются одновременно и бесшовно.

В этой главе мы много будем говорить о *работоспособных процессах* (Runnable process). Работоспособный процесс – это, в первую очередь, незаблокированный процесс. Процессы, взаимодействующие с пользователями, выполняющие обширный ввод-вывод или отвечающие на события ввода-вывода или сетевые события, обычно проводят много времени в заблокированном состоянии, ожидая, пока станут доступными ресурсы, и на протяжении этих длинных периодов (длинных по сравнению со временем, которое необходимо на выполнение машинных инструкций) они не работоспособны. У работоспособного процесса также должна оставаться хотя бы часть его *кванта времени* (times-

lice) — интервала времени, в течение которого планировщик разрешил ему выполнятся. Ядро помещает все работоспособные процессы в *список выполнения* (run list). Как только квант времени процесса заканчивается, он удаляется из этого списка и не считается снова работоспособным до тех пор, пока у всех остальных работоспособных процессов также не закончатся их кванты времени.

Если работоспособный процесс только один (или вообще ни одного), то работа планировщика процессов тривиальна. Планировщик подтверждает свою значимость, однако когда работоспособных процессов больше, чем процессоров. В подобных ситуациях, очевидно, некоторые процессы должны выполнятьсь, а другие — находиться в состоянии ожидания. Принятие решения, какие процессы выполнять, когда и в течение какого времени, — это фундаментальная ответственность планировщика.

Операционная система на однопроцессорной машине является *многозадачной* (multitasking), если она может чередовать выполнение нескольких процессов, создавая иллюзию того, что одновременно работает более одного процесса. На многопроцессорных машинах многозадачная операционная система фактически позволяет процессам выполняться параллельно, на разных процессорах. Немногозадачная операционная система, такая, как DOS, может одновременно выполнять только одно приложение.

Существует два варианта многозадачных операционных систем: *совместные* или *кооперативные* (cooperative), и *вытесняющие* или *приоритетные* (preemptive). В Linux реализована вторая форма многозадачности, в которой планировщик решает, когда один процесс должен остановиться, а другой возобновиться. Мы называем акт приостановки работающего процесса для того, чтобы освободить место для другого, *вытеснением* (preemption). И снова, длина промежутка времени, в течение которого процесс выполняется до того, как планировщик вытеснит его, называется *квантом времени процесса* (process timeslice) (это название появилось потому, что планировщик выделяет каждому работоспособному процессу «квант» времени процессора).

При кооперативной многозадачности, наоборот, процесс не останавливается до тех пор, пока он добровольно не решает прекратить выполнение. Мы называем акт добровольной приостановки процессом самого себя *уступкой* (yielding). В идеальном случае процессы часто уступают, но операционная система не может принуждать их к этому. Неаккуратно написанная программа может выполняться дольше оптимального времени и даже снижать эффективность всей системы. Из-за недостатков такого подхода современные операционные системы практически все работают на основе вытесняющей многозадачности, и Linux не исключение.

Планировщик процессов O(1), представленный в семействе ядер 2.5, является сердцевиной планирования в Linux¹. Алгоритм планирования в Linux

¹ Для любознательных читателей — планировщик процессов представляет собой самообслуживающийся механизм, который определен в файле `kernel/sched.c` в дереве ресурсов ядра.

обеспечивает вытесняющую многозадачность и поддержку нескольких процессоров, привязку процессов к процессорам, конфигурации *неоднородного доступа к памяти* (nonuniform memory access, NUMA), многопоточность, процессы реального времени и расстановку приоритетов пользователем.

Запись с О большим

$O(1)$ — читается как «О большое от единицы» — это пример записи с О большим, которая представляет сложность и масштабируемость алгоритма. Формально, если

$$f(x) \text{ — это } O(g(x)),$$

то

$$\exists c, x' \text{ такие, что } f(x) \leq c g(x), \forall x \geq x'.$$

Говоря обычным языком, значение некоторого алгоритма, f , всегда меньше или равно значению g , умноженному на некоторую произвольную константу, если только входное значение x больше некоторого базового значения x' . Это означает, что значение g больше или равно f ; g ограничивает f сверху.

$O(1)$, таким образом, подразумевает, что значение рассматриваемого алгоритма оценивается меньше некоторой константы, c . Вся эта пышность и торжественность, в действительности, несет одно важное обещание: планировщик процессов Linux всегда будет работать одинаково, независимо от числа процессов в системе. Это важно, так как на первый взгляд акт выбора нового процесса для запуска должен включать, как минимум, одно прохождение по спискам процессов, если не несколько. Когда использовались более простые планировщики (включая планировщики из ранних версий Linux), по мере того, как количество процессов в системе росло, подобные прохождения могли превращаться в «узкие места» системы. В лучшем случае циклы прохождения привносили неопределенность — отсутствие детерминизма — в процесс планирования.

Планировщик Linux, работающий в реальном времени независимо ни от каких факторов, не создает подобных узких мест.

Кванты времени

Квант времени, который Linux выделяет каждому процессу, — это важная переменная в общем поведении и эффективности системы. Если кванты времени слишком велики, то процессам приходится долго ждать между интервалами, когда им дозволяется выполняться, что снижает впечатление одновременного исполнения. Пользователя может разочаровать ощущение задержка. И наоборот, если кванты времени слишком малы, значительная доля времени системы тратится на переключение с одного приложения на другое и преимущества со средоточенности во времени исчезают.

Поэтому определить идеальную длину кванта времени нелегко. Некоторые операционные системы дают процессам большие кванты времени, надеясь мак-

симиизировать пропускную способность и общую производительность. Другие операционные системы предоставляют очень маленькие кванты времени в надежде обеспечить высокую интерактивность. Как мы увидим, Linux пытается достичь лучшего в обоих мирах, динамически выделяя процессам кванты времени.

Обратите внимание, что процессу не обязательно тратить весь свой квант за один раз. Процесс, которому предоставлен квант времени длиной 100 миллисекунд, может выполняться в течение 20 миллисекунд, а затем блокироваться в ожидании какого-то ресурса, например ввода с клавиатуры. Тогда планировщик временно удаляет этот процесс из списка работоспособных процессов. Когда заблокированный ресурс становится доступным — когда буфер клавиатуры заполняется какими-то данными — планировщик будет процесс. Процесс затем может продолжать выполнение до тех пор, пока не потратит оставшиеся 80 миллисекунд или пока снова не заблокируется в ожидании ресурса.

Процессы, ограниченные вводом-выводом, и процессы, ограниченные процессором

Процессы, постоянно расходующие все доступные им кванты времени, считаются ограниченными процессором. Такие процессы постоянно требуют процессорного времени и тратят все время, предоставляемое им планировщиком. Простейший пример — бесконечный цикл. Прочие примеры включают научные вычисления, математические вычисления и обработку изображений.

С другой стороны, процессы, которые проводят больше времени в заблокированном состоянии, ожидая каких-то ресурсов, чем в состоянии выполнения, называются ограниченными вводом-выводом. Такие процессы часто совершают вызовы ввода-вывода и ожидают их завершения, блокируются в ожидании ввода с клавиатуры или ждут, когда пользователь переместит мышь. Примеры ограниченных вводом-выводом приложений включают файловые утилиты, которые мало чего делают, за исключением отправки системных вызовов, запрашивающих у ядра операции ввода-вывода, например `cp` или `mv`, и многие приложения с графическим интерфейсом пользователя, которым приходится подолгу ожидать ввода со стороны пользователя.

Приложения, ограниченные процессором и вводом-выводом, отличаются в том смысле, что для них оптимальными являются разные варианты поведения планировщика. Приложениям, ограниченным процессором, необходимы как можно более длинные кванты времени, позволяющие максимизировать коэффициент попадания в кеш (за счет сосредоточенности во времени) и выполнять работу как можно быстрее. В противоположность этому процессам, ограниченным вводом-выводом, не обязательно нужны большие кванты времени, так как обычно они выполняются в течение очень коротких промежутков времени, а затем отправляют запрос ввода-вывода или блокируются в ожидании какого-то ресурса ядра. Для них, однако, полезно постоянное внимание со стороны планировщика. Чем быстрее такое приложение сможет возобновить работу после блокировки и отправить новые запросы ввода-вывода, тем эффективнее оно будет использовать аппаратное обеспечение системы. Помимо этого,

если приложение ожидает ввода со стороны пользователя, то чем быстрее ему предоставляется квант времени, тем лучшее впечатление создается у пользователя относительно бесшовного выполнения.

Жонглировать требованиями процессов, ограниченных процессором и вводом-выводом, непросто. Планировщик Linux пытается идентифицировать и обеспечивать приоритетную обработку приложениям, ограниченным вводом-выводом: программы, выполняющие большое количество операций ввода-вывода, получают более высокий приоритет, а сильно ограниченные процессором приложения обслуживаются во вторую очередь.

В действительности большинство приложений представляет собой смесь — они ограничиваются и вводом-выводом, и процессором. Кодировщики и декодировщики звука и видео — это хороший пример приложений, доказывающих невозможность жесткой категоризации. Многие игры также довольно сложны. Не всегда возможно идентифицировать принадлежность каждого конкретного приложения, и в разные моменты времени один и тот же процесс может вести себя по-разному.

Вытесняющее планирование

Когда у процесса заканчивается его квант времени, ядро приостанавливает процесс и начинает выполнять другой процесс. Если в системе не остается работоспособных процессов, ядро берет набор процессов с закончившимися квантами, снова предоставляет им кванты и начинает заново выполнять эти процессы. Таким образом, все процессы в конечном итоге получают возможность поработать, даже если в системе существуют процессы с более высоким приоритетом, — процессам с низким приоритетом просто приходится ждать, пока у приоритетных процессов не закончатся кванты времени или они не заблокируются. Такое поведение формулирует важное неявное правило планирования в Unix: все процессы должны постепенно выполняться.

Если в системе не остается никаких работоспособных процессов, ядро начинает «выполнять» *процесс бездействия* (idle process). В действительности процесс бездействия — это вовсе не процесс и он не выполняется (до истощения батарей). Процесс бездействия — это специальная процедура, которую ядро исполняет для упрощения алгоритма планировщика и для облегчения учета. Время бездействия — это просто время, потраченное на выполнение процесса бездействия.

Если выполняется какой-то процесс и в это время становится работоспособным процесс с более высоким приоритетом (например, до этого он был заблокирован в ожидании ввода с клавиатуры, а пользователь только что ввел слово), то текущий процесс мгновенно приостанавливается, а ядро переключается на приоритетный процесс. Таким образом, в системе никогда не может быть работоспособного, но неработающего процесса с приоритетом выше, чем у выполняющегося в данный момент процесса. Работающий процесс — это всегда работоспособный процесс с наивысшим приоритетом в системе.

Управление потоками

Потоки — это единицы выполнения процесса. В каждом процессе есть, по крайней мере, один поток. У каждого потока собственная виртуализация процессора: собственный набор регистров, указатель команд и состояние процессора. Хотя у большинства процессов только один поток, процессы могут обладать множеством потоков, выполняющих различные задачи, но совместно использующих одно и то же адресное пространство (и, следовательно, одну и ту же динамическую память, отображенные файлы, конечную программу и т. д.), список открытых файлов и другие ресурсы ядра.

У ядра Linux интересное и уникальное представление о потоках. Фактически, у ядра нет такой концепции. Для ядра Linux все потоки являются уникальными процессами. Проще говоря, нет никакого различия между двумя несвязанными процессами и двумя потоками внутри одного процесса. Ядро просто рассматривает потоки как процессы, совместно использующие ресурсы. То есть ядро считает процесс, состоящий из двух потоков, двумя различными процессами, совместно использующими набор ресурсов ядра (адресное пространство, список открытых файлов и т. д.).

Многопоточное программирование (multithreaded programming) — это искусство программирования на основе потоков. Самый распространенный API Linux для программирования с потоками — это API, стандартизированный в IEEE Std 1003.1c-1995 (стандарт POSIX 1995 или POSIX.1c). Разработчики часто называют библиотеку, в которой реализуется этот API, `pthreads`. Программирование с потоками — сложная тема, и API `pthreads` — очень большой и трудный интерфейс. Поэтому в данной книге мы не будем касаться `pthreads`. Вместо этого сфокусируемся на интерфейсах, на которых построена библиотека `pthreads`.

Уступка процессора

Хотя Linux — это многозадачная операционная система с вытеснением, в ней также предусмотрен системный вызов, позволяющий процессам явно уступать очередь выполнения и заставляющий планировщик выбирать новый процесс:

```
#include <sched.h>
```

```
int sched_yield (void);
```

Вызов `sched_yield()` приводит к приостановке выполняющегося в данный момент процесса, после чего планировщик процессов выбирает, какой процесс будет выполняться следующим, точно так же, как если бы ядро само вытеснило текущий процесс в пользу нового. Обратите внимание, что, если никаких других работоспособных процессов нет, как часто бывает, уступающий процесс сразу же продолжает выполняться. Из-за такой неопределенности в сочетании с общей уверенностью в наличии лучшего выбора используется данный системный вызов не часто.

В случае успеха вызов возвращает значение 0; в случае ошибки он возвращает -1 и присваивает переменной `errno` подходящий код ошибки. В Linux и, более чем вероятно, в большинстве других операционных систем вызов `sched_yield()` не может завершаться ошибкой и поэтому всегда возвращает значение 0. Внимательный программист тем не менее может все же проверять возвращаемое значение:

```
if (sched_yield ( ))  
    perror ("sched_yield").
```

Правомерное использование

На практике существует очень мало (если вообще хоть сколько-нибудь) вариантов правомерного использования вызова `sched_yield()` в правильной многозадачной системе с вытеснением, такой, как Linux. Ядро полностью справляется с принятием оптимальных и самых эффективных решений о планировании — определенно, ядро лучше вооружено, чем любое отдельное приложение, и ему намного удобнее решать, что нужно вытеснять и когда. Это именно то, почему разработчики операционных систем отказались от кооперативной многозадачности в пользу вытесняющей.

Почему же тогда вообще существует этот системный вызов, говорящий «измените мне расписание»? Ответ лежит в приложениях, которым приходится ожидать внешние события, которые могут вызываться пользователем, аппаратным компонентом или другим процессом. Например, если одному процессу приходится ждать другой, первым приходит на ум решение «просто уступить процессор, пока другой процесс не завершится». В качестве примера: реализация безыскусного потребителя в паре потребитель/производитель может быть приблизительно такой:

```
/* потребитель .. */  
do {  
    while (producer_not_ready ( ))  
        sched_yield ( );  
    process_data ( );  
} while (!time_to_quit ( ));
```

К счастью, программисты в Unix обычно не пишут такой код. Программы для Unix приводятся в движение событиями и чаще всего включают какой-либо поддерживающий блокировку механизм (например, конвейер) между потребителем и производителем, а не используют `sched_yield()`. В таком случае потребитель считывает данные из конвейера, блокируясь по необходимости, пока данные не станут доступными. Производитель, в свою очередь, записывает в конвейер, как только у него появляются свежие данные. Это снимает ответственность за координирование с процессом из пользовательского пространства, перекладывая ее на ядро, умеющее оптимально управлять ситуацией: приостанавливать процессы и пробуждать их по мере необходимости. В целом, программы Unix должны основываться на решениях, приводимых в движение событиями и полагающихся на дескрипторы файла с возможностью блокировки.

До недавних пор только одна ситуация мучительно требовала использования вызова `sched_yield()`: блокировка потоков в пользовательском пространстве. Когда поток пытался захватить блокировку, которая уже принадлежала другому потоку, этот поток уступал процессор до тех пор, пока блокировка не становилась доступной для него. С учетом того, что ядро не поддерживало блокировки в пользовательском пространстве, такой подход был самым простым и эффективным. К счастью, реализация потоков в современных версиях Linux (новая библиотека управления потоками New POSIX Threading Library или NPTL) привнесла в наш мир оптимальное решение, использующее *фьютексы* (*futex*), которые обеспечивают поддержку блокировок в пользовательском пространстве со стороны ядра.

Еще один вариант применения `sched_yield()` — это «любезная игра»: сильно нагружающие процессор программы могут периодически вызывать `sched_yield()`, пытаясь минимизировать свое воздействие на систему. Благородная по сути, эта стратегия тем не менее имеет два недостатка. Во-первых, ядро умеет принимать глобальные решения о планировании намного лучше, чем отдельные процессы, и, следовательно, ответственность за обеспечение плавного функционирования системы должна лежать на планировщике процессов, а не на процессах. В этом смысле планировщик старается вознаграждать приложения, сильно зависящие от ввода-вывода, и наказывать приложения, оказывающие большое давление на процессор. Во-вторых, смягчение нагрузки, оказываемой зависящими от процессора приложениями, и предоставление возможности свободно дышать другим приложениям — это ответственность пользователя, а не отдельных приложений. Пользователь может выражать свои относительные предпочтения касательно производительности приложений при помощи команды оболочки `nice`, о которой мы поговорим далее в этой главе.

Уступка, прошлое и настоящее

До появления ядра Linux 2.6 вызов `sched_yield()` мог оказывать лишь небольшой эффект. Если был доступен другой работоспособный процесс, ядро переключалось на него и помещало вызывающий процесс в хвост списка работоспособных процессов. Сразу же после этого ядро меняло расписание выполнения для вызывающего процесса. В весьма вероятном случае отсутствия других работоспособных процессов вызывающий процесс просто продолжал выполняться.

В ядре 2.6 это поведение было изменено. На сегодняшний день алгоритм таков:

1. Данный процесс является процессом реального времени? Если да, то поместить его в конец списка работоспособных процессов и вернуть результат (старое поведение). Если нет, то перейти к следующему шагу. (Подробнее о процессах реального времени говорится в разделе «Системы реального времени» далее в этой главе.)
2. Удалить данный процесс из списка работоспособных процессов и поместить его в список просроченных процессов. Это подразумевает, что все работо-

способные процессы должны выполниться и исчерпать свои кванты времени и только тогда вызывающий процесс, а также прочие просроченные процессы смогут возобновить работу.

3. Запланировать выполнение следующего работоспособного процесса в списке.

Таким образом, суммарным результатом вызова `sched_yield()` является та-
кая же ситуация, как если бы у процесса просто закончился его квант времени.
Такое поведение отличается от поведения предыдущих версий ядра, где эффект
`sched_yield()` был мягче (приблизительно таким: «если другой процесс готов
и ожидает, то повыполнит его немного, а потом вернись ко мне»).

Одной из причин такого изменения было предотвращение патологического
случая «пинг-понга». Представьте себе два процесса, A и B, и оба они вызывают
`sched_yield()`. Предположим, что это единственныe два работоспособных
процесса (могут быть и другие процессы, способные выполняться, но ни
одного с ненулевым квантом времени). Если бы работало старое поведение
`sched_yield()`, то результатом было бы то, что ядро по очереди заставляло бы
выполняться оба процесса, но они отвечали бы на это: «Нет, пусть кто-то дру-
гой работает!» Это продолжалось бы до тех пор, пока у обоих процессов не за-
кончились бы кванты времени. Если попытаться нарисовать диаграмму выбора
процессов планировщиком процессов, то она была бы такой: «A, B, A, B, A, B...»
и так далее — отсюда название «эффект пинг-понга».

Новое поведение не позволяет такому происходить. Как только процесс A
пытается уступить процессор, планировщик удаляет его из списка работоспо-
собных процессов. Аналогично, как только процесс B делает такую же попытку,
планировщик также удаляет его из списка работоспособных процессов. Плани-
ровщик не будет пытаться выполнять процессы A и B до тех пор, пока не оста-
нется никаких других работоспособных процессов, предотвращая эффект пинг-
понга и позволяя прочим процессам получить заслуженную долю процессорно-
го времени.

Следовательно, стараясь уступить процессор, процесс должен действительно
желать этого!

Приоритеты процессов

ПРИМЕЧАНИЕ

Обсуждение в этом разделе касается обычных процессов, а не процессов реального времени. По-
следние требуют других критериев планирования и отдельной системы приоритетов. Мы поговорим
о вычислениях в реальном времени далее в этой главе.

Linux не планирует процессы абы как. Всем приложениям назначаются *приоритеты* (*priority*), влияющие на то, когда их процессы выполняются и как
долго. Исторически в Unix эти приоритеты называются *значениями любезности* (*nice value*), так как идея, лежащая в основе их применения, заключается в том,
чтобы «быть любезным» к остальным процессам в системе, понижая приоритет

процесса и позволяя прочим процессам потреблять больше процессорного времени системы.

Значение любезности диктует, когда процесс будет выполняться. Linux планирует выполнение работоспособных процессов в порядке убывания: от самого высокого приоритета к низшему. Процесс с более высоким приоритетом выполняется раньше процесса с более низким. Значение любезности также диктует размер кванта времени для процесса.

Допустимые значения любезности – от -20 до 19 включительно, значение по умолчанию равно 0. Это может сбить с толку, но чем ниже значение любезности процесса, тем выше его приоритет и больше квант времени; и наоборот, чем больше значение, тем ниже приоритет процесса и меньше квант времени. Увеличение значения любезности, таким образом, – это «любезный» поступок по отношению к остальным составляющим системы. Такая числовая перестановка создает определенный беспорядок. Когда мы говорим, что у процесса «высокий приоритет», мы имеем в виду, что он быстрее выбирается для выполнения и может работать дольше, чем низкоприоритетные процессы, но у такого процесса значение любезности меньше.

Системный вызов nice()

В Linux предусмотрено несколько системных вызовов для получения и установки значения любезности процесса. Самый простой из них – это nice():

```
#include <unistd.h>
```

```
int nice (int inc);
```

Успешный вызов nice() увеличивает значение любезности процесса на значение inc и возвращает обновленное значение. Только процесс с характеристикой CAP_SYS_NICE (фактически, процесс, владельцем которого является пользователь root) может передавать отрицательное значение для параметра inc, уменьшая свое значение любезности и, таким образом, повышая приоритет. Следовательно, процессы, не принадлежащие пользователю root, имеют право только понижать свой приоритет (увеличивая значение любезности).

В случае ошибки вызов nice() возвращает значение -1. Однако, так как nice() возвращает новое значение любезности для процесса, -1 также может возвращаться в случае успеха. Для того чтобы различать эти две ситуации, можно присваивать переменной errno нулевое значение перед вызовом и снова проверять его после изменения значения любезности. Например:

```
int ret;
```

```
errno = 0;
ret = nice (10);           /* увеличение значения любезности на 10 */
if (ret == -1 && errno != 0)
    perror ("nice");
else
    printf ("значение любезности теперь равно %d\n", ret);
```

Linux возвращает только один код ошибки, EPERM, указывающий, что вызывающий процесс попытался увеличить свой приоритет (при помощи отрицательного значения `inc`), но не обладает характеристикой CAP_SYS_NICE. Прочие системы также возвращают код EINVAL, когда прибавление `inc` выводит значение любезности за пределы допустимого диапазона, но в Linux так не делается. Вместо этого Linux по необходимости просто тихо округляет недопустимые значения `inc` вверх или вниз, ограничиваясь разрешенным диапазоном.

Передать 0 в качестве аргумента `inc` — это простой способ узнать текущее значение любезности:

```
printf ("текущее значение любезности равно %d\n", nice (0));
```

Часто процессу бывает необходимо установить абсолютное значение любезности, а не увеличить его относительно текущего. Это можно сделать при помощи такого кода:

```
int ret, val;  
  
/* получение текущего значения любезности */  
val = nice (0);  
  
/* нам нужно значение любезности, равное 10 */  
val = 10 - val;  
errno = 0;  
ret = nice (val);  
if (ret == -1 && errno != 0)  
    perror ("nice");  
else  
    printf ("текущее значение любезности равно %d\n", ret);
```

Системные вызовы `getpriority()` и `setpriority()`

Предпочтительное решение — применять системные вызовы `getpriority()` и `setpriority()`, обеспечивающие более высокую степень контроля, но также и более сложные в использовании.

```
#include <sys/time.h>  
#include <sys/resource.h>  
  
int getpriority (int which, int who);  
int setpriority (int which, int who, int prio);
```

Эти вызовы работают на процессе, группе процессов или пользователе, в зависимости от значений параметров `which` и `who`. Параметр `which` может принимать значения PRIO_PROCESS, PRIO_PGRP и PRIO_USER, а при помощи параметра `who` указывается идентификатор процесса, идентификатор группы процессов или идентификатор пользователя соответственно. Если значение `who` равно 0, то вызов работает на текущем идентификаторе процесса, идентификаторе группы процессов или идентификаторе пользователя, соответственно.

Вызов `getpriority()` возвращает наибольшее значение приоритета (наименьшее значение любезности) среди указанных процессов. Вызов `setpriority()` устанавливает в качестве приоритета всех указанных процессов значение `prio`. Как и в случае системного вызова `nice()`, только процесс, обладающий характе-

ристикой CAP_SYS_NICE, имеет право повышать приоритет процесса (уменьшать числовое значение любезности). Помимо этого, только процесс с такой характеристикой может повысить или понизить приоритет процесса, которым вызывающий пользователь не владеет.

Как и вызов nice(), вызов getpriority() возвращает значение -1 в случае ошибки. Так как это также может быть успешным возвращаемым значением, рекомендуется сбрасывать переменную errno перед вызовом, чтобы иметь возможность обрабатывать условия ошибки. С вызовом setpriority() таких проблем не связано: setpriority() всегда возвращает 0 в случае успеха и -1 в случае ошибки.

Следующий код возвращает значение приоритета текущего процесса:

```
int ret;
```

```
ret = getpriority (PRIO_PROCESS, 0);  
printf ("значение любезности равно %d\n", ret).
```

Следующий код устанавливает в качестве приоритета всех процессов в текущей группе процессов значение 10:

```
int ret;
```

```
ret = setpriority (PRIO_PGRP, 0, 10);  
if (ret == -1)  
    perror ("setpriority");
```

В случае ошибки обе функции присваивают переменной errno одно из следующих значений:

EACCES

Процесс попытался повысить приоритет указанного процесса, но не обладает характеристикой CAP_SYS_NICE (только для setpriority()).

EINVAL

Значение параметра which не равно PRIO_PROCESS, PRIO_PGRP или PRIO_USER.

EPERM

Действительный идентификатор пользователя указанного процесса не совпадает с действительным идентификатором пользователя вызывающего процесса, и вызывающий процесс не обладает характеристикой CAP_SYS_NICE (только для setpriority()).

ESRCH

Не найдено ни одного процесса, отвечающего критериям параметров which и who.

Приоритеты ввода-вывода

Помимо приоритетов планирования, Linux позволяет процессам указывать также *приоритеты ввода-вывода* (I/O priority). Это значение влияет на относительный приоритет запросов ввода-вывода, делаемых процессами. Планировщик ввода-вывода ядра (о котором говорилось в главе 4) обслуживает запросы,

приходящие от процессов с более высокими приоритетами ввода-вывода, быстрее, чем запросы процессов с низкими приоритетами ввода-вывода.

По умолчанию планировщики ввода-вывода используют для определения приоритета ввода-вывода значение любезности. Следовательно, установка значения любезности автоматически меняет приоритет ввода-вывода. Однако ядро Linux дополнительно предоставляет два системных вызова для явной установки и извлечения значения приоритета ввода-вывода независимо от значения любезности:

```
int ioprio_get (int which, int who)
int ioprio_set (int which, int who, int ioprio)
```

К сожалению, ядро пока что не экспортирует эти системные вызовы, и в glibc не предоставляются способы доступа из пользовательского пространства. Без поддержки glibc использовать их в лучшем случае затруднительно. Помимо этого, когда и если поддержка glibc все-таки появится, интерфейсы могут отличаться от системных вызовов. До того времени, как такая поддержка станет доступной, остается только два переносимых способа манипулирования приоритетом ввода-вывода процесса: при помощи значения любезности или при помощи такой утилиты, как ionice, входящей в пакет util-linux¹.

Не все планировщики ввода-вывода поддерживают приоритеты ввода-вывода. В частности, планировщик ввода-вывода CFQ (Complete Fair Queueing) поддерживает их; в настоящее время он единственный, остальные стандартные планировщики не работают с приоритетами. Если текущий планировщик ввода-вывода не поддерживает приоритеты ввода-вывода, они просто игнорируются.

Привязка процессов к процессорам

Linux поддерживает наличие нескольких процессоров в одной системе. Не считая процесс загрузки, большая часть работы по поддержке нескольких процессоров лежит на плечах планировщика процессов. На машине с *симметричной многопроцессорной обработкой* (symmetric multiprocessing, SMP) планировщик процессов должен принимать решения, какие процессы должны выполняться на каких процессорах. Эта ответственность влечет за собой две сложности: планировщик должен стремиться к полной загрузке всех процессоров в системе, так как это неэффективная ситуация, когда процессор находится в бездействии, пока какой-то процесс ожидает выполнения.

При этом, после того как процесс один раз распределяется на один из процессоров, планировщик также должен пытаться распределять его на тот же процессор в будущем. Это удобно, так как *перемещение* (migrating) процесса с одного процессора на другой несет с собой большие затраты.

Самые крупные из этих затрат связаны с *эффектами кэширования* (cache effects), сопутствующими переносу. Вследствие конструкции современных си-

¹ Пакет util-linux можно найти по адресу: <http://www.kernel.org/pub/linux/utils/util-linux>. Он лицензирован GNU General Public License v2.

стем SMP кэши, привязанные к каждому из процессоров, разделены и самостоятельны. Это означает, что данные, находящиеся в кэше одного процессора, не принадлежат кэшу другого. Таким образом, если процесс перемещается на другой процессор и записывает в память новые данные, данные в кэше старого процессора могут устаревать. Если полагаться на эти данные, это может приводить к повреждению информации. Для предотвращения этого кэши делают недействительными данные друг друга, как только они кэшируют новый фрагмент памяти. Следовательно, каждый конкретный фрагмент данных в каждый момент времени может находиться в одном, и только одном, кэше процессора (в предположении, что эти данные вообще кэшируются). Когда процесс перемещается с одного процессора на другой, возникает два вида затрат: кэшированные данные становятся недоступными для перемещенного процесса, а данные в кэше исходного процессора необходимо сделать недействительными. Из-за этих затрат планировщики процессов стараются привязывать процессы к определенным процессорам на как можно более длительные интервалы времени.

Две цели планировщика, конечно же, потенциально могут конфликтовать друг с другом. Если с одним процессором связано намного больше процессов, чем с другим, или, что еще хуже, если один процессор загружен, пока другой бездействует, имеет смысл перенаправить некоторые процессы на менее загруженный процессор. Принятие решений, когда стоит перемещать процессы в ответ на подобную несбалансированность, называется *балансировкой нагрузки* (*load balancing*) и имеет огромную важность для производительности SMP-машин.

Привязка процессов к процессорам (*processor affinity*) обозначает вероятность того, что процесс будет постоянно распределяться на один и тот же процессор. Термин *мягкая привязка* (*soft affinity*) относится к естественной склонности планировщика продолжать распределять процесс на тот же процессор. Как я уже сказал ранее, это стоящее свойство. Планировщик Linux старается распределять одни и те же процессы на одни и те же процессоры как можно дольше, перенося процессы с одного процессора на другой только в ситуациях критической несбалансированности загрузки. Это позволяет планировщику минимизировать эффекты кэширования от переноса, все так же гарантируя, что все процессоры в системе будут равномерно загружены.

Иногда, однако, пользователь или приложение желает принудительно обеспечить связку процесс—процессор. Часто бывает, что процесс очень чувствителен к кэшу и для него желательно оставаться все время на одном и том же процессоре. Сцепление процесса с определенным процессором и поддержание этой цепочки ядром называется установкой *жесткой привязки* (*hard affinity*).

Системные вызовы `sched_getaffinity()` и `sched_setaffinity()`

Процессы наследуют привязку к процессорам от своих предков и по умолчанию могут выполняться на любом процессоре. Linux предоставляет два системных

вызова, позволяющих узнавать и определять жесткую привязку процессов к процессорам:

```
#define _GNU_SOURCE
#include <sched.h>
typedef struct cpu_set_t:
size_t CPU_SETSIZE;

void CPU_SET (unsigned long cpu, cpu_set_t *set);
void CPU_CLR (unsigned long cpu, cpu_set_t *set);
int CPU_ISSET (unsigned long cpu, cpu_set_t *set);
void CPU_ZERO (cpu_set_t *set);

int sched_setaffinity (pid_t pid, size_t setsiz,
                      const cpu_set_t *set);

int sched_getaffinity (pid_t pid, size_t setsiz,
                      cpu_set_t *set);
```

Вызов `sched_getaffinity()` позволяет узнать привязку к процессору процесса с идентификатором `pid` и сохраняет результат в специальном типе `cpu_set_t`, доступ к которому осуществляется через специальный макрос. Если значение `pid` равно 0, то вызов извлекает информацию о привязке для текущего процесса. Параметр `setsiz` — это размер типа `cpu_set_t`, который может использоваться библиотекой `glibc` для обеспечения совместимости с будущими изменениями размера данного типа. В случае успеха вызов `sched_getaffinity()` возвращает значение 0; в случае ошибки он возвращает -1 и устанавливает переменную `errno`. Рассмотрим пример:

```
cpu_set_t set;
int ret, i;

CPU_ZERO (&set);
ret = sched_getaffinity (0, sizeof (cpu_set_t), &set);
if (ret == -1)
    perror ("sched_getaffinity");

for (i = 0; i < CPU_SETSIZE; i++) {
    int cpu;

    cpu = CPU_ISSET (i, &set);
    printf ("cpu=%i is %s\n",
           i, cpu ? "set" : "unset");
}
```

Перед вызовом мы используем `CPU_ZERO`, чтобы обнулить все биты в наборе `set`. Затем мы проходим по набору от 0 до `CPU_SETSIZE`. Обратите внимание, что `CPU_SETSIZE` — это не размер набора и это значение *никогда* нельзя передавать в качестве параметра `setsiz`. Это число процессоров, которые может потенциально представлять набор. Так как текущая реализация представляет каждый процессор одним битом, `CPU_SETSIZE` намного больше `sizeof(cpu_set_t)`. Мы используем `CPU_ISSET` для проверки, привязан или не привязан текущий систем-

ный процессор, *i*, к процессу. Этот макрос возвращает 0 в случае отсутствия привязки и ненулевое значение, если процессор привязан к процессу.

Устанавливаются только процессоры, физически присутствующие в системе. Таким образом, выполнение этого фрагмента кода в системе с двумя процессорами даст такой результат:

```
cpu=0 is set  
cpu=1 is set  
cpu=2 is unset  
cpu=3 is unset  
  
cpu=1023 is unset
```

Как показывает этот вывод, значение CPU_SETSIZE (для которого нумерация начинается с нуля) в настоящее время равно 1024.

Нас интересуют только процессоры 0 и 1, так как это единственные физические процессоры в данной системе. Предположим, что нам нужно удостовериться, что наш процесс будет выполняться только на процессоре 0, но никогда на процессоре 1. Следующий код позволяет сделать это:

```
cpu_set_t set:  
int ret, i;  
  
CPU_ZERO (&set);      /* очистить набор процессоров */  
CPU_SET (0, &set);    /* разрешить процессор #0 */  
CPU_CLR (1, &set);    /* запретить процессор #1 */  
ret = sched_setaffinity (0, sizeof (cpu_set_t), &set);  
if (ret == -1)  
    perror ("sched_setaffinity");  
  
for (i = 0; i < CPU_SETSIZE; i++) {  
    int cpu;  
  
    cpu = CPU_ISSET (i, &set);  
    printf ("%s\n",  
           cpu ? "set" : "unset");  
}
```

Мы начинаем, как всегда, с обнуления набора при помощи CPU_ZERO. После этого мы устанавливаем процессор 0 при помощи CPU_SET и сбрасываем (очищаем) процессор 1 при помощи CPU_CLR. Операция CPU_CLR здесь, в действительности, не нужна, так как мы обнулили весь набор, но она добавлена для полноты примера.

Если выполнить последний пример в той же двухпроцессорной системе, то результат будет немного отличаться:

```
cpu=0 is set  
cpu=1 is unset  
cpu=2 is unset  
  
cpu=1023 is unset
```

Теперь мы видим, что процессор 1 не установлен. Данный процесс будет выполняться только на процессоре 0, независимо ни от чего.

Возможны четыре значения переменной `errno`:

EFAULT

Переданный вызову указатель находится за пределами адресного пространства процесса или недействителен по другой причине.

EINVAL

В данном случае в системе не было физически включено ни одного процессора из набора `set` (только для `sched_setaffinity()`) или `setsz` меньше размера внутренней структуры данных ядра, представляющей набор процессоров.

EPERM

Процесс, указанный при помощи `pid`, не принадлежит текущему действительному идентификатору пользователя вызывающего процесса, и процесс не обладает характеристикой `CAP_SYS_NICE`.

ESRCH

Не было обнаружено ни одного процесса, связанного с `pid`.

Системы реального времени

В вычислительной технике термин *реальное время* (real-time) очень часто бывает источником неразберихи и непонимания. Система является системой «реального времени», если она подчиняется *операционным предельным срокам* (operational deadline): минимальным обязательным для соблюдения интервалам между возбудителями и реакциями. Знакомая система реального времени — это антиблокировочная тормозная система (antilock braking system, ABS), которая встраивается практически во все современные автомобили. В этой системе, когда нажимается тормоз, компьютер регулирует давление тормоза, много раз в секунду применяя максимальное давление и отпуская тормоз. Это предотвращает «блокирование» колес, которое снижает эффективность торможения или даже может отправлять автомобиль в неконтролируемый занос. В подобной системе операционные предельные сроки определяют, насколько быстро система должна реагировать на «заблокированное» состояние колеса и насколько быстро система может применять тормозное давление.

Большинство современных операционных систем, в том числе Linux, предлагаю некоторую степень поддержки процессов реального времени.

Жесткие и мягкие системы реального времени

Системы реального времени бывают двух видов: жесткие и мягкие. *Жесткая система реального времени* (hard real-time system) требует абсолютно точного соблюдения операционных предельных сроков. Превышение предельного срока означает сбой, и это для нее серьезная неполадка. *Мягкая система реального*

времени (soft real-time system), с другой стороны, не считает превышение предельного срока критической ошибкой.

Жесткие приложения реального времени легко распознать: несколько примеров включают антиблокировочные тормозные системы, военные системы вооружения, медицинские устройства и приложения обработки сигналов. Мягкие приложения реального времени идентифицировать не так просто. Один очевидный член этой группы – это приложения обработки видео: если лимиты превышаются, пользователи замечают падение качества, но пропуск нескольких кадров не критичен.

На многие другие приложения накладываются временные ограничения, которые в случае нарушения приводят к ухудшению впечатлений пользователя от работы с приложением. Сразу же приходят на ум мультимедийные приложения, игры и сетевые программы. А что насчет текстового редактора? Если программа не может достаточно быстро реагировать на нажатие клавиш, то она производит плохое впечатление, что раздражает и злит пользователя. Является ли текстовый редактор мягким приложением реального времени? Определенно, когда разработчики писали приложение, они осознавали, что оно должно будет своевременно отвечать на нажатия клавиш. Но считается ли это операционным предельным сроком? Невозможно четко определить границу вокруг множества мягких приложений реального времени.

В противоположность распространенному мнению, система реального времени не обязательно работает быстро. Действительно, если взять сравнимое аппаратное обеспечение, то система реального времени будет, вероятно, функционировать медленнее, чем *система модельного времени* (nonreal-time system), например, из-за увеличения нагрузки вследствие необходимости поддерживать процессы реального времени. Схожим образом различие между жесткими и мягкими системами реального времени не зависит от величины операционных предельных сроков. Ядерный реактор перегреется, если система SCRAM не опустит стержни регулирования мощности реактора в течение нескольких секунд после распознания избыточного нейтронного потока. Это жесткая система реального времени с большим (для мира машинных вычислений) операционным предельным сроком. И наоборот, видеопроигрыватель может пропустить кадр или споткнуться при проигрывании звука, если приложение не заполнит заново буфер воспроизведения в течение 100 миллисекунд. Это мягкая система реального времени с небольшим операционным предельным сроком.

Задержка, дрожание и предельные сроки

Задержка (latency) обозначает период между срабатыванием раздражителя и реализацией ответного действия. Если задержка меньше или равна предельному сроку, то система работает правильно. Во многих жестких системах реального времени операционный предельный срок равен задержке – система обрабатывает раздражители в течение фиксированных интервалов в точные моменты времени. В мягких системах реального времени момент ответа может быть

менее точным и значение задержки может варьироваться — целью является исключительно реагирование до наступления предельного срока.

Очень часто бывает сложно измерить задержку, так как ее вычисление требует знания момента, когда произошло раздражение. Возможность связывать временную метку с раздражением, однако, зачастую снижает способность реагировать на него. Таким образом, обычно измеряется не задержка, а вариация во времени между ответами. Вариация длины интервалов между последовательными событиями — это *дрожание* (jitter), а не задержка.

Например, предположим, что раздражение создается каждые 10 миллисекунд. Для того чтобы измерить производительность нашей системы, мы могли бы присваивать метки времени ответам системы, проверяя, что те происходят каждые десять миллисекунд. Отклонение от этой цели, однако, это не задержка, а дрожание. То, что мы измеряем, — это вариация между последовательными ответами. Не зная, когда осуществляется раздражение, мы не можем знать точную разницу во времени между раздражением и реакцией. Даже если мы знаем, что интервал между раздражениями составляет 10 миллисекунд, нам не известно, когда произошло *первое* раздражение. Это может звучать неожиданно, но во многих попытках измерить задержку делается эта ошибка и они выдают в результате значение дрожания, а не задержки. Не сомневайтесь, дрожание — полезный показатель, и измерение его значения также может пригодиться. Но тем не менее нужно называть вещи своими именами!

Жесткие системы реального времени часто демонстрируют очень низкий уровень дрожания, так как они реагируют на раздражение после — а не в течение — определенного интервала. Целью подобных систем является нулевое значение дрожания и задержка, равная предельному сроку. Если задержка превышает предельный срок, это создает сбой в системе.

Мягкие системы реального времени допускают большее дрожание. В этих системах время ответа идеально вписывается в предельный срок и часто наступает намного раньше, хотя иногда и позже него. Дрожание, таким образом, является отличной заменой задержке в качестве показателя производительности.

Поддержка процессов реального времени в Linux

Linux предоставляет приложениям мягкую поддержку приложений реального времени через семейство системных вызовов, определенных в стандарте IEEE Std 1003.1b–1993 (часто называемом просто POSIX 1993 или POSIX.1b).

Говоря технически, стандарт POSIX не диктует жесткость или мягкость поддержки реального времени. В действительности все, что делает стандарт POSIX, — это описывает несколько политик планирования, учитывающих приоритеты. Какой тип временных ограничений реализует операционная система в соответствии с этими политиками — дело разработчиков операционной системы.

С течением времени ядро Linux эволюционировало, и поддержка процессов реального времени становилась все совершеннее, обеспечивая более низкую задержку, более стабильное дрожание и при этом отсутствие снижения производительности системы. Во многом это происходило благодаря тому, что умень-

шение задержки приносило преимущества многим классам приложений, таким, как настольные процессы и процессы, ограниченные вводом-выводом, а не только приложениям реального времени. Усовершенствования также можно приписать успеху Linux в плане встроенных систем и систем реального времени.

К сожалению, многие из этих встроенных модификаций и модификаций реального времени, которые были внедрены в ядро Linux, существуют только в индивидуальных решениях Linux за пределами основного официального ядра. Некоторые из подобных модификаций обеспечивают дальнейшее снижение задержки и даже позволяют добиваться жесткой реализации процессов реального времени. В следующих разделах обсуждаются только официальные интерфейсы ядра и поведение основного ядра. К счастью, большинство относящихся к реальному времени модификаций продолжают использовать интерфейсы POSIX. Следовательно, последующее обсуждение также распространяется на модифицированные системы.

Политики планирования и приоритеты в Linux

Поведение планировщика по отношению к процессу зависит от *политики планирования* (scheduling policy) процесса, также называемой классом планирования. Помимо обычной политики, по умолчанию в Linux есть еще две политики планирования реального времени. Каждую политику представляют макросы препроцессора из файла заголовка `<sched.h>`: `SCHED_FIFO`, `SCHED_RR` и `SCHED_OTHER`.

Каждому процессу соответствует *статический приоритет* (static priority), не связанный с его значением любезности. Для обычных приложений приоритет всегда равен 0. Для процесса реального времени он может принимать значения от 1 до 99 включительно. Планировщик всегда выбирает для выполнения процесс с наивысшим приоритетом (то есть процесс с наибольшим числовым значением статического приоритета). Если в данный момент выполняется процесс со статическим приоритетом 50 и становится работоспособным процесс с приоритетом 51, то планировщик немедленно приостанавливает выполняющийся процесс и переключается на ставший доступным работоспособный процесс. И наоборот, если работает процесс с приоритетом 50, а становится работоспособным процесс с приоритетом 49, то планировщик не запускает второй до тех пор, пока процесс с приоритетом 50 не блокируется, то есть пока он не перестает быть работоспособным.

Так как у обычных процессов приоритет равен 0, любой работоспособный процесс реального времени всегда подавляет обычные процессы и выполняется в первую очередь.

Политика «первый вошел — первый вышел»

Класс FIFO (first in, first out — первый вошел, первый вышел) — это очень простая политика реального времени без квантов времени. Процесс из класса FIFO продолжает выполняться до тех пор, пока не появляется другой работоспособный процесс с более высоким приоритетом. Класс FIFO представляется макросом `SCHED_FIFO`.

Так как в этой политике кванты времени отсутствуют, ее правила работы довольно просты:

- работоспособный процесс из класса FIFO всегда выполняется, если у него наивысший приоритет в системе. В частности, как только процесс из класса FIFO становится работоспособным, он сразу же подавляет обычный процесс и начинает выполняться;
- процесс из класса FIFO продолжает выполняться до тех пор, пока он не заблокируется или не вызовет `sched_yield()` или же пока не станет работоспособным какой-то другой процесс с более высоким приоритетом;
- когда процесс из класса FIFO блокируется, планировщик удаляет его из списка работоспособных процессов. Когда он становится работоспособным в очередной раз, он вставляется в конец списка процессов с таким же приоритетом. Это означает, что он не начнет снова выполнять до тех пор, пока не закончат работу все остальные процессы с более высоким *или равным* приоритетом;
- когда процесс из класса FIFO вызывает `sched_yield()`, планировщик перемещает его в конец списка процессов с таким же приоритетом. Это означает, что он не продолжит работать до тех пор, пока не завершатся все остальные процессы с равным приоритетом. Если вызывающий процесс – единственный процесс с таким значением приоритета, то `sched_yield()` не оказывает никакого эффекта;
- когда процесс с более высоким приоритетом подавляет процесс из класса FIFO, этот приостановленный процесс остается на том же месте в списке процессов, где и был с учетом его приоритета. Это означает, что, как только процесс с более высоким приоритетом завершается, приостановленный процесс из класса FIFO продолжает выполнение;
- когда процесс присоединяется к классу FIFO или когда значение статического приоритета процесса меняется, он помещается в голову списка процессов с таким же значением приоритета. Следовательно, только что получивший приоритет процесс из класса FIFO может приостановить выполняющийся процесс с таким же приоритетом.

По существу, мы можем сказать, что процессы из класса FIFO всегда выполняются так долго, как они этого хотят, если, конечно, у них самый высокий приоритет в системе. Самые интересные правила касаются того, что происходит с процессами из класса FIFO с одним и тем же значением приоритета.

Карусельная политика

Класс RR (round-robin, карусель) идентичен классу FIFO, за исключением того, что он включает дополнительные правила на случай процессов с одинаковыми приоритетами. Этот класс представляет макрос `SCHED_RR`.

Планировщик выделяет каждому процессу из класса RR квант времени. Когда у процесса заканчивается его квант, планировщик перемещает его в конец списка процессов с таким же приоритетом. Таким образом, процессы из класса

RR, имеющие одинаковые приоритеты, выполняются в своем тесном кругу в «карусельном режиме». Если процесс с данным значением приоритета только один, то класс RR идентичен классу FIFO. В таком случае, когда его квант времени истекает, процесс просто возобновляется.

Можно считать процесс из класса RR идентичным процессу из класса FIFO, за исключением того, что он прекращает выполняться, когда истекает его квант времени, и перемещается в этот момент в конец списка работоспособных процессов с таким же приоритетом.

Выбор между SCHED_FIFO и SCHED_RR – это исключительно вопрос желаемого поведения процессов с одинаковыми приоритетами. Кванты времени в классе RR имеют смысл только для равных по приоритету процессов. Процессы из класса FIFO продолжают выполняться без помех; выполнение процессов из класса RR планируется на каждом уровне приоритета. В обоих случаях процессы с более низкими приоритетами не могут начать работать до тех пор, пока не завершатся все процессы с высокими приоритетами.

Обычная политика

Макрос SCHED_OTHER представляет стандартную политику планирования, класс модельного времени по умолчанию. У всех процессов из обычного класса значение статического приоритета равно 0. Следовательно, любой работоспособный процесс из класса FIFO или RR может подавить и приостановить работоспособный процесс из обычного класса.

Для определения приоритетов процессов внутри обычного класса планировщик использует значение любезности, о котором говорилось выше. Значение любезности никак не влияет на статический приоритет, который остается равным нулю.

Пакетная политика планирования

Макрос SCHED_BATCH представляет *пакетную политику планирования* (batch scheduling policy) или *политику планирования бездействия* (idle scheduling policy). Его поведение в каком-то роде противоположно политикам реального времени: процессы из этого класса выполняются только тогда, когда других работоспособных процессов в системе нет; они не запускаются даже тогда, когда у прочих процессов просто заканчиваются их кванты времени. Это отличается от поведения процессов с самыми высокими значениями любезности (то есть самых низкоприоритетных процессов), которые в конечном итоге тоже начинают выполняться, когда у высокоприоритетных процессов истекают кванты времени.

Установка политики планирования в Linux

Процессы могут манипулировать политикой планирования в Linux при помощи системных вызовов `sched_getscheduler()` и `sched_setscheduler()`:

```
#include <sched.h>

struct sched_param {  
    /* ... */
```

```

int sched_priority;
/* ... */
};

int sched_getscheduler (pid_t pid);
int sched_setscheduler (pid_t pid,
                      int policy,
                      const struct sched_param *sp);

```

Успешный вызов `sched_getscheduler()` возвращает политику планирования для процесса, указанного при помощи параметра `pid`. Если `pid` равен 0, то вызов возвращает политику планирования для вызывающего процесса. Целое число, определенное в файле заголовка `<sched.h>`, представляет политику планирования: политике «первый вошел — первый вышел» соответствует значение `SCHED_FIFO`; карусельная политика — это `SCHED_RR`, а обычная политика — это `SCHED_OTHER`. В случае ошибки вызов возвращает значение `-1` (которое никогда не может отражать допустимую политику планирования) и соответствующим образом устанавливает переменную `errno`.

Использовать этот вызов просто:

```

int policy;

/* узнаем нашу политику планирования */
policy = sched_getscheduler (0);

switch (policy) {
case SCHED_OTHER:
    printf ("Обычная политика\n");
    break;
case SCHED_RR:
    printf ("Карусельная политика\n");
    break;
case SCHED_FIFO:
    printf ("Политика FIFO\n");
    break;
case -1:
    perror ("sched_getscheduler");
    break;
default:
    fprintf (stderr, "Неизвестная политика!\n");
}

```

Вызов `sched_setscheduler()` устанавливает для процесса, указанного при помощи параметра `pid`, политику планирования `policy`. Все параметры, связанные с политикой, устанавливаются через аргумент `sp`. Если `pid` равен 0, то устанавливаются политика и параметры для вызывающего процесса. В случае успеха вызов возвращает значение 0. В случае ошибки он возвращает `-1` и присваивает переменной `errno` соответствующее значение.

Допустимые поля структуры `sched_param` зависят от того, какие политики планирования поддерживаются в операционной системе. Политики `SCHED_RR` и `SCHED_FIFO` требуют только одно поле, `sched_priority`, представляющее статический приоритет. Для политики `SCHED_OTHER` никакие поля не используются,

а возможные будущие политики планирования могут использовать новые поля. Таким образом, в переносимых и допустимых программах не следует делать предположения относительно компоновки структуры.

Установить политику планирования и параметры процесса просто:

```
struct sched_param sp = { .sched_priority = 1 };
int ret;

ret = sched_setscheduler (0, SCHED_RR, &sp);
if (ret == -1) {
    perror ("sched_setscheduler");
return 1;
}
```

В этом фрагменте для вызывающего процесса устанавливается карусельная политика планирования и ему присваивается статический приоритет, равный 1. Мы предполагаем, что единица представляет собой допустимый приоритет, но технически это не всегда может быть так. Подробнее о том, как находить допустимый диапазон приоритетов для конкретных политик, говорится в следующем разделе.

Установка любой политики планирования, за исключением SCHED_OTHER, требует наличия характеристики CAP_SYS_NICE. Следовательно, процессы реального времени обычно выполняются от имени пользователя root. Начиная с версии ядра 2.6.12, лимит ресурсов RLIMIT_RTPRIO позволяет пользователям, отличным от пользователя root, устанавливать политики реального времени до определенной границы приоритета.

Коды ошибок

В случае ошибки возможны четыре значения переменной errno:

EFAULT

Указатель sp указывает на недопустимую или недостижимую область памяти.

EINVAL

Параметр policy содержит недопустимую политику планирования или значение параметра sp не имеет смысла для данной политики (только sched_setscheduler()).

EPERM

У вызывающего процесса нет необходимых характеристик.

ESRCH

Значение pid не соответствует никакому выполняющемуся процессу.

Установка параметров планирования

Определенные в стандарте POSIX интерфейсы sched_getparam() и sched_setparam() позволяют извлечь или определить параметры для уже установленной политики планирования:

```
#include <sched.h>

struct sched_param {
    /* ... */
    int sched_priority;
    /* ... */
};

int sched_getparam (pid_t pid, struct sched_param *sp);

int sched_setparam (pid_t pid, const struct sched_param *sp);
```

Интерфейс `sched_getscheduler()` возвращает только политику планирования, но не связанные с ней параметры. Вызов `sched_getparam()` возвращает через аргумент `sp` параметры планирования для процесса с идентификатором `pid`:

```
struct sched_param sp;
int ret;

ret = sched_getparam (0, &sp);
if (ret == -1) {
    perror ("sched_getparam");
    return 1;
}

printf ("Приоритет равен %d\n", sp.sched_priority);
```

Если значение параметра `pid` равно 0, то вызов возвращает параметры для вызывающего процесса. В случае успеха вызов возвращает 0. В случае ошибки он возвращает `-1` и присваивает переменной `errno` соответствующее значение.

Так как вызов `sched_setscheduler()` также позволяет установить параметры планирования, вызов `sched_setparam()` полезен только для дальнейшей модификации параметров:

```
struct sched_param sp;
int ret;

sp.sched_priority = 1;
ret = sched_setparam (0, &sp);
if (ret == -1) {
    perror ("sched_setparam");
    return 1;
}
```

В случае успеха для процесса с идентификатором `pid` устанавливаются параметры планирования в соответствии со значением `sp`, а вызов возвращает значение 0. В случае ошибки вызов возвращает `-1` и присваивает переменной `errno` соответствующее значение.

Если выполнить два предыдущих фрагмента кода по порядку, то мы получим следующий результат:

Приоритет равен 1

В этом примере снова подразумевается, что значение 1 является допустимым значением приоритета. Это действительно так, но в переносимых прило-

жениях необходимо выполнять дополнительную проверку. Через мгновение я покажу вам, как проверить диапазон разрешенных значений приоритета.

Коды ошибок

В случае ошибки возможны четыре значения переменной `errno`:

`EFAULT`

Указатель `sp` указывает на недопустимую или недостижимую область памяти.

`EINVAL`

Значение параметра `sp` не имеет смысла для данной политики (только для `sched_setparam()`).

`EPERM`

У вызывающего процесса отсутствуют необходимые характеристики.

`ESRCH`

Значение `pid` не соответствует ни одному выполняющемуся процессу.

Определение диапазона допустимых значений приоритета

В предыдущих примерах мы передавали системным вызовам, отвечающим за настройку планирования, жестко закодированные значения приоритета. POSIX не дает никакой гарантии относительно того, какие приоритеты планирования могут существовать в конкретной системе, за исключением того, что между минимальным и максимальным значениями должно быть, как минимум, еще 32 значения приоритета. Как упоминалось ранее в разделе «Политики планирования и приоритеты в Linux», в Linux для двух политик планирования реального времени предусмотрен диапазон значений от 1 до 99 включительно. В правильной переносимой программе обычно реализуется собственный диапазон значений приоритета, который отображается на диапазон операционной системы. Например, если вы хотите выполнять процессы на четырех различных уровнях приоритета реального времени, можно динамически определять диапазон приоритетов и выбирать четыре значения.

В Linux предоставляется два системных вызова, позволяющих узнавать диапазон допустимых значений приоритета. Один возвращает минимальное значение, а второй максимальное:

```
#include <sched.h>

int sched_get_priority_min (int policy);

int sched_get_priority_max (int policy);
```

В случае успеха вызов `sched_get_priority_min()` возвращает минимальное, а вызов `sched_get_priority_max()` — максимальное допустимое значение приоритета для политики планирования, указанной при помощи параметра `policy`. Оба вызова возвращают 0. В случае ошибки оба системных вызова возвращают

значение **-1**. Единственная возможная ошибка — это недопустимое значение параметра **policy**, которому соответствует значение переменной **errno EINVAL**.

Использовать вызовы просто:

```
int min, max;

min = sched_get_priority_min (SCHED_RR);
if (min == -1) {
    perror ("sched_get_priority_min");
    return 1;
}

max = sched_get_priority_max (SCHED_RR);
if (max == -1) {
    perror ("sched_get_priority_max");
    return 1;
}

printf ("Диапазон приоритетов для SCHED_RR равен %d - %d\n", min, max);
```

В стандартной системе Linux выполнение этого фрагмента кода даст такой результат:

Диапазон приоритетов для SCHED_RR равен 1-99

Как говорилось выше, чем больше числовое значение приоритета, тем выше приоритет. Для того чтобы присвоить процессу самый высокий приоритет для его политики планирования, можно применить следующий код:

```
/*
 * set_highest_priority – устанавливает для процесса с идентификатором pid
 * значение приоритета планирования, равное максимальному приоритету для
 * его текущей политики планирования. Если pid равен 0, то устанавливается
 * приоритет текущего процесса.
 *
 * Возвращает 0 в случае успеха.
 */
int set_highest_priority (pid_t pid)
{
    struct sched_param sp;
    int policy, max, ret;

    policy = sched_getscheduler (pid);
    if (policy == -1)
        return -1;

    max = sched_get_priority_max (policy);
    if (max == -1)
        return -1;

    memset (&sp, 0, sizeof (struct sched_param));
    sp.sched_priority = max;
    ret = sched_setparam (pid, &sp);

    return ret;
}
```

В программах обычно извлекается минимальное и максимальное значения приоритета для данной системы, а затем для назначения необходимый приоритетов используются единичные приращения (например, max-1, max-2 и т. д.).

Интерфейс `sched_rr_get_interval()`

Как обсуждалось ранее, процессы SCHED_RR ведут себя так же, как и процессы SCHED_FIFO, за исключением того, что планировщик назначает этим процессам кванты времени. Когда у процесса SCHED_RR заканчивается его квант, планировщик переносит его в конец списка выполнения, соответствующего текущему приоритету процесса. Таким образом, все процессы SCHED_RR с одинаковым приоритетом выполняются в карусельном режиме. Процессы с более высоким приоритетом (и процессы SCHED_FIFO с таким же и более высоким приоритетом) всегда могут приостановить работающий процесс SCHED_RR, независимо от того, израсходован его квант времени или еще нет.

В POSIX определяется интерфейс, позволяющий узнать длину кванта времени для конкретного процесса:

```
#include <sched.h>

struct timespec {
    time_t tv_sec;      /* секунды */
    long tv_nsec;        /* наносекунды */
};

int sched_rr_get_interval (pid_t pid, struct timespec *tp);
```

Успешный вызов функции с ужасным именем `sched_rr_get_interval()` сохраняет в структуре `timespec`, на которую указывает параметр `tp`, значение, соответствующее кванту времени, выделенному процессу с идентификатором `pid`, и возвращает 0. В случае ошибки она возвращает -1 и присваивает переменной `errno` соответствующее значение.

Согласно POSIX, эта функция обязана работать только с процессами SCHED_RR. В Linux, однако, она позволяет узнать длину кванта времени для любого процесса. В переносимых приложениях следует предполагать, что она предназначена только для процессов из карусельного класса; в программах, предназначенных исключительно для Linux, ее можно использовать по необходимости. Пример использования функции:

```
struct timespec tp;
int ret;

/* узнаем длину кванта времени текущего задания */
ret = sched_rr_get_interval (0, &tp);
if (ret == -1) {
    perror ("sched_rr_get_interval");
    return 1;
}

/* преобразуем секунды и наносекунды в миллисекунды */
printf ("Наш квант времени равен %.2lf миллисекунд\n".
(tp.tv_sec * 1000.0f) + (tp.tv_nsec / 1000000.0f));
```

Если процесс выполняется в классе FIFO, то значения tv_sec и tv_nsec оба равны 0, обозначая бесконечность.

Коды ошибок

В случае ошибки возможны три значения переменной errno:

EFAULT

Память, на которую указывает tp, недопустима или недоступна.

EINVAL

Значение pid недопустимо (например, оно меньше нуля).

ESRCH

Значение 0 допустимо, но относится к несуществующему процессу.

Предосторожности при использовании процессов реального времени

Из-за природы процессов реального времени разработчикам нужно соблюдать осторожность при разработке и отладке подобных программ. Если программа реального времени начнет предпринимать необдуманные действия, система может полностью перестать реагировать. Любой привязанный к процессору цикл в программе реального времени, то есть любой не блокирующийся фрагмент кода, будет продолжать выполняться бесконечно, пока не станет работоспособным какой-либо процесс реального времени с более высоким приоритетом.

Следовательно, разработка программ реального времени требует внимания и осторожности. Такие программы безраздельно властвуют на машине и могут с легкостью сломать систему. Вот несколько советов и рекомендаций:

- помните, что любой привязанный к процессору цикл будет работать до завершения без перерыва, если в системе нет другого процесса с более высоким приоритетом. Если цикл бесконечный, система просто перестанет отвечать;
- так как процессы реального времени выполняются ценой всего остального в системе, к их разработке необходимо подходить с особым вниманием. Помните, что необходимо стараться не сжирать оставшееся процессорное время;
- будьте очень осторожны с активным ожиданием. Если процесс реального времени активно ожидает ресурс, удерживаемый процессом с более низким приоритетом, то он это будет делать бесконечно;
- разрабатывая процесс реального времени, держите терминал открытym, чтобы он выполнялся как процесс реального времени с более высоким приоритетом, чем у разрабатываемого процесса. В случае крайней необходимости терминал будет реагировать на ваши действия и позволит убить пошедший вразнос процесс реального времени (пока терминал бездействует, ожидая ввода с клавиатуры, он не мешает работе прочих процессов реального времени);

- утилита `chrt`, входящая в состав пакета инструментов `util-linux`, значительно упрощает проверку и установку атрибутов реального времени для других процессов. С помощью этой утилиты можно с легкостью запускать произвольные программы в классе планирования реального времени, такие, как вышеупомянутый терминал, или менять приоритеты реального времени существующих приложений.

Детерминизм

Процессы реального времени — большие поклонники детерминизма. В вычислениях в реальном времени действие является *детерминированным* (*deterministic*), если на одном и том же входе оно всегда выдает один и тот же результат за одно и то же время. Современные компьютеры представляют собой воплощение недетерминированности: несколько уровней кэша (попадания и промахи в которых случаются непредсказуемо), несколько процессоров, разбиение памяти на страницы, подкачка и многозадачность привносят хаос в любые оценки длительности каждого конкретного действия. Определенно, мы достигли точки, когда практически любое действие (за исключением доступа к жесткому диску) выполняется «невероятно быстро», но одновременно в современных системах очень трудно точно определять, сколько будет длиться каждая операция.

Приложения реального времени часто пытаются ограничить непредсказуемость в целом и задержки в худшем случае в частности. В следующих разделах обсуждаются два применяющихся для этого метода.

Доаварийная запись данных и фиксирование памяти

Представьте себе такую картину: срабатывает аппаратное прерывание от специализированного монитора приближения межконтинентальных баллистических ракет, и драйвер устройства быстро копирует данные с оборудования в ядро. Драйвер замечает, что процесс находится в спящем режиме, заблокированный в ожидании данных от узла устройств ядра. Драйвер приказывает ядру разбудить процесс. Ядро, видя, что этот процесс выполняется с политикой планирования реального времени и высоким приоритетом, немедленно приостанавливает работающий в данный момент процесс и входит в режим перегрузки с целью немедленно начать выполнять процесс реального времени. Планировщик запускает процесс реального времени и переключает контекст в адресное пространство этого процесса. Теперь процесс выполняется. Все это заняло 0,3 миллисекунды, что отлично укладывается в допустимый период задержки для худшего случая, равный 1 миллисекунде.

Теперь, уже в пользовательском пространстве, процесс реального времени замечает приближающуюся ракету и начинает обрабатывать ее траекторию. Вычислив баллистику, процесс реального времени инициирует развертывание системы противоракетной обороны. Прошла всего лишь 0,1 миллисекунды — достаточно для того, чтобы организовать ответ системы ПВО и сохранить

жизни. Но — о, нет! — код системы ПВО был сброшен на диск. Происходит ошибка обращения к отсутствующей странице, процессор переключается обратно в режим ядра, и ядро инициирует ввод-вывод с жестким диском, чтобы восстановить сброшенные на диск данные. Планировщик переводит процесс в спящий режим до того момента, когда ошибка отсутствия страницы будет исправлена. Проходит несколько секунд. Слишком поздно.

Очевидно, что разбиение на страницы и подкачка страниц привносят в систему довольно недетерминированное поведение, способное обращать процесс реального времени в хаос. Для предотвращения подобной катастрофы приложения реального времени часто «блокируют» или «жестко фиксируют» все страницы из своего адресного пространства в физической памяти, выполняя *доаварийную запись* (prefaulting) их в память и не позволяя системе сбрасывать страницы обратно на диск. Как только страницы фиксируются в памяти, ядро больше никогда не сбрасывает их на диск. Никакие операции доступа к этим страницам не приводят к ошибкам обращения к несуществующим страницам. Большинство приложений реального времени фиксируют некоторые или все свои страницы в физической памяти.

В Linux предусмотрены интерфейсы и для доаварийной записи, и для фиксирования данных. В главе 4 обсуждаются интерфейсы для доаварийной записи данных в физическую память. В главе 8 рассматриваются интерфейсы для фиксирования данных в физической памяти.

Привязка к процессору и процессы реального времени

Вторая проблема приложений реального времени — это многозадачность. Хотя ядро Linux работает на основе вытесняющей многозадачности, его планировщик не всегда способен мгновенно приостановить один процесс для того, чтобы позволить выполниться другому. Иногда бывает так, что выполняющийся в данный момент процесс работает внутри критической области ядра и планировщик не в состоянии приостановить его до тех пор, пока он не покинет эту область. Если ожидающий выполнения процесс — это процесс реального времени, такая задержка может быть недопустимой и приводить к быстрому превышению операционного предельного срока.

Итак, многозадачность добавляет в систему недетерминированность, по своей природе сходную с непредсказуемостью, сопутствующую использованию страниц памяти. Решение относительно многозадачности точно такое же, как для страниц памяти, — избегать ее. Конечно же, велики шансы, что вам просто не удастся избавиться от всех прочих процессов в системе. Если бы это было возможно в вашей среде, то, вероятно, вам вообще не понадобилась бы операционная система Linux — простой пользовательской операционной системы было бы достаточно. Если же в вашей системе несколько процессоров, то один или несколько из них можно выделить исключительно для процесса или процессов реального времени. Фактически, можно оградить процессы реального времени от многозадачности.

Мы обсуждали системные вызовы для манипулирования привязкой процессов к процессорам ранее в этой главе. Возможный вариант оптимизации для

приложений реального времени — зарезервировать один процессор для процессов реального времени и позволить всем остальным процессам делить между собой время второго процессора.

Самый простой способ сделать это — модифицировать программу init в Linux, SysVinit¹, чтобы до того, как начинать процесс загрузки, она делала что-то подобное:

```
cpu_set_t set;
int ret;

CPU_ZERO (&set);          /* очистить набор процессоров */
ret = sched_getaffinity (0, sizeof (cpu_set_t), &set);
if (ret == -1) {
    perror ("sched_getaffinity");
    return 1;
}

CPU_CLR (1, &set);        /* запретить использование процессора #1 */
ret = sched_setaffinity (0, sizeof (cpu_set_t), &set);
if (ret == -1) {
    perror ("sched_setaffinity");
    return 1;
}
```

Код в этом фрагменте захватывает текущий набор разрешенных процессоров программы init, который, как ожидается, включает все существующие процессоры. После этого он удаляет один процессор, процессор с номером 1, из этого набора и обновляет список разрешенных процессоров.

Так как набор разрешенных процессоров наследуется потомками от предков, а процесс init — это прародитель всех процессов в системе, все они выполняются с учетом данного набора разрешенных процессоров. Следовательно, никакие процессы никогда не смогут выполняться на процессоре 1.

После этого нужно модифицировать процесс реального времени, разрешив ему выполняться только на процессоре 1:

```
cpu_set_t set;
int ret;

CPU_ZERO (&set);          /* очистить набор процессоров */
CPU_CLR (1, &set);        /* запретить процессор #1 */
ret = sched_setaffinity (0, sizeof (cpu_set_t), &set);
if (ret == -1) {
    perror ("sched_setaffinity");
    return 1;
}
```

Таким образом, процесс реального времени будет работать только на процессоре 1, а все остальные процессы — только на других процессорах.

¹ Исходный код SysVinit можно найти на ftp-сервере: <ftp://ftp.cistron.nl/pub/people/miquels/sysvinit>. Он лицензирован согласно лицензии GNU General Public License v2.

Лимиты ресурсов

Ядро Linux накладывает на процессы несколько *лимитов ресурсов* (resource limit). Эти лимиты определяют жесткие максимальные объемы ресурсов ядра, которые разрешается использовать процессу, — число открытых файлов, страниц в памяти, ожидающих сигналов и т. д. Лимиты строго соблюдаются; ядро не допускает никаких действий, которые могут привести к превышению жесткого лимита потребления ресурсов процессом. Например, если после открытия файла у процесса окажется больше открытых файлов, чем разрешается применимым лимитом ресурсов, то вызов `open()` завершится ошибкой¹.

В Linux предоставляется два системных вызова для манипулирования лимитами ресурсов. В стандарте POSIX определяются оба интерфейса, но Linux поддерживает еще несколько лимитов в дополнение к тем, которые диктуются стандартом. Лимиты можно проверять при помощи системного вызова `getrlimit()` и устанавливать при помощи `setrlimit()`:

```
#include <sys/time.h>
#include <sys/resource.h>

struct rlimit {
    rlim_t rlim_cur; /* soft limit */
    rlim_t rlim_max; /* hard limit */
};

int getrlimit (int resource, struct rlimit *rlim);
int setrlimit (int resource, const struct rlimit *rlim);
```

Целочисленные константы, такие, как `RLIMIT_CPU`, представляют ресурсы. Структура `rlimit` представляет фактические лимиты. В структуре определяется два максимальных значения: *мягкий лимит* (soft limit) и *жесткий лимит* (hard limit). Ядро соблюдает мягкие лимиты ресурсов для процессов, но процесс может свободно менять свой мягкий лимит, устанавливая любое значение от 0 и до жесткого лимита включительно. Процесс, не имеющий характеристики `CAP_SYS_RESOURCE` (то есть процесс, не принадлежащий пользователю `root`), может только понижать свой жесткий лимит. Непrivилегированному процессу никогда не удастся повысить свой жесткий лимит, даже поднять его до существовавшего ранее более высокого значения — уменьшение жесткого лимита невозможно. Привилегированному процессу разрешается устанавливать в качестве жесткого лимита любое допустимое значение.

То, что в действительности представляют собой лимиты, зависит от конкретного ресурса. Например, если значение параметра `resource` равно `RLIMITFSIZE`, то лимит относится к максимальному размеру (в байтах) файла, который может создать процесс. В данном случае, если значение `rlim_cur` равно 1024, процесс не может создавать файлы размером больше килобайта или увеличивать их за пределами этого значения.

¹ В этом случае вызов присвоит переменной `errno` значение `EMFILE`, указывающее, что процесс достиг лимита ресурсов, определяющего максимальное количество открытых файлов. Системный вызов `open()` обсуждается в главе 2.

У всех лимитов ресурсов есть два специальных значения: 0 и бесконечность. Первое вообще отключает использование данного ресурса. Например, если значение `RLIMIT_CORE` равно 0, то ядро не может создавать файл ядра. И наоборот, второе специальное значение снимает любые лимиты, наложенные на ресурс. Внутри ядра бесконечность обозначается специальным значением `RLIM_INFINITY`, которое равно -1 (это может вызывать некоторую неразбериху, так как -1 — это также возвращаемое значение, указывающее на ошибку). Если `RLIMIT_CORE` равно бесконечности, то ядро может создавать файлы ядра любого размера.

Функция `getrlimit()` помещает текущий жесткий и мягкий лимиты для ресурса, указанного при помощи параметра `resource`, в структуру, на которую указывает параметр `rlim`. В случае успеха вызов возвращает значение 0. В случае ошибки он возвращает -1 и соответствующим образом устанавливает переменную `errno`.

Соответственно, функция `setrlimit()` устанавливает для ресурса `resource` жесткий и мягкий лимиты, используя для этого значения из структуры `rlim`. В случае успеха вызов возвращает 0, а ядро обновляет лимиты ресурса в соответствии с запросом. В случае ошибки вызов возвращает -1 и устанавливает переменную `errno`.

Лимиты

В Linux в настоящий момент существует 15 лимитов ресурсов:

`RLIMIT_AS`

Ограничивает максимальный размер адресного пространства процесса в байтах. Попытки увеличить размер адресного процесса за пределы этого лимита — при помощи таких вызовов, как `mmap()` и `brk()` — завершаются ошибкой с кодом `ENOMEM`. Если стек процесса, который по необходимости автоматически увеличивается, выходит за пределы этого лимита, ядро отправляет процессу сигнал `SIGSEGV`. Обычно этот лимит равен `RLIM_INFINITY`.

`RLIMIT_CORE`

Диктует максимальный размер файлов ядра в байтах. Если он не равен нулю, то файлы ядра, размер которых превышает этот лимит, усекаются до максимального допустимого размера. Если данный лимит равен нулю, то файлы ядра никогда не создаются.

`RLIMIT_CPU`

Диктует максимальное процессорное время, которое процесс может потратить, в секундах. Если процесс выполняется дольше, чем определено лимитом, то ядро отправляет ему сигнал `SIGXCPU`, который процессы могут захватывать и обрабатывать. Портативные программы при получении этого сигнала должны завершаться, так как в стандарте POSIX не определяется, какое действие ядро может предпринять далее. Некоторые системы могут прерывать процесс, если он продолжает выполняться. Linux, однако, не запрещает процессу работать, но продолжает отправлять сигналы `SIGXCPU`.

с интервалом в одну секунду. Как только процесс достигает жесткого лимита, ему отправляется сигнал SIGKILL и процесс убивается.

RLIMIT_DATA

Управляет максимальным размером сегмента данных и кучи процесса в байтах. Попытки расширить сегмент данных за пределы этого значения при помощи вызова `brk()` завершаются ошибкой с кодом ENOMEM.

RLIMITFSIZE

Указывает максимальный размер файла, который процесс может создать, в байтах. Если процесс увеличивает файл за пределы этого значения, то ядро отправляет процессу сигнал SIGXFSZ. По умолчанию данный сигнал прекращает процесс. Процесс может, однако, захватить и обработать этот сигнал, и тогда ставший его причиной системный вызов завершится ошибкой и вернет значение EFBIG.

RLIMITLOCKS

Управляет максимальным числом файловых блокировок, которые может удерживать процесс (подробнее о блокировках файлов говорится в главе 7). Как только этот лимит достигается, все остальные попытки захватить дополнительные файловые блокировки должны завершаться ошибкой и возвращать ошибку ENOLCK. В ядре Linux 2.4.25, однако, эта функциональность была устранена. В текущих версиях ядра данный лимит можно устанавливать, но он не оказывает никакого эффекта.

RLIMITMEMLOCK

Определяет максимальное число байтов памяти, которое процесс без характеристики CAP_SYS_IPC (фактически, процесс, не принадлежащий пользователю root) может зафиксировать при помощи `mlock()`, `mlockall()` или `shmctl()`. Если этот лимит превышается, то перечисленные вызовы завершаются ошибкой и возвращают код EPERM. На практике реальный лимит округляется вниз до целого числа страниц. Процессы, обладающие характеристикой CAP_SYS_IPC, могут фиксировать любое число страниц в памяти, и на них данный лимит не распространяется. До появления версии ядра 2.6.9 этот лимит определял максимальный объем памяти, который мог зафиксировать процесс с характеристикой CAP_SYS_IPC, а непривилегированные процессы вообще не могли фиксировать никакие страницы. Данный лимит не является частью стандарта POSIX; он впервые появился в BSD.

RLIMITMSGQUEUE

Определяет максимальное число байтов, которое пользователь может выделить для очередей сообщений POSIX. Если только что созданная очередь сообщений может превысить этот лимит, то `mq_open()` завершается ошибкой и возвращает код ENOMEM. Данный лимит не является частью стандарта POSIX; он был добавлен в версии ядра 2.6.8 и унаследован для Linux.

RLIMITNICE

Указывает максимальное значение, до которого процесс может понизить свое значение любезности (повысить приоритет). Как говорилось ранее в этой

главе, обычно процессы могут только увеличивать свои значения любезности (понижать приоритет). Данный лимит позволяет администратору определять максимальный уровень (минимальное значение любезности), до которого процессам разрешается поднимать свой приоритет. Так как значения любезности могут быть отрицательными, ядро интерпретирует это значение как $20 - rlim_{cur}$. Таким образом, если установлен лимит, равный 40, то процессу разрешается понизить свое значение любезности максимум до -20 (это будет высочайший для него приоритет). Данный лимит впервые появился в версии ядра 2.6.12.

RLIMIT_NOFILE

Это значение на единицу превышает максимальное число дескрипторов файла, которое процесс может держать в открытом состоянии. Попытки пре-взойти данный лимит приводят к ошибке, а соответствующий системный вызов возвращает значение `EMFILE`. Этот лимит также можно указывать при помощи имени `RLIMIT_OFILE`, которое используется в BSD.

RLIMIT_NPROC

Определяет максимальное число процессов, которое у пользователя может выполняться в системе в каждый момент времени. Попытки превысить этот лимит приводят к ошибке, а системный вызов `fork()` возвращает код `EAGAIN`. Данный лимит не входит в стандарт POSIX; впервые он был представлен в BSD.

RLIMIT_RSS

Указывает максимальное число страниц, которое может постоянно храниться в памяти для данного процесса (также носит название *размер резидентной части* – resident set size, RSS). Данный лимит соблюдался только в первых версиях ядра 2.4. Текущие версии ядра позволяют устанавливать его, но он не соблюдается. Данный лимит не входит в стандарт POSIX; впервые он был представлен в BSD.

RLIMIT_RTPRIO

Указывает максимальный уровень приоритета реального времени, который может запросить процесс, не обладающий характеристикой `CAP_SYS_NICE` (то есть процесс, не принадлежащий пользователю `root`). Обычно непривилегированным процессам не разрешается запрашивать никакие классы планирования реального времени. Данный лимит не является частью стандарта POSIX; он был добавлен в версии ядра 2.6.12 и уникален для Linux.

RLIMIT_SIGPENDING

Определяет максимальное число сигналов (стандартных и реального времени), которые могут находиться в очереди для данного пользователя. Попытки добавить в очередь другие сигналы завершаются ошибкой, а системные вызовы, такие, как `sigqueue()`, возвращают код `EAGAIN`. Обратите внимание, что, независимо от этого лимита, всегда можно поместить в очередь один экземпляр еще не поставленного в очередь сигнала. Это означает, что всегда есть возможность доставить процессу сигнал `SIGKILL` или `SIGTERM`. Данный лимит не является частью стандарта POSIX; он уникален для Linux.

RLIMIT_STACK

Определяет максимальный размер стека процесса в байтах. Превышение данного лимита приводит к отправке сигнала SIGSEGV.

Ядро хранит лимиты отдельно для каждого пользователя. Другими словами, у всех процессов, которые выполняются от имени одного пользователя, одинаковые мягкие и жесткие лимиты для каждого конкретного ресурса. Сами лимиты, однако, могут описывать разрешенные значения для отдельных процессов (а не для отдельных пользователей). Например, ядро хранит лимит RLIMIT_NOFILE отдельно для каждого пользователя; его значение по умолчанию равно 1024. Этот лимит, однако, определяет максимальное число файлов, которые может открыть *каждый процесс*, а не общее разрешенное количество открытых файлов для пользователя. Обратите внимание: это не означает, что данный лимит можно настроить индивидуально для каждого из процессов пользователя, — если один процесс меняет мягкий лимит RLIMIT_NOFILE, то это изменение распространяется на все процессы, которыми владеет данный пользователь.

Лимиты по умолчанию

Доступные процессу лимиты по умолчанию зависят от трех переменных: первоначальный мягкий лимит, первоначальный жесткий лимит и ваш системный администратор. Ядро диктует первоначальные значения жесткого и мягкого лимитов; они перечислены в табл. 6.1. Ядро устанавливает эти лимиты для процесса `init`, и, так как потомки наследуют лимиты своих родителей, все последующие процессы также получают мягкий и жесткий лимиты процесса `init`.

Лимиты по умолчанию могут меняться в двух ситуациях:

- любой процесс может свободно увеличивать мягкий лимит до любого значения от 0 до жесткого лимита, а также уменьшать жесткий лимит. Потомки наследуют обновленные значения лимита во время ветвления;
- привилегированному процессу разрешается устанавливать любое значение жесткого лимита. Потомки наследуют обновленные лимиты во время ветвления.

Маловероятно, что корневой процесс (процесс, принадлежащий пользователю `root`) в обычной генеалогии процессов будет менять значения жестких лимитов. Следовательно, первый вариант является более вероятным источником модификаций лимитов, чем второй. Действительно, фактические значения лимитов для процесса обычно устанавливаются оболочкой пользователя, которую системный администратор может настроить так, чтобы предоставлять различные лимиты. В оболочке Bourne-again (`bash`), например, администратор делает это при помощи команды `ulimit`. Обратите внимание, что администратор не обязательно уменьшает значения; он также может повышать мягкие лимиты до жестких, давая пользователям более разумные лимиты по умолчанию. Это часто делается с `RLIMIT_STACK`, для которого во многих системах устанавливается значение `RLIMIT_INFINITY`.

Таблица 6.1. Мягкие и жесткие лимиты ресурсов по умолчанию

Лимит ресурса	Мягкий лимит	Жесткий лимит
RLIMIT_AS	RLIMIT_INFINITY	RLIMIT_INFINITY
RLIMIT_CORE	0	RLIMIT_INFINITY
RLIMIT_CPU	RLIMIT_INFINITY	RLIMIT_INFINITY
RLIMIT_DATA	RLIMIT_INFINITY	RLIMIT_INFINITY
RLIMIT_FSIZE	RLIMIT_INFINITY	RLIMIT_INFINITY
RLIMIT_LOCKS	RLIMIT_INFINITY	RLIMIT_INFINITY
RLIMIT_MEMLOCK	8 страниц	8 страниц
RLIMIT_MSGQUEUE	800 Кбайт	800 Кбайт
RLIMIT_NICE	0	0
RLIMIT_NOFILE	1024	1024
RLIMIT_NPROC	0 (подразумевает отсутствие лимита)	0 (подразумевает отсутствие лимита)
RLIMIT_RSS	RLIMIT_INFINITY	RLIMIT_INFINITY
RLIMIT_RTPRIO	0	0
RLIMIT_SIGPENDING	0	0
RLIMIT_STACK	8 Мбайт	RLIMIT_INFINITY

Установка и проверка лимитов

Закончив изучение разнообразных лимитов ресурсов, можно обратиться к способам проверки и установки этих лимитов. Код, позволяющий узнать лимит ресурса, довольно прост:

```
struct rlimit rlim;
int ret;

/* получение лимитов для размера файлов ядра */
ret = getrlimit (RLIMIT_CORE, &rlim);
if (ret == -1) {
    perror ("getrlimit");
    return 1;
}

printf ("Лимиты RLIMIT_CORE: мягкий=%ld жесткий=%ld\n",
       rlim.rlim_cur, rlim.rlim_max);
```

Если скомпилировать этот фрагмент в составе большой программы и выполнить его, то мы получим следующий результат:

Лимиты RLIMIT_CORE: мягкий=0 жесткий=-1

Мягкий лимит равен нулю, а жесткий лимит равен бесконечности (-1 обозначает RLIM_INFINITY). Таким образом, можно установить новый мягкий лимит, присвоив ему любое значение. В следующем примере определяется максимальный размер файла ядра, равный 32 Мбайт:

```
struct rlimit rlim;
int ret;

rlim.rlim_cur = 32 * 1024 * 1024; /* 32 Мбайт */
rlim.rlim_max = RLIM_INFINITY; /* не менять */
ret = setrlimit (RLIMIT_CORE, &rlim);
if (ret == -1) {
    perror ("setrlimit");
    return 1;
}
```

Коды ошибок

В случае ошибки возможны три кода errno:

EFAULT

Память, на которую указывает rlim, недопустима или недоступна.

EINVAL

Значение, указанное при помощи resource, недопустимо, или rlim.rlim_cur больше rlim.rlim_max (только для setrlimit()).

EPERM

Вызывающий не обладает характеристикой CAP_SYS_RESOURCE, но сделал попытку поднять жесткий лимит.

7

Управление файлами и каталогами

В главах 2, 3 и 4 мы рассмотрели множество подходов к файловому вводу-выводу. В этой главе мы снова будем заниматься файлами, на этот раз фокусируясь не на чтении или записи файлов, а на манипулировании и управлении ими и их метаданными.

Файлы и их метаданные

Как я рассказывал в главе 1, каждый файл представляется структурой inode, которой присваивается адрес при помощи уникального для файловой системы числового значения, называемого *номером inode* (inode number). Inode – это и физический объект, находящийся на диске в Unix-подобной файловой системе, и концептуальная сущность, представляемая структурой данных в ядре Linux. В inode хранятся *метаданные* (metadata), связанные с файлом, такие, как разрешения на доступ к файлу, временная отметка, указывающая время последнего доступа к файлу, владелец, группа и размер файла, а также местоположение данных файла.

Получить номер inode для файла можно, добавив флаг `-i` к команде `ls`:

```
$ ls -i
1689459 Kconfig    1689461 main.c      1680144 process.c  1689464 swsus.c
1680137 Makefile   1680141 pm.c       1680145 smp.c     1680149 user.c
1680138 console.c 1689462 power.h   1689463 snapshot.c
1689460 disk.c     1680143 poweroff.c 1680147 swap.c
```

Этот вывод демонстрирует, что, например, у файла `disk.c` номер inode равен 1689460. В данной конкретной файловой системе ни у какого другого файла нет такого же номера inode. В другой файловой системе, однако, такие гарантии дать невозможно.

Семейство stat

В Unix есть семейство функций, предназначенных для извлечения метаданных файла:

```
#include <sys/types.h>
#include <sys/stat.h>
#include <unistd.h>

int stat (const char *path, struct stat *buf);
int fstat (int fd, struct stat *buf);
int lstat (const char *path, struct stat *buf);
```

Каждая из этих функций возвращает информацию о файле. `stat()` предоставляет информацию о файле, указанном при помощи пути в параметре `path`, тогда как `fstat()` возвращает сведения о файле, представляющем дескриптором файла `fd`. `Lstat()` идентична `stat()`, за исключением того, что, если передать ей символьическую ссылку, то `lstat()` вернет информацию о самой ссылке, а не о целевом файле.

Каждая из этих функций сохраняет информацию в структуре `stat`, которую ей предоставляет пользователь. Сама структура определяется в файле заголовка `<bits/stat.h>`, который включается через файл заголовка `<sys/stat.h>`:

```
struct stat {
    dev_t st_dev;          /* идентификатор устройства, на котором */
                           /* хранится файл */
    ino_t st_ino;          /* номер inode */
    mode_t st_mode;        /* разрешения */
    nlink_t st_nlink;      /* число жестких ссылок */
    uid_t st_uid;          /* идентификатор пользователя владельца */
    gid_t st_gid;          /* идентификатор группы владельца */
    dev_t st_rdev;          /* идентификатор устройства */
                           /* (для специальных файлов) */
    off_t st_size;          /* общий размер в байтах */
    blksize_t st_blksize;   /* размер блока для ввода-вывода */
                           /* в файловой системе */
    blkcnt_t st_blocks;     /* количество выделенных блоков */
    time_t st_atime;        /* время последнего доступа */
    time_t st_mtime;        /* время последней модификации */
    time_t st_ctime;        /* время последнего изменения статуса */
};
```

Подробнее о полях:

- поле `st_dev` описывает узел устройства, на котором располагается файл (далее мы поговорим об узлах устройств). Если файл не поддерживается устройством, например находится на подмонтированном ресурсе NFS, то это значение равно 0;
- в поле `st_ino` находится номер inode для данного файла;
- поле `st_mode` содержит байты режима файла. О байтах режима и разрешениях рассказывалось в главах 1 и 2;
- в поле `st_nlink` хранится число жестких ссылок, указывающих на файл. У каждого файла есть, по крайней мере, одна жесткая ссылка;
- поле `st_uid` предоставляет идентификатор пользователя, владеющего файлом;
- поле `st_gid` предоставляет идентификатор группы, владеющей файлом;
- если файл является узлом устройства, то поле `st_rdev` описывает устройство, которое данный файл представляет;

- в поле `st_size` хранится размер файла в байтах;
- поле `st_blksize` содержит предпочтительный размер блока для эффективного файлового ввода-вывода. Это значение (или кратное ему значение) представляет оптимальный размер блока для ввода-вывода с пользовательской буферизацией (см. главу 3);
- в поле `st_blocks` находится число блоков файловой системы, связанных с данным файлом. Это значение меньше значения в поле `st_size`, если в файле есть дыры (то есть если это разреженный файл);
- поле `st_atime` содержит *время последнего доступа к файлу* (last file access time). Это момент, когда было выполнено самое последнее обращение к файлу (например, при помощи вызова `read()` или `execle()`);
- поле `st_mtime` содержит *время последней модификации файла* (last file modification time), то есть момент последней записи в файл;
- поле `st_ctime` содержит *время последнего изменения файла* (last file change time). Часто его считают временем создания файла, которое в Linux и других Unix-подобных системах не сохраняется. В действительности это поле описывает момент, когда в последний раз менялись метаданные файла (например, владелец или разрешения).

В случае успеха все три вызова возвращают значение 0 и записывают метаданные файла в предоставленную им структуру `stat`. В случае ошибки они возвращают -1 и присваивают переменной `errno` одно из следующих значений:

EACCESS

У вызывающего процесса отсутствуют разрешения на поиск для одного из каталогов в пути `path` (только для функций `stat()` и `lstat()`).

EBADF

Недопустимое значение `fd` (только для функции `fstat()`).

EFAULT

Параметр `path` или `buf` содержит недопустимый указатель.

ELOOP

Параметр `path` включает слишком много символьических ссылок (только для функций `stat()` и `lstat()`).

ENAMETOOLONG

Слишком длинное значение параметра `path` (только для функций `stat()` и `lstat()`).

ENOENT

Компонент пути `path` не существует (только для функций `stat()` и `lstat()`).

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOTDIR

Компонент пути `path` не является каталогом (только для функций `stat()` и `lstat()`).

В следующей программе функция stat() используется для извлечения размера файла, указанного в командной строке:

```
#include <sys/types.h>
#include <sys/stat.h>
#include <unistd.h>
#include <stdio.h>

int main (int argc, char *argv[])
{
    struct stat sb;
    int ret;

    if (argc < 2) {
        fprintf (stderr,
                 "использование: %s <файл>\n", argv[0]);
        return 1;
    }

    ret = stat (argv[1], &sb);
    if (ret) {
        perror ("stat");
        return 1;
    }

    printf ("Размер %s равен %ld байтов\n",
            argv[1], sb.st_size);

    return 0;
}
```

Результат выполнения этой программы на ее исходном файле:

```
$ ./stat stat.c
Размер stat.c равен 392 байт
```

В следующем фрагменте, в свою очередь, функция fstat() используется для проверки того, находится ли уже открытый файл на физическом устройстве (в противоположность сетевому ресурсу):

```
/*
 * is_on_physical_device – возвращает положительное целое число,
 * если файл с дескриптором fd находится на физическом устройстве.
 * 0 – если файл находится на нефизическом или виртуальном
 * устройстве (например, подмонтированном ресурсе NFS) и
 * -1 в случае ошибки.
 */
int is_on_physical_device (int fd)
{
    struct stat sb;
    int ret;

    ret = fstat (fd, &sb);
    if (ret) {
        perror ("fstat");
        return -1;
    }

    return gnu_dev_major (sb.st_dev);
}
```

Разрешения

Хотя вызовы `stat` можно использовать для получения значений разрешений для данного файла, для установки этих значений используются два других системных вызова:

```
#include <sys/types.h>
#include <sys/stat.h>

int chmod (const char *path, mode_t mode);
int fchmod (int fd, mode_t mode);
```

Оба вызова `chmod()` и `fchmod()` устанавливают для файла разрешения, указанные при помощи параметра `mode`. В вызове `chmod()` параметр `path` содержит относительный или абсолютный путь к модифицируемому файлу. Для вызова `fchmod()` файл указывается при помощи дескриптора `fd`.

Допустимые значения параметра `mode`, представляемого непрозрачным целочисленным типом `mode_t`, те же, что и значения, возвращаемые в поле `st_mode` структуры `stat`. Хотя это простые целочисленные значения, они могут иметь разный смысл в разных реализациях Unix. Следовательно, в POSIX определяется набор констант, представляющих разнообразные разрешения (подробнее об этом говорится в разделе «Разрешения новых файлов» главы 2). Для формирования допустимых значений параметра `mode` эти константы можно объединять между собой операцией двоичного ИЛИ. Например, (`S_IRUSR | S_IRGRP`) устанавливает для файла разрешения, предоставляющие возможность чтения файла его владельцу и группе.

Чтобы иметь возможность изменить разрешения файла, действительный идентификатор процесса,зывающего `chmod()` или `fchmod()`, должен совпадать с идентификатором владельца файла или же процесс должен обладать характеристикой `CAP_FOWNER`.

В случае успеха оба вызова возвращают значение 0. В случае ошибки оба возвращают -1 и присваивают переменной `errno` один из следующих кодов ошибки:

`EACCES`

Узывающему процессу нет разрешений на поиск для компонента пути `path` (только для `chmod()`).

`EBADF`

Недопустимый дескриптор файла `fd` (только для `fchmod()`).

`EFAULT`

Параметр `path` содержит недопустимый указатель (только для `chmod()`).

`EIO`

В файловой системе произошла внутренняя ошибка ввода-вывода. Это очень плохая ситуация, которая может указывать на повреждение диска или файловой системы.

`ELoop`

Во время разрешения пути `path` ядро встретило слишком много симвлических ссылок (только для `chmod()`).

ENAMETOOLONG

Слишком длинное значение параметра `path` (только для `chmod()`).

ENOENT

Путь `path` не существует (только для параметра `chmod()`).

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOTDIR

Компонент пути `path` не является каталогом (только для `chmod()`).

EPERM

Действительный идентификатор вызывающего процесса не соответствует владельцу файла, и у процесса отсутствует характеристика `CAP_FOWNER`.

EROFS

Файл находится в файловой системе, доступной только для чтения.

В следующем фрагменте кода для файла `map.png` устанавливаются разрешения, позволяющие владельцу считывать и записывать файл:

```
int ret;  
  
/*  
 * Установка для файла map.png в текущем каталоге разрешений на  
 * считывание и запись для владельца файла. Эквивалентно команде  
 * chmod 600 ./map.png.  
 */  
ret = chmod("./map.png", S_IRUSR | S_IWUSR);  
if (ret)  
    perror("chmod");
```

Следующий фрагмент кода выполняет ту же самую задачу, считая, что дескриптор `fd` представляет открытый файл `map.png`:

```
int ret;  
  
/*  
 * Установка для файла с дескриптором fd разрешений на чтение  
 * и запись владельцем файла.  
 */  
ret = fchmod(fd, S_IRUSR | S_IWUSR);  
if (ret)  
    perror("fchmod");
```

Оба вызова — `chmod()` и `fchmod()` — доступны во всех современных системах Unix. Стандарт POSIX требует реализации первого вызова, а второй считается необязательным.

Владение

В структуре `stat` в полях `st_uid` и `st_gid` хранятся владелец и группа файла соответственно. Три системных вызова позволяют пользователю менять эти два значения:

```
#include <sys/types.h>
#include <unistd.h>

int chown (const char *path, uid_t owner, gid_t group);
int lchown (const char *path, uid_t owner, gid_t group);
int fchown (int fd, uid_t owner, gid_t group);
```

Вызовы `chown()` и `lchown()` определяют владение файлом, указанным при помощи пути в параметре `path`. Они работают одинаково для всех файлов, за исключением символических ссылок. В этом случае первый вызов проходит по символьским ссылкам и меняет настройки владения целевого файла, а не самой ссылки, а `lchown()` не проходит по символьским ссылкам и, таким образом, меняет настройки владения самого файла символьской ссылки. Вызов `fchown()` устанавливает параметры владения для файла, представляемого дескриптором файла `fd`.

В случае успеха все три вызова делают владельцем файла пользователя `owner`, группой файла группу `group` и возвращают значение 0. Если параметр `owner` или `group` равен -1, то соответствующее значение не устанавливается. Только процесс, владеющий характеристикой `CAP_CHOWN` (обычно это процесс, принадлежащий пользователю `root`), может менять владельца файла. Владелец файла может сменить группу файла на любую другую, членом которой он является; процессы с характеристикой `CAP_CHOWN` имеют право менять группу файла на любое другое значение.

В случае ошибки вызовы возвращают -1 и присваивают переменной `errno` одно из следующих значений:

EACCES

У вызывающего процесса отсутствуют разрешения на поиск для компонента пути `path` (только для `chown()` и `lchown()`).

EBADF

Недопустимое значение `fd` (только для `fchown()`).

EFAULT

Недопустимое значение `path` (только для `chown()` и `lchown()`).

EIO

Внутренняя ошибка ввода-вывода (очень плохо).

ELOOP

Во время разрешения пути `path` ядро встретило слишком много символьских ссылок (только для `chown()` и `lchown()`).

ENAMETOOLONG

Слишком длинное значение `path` (только для `chown()` и `lchown()`).

ENOENT

Файл не существует.

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOTDIR

Компонент пути path не является каталогом (только для chown() и lchown()).

EPERM

У вызывающего процесса отсутствуют необходимые права для изменения владельца или группы в соответствии с запросом.

EROFS

Файловая система доступна только для чтения.

Следующий фрагмент кода меняет группу файла manifest.txt в текущем рабочем каталоге на группу officers. Для того чтобы это действие удалось, вызывающий пользователь должен либо обладать характеристикой CAP_CHOWN, либо быть пользователем kidd в группе officers:

```
struct group *gr;
int ret;

/*
 * getgrnam( ) возвращает информацию о группе, принимая в качестве
 * параметра ее имя.
 */
gr = getgrnam ("officers");
if (!gr) {
    /* вероятно, недопустимая группа */
    perror ("getgrnam");
    return 1;
}

/* установка для файла manifest.txt группы officers */
ret = chmod ("manifest.txt", -1, gr->gr_gid);
if (ret)
    perror ("chmod");
```

Перед вызовом файл принадлежит группе crew:

```
$ ls -l
-rw-r--r- 1 kidd crew 13274 May 23 09:20 manifest.txt
```

После вызова файл принадлежит только группе officers:

```
$ ls -l
-rw-r--r- 1 kidd officers 13274 May 23 09:20 manifest.txt
```

Владелец файла, пользователь kidd, не изменился, так как в коде для параметра uid было передано значение –1.

Эта функция устанавливает для файла, представляемого дескриптором fd, пользователя-владельца и группу-владельца root:

```
/*
 * make_root_owner – меняет владельца и группу файла, указанного
 * при помощи параметра fd, на root. Возвращает 0 в случае успеха
 * и -1 в случае ошибки.
 */
int make_root_owner (int fd)
{
    int ret;
```

```
/* 0 – это gid и uid для root */
ret = fchown (fd, 0, 0);
if (ret)
    perror ("fchown");

return ret.
}
```

Вызывающий процесс должен обладать характеристикой `CAP_CHOWN`. Что касается характеристик, обычно это означает, что вызывающий процесс должен принадлежать пользователю `root`.

Расширенные атрибуты

Расширенные атрибуты (extended attributes), также называемые `xattr`, обеспечивают механизм для создания постоянных связок между файлами и парами из ключа и значения. В этой главе мы уже обсудили всевозможные метаданные, включающие ключ и значение и относящиеся к файлам: размер файла, владелец, время последней модификации и т. д. Расширенные атрибуты позволяют существующим файловым системам поддерживать новые возможности, которые не были предусмотрены исходно, например обязательное управление доступом в целях соблюдения безопасности. Такими интересными расширенными атрибутами делает то, что приложения из пользовательского пространства имеют возможность по своему усмотрению создавать, считывать и записывать пары из ключа и значения.

Расширенные атрибуты *не связаны с файловой системой* (filesystem-agnostic) в том смысле, что приложения для манипулирования атрибутами применяют стандартный интерфейс. Этот интерфейс не является уникальной особенностью какой-либо файловой системы. Таким образом, в приложениях можно применять расширенные атрибуты, не задумываясь о том, в какой файловой системе хранятся файлы или как внутри файловой системы представляются ключи и значения. И все же сама реализация расширенных атрибутов сильно зависит от файловой системы. В разных файловых системах расширенные атрибуты хранятся по-разному, но ядро скрывает эти различия, абстрактно представляя их при помощи интерфейса расширенных атрибутов.

Например, в файловой системе ext3 расширенные атрибуты файла хранятся в пустом пространстве внутри inode файла¹. Благодаря этому чтение расширенных атрибутов осуществляется очень быстро. Так как блок файловой системы, содержащий inode, считывается с диска в память, как только приложение обращается к файлу, расширенные атрибуты также «автоматически» попадают

¹ Конечно, до тех пор, пока в inode остается пространство. После этого файловая система ext3 начинает записывать расширенные атрибуты в дополнительные блоки файловой системы. В старых версиях ext3 отсутствует эта возможность хранить расширенные атрибуты внутри inode.

в память, где к ним можно обращаться, не создавая никакую дополнительную нагрузку на систему.

Прочие файловые системы, такие, как FAT и minixfs, не поддерживают расширенные файловые атрибуты. Эти файловые системы возвращают код ошибки ENOTSUP, если попытаться выполнить на их файлах операции, включающие расширенные атрибуты.

Ключи и значения

Уникальный *ключ* (key) идентифицирует каждый расширенный атрибут. Ключ должен быть в формате UTF-8. Каждый ключ имеет вид `пространство_имен.атрибут`. Каждый ключ должен быть полностью уточнен, то есть начинаться с допустимого пространства имен, за которым должна следовать точка. Пример допустимого имени ключа: `user.mime_type` — этот ключ принадлежит пространству имен `user`, а имя атрибута у него равно `mime_type`.

Ключ может быть *определенным* (defined) или *неопределенным* (undefined). Если ключ определен, то его значение может быть пустым или непустым. Это означает, что между неопределенным ключом и определенным ключом, с которым не связано значение, существует определенное различие. Как вы узнаете далее, для удаления ключей необходимо использовать особый интерфейс (недостаточно просто назначить ключу пустое значение).

Значение, связанное с ключом, если оно не пусто, может представлять собой произвольный массив байтов. Так как значение не обязательно должно быть строкой, оно не должно завершаться нулем, хотя это определенно имеет смысл, если в качестве значения ключа вы используете строку C. Так как значения не обязательно завершаются нулем, все операции на расширенных атрибутах должны знать и учитывать размер значения. При считывании атрибута размер предоставляет ядро; при записи атрибута вы должны самостоятельно указывать размер.

ЛУЧШИЙ СПОСОБ ХРАНЕНИЯ ТИПОВ MIME

Диспетчеры файлов с графическим интерфейсом пользователя, такие, как Nautilus от GNOME, демонстрируют разное поведение для разных типов файлов: предлагают уникальные пиктограммы, разнообразное поведение по умолчанию при щелчке мышью, специальные списки операций над файлами и т. д. Для того чтобы уметь делать это, диспетчерам файлов необходимо знать формат каждого файла. Для определения формата файловых систем, подобные Windows, просто проверяют расширение файла. По причинам, связанным как с традициями, так и с безопасностью, системы Unix обычно изучают файл и интерпретируют его тип. Этот процесс называется анализом типа MIME (MIME type sniffing).

Некоторые диспетчеры файлов генерируют эту информацию на лету; другие создают ее единожды и сохраняют в кэше. Кэширующие диспетчеры обычно помещают эти сведения в собственную базу данных. Диспетчеру файлов приходится поддерживать базу данных, синхронизируя ее с файлами, которые могут меняться без его ведома. Гораздо лучший подход — просто избавиться от собственной базы данных и хранить подобные метаданные в расширенных атрибутах: их проще поддерживать, доступ к ним быстрее и они всегда доступны для любого приложения.

Linux не накладывает никаких ограничений на число ключей, длину ключей, размер значения или общее пространство, занимаемое всеми ключами и значениями, связанными с файлом. Однако в файловых системах существуют опре-

деленные практические лимиты. Они обычно ограничивают общий объем всех ключей и значений для каждого файла.

В ext3, например, все расширенные атрибуты для каждого файла должны умещаться в пустое пространство внутри inode файла и занимать не больше одного блока файловой системы (в более старых версиях ext3 существовало только ограничение на один блок файловой системы, без дополнительного хранилища внутри inode). На практике это эквивалентно лимиту от 1 до 8 Кбайт на файл, в зависимости от размера блоков в файловой системе. В XFS, в противоположность этому, нет никаких практических лимитов. Даже в ext3 такие лимиты в действительности не представляют проблемы, так как большинство ключей и значений — это короткие текстовые строки. Тем не менее не стоит забывать о них — подумайте дважды, прежде чем сохранить всю историю управления версиями проекта в расширенных атрибутах файла!

Пространства имен расширенных атрибутов

Пространства имен, связанные с расширенными атрибутами, — это больше, чем просто организационные инструменты. Ядро реализует разные политики доступа в зависимости от пространства имен.

В Linux в настоящее время определяются четыре пространства имен расширенных атрибутов и в будущем могут появиться дополнительные. Существуют следующие пространства имен:

system

Пространство имен system используется для реализации функций ядра, применяющих расширенные атрибуты, такие, как *справки управления доступом* (access control list, ACL). Пример расширенного атрибута из этого пространства имен — system.posix_acl_access. Может ли пользователь читать или записывать эти атрибуты, — зависит от используемого модуля безопасности. Предполагайте худший случай, когда ни один пользователь (в том числе, root) не может даже считывать эти атрибуты.

security

Пространство имен security применяется для реализации модулей безопасности, таких, как SELinux. Могут ли приложения из пользовательского пространства обращаться к этим атрибутам, — зависит, опять же, от используемого модуля безопасности. По умолчанию все процессы могут считывать эти атрибуты, но только процессы с характеристикой CAP_SYS_ADMIN способны записывать их.

trusted

В пространстве имен trusted в пользовательском пространстве хранится закрытая информация. Только процессы с характеристикой CAP_SYS_ADMIN могут считывать или записывать эти атрибуты.

user

Пространство имен user — это стандартное пространство имен, которое используется обычными процессами. Ядро управляет доступом к этому про-

странству при помощи обычных битов разрешений для файлов. Чтобы иметь возможность читать значения из существующих ключей, у процесса должен быть доступ на чтение к соответствующим файлам. Для того чтобы создать новый ключ или записать значение в существующий ключ, процессу нужен доступ на запись к соответствующему файлу. Расширенные атрибуты в пользовательском пространстве имен можно назначать только обычным файлам, но не символьским ссылкам или файлам устройств. При разработке в пользовательском пространстве приложения, использующего расширенные атрибуты, вы, вероятнее всего, будете применять именно это пространство имен.

Операции над расширенными атрибутами

В стандарте POSIX определяются четыре операции, которые приложения могут выполнять над расширенными атрибутами файла:

- если известны файл и ключ, можно восстановить соответствующее значение атрибута;
- если известны файл, ключ и значение атрибута, то можно назначить это значение известному ключу;
- если известны файл, можно восстановить список всех ключей расширенных атрибутов, связанных с файлом;
- если известен файл и ключ, можно удалить этот расширенный атрибут из файла.

Для каждой операции в POSIX предусмотрено три системных вызова:

- версия, работающая с путем к файлу; если путь относится к символьской ссылке, то операция выполняется на цели данной ссылки (обычное поведение);
- версия, работающая с путем к файлу; если путь относится к символьской ссылке, то операция выполняется над самой ссылкой (стандартный l-вариант системного вызова);
- версия, работающая с дескриптором файла (стандартный f-вариант).

В следующих подразделах мы рассмотрим все варианты.

Извлечение расширенного атрибута

Самая простая операция — это возвращение значения расширенного атрибута для данного файла, когда известен ключ:

```
#include <sys/types.h>
#include <attr/xattr.h>

ssize_t getxattr (const char *path, const char *key,
                  void *value, size_t size);
ssize_t lgetxattr (const char *path, const char *key,
                   void *value, size_t size);
ssize_t fgetxattr (int fd, const char *key,
                   void *value, size_t size);
```

Успешный вызов `getxattr()` сохраняет расширенный атрибут с именем `key` для файла `path` в предоставленном пользователем буфере `value`, длина которого составляет `size` байтов. Вызов возвращает фактический размер значения.

Если значение параметра `size` равно 0, то вызов возвращает размер значения, не сохраняя само значение в буфере `value`. Таким образом, можно передать 0, чтобы определить правильный размер буфера, в котором позже будет сохранено значение ключа. Такая возможность упрощает в приложениях выделение и изменение размера буфера.

`lgetxattr()` работает так же, как `getxattr()`, если только `path` — не символьическая ссылка; в таком случае этот вызов возвращает расширенные атрибуты для самой ссылки, а не для цели, на которую она указывает. Вспомните, в предыдущем разделе говорилось, что атрибуты в пользовательском пространстве невозможно применять к символьическим ссылкам и поэтому данный вызов используется редко.

Вызов `fgetxattr()` работает с дескриптором файла `fd`; во всех остальных аспектах он аналогичен вызову `getxattr()`.

В случае ошибки все три вызова возвращают значение `-1` и присваивают переменной `errno` одно из следующих значений:

EACCES

У вызывающего процесса отсутствуют разрешения на поиск для одного компонента в пути `path` (только для `getxattr()` и `lgetxattr()`).

EBADF

Недопустимое значение параметра `fd` (только для `fgetxattr()`).

EFAULT

Параметр `path`, `key` или `value` содержит недопустимый указатель.

ELOOP

Путь `path` включает слишком много символьических ссылок (только для `getxattr()` и `lgetxattr()`).

ENAMETOOLONG

Слишком длинное значение `path` (только для `getxattr()` и `lgetxattr()`).

ENOATTR

Атрибут `key` не существует или у процесса нет доступа к атрибуту.

ENOENT

Компонент пути `path` не существует (только для `getxattr()` и `lgetxattr()`).

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOTDIR

Компонент пути `path` не является каталогом (только для `getxattr()` и `lgetxattr()`).

ENOTSUP

Файловая система, в которой располагается путь `path` или файл с дескриптором `fd`, не поддерживает расширенные атрибуты.

ERANGE

Размер `size` слишком мал, чтобы вместить значение ключа `key`. Как говорилось выше, вызов можно повторить, присвоив параметру `size` значение 0; будет возвращено значение, указывающее требуемый размер буфера, что позволит установить правильный размер буфера.

Установка расширенного атрибута

Три следующих системных вызова устанавливают указанный расширенный атрибут:

```
#include <sys/types.h>
#include <attr/xattr.h>

int setxattr (const char *path, const char *key,
              const void *value, size_t size, int flags);
int lsetxattr (const char *path, const char *key,
               const void *value, size_t size, int flags);
int fsetxattr (int fd, const char *key,
               const void *value, size_t size, int flags);
```

Успешный вызов `setxattr()` устанавливает расширенный атрибут `key` для файла `path`, присваивая атрибуту значение `value`, длина которого составляет `size` байтов. Поле `flags` модифицирует поведение вызова. Когда значение `flags` равно `XATTR_CREATE` (создание атрибута), вызов завершается ошибкой, если расширенный атрибут уже существует. Когда значение `flags` равно `XATTR_REPLACE` (замена атрибута), вызов завершается ошибкой, если расширенный атрибут еще не существует. Поведение по умолчанию, которое реализуется, когда значение `flags` равно 0, допускает и создание, и замену. Независимо от значения `flags`, никакие другие ключи, за исключением `key`, не меняются.

Вызов `lsetxattr()` работает так же, как и `setxattr()`, если только `path` не содержит символьскую ссылку; в этом случае вызов устанавливает расширенные атрибуты для самой ссылки, а не для файла, являющегося целью ссылки. Вспомните, что атрибуты в пользовательском пространстве имен невозможно применять к символьским ссылкам, поэтому данный вызов применяется редко.

Вызов `fsetxattr()` работает с дескриптором файла `fd`; в остальном его поведение аналогично `setxattr()`.

В случае успеха все три вызова возвращают значение 0; в случае ошибки вызовы возвращают -1 и присваивают переменной `errno` одно из следующих значений:

EACCESS

У вызывающего процесса отсутствуют разрешения на поиск для компонента пути `path` (только для `setxattr()` и `lsetxattr()`).

EBADF

Недопустимое значение fd (только для fsetxattr()).

EDQUOT

Предельный размер квоты не позволяет использовать все пространство, необходимое запрошенной операции.

EEXIST

Значение flags равно XATTR_CREATE, но атрибут с ключом key у данного файла уже существует.

EFAULT

Параметр path, key или value содержит недопустимый указатель.

EINVAL

Недопустимое значение flags.

ELOOP

Путь path содержит слишком много символьических ссылок (только для setxattr() и lsetxattr()).

ENAMETOOLONG

Слишком длинное значение path (только для setxattr() и lsetxattr()).

ENOATTR

Значение flags равно XATTR_REPLACE, но у данного файла отсутствует атрибут с ключом key.

ENOENT

Компонент пути path не существует (только для setxattr() и lsetxattr()).

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOSPC

В файловой системе недостаточно пространства для записи расширенного атрибута.

ENOTDIR

Компонент пути path не является каталогом (только для setxattr() и lsetxattr()).

ENOTSUP

Файловая система, в которой находится файл path или файл с дескриптором fd, не поддерживает расширенные атрибуты.

Перечисление расширенных атрибутов файла

Следующие три системных вызова предназначены для перечисления набора ключей расширенных атрибутов, связанных с данным файлом:

```
#include <sys/types.h>
#include <attr/xattr.h>
```

```
ssize_t listxattr (const char *path,
                   char *list, size_t size);
ssize_t llistxattr (const char *path,
                    char *list, size_t size);
ssize_t flistxattr (int fd,
                    char *list, size_t size);
```

Успешный вызов `listxattr()` возвращает список ключей расширенных атрибутов, связанных с файлом, который указан при помощи параметра `path`. Список сохраняется в буфере `list`, длина которого составляет `size` байтов. Системный вызов возвращает фактический размер списка в байтах.

Каждый ключ расширенного атрибута, который возвращается в буфере `list`, завершается символом нуля, поэтому список выглядит приблизительно так:

```
"user.md5_sum\0user.mime_type\0system.posix_acl_default\0"
```

Таким образом, хотя каждый ключ представляет собой обычную строку C, завершающуюся нулем, для того чтобы пройти по списку ключей, необходимо знать длину всего списка (которую можно получить из возвращаемого значения вызова). Чтобы выяснить, буфер какого размера потребуется выделить для списка, выполните одну из этих функций, передав ей в качестве параметра `size` значение 0; в этом случае функция вернет фактическую длину полного списка ключей. Как и в случае `getattr()`, такая функциональность очень удобна для выделения памяти или изменения размера буфера.

`llistxattr()` работает так же, как и `listxattr()`, за исключением случая, когда `path` представляет собой символьическую ссылку. В такой ситуации данный вызов перечисляет ключи расширенных атрибутов, связанные с самой ссылкой, а не с файлом, на который ссылка нацелена. Вспомните, что в пользовательском пространстве имен нельзя применять атрибуты к символьическим ссылкам, поэтому данный вызов используется редко.

Вызов `flistxattr()` работает с дескриптором файла `fd`; во всем остальном его поведение абсолютно аналогично `listxattr()`.

В случае ошибки все три вызова возвращают значение `-1` и присваивают переменной `errno` один из следующих кодов ошибки:

EACCESS

У вызывающего процесса отсутствуют разрешения на поиск для компонента пути `path` (только для `listxattr()` и `llistxattr()`).

EBADF

Недопустимое значение `fd` (только для `flistxattr()`).

EFAULT

Параметр `path` или `list` содержит недопустимый указатель.

ELOOP

Путь `path` содержит слишком много символьических ссылок (только для `listxattr()` и `llistxattr()`).

ENAMETOOLONG

Слишком большое значение `path` (только для `listxattr()` и `llistxattr()`).

ENOENT

Компонент пути path не существует (только для `listxattr()` и `llistxattr()`).

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOTDIR

Компонент пути path не является каталогом (только для `listxattr()` и `llistxattr()`).

ENOTSUP

Файловая система, которой принадлежит файл path или файл с дескриптором fd, не поддерживает расширенные атрибуты.

ERANGE

Значение size не равно нулю, но недостаточно велико, чтобы вместить полный список ключей. Можно повторить вызов, передав в качестве параметра size значение 0, чтобы выяснить фактический размер списка. После этого следует переопределить значение size и повторить системный вызов.

Удаление расширенного атрибута

Наконец, следующие три системных вызова удаляют указанный ключ из указанного файла:

```
#include <sys/types.h>
#include <attr/xattr.h>

int removexattr (const char *path, const char *key);
int lremovexattr (const char *path, const char *key);
int fremovexattr (int fd, const char *key);
```

Успешный вызов `removexattr()` удаляет ключ key расширенного атрибута из файла path. Вспомните, что существует различие между неопределенным ключом и определенным ключом с пустым значением (значением нулевой длины).

Вызов `lremovexattr()` работает так же, как и `removexattr()`, за исключением случаев, когда path содержит символьическую ссылку. Тогда вызов удаляет ключ расширенного атрибута, связанный с самой ссылкой, а не с файлом, на который ссылка нацелена. Вспомните, что в пользовательском пространстве имен атрибуты нельзя применять к символьским ссылкам, поэтому данный вызов используется редко.

Вызов `fremovexattr()` работает с дескриптором файла fd; в остальном он аналогичен вызову `removexattr()`.

В случае успеха все три системных вызова возвращают значение 0. В случае ошибки все три возвращают -1 и присваивают переменной errno одно из следующих значений:

EACCESS

У вызывающего процесса отсутствуют разрешения на поиск для компонента пути path (только для `removexattr()` и `lremovexattr()`).

EBADF

Недопустимое значение fd (только для `fremovexattr()`).

EFAULT

Параметр `path` или `key` содержит недопустимый указатель.

ELOOP

В пути `path` содержится слишком много символьических ссылок (только для `removexattr()` и `lremovexattr()`).

ENAMETOOLONG

Слишком большое значение `path` (только для `removexattr()` и `lremovexattr()`).

ENOATTR

Для данного файла не существует атрибута с ключом `key`.

ENOENT

Компонент пути `path` не существует (только для `removexattr()` и `lremovexattr()`).

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOTDIR

Компонент пути `path` не является каталогом (только для `removexattr()` и `lremovexattr()`).

ENOTSUP

Файловая система, в которой расположен файл `path` или файл с дескриптором `fd`, не поддерживает расширенные атрибуты.

Каталоги

В Unix *каталог* (*derictory*) — это очень простая концепция: каталог всего лишь содержит список имен файлов, с каждым из которых сопоставлен номер `inode`. Имена называются *записями каталога* (*directory entry*), а каждая пара из имени и номера `inode` — это *ссылка* (*link*). Содержимое каталога — то, что пользователь видит в результате выполнения команды `ls`, — это перечисление всех имен файла в этом каталоге. Когда пользователь пытается открыть файл в каталоге, ядро находит имя файла в списке данного каталога и извлекает соответствующий номер `inode`. После этого ядро передает номер `inode` файловой системе, которая использует его для поиска физического местоположения файла на устройстве.

Каталоги могут включать не только файлы, но и другие каталоги. *Подкаталог* (*subdirectory*) — это каталог внутри другого каталога. С учетом этого определения все каталоги являются подкаталогами каких-то *родительских каталогов* (*parent directory*), за исключением каталога в корне дерева файловой системы, `/`. Неудивительно, что этот каталог называется *корневым каталогом*.

(root directory) (но не путайте его с домашним каталогом пользователя root, /root).

Путь к файлу (pathname) состоит из имени файла и одного или нескольких его родительских каталогов. Абсолютный путь к файлу (absolute pathname) – это путь, начинающийся с корневого каталога, например /usr/bin/sextant. Относительный путь к файлу (relative pathname) – это путь, начинающийся с каталога, отличного от корневого, например bin/sextant. Для того чтобы использовать подобный путь, операционной системе необходимо знать, по отношению к какому каталогу он указан. В качестве начальной точки используется текущий рабочий каталог (о котором речь пойдет в следующем разделе).

Имена файлов и каталогов могут содержать любые символы, за исключением /, который разделяет каталоги в пути, а также пустого значения (null), который завершает путь к файлу. И все-таки стандартной практикой является ограничение набора допустимых символов: обычно в названиях файлов и каталогов используются только допустимые печатаемые символы из текущих национальных настроек, или даже только символы ASCII. Так как ни ядро, ни библиотека С не принуждают к этому, разработчики приложений свободны самостоятельно принимать решения об использовании только допустимых печатаемых символов.

В старых системах Unix имена файлов ограничивались 14 символами. Сегодня все современные файловые системы Unix допускают имена файлов длиной до 255 байт¹. Многие файловые системы в Linux позволяют использовать имена файлов еще большей длины².

В каждом каталоге есть два специальных каталога, . и .. (они называются точка и точка-точка). Каталог точка представляет собой ссылку на сам данный каталог. Каталог точка-точка – это ссылка на родительский каталог данного каталога. Например, /home/kidd/gold/.. – это то же самое, что и /home/kidd. Каталоги точка и точка-точка в корневом каталоге указывают на самое себя: то есть /./ и ../ – это один и тот же каталог. Технически, можно сказать, что даже корневой каталог является подкаталогом – в данном случае самого себя.

Текущий рабочий каталог

У каждого процесса есть текущий каталог, который процесс первоначально наследует от своего родительского процесса. Этот каталог называется *текущим рабочим каталогом* процесса (current working directory, cwd). Текущий рабочий

¹ Обратите внимание, что ограничение составляет 255 байт, а не 255 символов. Многобайтовые символы, очевидно, занимают более одного из этих 255 байт.

² Конечно же, в старых файловых системах, которые Linux поддерживает для обеспечения обратной совместимости, таких, как FAT, до сих пор действуют собственные ограничения. В случае FAT имя файла может включать только восемь символов, за которыми должны следовать точка и еще три символа. Следует заметить, что принудительное использование точки в качестве специального символа в файловой системе – это не самый умный шаг.

каталог представляет собой начальную точку, основываясь на которой ядро разрешает относительные пути к файлам. Например, если текущий рабочий каталог процесса — это `/home/blackbeard` и этот процесс пытается открыть файл `parrot.jpg`, то ядро для этого ищет файл `/home/blackbeard/parrot.jpg`. Однако если процесс попытается открыть файл `/usr/bin/mast`, то ядро действительно откроет файл `/usr/bin/mast` — текущий рабочий каталог никак не влияет на абсолютные пути к файлам (то есть пути, начинающиеся с косой черты).

Процесс может узнать и изменить свой текущий рабочий каталог.

Проверка текущего рабочего каталога

Предпочтительный способ получения текущего рабочего каталога — системный вызов `getcwd()`, стандартизованный в POSIX:

```
#include <unistd.h>
```

```
char * getcwd (char *buf, size_t size);
```

Успешный вызов `getcwd()` копирует текущий рабочий каталог в форме абсолютного пути в буфер длиной `size`, на который указывает параметр `buf`, и возвращает указатель на `buf`. В случае ошибки вызов возвращает значение `NULL` и присваивает переменной `errno` одно из следующих значений:

`DEFAULT`

Параметр `buf` содержит недопустимый указатель.

`EINVAL`

Значение `size` равно 0, но значение `buf` не равно `NULL`.

`ENOENT`

Текущий рабочий каталог более не действителен. Это может происходить при удалении текущего рабочего каталога.

`ERANGE`

Значение `size` слишком мало, чтобы можно было сохранить текущий рабочий каталог в буфере `buf`. Необходимо выделить буфер большего размера и повторить попытку.

Пример использования вызова `getcwd()`:

```
char cwd[BUF_LEN];  
  
if (!getcwd (cwd, BUF_LEN)) {  
    perror ("getcwd");  
    exit (EXIT_FAILURE);  
}  
  
printf ("cwd = %s\n", cwd);
```

Стандарт POSIX диктует, что в случае, когда значение `buf` равно `NULL`, поведение вызова `getcwd()` не определено. Однако библиотека C в Linux в этом случае выделяет буфер длиной `size` байтов и записывает туда текущий рабочий каталог. Если значение `size` равно 0, то библиотека C выделяет буфер достаточно

большого размера, чтобы сохранить текущий рабочий каталог. После этого, когда содержимое буфера приложению больше не нужно, приложение ответственно за его освобождение при помощи вызова `free()`. Так как подобное поведение уникально для Linux, в приложениях, имеющих своей целью хорошую переносимость на разные платформы или строгое следование стандарту POSIX, полагаться на данную функциональность нельзя. Однако она делает использование вызова очень простым. Вот пример:

```
char *cwd:  
  
 cwd = getcwd (NULL, 0);  
 if (!cwd) {  
     perror ("getcwd");  
     exit (EXIT_FAILURE);  
 }  
  
 printf ("cwd = %s\n", cwd);  
  
 free (cwd);
```

В библиотеке C в Linux также есть функция `get_current_dir_name()`, работающая так же, как и `getcwd()`, если передавать этому вызову параметры `buf` и `size` со значениями `NULL` и `0` соответственно:

```
#define _GNU_SOURCE  
#include <unistd.h>  
  
char * get_current_dir_name (void);
```

Таким образом, следующий фрагмент кода делает то же самое, что и предыдущий:

```
char *cwd:  
  
 cwd = get_current_dir_name ( );  
 if (!cwd) {  
     perror ("get_current_dir_name");  
     exit (EXIT_FAILURE);  
 }  
  
 printf ("cwd = %s\n", cwd);  
  
 free (cwd);
```

В старых системах BSD использовался вызов `getwd()`, который в Linux присутствует для обеспечения обратной совместимости:

```
#define _XOPEN_SOURCE_EXTENDED /* or _BSD_SOURCE */  
#include <unistd.h>
```

```
char * getwd (char *buf);
```

Вызов `getwd()` копирует текущий рабочий каталог в буфер `buf`, длина которого должна составлять, как минимум, `PATH_MAX`. Вызов возвращает `buf` в случае успеха и `NULL` в случае ошибки, например:

```
char cwd[PATH_MAX];  
  
if (!getwd (cwd)) {  
    perror ("getwd");  
    exit (EXIT_FAILURE);  
}  
  
printf ("cwd = %s\n", cwd);
```

По причинам, связанным с переносимостью и безопасностью, не следует применять в приложениях вызов getwd() — лучше использовать getcwd().

Изменение текущего рабочего каталога

Когда пользователь впервые входит в систему, процесс входа в систему устанавливает для него в качестве текущего рабочего каталога его домашний каталог, указанный в файле /etc/passwd. В некоторых случаях у процесса может возникать необходимость изменить свой текущий рабочий каталог. Например, оболочка делает это, когда пользователь выполняет команду cd.

В Linux есть два системных вызова, предназначенные для изменения текущего рабочего каталога: один принимает полный путь к каталогу, а второй — дескриптор файла, представляющий открытый каталог:

```
#include <unistd.h>  
  
int chdir (const char *path);  
int fchdir (int fd);
```

Вызов chdir() меняет текущий рабочий каталог на каталог, указанный при помощи параметра path, который может содержать либо абсолютный, либо относительный путь. Подобным образом fchdir() меняет текущий рабочий каталог на каталог, представленный дескриптором файла fd, который должен быть открыт для этого каталога. В случае успеха оба вызова возвращают значение 0. В случае ошибки оба возвращают -1.

Если происходит ошибка, chdir() также присваивает переменной errno одно из следующих значений:

EACCES

У вызывающего процесса отсутствуют разрешения на поиск для компонента пути path.

EFAULT

Параметр path содержит недопустимый указатель.

EIO

Произошла внутренняя ошибка ввода-вывода.

ELOOP

Во время разрешения пути path ядру встретилось слишком много символьических ссылок.

ENAMETOOLONG

Слишком длинное значение path.

ENOENT

Каталог, на который указывает параметр *path*, не существует.

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOTDIR

Один или несколько компонентов пути *path* не являются каталогами.

Fchdir() может присваивать переменной *errno* следующие значения:

EACCESS

У вызывающего процесса отсутствуют разрешения на поиск в каталоге с дескриптором файла, указанным при помощи параметра *fd* (то есть бит исполнения не установлен). Такая ситуация случается, когда каталог более высокого уровня доступен для чтения, но не для исполнения, вызов *open()* выполняется успешно, но *fchdir()* завершается ошибкой.

EBADF

Значение параметра *fd* не является открытым дескриптором файла.

В зависимости от файловой системы оба вызова могут использовать другие значения ошибки.

Эти системные вызовы влияют только на текущий выполняющийся процесс. В Unix не существует механизма изменения текущего рабочего каталога другого процесса. Таким образом, команда *cd* в оболочках не может быть отдельным процессом (как большинство команд), который просто выполняет *chdir()* на первом аргументе в командной строке, а затем завершается. Нет, *cd* – это специальная встроенная команда, заставляющая саму оболочку вызывать *chdir()*, меняя собственный текущий рабочий каталог.

Команда *getcwd()* чаще всего применяется для сохранения текущего рабочего каталога, чтобы процесс мог позже вернуться в него. Например:

```
char *swd;
int ret;

/* сохранение текущего рабочего каталога */
swd = getcwd (NULL, 0);
if (!swd) {
    perror ("getcwd");
    exit (EXIT_FAILURE);
}

/* переход в другой каталог */
ret = chdir (some_other_dir);
if (ret) {
    perror ("chdir");
    exit (EXIT_FAILURE);
}

/* какая-то работа в новом каталоге... */

/* возврат в сохраненный каталог */
```

```

ret = chdir (swd);
if (ret) {
    perror ("chdir");
    exit (EXIT_FAILURE);
}

free (swd);

```

Лучше, однако, открывать текущий каталог вызовом `open()`, чтобы потом возвращаться в него при помощи `fchdir()`. Этот подход быстрее, так как ядро не хранит полный путь текущего рабочего каталога в памяти; хранится только структура `inode`. Следовательно, когда пользователь вызывает `getcwd()`, ядру необходимо формировать путь, проходя по структуре каталогов. В противоположность этому, открытие текущего рабочего каталога обходится дешевле, потому что в этом случае у ядра уже есть его `inode`, а удобный для восприятия человеком путь к каталогу для открытия файла не требуется. Этот подход реализуется в следующем фрагменте кода:

```

int swd_fd;

swd_fd = open ("..", O_RDONLY);
if (swd_fd == -1) {
    perror ("open");
    exit (EXIT_FAILURE);
}

/* переход в другой каталог */
ret = chdir (some_other_dir);
if (ret) {
    perror ("chdir");
    exit (EXIT_FAILURE);
}

/* какая-то работа в новом каталоге... */

/* возврат в сохраненный каталог*/
ret = fchdir (swd_fd);
if (ret) {
    perror ("fchdir");
    exit (EXIT_FAILURE);
}

/* закрытие дескриптора fd каталога */
ret = close (swd_fd);
if (ret) {
    perror ("close");
    exit (EXIT_FAILURE);
}

```

Именно так в оболочках реализуется кэширование предыдущего каталога (например, для `cd` — в `bash`).

Процесс, которому неважен его текущий рабочий каталог, например демон, обычно устанавливает в качестве этого значения корневой каталог, применяя

вызов `chdir("/")`. Приложение, взаимодействующее с пользователем и его данными, такое, как текстовый редактор, чаще всего устанавливает в качестве своего текущего рабочего каталога домашний каталог пользователя или специальный каталог для хранения документов. Так как текущие рабочие каталоги имеют смысл только в контексте относительных путей к файлам, рабочий каталог имеет наибольшую практическую пользу для утилит командной строки, которые пользователь вызывает из оболочки.

Создание каталогов

В Linux есть только один системный вызов, стандартизированный POSIX, для создания новых каталогов:

```
#include <sys/stat.h>
#include <sys/types.h>

int mkdir (const char *path, mode_t mode);
```

Успешный вызов `mkdir()` создает каталог `path` (параметр может содержать абсолютный или относительный путь) с битами разрешения `mode` (как они модифицированы текущим значением `umask`) и возвращает значение 0.

Текущее значение `umask` модифицирует обычным способом аргумент `mode`, а также любые биты режима, уникальные для операционной системы: в Linux биты разрешения для создаваемых каталогов составляют значение (`mode & ~umask & 01777`). Другими словами, `umask` для процесса определяет ограничения, которые вызов `mkdir()` переопределить не может. Если для родительского каталога текущего каталога бит `set group ID (sgid)` установлен или если файловая система подмонтирована с семантикой групп BSD, то новый каталог наследует принадлежность к группам своего предка. В противном случае на каталог распространяется действительный идентификатор группы процесса.

В случае ошибки `mkdir()` возвращает значение -1 и присваивает переменной `errno` одно из следующих значений:

EACCESS

У текущего процесса нет прав на запись в родительский каталог, или один или несколько компонентов в пути `path` недоступны для поиска.

EEXIST

Путь `path` уже существует (и не обязательно является каталогом).

EFAULT

Параметр `path` содержит недопустимый указатель.

ELOOP

Во время разрешения пути `path` ядру встретилось слишком много символьических ссылок.

ENAMETOOLONG

Слишком длинное значение `path`.

ENOENT

Компонент пути *path* не существует или является символьской ссылкой, указывающей на несуществующий объект.

ENOMEM

Недостаточно памяти ядра для выполнения данного запроса.

ENOSPC

На устройстве, где находится *path*, недостаточно пространства либо превышена дисковая квота для пользователя.

ENOTDIR

Один или несколько компонентов пути *path* не являются каталогами.

EPERM

Файловая система, которой принадлежит *path*, не поддерживает создание каталогов.

EROFS

Файловая система, которой принадлежит *path*, подмонтирована с доступом только на чтение.

Удаление каталогов

В противоположность `mkdir()`, стандартизованный в POSIX системный вызов `rmdir()` удаляет каталог из иерархии файловой системы:

```
#include <unistd.h>
```

```
int rmdir (const char *path);
```

В случае успеха `rmdir()` удаляет каталог *path* из файловой системы и возвращает значение 0. Каталог, указанный при помощи параметра *path*, должен быть пустым, за исключением каталогов точка и точка-точка. Системного вызова, реализующего эквивалент рекурсивного удаления, как `rm -r`, не существует. Если вам необходим подобный инструмент, то нужно вручную выполнить проход по файловой системе в глубину, удаляя все файлы и каталоги, начиная с листьев, и продвигаясь выше к корню системы; вызов `rmdir()` можно использовать на всех этапах для удаления каталогов после того, как все файлы из каталога уже удалены.

В случае ошибки `rmdir()` возвращает -1 и присваивает переменной `errno` одно из следующих значений:

EACCESS

Запрещен доступ на запись в родительский каталог каталога *path* или одна или несколько составляющих *path* недоступны для поиска.

EBUSY

Каталог *path* в данный момент используется системой и удалить его невозможно. В Linux это может произойти только в том случае, если *path* являет-

ся точкой монтирования или корневым каталогом (корневые каталоги не обязаны быть точками монтирования — спасибо chroot()!).

EFAULT

Параметр *path* содержит недопустимый указатель.

EINVAL

Последним компонентом пути *path* является каталог точки.

ELOOP

При разрешении пути *path* ядру встретилось слишком много символьических ссылок.

ENAMETOOLONG

Слишком длинное значение *path*.

ENOENT

Компонент пути *path* не существует или представляет собой символьскую ссылку, указывающую на несуществующий объект.

ENOMEM

Недостаточно памяти ядра для выполнения данного запроса.

ENOTDIR

Один или несколько компонентов пути *path* не являются каталогами.

ENOTEMPTY

Каталог *path* содержит другие записи, кроме специальных каталогов точки и точка-точка.

EPERM

Для каталога, являющегося предком каталога *path*, установлен бит закрепления в памяти (sticky bit, S_ISVTX), но действительный идентификатор пользователя процесса не совпадает ни с идентификатором пользователя указанного родителя, ни с идентификатором самого каталога *path*, а также процесс не обладает характеристикой CAP_FOWNER. Также возможно, что файловая система, которой принадлежит *path*, не поддерживает удаление каталогов.

EROFS

Файловая система, которой принадлежит *path*, подмонтирована в режиме только для чтения.

Использовать этот системный вызов просто:

int ret;

```
/* удаление каталога /home/barbary/maps */
ret = rmdir ("/home/barbary/maps");
if (ret)
    perror ("rmdir");
```

Чтение содержимого каталога

В стандарте POSIX определяется семейство функций, предназначенных для считывания содержимого каталогов, то есть получения списка файлов, находящихся в определенном каталоге. Эти функции полезны для реализации команды `ls` или графического диалогового окна сохранения файла, для выполнения операций на всех файлах в каталоге или для поиска в каталоге файлов по определенному шаблону.

Для того чтобы начать читать содержимое каталога, нужно создать *поток каталога* (directory stream), представляемый объектом DIR:

```
#include <sys/types.h>
#include <dirent.h>

DIR * opendir (const char *name);
```

Успешный вызов `opendir()` создает поток каталога, который представляет каталог, указанный при помощи параметра `name`.

Поток каталога — это всего лишь дескриптор файла, представляющий открытый каталог, некоторое количество метаданных и буфер для записи содержимого каталога. Следовательно, имея поток каталога, можно получить дескриптор файла для соответствующего каталога:

```
#define _BSD_SOURCE /* или _SVID_SOURCE */
#include <sys/types.h>
#include <dirent.h>

int dirfd (DIR *dir);
```

Успешный вызов `dirfd()` возвращает дескриптор файла, соответствующий потоку каталога `dir`. В случае ошибки вызов возвращает значение `-1`. Так как функции потока каталога используют данный дескриптор для своей работы, в программах допускается только применение вызовов, не меняющих позицию в файле. `Dirfd()` — это расширение BSD, которое не стандартизовано в POSIX; если вы стремитесь к тому, чтобы ваши программы в точности соответствовали POSIX, то вам следует избегать использования этого вызова.

Чтение из потока каталога

После того как поток каталога создан при помощи вызова `opendir()`, программа может начинать читать записи из каталога. Для этого используется вызов `readdir()`, возвращающий одну за другой записи из указанного объекта DIR:

```
#include <sys/types.h>
#include <dirent.h>

struct dirent * readdir (DIR *dir);
```

Успешный вызов `readdir()` возвращает очередную запись в каталоге, указанном при помощи параметра `dir`. Запись каталога представляется структурой `dirent`. Она определена в Linux в файле заголовка `<dirent.h>` и выглядит следующим образом:

```
struct dirent {  
    ino_t d_ino;           /* номер inode */  
    off_t d_off;          /* смещение к следующей записи dirent */  
    unsigned short d_reclen; /* длина данной записи */  
    unsigned char d_type; /* тип файла */  
    char d_name[256];     /* имя файла */  
};
```

В стандарте POSIX обязательным является только поле `d_name`, то есть имя файла в каталоге. Прочие поля не обязательны или уникальны для Linux. В приложениях, стремящихся к хорошей переносимости на другие системы или к полному соблюдению POSIX, следует обращаться только к полю `d_name`.

Приложения последовательно вызывают `readdir()`, извлекая один за другим файлы из каталога, до тех пор пока еще обнаруживаются новые файлы или пока записи каталога не заканчиваются. В последнем случае `readdir()` возвращает значение `NULL`.

В случае ошибки вызов `readdir()` также возвращает `NULL`. Для того чтобы отличать эту ошибку от ситуации, когда в каталоге просто заканчиваются файлы, перед каждым вызовом `readdir()` следует присваивать переменной `errno` значение 0, а после вызова проверять и возвращаемое значение, и значение `errno`. Единственное значение `errno`, которое может установить `readdir()`, — это `EBADF`, указывающее на недопустимое значение `dir`. Однако во многих приложениях проверка ошибок не выполняется и предполагается, что `NULL` обозначает только случай, когда в каталоге заканчиваются файлы.

Закрытие потока каталога

Для того чтобы закрыть поток каталога, открытый при помощи вызова `opendir()`, используйте вызов `closedir()`:

```
#include <sys/types.h>  
#include <dirent.h>
```

```
int closedir (DIR *dir);
```

Успешный вызов `closedir()` закрывает поток каталога, переданный при помощи параметра `dir`, а также соответствующий дескриптор файла и возвращает значение 0. В случае ошибки функция возвращает -1 и присваивает переменной `errno` значение `EBADF` — единственный возможный код ошибки, указывающий, что `dir` не является открытым потоком каталога.

В следующем фрагменте кода реализуется функция `find_file_in_dir()`, в которой вызов `readdir()` применяется для поиска в указанном каталоге указанного имени файла. Если файл в каталоге существует, функция возвращает 0. В противном случае она возвращает ненулевое значение:

```
/*  
 * find_file_in_dir — выполняет поиск в каталоге path  
 * и ищет файл с именем file.  
 *  
 * Возвращает 0, если file существует в path, и ненулевое  
 * значение в противном случае.  
 */
```

```

*/
int find_file_in_dir (const char *path, const char *file)
{
    struct dirent *entry,
    int ret = 1;
    DIR *dir;

    dir = opendir (path);

    errno = 0;
    while ((entry = readdir (dir)) != NULL) {
        if (!strcmp(entry->d_name, file)) {
            ret = 0;
            break;
        }
    }

    if (errno && !entry)
        perror ("readdir");

    closedir (dir);
    return ret;
}

```

Системные вызовы для чтения содержимого каталогов

Функции для чтения содержимого каталогов, о которых речь шла выше, стандартизованы в POSIX и предоставляются библиотекой С. Внутри системы эти функции используют один из двух системных вызовов, `readdir()` или `getdents()`, которые я добавляю для полноты обсуждения:

```

#include <unistd.h>
#include <linux/types.h>
#include <linux/dirent.h>
#include <linux/unistd.h>
#include <errno.h>

/*
 * Не определены в пользовательском пространстве: необходимо для доступа
 * использовать макрос _syscall3( ).
*/
int readdir (unsigned int fd,
             struct dirent *dirp,
             unsigned int count);

int getdents (unsigned int fd,
              struct dirent *dirp,
              unsigned int count);

```

Вы не должны использовать эти системные вызовы! Они бессмысленны и непереносимы. Вместо этого в приложениях из пользовательского пространства следует применять системные вызовы `opendir()`, `readdir()` и `closedir()` из библиотеки С.

Ссылки

Обсуждая каталоги, я упоминал, что каждое соответствие между именем и номером inode в каталоге называется *ссылкой* (*link*). Учитывая это простое определение — ссылка представляет собой всего лишь имя в списке (каталоге), указывающее на inode, — вас может удивлять, почему возможно существование нескольких ссылок на одну и ту же структуру inode. Это означает, что на одну структуру inode (и, следовательно, на один файл) могут существовать ссылки и от, скажем, файла /etc/customs, и от /var/run/ledger.

Но это действительно так и есть один подвох: так как ссылки привязаны к номерам inode, а номера inode уникальны в каждой файловой системе, оба файла — /etc/customs и /var/run/ledger — должны принадлежать одной и той же файловой системе. В пределах одной файловой системы на каждый файл может указывать множество ссылок. Единственное ограничение относится к размеру целочисленного типа данных, который используется для хранения числа ссылок. Среди нескольких ссылок ни одна не является «оригинальной» или «основной». У всех ссылок одинаковый статус, и все они указывают на один и тот же файл.

Мы называем этот тип ссылок *жесткими ссылками* (*hard link*). У файла может быть ни одной, одна или несколько ссылок. У большинства файлов счетчик ссылок равен единице, то есть на них указывает только одна запись каталога, но у некоторых файлов может быть две и более ссылки. Файлы, для которых счетчик ссылок равен 0, не имеют соответствующих записей каталога в файловой системе. Когда счетчик ссылок файла достигает нуля, файл помечается как освобожденный, а его дисковые блоки становятся доступными для повторного использования¹. Однако подобный файл остается в файловой системе до тех пор, пока он открыт у какого-то процесса. Как только все процессы закрывают этот файл, файл удаляется.

В ядро Linux такое поведение реализуется с использованием счетчика ссылок и счетчика пользователей. *Счетчик пользователей* (*usage count*) — это общее число экземпляров, где данный файл остается открытым. Файл не удаляется из файловой системы до тех пор, пока оба счетчика — счетчик ссылок и счетчик пользователей — не становятся равными нулю.

Еще один тип ссылок, *символические ссылки* (*symbolic link*), представляет собой не отображение на уровне файловой системы, а высокоуровневый указатель, который интерпретируется во время выполнения. Подобные ссылки могут охватывать разные файловые системы — мы поговорим о них ниже.

¹ Поиск файлов со счетчиком ссылок, равным нулю, но чьи блоки помечены как занятые, — это основная задача fsck, утилиты контроля над файловой системой. Подобное условие может возникать, когда файл удаляется, но остается открытым и в системе происходит аварийный сбой до того, как файл закрывается. Ядро не получает возможности пометить блоки файловой системы как свободные, из-за чего возникают несоответствия. Протоколирование файловых систем устраняет такой тип ошибок.

Жесткие ссылки

Системный вызов `link()`, один из первоначальных системных вызовов Unix, теперь стандартизованный в POSIX, создает новую ссылку для существующего файла:

```
#include <unistd.h>
```

```
int link (const char *oldpath, const char *newpath);
```

Успешный вызов `link()` создает для пути `newpath` новую ссылку, указывающую на существующий файл `oldpath`, и возвращает значение 0. После завершения оба пути, `oldpath` и `newpath`, указывают на один и тот же файл, то есть, фактически, невозможно сказать, какая из ссылок была «исходной».

В случае ошибки вызов возвращает значение -1 и присваивает переменной `errno` одно из следующих значений:

EACCESS

У вызывающего процесса отсутствуют разрешения на поиск для компонента пути `oldpath` или у вызывающего процесса нет разрешений на запись в каталог, содержащий `newpath`.

EXEXIST

Путь `newpath` уже существует — системный вызов `link()` не переопределяет существующие записи каталога.

EFAULT

Параметр `oldpath` или `newpath` содержит недопустимый указатель.

EIO

Произошла внутренняя ошибка ввода-вывода (это плохо!).

ELOOP

При разрешении пути `oldpath` или `newpath` встретилось слишком много символических ссылок.

EMLINK

На структуру `inode`, на которую указывает путь `oldpath`, уже указывает максимально допустимое количество ссылок.

ENAMETOOLONG

Слишком длинное значение `oldpath` или `newpath`.

ENOENT

Компонент `oldpath` или `newpath` не существует.

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOSPC

На устройстве, которому принадлежит `newpath`, недостаточно свободного пространства для новой записи каталога.

ENOTDIR

Компонент oldpath или newpath не является каталогом.

EPERM

Файловая система, которой принадлежит newpath, не допускает создание новых жестких ссылок или oldpath является каталогом.

EROFS

Newpath принадлежит файловой системе, доступной только для чтения.

EXDEV

Пути newpath и oldpath не принадлежат одной и той же подмонтированной системе. (Linux допускает монтирование одной файловой системы в нескольких точках, но даже в этом случае нельзя создавать жесткие ссылки между точками монтирования.)

В следующем примере создается новая запись каталога, pirate, которая со-поставляется с тем же номером inode (и, следовательно, с тем же файлом), что и существующий файл privateer. Оба файла принадлежат каталогу /home/kidd:

```
int ret;
```

```
/*
 * создание новой записи каталога.
 * /home/kidd/privateer, указывающей на ту же
 * структуру inode, что и /home/kidd/pirate
 */
ret = link ("/home/kidd/privateer", "/home/kidd/pirate");
if (ret)
    perror ("link");
```

Символические ссылки

Символические ссылки (symbolic link), также известные как *симсылки* (symlink) и *мягкие ссылки* (soft links), схожи с жесткими ссылками в том, что оба типа ссылок указывают на файлы в файловой системе. Однако символические ссылки отличаются тем, что они существуют не просто в форме дополнительных записей каталога, а представляют собой специальный тип файлов. В таком специальном файле содержится путь к *другому* файлу, называемому *целью* (target) символической ссылки. Во время выполнения ядро на лету заменяет путь к файлу символической ссылки путем к цели символической ссылки (если только в программе не используются разнообразные l-версии системных вызовов, такие, как lstat(), работающие на самом файле ссылки, а не ее цели). Таким образом, тогда как одну жесткую ссылку невозможно отличить от другой жесткой ссылки на тот же файл, различия между символической ссылкой и ее целевым файлом очевидны.

Символическая ссылка может быть относительной или абсолютной. Она также может содержать особый каталог точка, о котором говорилось выше, ссылаясь на каталог, которому принадлежит символическая ссылка, или каталог точка-точка, ссылаясь на родительский каталог своего каталога.

В отличие от жестких ссылок, мягкие ссылки могут охватывать разные файловые системы. В действительности они могут указывать куда угодно! Символические ссылки могут указывать на существующие файлы (распространенная практика) и на несуществующие тоже. Второй тип ссылок называется *повисшими символическими ссылками* (dangling symlink). Иногда повисшие ссылки являются нежелательным результатом, например удаления целевого файла без учета символьских ссылок на него, но иногда они создаются намеренно.

Системный вызов, предназначенный для создания символьской ссылки, очень похож на своего двоюродного брата для определения жесткой ссылки:

```
#include <unistd.h>
```

```
int symlink (const char *oldpath, const char *newpath);
```

Успешный вызов `symlink()` создает символьскую ссылку `newpath`, указывающую на цель `oldpath`, и возвращает значение 0.

В случае ошибки `symlink()` возвращает -1 и присваивает переменной `errno` одно из следующих значений:

EACCES

У вызывающего процесса отсутствует разрешение на поиск для компонента пути `oldpath` или у вызывающего процесса нет разрешения на запись в каталог, содержащий `newpath`.

EEXIST

Файл `newpath` уже существует — системный вызов `symlink()` не переписывает существующую запись каталога.

EFAULT

Параметр `oldpath` или `newpath` содержит недопустимый указатель.

EIO

Произошла внутренняя ошибка ввода-вывода (это плохо!).

ELOOP

При разрешении пути `oldpath` или `newpath` встретилось слишком много символьских ссылок.

EMLINK

В структуре `inode`, на которую указывает `oldpath`, уже достигнуто максимальное число указывающих на нее ссылок.

ENAMETOOLONG

Слишком длинное значение `oldpath` или `newpath`.

ENOENT

Компонент пути `oldpath` или `newpath` не существует.

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOSPC

На устройстве, которому принадлежит newpath, недостаточно места для создания новой записи каталога.

ENOTDIR

Компонент пути oldpath или newpath не является каталогом.

EPRM

Файловая система, которой принадлежит newpath, не допускает создания новых символьических ссылок.

EROFS

Путь newpath принадлежит файловой системе, доступной только для чтения.

Следующий фрагмент кода аналогичен предыдущему, но здесь файл /home/kidd/pirate создается как символьическая ссылка (в противоположность жесткой ссылке) на /home/kidd/privateer:

```
int ret;  
  
/*  
 * создание символьической ссылки.  
 * /home/kidd/privateer, указывающей  
 * на /home/kidd/pirate  
 */  
ret = symlink ("/home/kidd/privateer", "/home/kidd/pirate");  
if (ret)  
    perror ("symlink");
```

Удаление ссылки

В противоположность созданию ссылки, отсоединение означает удаление из файловой системы пути к файлу. С этой задачей справляется единственный системный вызов, `unlink()`:

```
#include <unistd.h>  
  
int unlink (const char *pathname);
```

Успешный вызов `unlink()` удаляет из файловой системы путь pathname и возвращает значение 0. Если это имя было последней ссылкой на файл, то файл также удаляется из файловой системы. Если же этот файл открыт у какого-то процесса, то ядро не удаляет файл до тех пор, пока процесс не закроет файл. Когда все процессы закрывают данный файл, он удаляется.

Если pathname относится к символьической ссылке, то разрушается сама ссылка, а целевой файл не удаляется.

Если pathname относится к другому типу специального файла, например это файл устройства, конвейера FIFO или сокета, то данный специальный файл удаляется из файловой системы, но процессы, у которых он еще открыт, могут продолжать использовать его.

В случае ошибки `unlink()` возвращает -1 и присваивает переменной errno один из следующих кодов ошибки:

EACCESS

У вызывающего процесса нет разрешения на запись в родительский каталог файла pathname или у вызывающего процесса нет разрешения на поиск для компонента пути pathname.

EFAULT

Параметр pathname содержит недопустимый указатель.

EIO

Произошла ошибка ввода-вывода (это плохо!).

EISDIR

Путь pathname ссылается на каталог.

ELOOP

При прохождении pathname встретилось слишком много символьических ссылок.

ENAMETOOLONG

Слишком длинное значение pathname.

ENOENT

Компонент пути pathname не существует.

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOTDIR

Компонент пути pathname не является каталогом.

EPERM

Система не позволяет отсоединять файлы.

EROFS

Путь pathname принадлежит файловой системе, доступной только для чтения.

Системный вызов `unlink()` не поддерживает удаление каталогов. Для этого в приложениях следует использовать вызов `rmdir()`, о котором говорилось выше в разделе «Удаление каталогов».

Для упрощения бессмысленного разрушения файлов любого типа в языке С предусмотрена функция `remove()`:

```
#include <stdio.h>
```

```
int remove (const char *path);
```

Успешный вызов `remove()` удаляет файл path из файловой системы и возвращает значение 0. Если path — это файл, то `remove()` вызывает `unlink()`; если path — это каталог, то `remove()` вызывает `rmdir()`.

В случае ошибки `remove()` возвращает -1 и присваивает переменной `errno` один из допустимых кодов ошибки, определенных для `unlink()` и `rmdir()` соответственно.

Копирование и перемещение файлов

Две базовые задачи манипулирования файлами — это копирование и перемещение файлов, и для их выполнения часто применяют команды `cp` и `mv`. На уровне файловой системы *копирование* (*copying*) — это акт дублирования содержимого указанного файла в новом местоположении с новым путем к файлу. Это не то же самое, что создание новой жесткой ссылки на файл, так как изменение одного файла не влияет на второй, то есть после копирования существуют две совершенно разные копии файла под (по крайней мере) двумя записями каталога. *Перемещение* (*moving*), в свою очередь, — это акт переименования записи каталога, соответствующей файлу. Это действие не приводит к созданию второй копии.

Копирование

Это может звучать удивительно, но в Unix нет системного или библиотечного вызова, упрощающего копирование файлов и каталогов. Вместо этого такие утилиты, как `cp`, и диспетчер файлов Nautilus от GNOME выполняют подобные задачи вручную.

Для того чтобы скопировать файл `src` в файл `dst`, нужно сделать следующие шаги:

1. Открыть файл `src`.
2. Открыть файл `dst` — создать его, если он еще не существует, или усечь до нулевой длины, если такой файл уже есть.
3. Прочитать фрагмент содержимого файла `src` в память.
4. Записать этот фрагмент в файл `dst`.
5. Продолжать до тех пор, пока все содержимое файла `src` не будет считано и записано в `dst`.
6. Закрыть файл `dst`.
7. Закрыть файл `src`.

При копировании каталога отдельный каталог и любые подкаталоги создаются при помощи системного вызова `mkdir()`; каждый файл копируется индивидуально.

Перемещение

В отличие от копирования, в Unix есть системный вызов для перемещения файлов. В стандарте ANSI C был представлен вызов, предназначенный только для файлов, а POSIX стандартизировал его и для файлов, и для каталогов:

```
#include <stdio.h>
```

```
int rename (const char *oldpath, const char *newpath);
```

Успешный вызов `rename()` переименовывает путь к файлу `oldpath` в путь `newpath`. Содержимое файла и его `inode` не меняются. И `oldpath`, и `newpath` долж-

ны принадлежать одной файловой системе¹; если это не так, то вызов завершается ошибкой. Такие утилиты, как `mv`, должны обрабатывать данный случай при помощи копирования и отсоединения.

В случае успеха системный вызов `rename()` возвращает значение 0, а к файлу, на который раньше указывал путь `oldpath`, после этого можно обращаться только при помощи пути `newpath`. В случае ошибки вызов возвращает -1, не меняет `oldpath` и `newpath` и присваивает переменной `errno` одно из следующих значений:

EACCES

У вызывающего процесса нет разрешения на запись в родительский каталог файла `oldpath` или `newpath`, разрешения на поиск для компонента пути `oldpath` или `newpath` или разрешения на запись для `oldpath`, если `oldpath` — это каталог. Последний случай может быть проблемой, так как вызову `rename()` придется обновлять каталог, когда `oldpath` является каталогом.

EBUSY

`oldpath` или `newpath` — это точка монтирования.

EFAULT

Параметр `oldpath` или `newpath` содержит недопустимый указатель.

EINVAL

Путь `newpath` является составляющей `oldpath` и, таким образом, переименование одного в другой сделает `oldpath` подкаталогом самого себя.

EISDIR

Путь `newpath` существует и является каталогом, но `oldpath` — не каталог.

ELOOP

При разрешении пути `oldpath` или `newpath` встретилось слишком много символьских ссылок.

EMLINK

На `oldpath` уже указывает максимальное число ссылок, или `oldpath` — это каталог, а на `newpath` уже указывает максимальное число ссылок.

ENAMETOOLONG

Слишком длинное значение `oldpath` или `newpath`.

ENOENT

Компонент `oldpath` или `newpath` не существует или является повисшей символьской ссылкой.

ENOMEM

Недостаточно памяти ядра для выполнения данного запроса.

¹ Хотя Linux позволяет монтировать устройство в нескольких точках в структуре каталогов, невозможно переименовать файл с одной из этих точек монтирования на другую, несмотря на то, что фактически обе точки обеспечиваются одним устройством.

ENOSPC

Недостаточно пространства на устройстве для выполнения данного запроса.

ENOTDIR

Компонент (за исключением потенциально последнего компонента) `oldpath` или `newpath` не является каталогом или `oldpath` — это каталог, а `newpath` существует и не является каталогом.

ENOTEMPTY

`Newpath` — это каталог и он не пуст.

EPERM

По крайней мере, один из указанных в аргументах путей существует, для родительского каталога установлен бит закрепления в памяти, действительный идентификатор пользователя вызывающего процесса не совпадает ни с идентификатором пользователя файла, ни с идентификатором пользователя его предка, и процесс не привилегированный.

EROFS

Файловая система помечена как доступная только для чтения.

EXDEV

`Oldpath` и `newpath` принадлежат разным файловым системам.

В табл. 7.1 перечисляются результаты перемещения разных типов файлов.

Таблица 7.1. Результаты перемещения различных типов файлов

Целью является файл	Целью является каталог	Целью является ссылка	Цель не существует
Источником является файл	На место цели записывается источник	Ошибка с кодом EISDIR	Файл переименовывается и цель переписывается
Источником является каталог	Ошибка с кодом ENOTDIR	Если цель пуста, то источник переименовывается и ему присваивается имя цели; в противном случае создается ошибка с кодом ENOTEMPTY	Каталог переименовывается, и цель переписывается
Источником является ссылка	Ссылка переименовывается и цель переписывается	Ошибка с кодом EISDIR	Ссылка переименовывается, и цель переписывается
Источник не существует	Ошибка с кодом ENOENT	Ошибка с кодом ENOENT	Ошибка с кодом ENOENT

Для всех этих случаев, независимо от типа, если источник и цель находятся в разных файловых системах, вызов завершается ошибкой и возвращает код ENODEV.

Узлы устройств

Узлы устройств — это специальные файлы, позволяющие приложениям взаимодействовать с драйверами устройств. Когда приложение выполняет обычные операции ввода-вывода Unix — открытие, закрытие, чтение, запись и т. д. — на узле устройства, ядро обрабатывает эти запросы не так, как обычные запросы файлового ввода-вывода. Вместо этого ядро передает их драйверу устройства. Драйвер устройства обрабатывает операцию ввода-вывода и возвращает результаты пользователю. Узлы устройств обеспечивают абстрагирование устройств, чтобы для написания приложений не нужно было знать особенности конкретных устройств или овладевать специальными интерфейсами. Действительно, узлы устройств — это стандартный механизм для доступа к аппаратному обеспечению в системах Unix. Сетевые устройства относятся к редкому исключению, и в истории Unix встречались возражения, что такое исключение все же ошибка. Конечно же, есть какая-то элегантная красота в манипулировании всем машинным оборудованием при помощи вызовов `read()`, `write()` и `mmap()`.

Как ядро идентифицирует драйвер устройства, которому следует передать запрос? Каждому узлу устройства назначается два числовых значения, называемых *старшим номером* (majorg number) и *младшим номером* (minorg number). Пара номеров, старший и младший, соответствует конкретному драйверу устройства, загруженному в ядро. Если у узла устройства старший и младший номера не соответствуют драйверу устройства в ядре, — что по некоторым причинам иногда случается, — то запрос `open()` на этом узле возвращает значение `-1` с кодом `errno ENODEV`. Мы говорим, что такие узлы устройств представляют несуществующие устройства.

Специальные узлы устройств

Во всех системах Linux присутствует несколько узлов устройств. Они являются частью среды разработки Linux и их наличие считается частью ABI Linux.

У *фиктивного устройства* (null device) старший номер равен 1, а младший номер равен 3. Оно живет по адресу `/dev/null`. Владельцем файла устройства должен быть пользователь `root`, и оно должно быть доступно для чтения и записи всем пользователям. Ядро бесшумно удаляет все запросы на запись, направляемые на это устройство. Все запросы на чтение для этого файла возвращают значение EOF (end-of-file, конец файла).

Нулевое устройство (zero device) живет по адресу `/dev/zero`; его старший номер равен 1, а младший — 5. Как и для фиктивного устройства, ядро бесшумно

удаляет все запросы на запись в нулевое устройство. Чтение с нулевого устройства возвращает бесконечный поток нулевых байтов.

Полное устройство (full device) со старшим номером 1 и младшим номером 7 располагается по адресу `/dev/full`. Как и нулевое устройство, полное устройство при чтении с него возвращает символы нуля (`\0`). Запросы на запись, однако, всегда вызывают ошибку ENOSPC, указывая, что устройство заполнено.

Эти устройства служат разнообразным целям. Они полезны для тестирования обработки в приложении сложных и проблемных ситуаций — например, заполнения файловой системы. Так как фиктивное и нулевое устройства игнорируют запись, то они также обеспечивают не создающий дополнительной нагрузки способ избавления от нежелательного ввода-вывода.

Генератор случайных чисел

Устройства ядра, генерирующие случайные числа, располагаются в файлах `/dev/random` и `/dev/urandom`. Старший номер для них равен 1, а младшие — 8 и 9 соответственно.

Генератор случайных чисел ядра собирает шум на драйверах устройств и других источниках, а ядро сцепляет и необратимо хэширует собранный шум. Результат сохраняется в *пуле энтропии* (entropy pool). Ядро также хранит приблизительную оценку числа битов энтропии в пуле.

Запросы на чтение с устройства `/dev/random` возвращают энтропию из пула. Результаты можно использовать в качестве начальных чисел для генераторов случайных чисел, для генерирования ключей и других задач, требующих криптографически надежной энтропии.

В теории, враг, способный получить достаточно данных из пула энтропии и успешно сломать необратимый хэш, мог бы узнать о состоянии остального пула энтропии. Хотя подобная атака сегодня кажется лишь теоретической возможностью — пока что неизвестно ни об одном подобном случае, — ядро реагирует на эту вероятность, уменьшая оценку объема энтропии в пуле с каждым запросом на запись. Когда значение оценки достигает нуля, чтение блокируется до тех пор, пока система не сгенерирует дополнительную энтропию и оценка энтропии не повысится до уровня, необходимого для завершения чтения.

Устройство `/dev/urandom` не обладает этим свойством; чтение с данного устройства завершается успешно, даже если оценка энтропии, сделанная ядром, недостаточно велика для завершения запроса. Так как только самые защищенные приложения — например, генераторы ключей для безопасного обмена данными в GNU Privacy Guard — чувствительны к криптографической надежности энтропии, в большинстве приложений следует использовать `/dev/urandom`, а не `/dev/random`. Операция чтения из последнего может потенциально заблокироваться на очень длительное время, если не будет происходить достаточного количества операций ввода-вывода, чтобы заполнить пул энтропии ядра. Такая ситуация не так уж редко встречается на бездисковых безголовых серверах.

Внеполосная коммуникация

Файловая модель Unix впечатляет. При помощи простых операций чтения и записи Unix абстрагирует практически любое мыслимое действие, которое можно произвести на объекте. Иногда, однако, у программистов возникает необходимость обращаться к файлам за пределами основного потока данных. Например, рассмотрим устройство порта последовательного интерфейса. Операция чтения с этого устройства позволяет читать данные с оборудования, подключенного к порту последовательного интерфейса; операция записи на устройство отправляет данные на оборудование. Как же процессу прочитать один из специальных штырьков состояния порта, например сигнал DTR (data terminal ready, терминал данных готов)? Или как процессу установить четность порта последовательного интерфейса?

Ответ прост — нужно использовать системный вызов `ioctl()`. `ioctl()` расшифровывается как I/O control — управление вводом-выводом. Этот системный вызов позволяет осуществлять внеполосную коммуникацию:

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

```
int ioctl (int fd, int request, ...);
```

Системный вызов требует два параметра:

`fd`

Дескриптор данного файла.

`request`

Специальное значение кода запроса, предопределенное и согласованное между ядром и процессом, которое обозначает, какую операцию нужно выполнить на файле, указанном при помощи параметра `fd`.

Также вызов может получать один или несколько нетипизированных необязательных параметров (обычно целочисленные значения без знака или указатели) и передавать их ядру.

В следующей программе запрос `CDROMEJECT` используется для выброса лотка носителя из устройства CD-ROM — привода для компакт-дисков, которое пользователь указывает как первый параметр командной строки программы. Таким образом, данная программа работает так же, как стандартная команда `eject`:

```
#include <sys/types.h>
#include <sys/stat.h>
#include <fcntl.h>
#include <sys/ioctl.h>
#include <unistd.h>
#include <linux/cdrom.h>
#include <stdio.h>

int main (int argc, char *argv[])
{
    int fd, ret;
```

```
if (argc < 2) {
    fprintf (stderr,
             "использование: %s <устройство>\n",
             argv[0]);
    return 1;
}

/*
 * Открывает устройство CD-ROM в режиме только для чтения. Флаг O_NONBLOCK
 * говорит ядру, что мы хотим открыть устройство, даже если
 * в приводе отсутствует носитель.
 */
fd = open (argv[1], O_RDONLY | O_NONBLOCK);
if (fd < 0) {
    perror ("open");
    return 1;
}

/* Отправляет команду выброса на устройство CD-ROM. */
ret = ioctl (fd, CDROMEJECT, 0);
if (ret) {
    perror ("ioctl");
    return 1;
}

ret = close (fd);
if (ret) {
    perror ("close");
    return 1;
}

return 0;
}
```

Запрос CDROMEJECT относится к возможностям драйвера устройства CD-ROM в Linux. Когда ядро получает запрос ioctl(), оно находит файловую систему (в случае реальных файлов) или драйвер устройства (в случае узлов устройств), ответственные за указанный дескриптор файла, и передает им запрос для дальнейшей обработки. В данной программе драйвер привода для компакт-дисков получает запрос и физически выбрасывает лоток.

Далее в этой главе мы рассмотрим пример использования ioctl() с дополнительным параметром, позволяющим возвращать информацию запрашивающему процессу.

Мониторинг событий файлов

В Linux существует интерфейс inotify, предназначенный для мониторинга файлов, например, позволяющий узнавать, когда они перемещаются,читываются, записываются или удаляются. Предположим, вы пишете графический диспетчер файлов, такой, как Nautilus от GNOME. Если файл копируется в каталог,

пока в Nautilus отображается его содержимое, представление каталога в диспетчере файлов становится неверным.

Одним из решений может быть постоянное считывание содержимого каталога, позволяющее распознавать изменения и обновлять отображение. Однако это создает периодическую нагрузку и не является элегантным решением. Кроме того, при этом постоянно существует состязание между моментами удаления или добавления файлов в каталог и моментами, когда диспетчер файлов считывает содержимое каталога.

Благодаря *inotify* ядро *проталкивает* (push) событие в приложение точно в тот момент, когда оно происходит. Как только файл удаляется, ядро может уведомить об этом диспетчер Nautilus. В ответ Nautilus может мгновенно убрать удаленный файл из графического представления каталога.

Многие другие приложения также зависят от событий файлов. Рассмотрим утилиту архивации или инструмент индексации данных. Интерфейс *inotify* позволяет этим программам работать в реальном времени: как только файл создается, удаляется или записывается, утилиты обновляют архив или индекс.

Интерфейс *inotify* пришел на смену *dnotify*, более старому механизму мониторинга файлов с громоздким интерфейсом на основе сигналов. В приложениях всегда следует применять *inotify*, а не *dnotify*. *Inotify*, впервые представленный в ядре 2.6.13, достаточно гибок и прост в использовании, так как те же операции, которые программы выполняют на обычных файлах, — в частности, *select()* и *poll()* — работают и с *inotify*. В этой книге мы рассмотрим только интерфейс *inotify*.

Инициализация *inotify*

Для того чтобы использовать интерфейс *inotify*, процесс сначала должен инициализировать его. Системный вызов *inotify_init()* инициализирует *inotify* и возвращает дескриптор файла, представляющий инициализированный экземпляр:

```
#include <inotify.h>

int inotify_init (void);
```

В случае ошибки *inotify_init()* возвращает -1 и присваивает переменной *errno* один из следующих кодов:

EMFILE

Достигнут лимит максимального числа экземпляров *inotify* для пользователя.

ENFILE

Достигнут лимит максимального числа дескрипторов файла для всей системы.

ENOMEM

Недостаточно памяти для выполнения данного запроса.

Давайте инициализируем `inotify`, чтобы применить этот интерфейс на последующих шагах:

```
int fd:  
  
fd = inotify_init ( );  
if (fd == -1) {  
    perror ("inotify_init");  
    exit (EXIT_FAILURE);  
}
```

Стражи

После того как процесс инициализирует `inotify`, он настраивает *стражи* (`watch`). Страж, представляемый дескриптором стража, — это стандартный путь в Unix, с которым связана маска стража, сообщающая ядру, в каких событиях заинтересован процесс: в чтении, записи или в том и другом.

`inotify` умеет наблюдать и за файлами, и за каталогами. Наблюдая за каталогом, `inotify` сообщает о событиях, относящихся к самому каталогу, а также ко всем файлам, принадлежащим данному каталогу (но не к файлам в подкаталогах наблюдаемого каталога — страж не поддерживает рекурсивность).

Добавление нового стража

Системный вызов `inotify_add_watch()` добавляет к экземпляру `inotify`, представляемому дескриптором файла `fd`, страж для события или событий файла или каталога `path`, указанных при помощи параметра `mask`:

```
#include <inotify.h>
```

```
int inotify_add_watch (int fd,  
                      const char *path,  
                      uint32_t mask);
```

В случае успеха вызов возвращает дескриптор нового стража. В случае ошибки `inotify_add_watch()` возвращает `-1` и присваивает переменной `errno` один из следующих кодов:

EACCESS

Запрещен доступ на чтение к файлу, указанному при помощи параметра `path`. Для того чтобы добавлять страж к файлу, вызывающий процесс должен иметь возможность читать этот файл.

EBADF

Дескриптор файла `fd` не соответствует допустимому экземпляру `inotify`.

EFAULT

Недопустимый указатель в параметре `path`.

EINVAL

Маска стража, `mask`, не содержит допустимые события.

ENOMEM

Недостаточно памяти для выполнения данного запроса.

ENOSPC

Достигнут лимит общего количества стражей *inotify* для пользователя.

Маски стражи

Маска стражи представляет собой одно или несколько событий, объединенных операцией бинарного ИЛИ. Эти события определяются в файле *<inotify.h>*:

IN_ACCESS

Файл был считан.

IN_MODIFY

В файл были записаны данные.

IN_ATTRIB

Были изменены метаданные файла (например, владелец, разрешения или расширенные атрибуты).

IN_CLOSE_WRITE

Был закрыт файл, который ранее был открыт для записи.

IN_CLOSE_NOWRITE

Был закрыт файл, который ранее не был открыт для записи.

IN_OPEN

Файл был открыт.

IN_MOVED_FROM

Из наблюдаемого каталога был перемещен файл.

IN_MOVED_TO

В наблюдаемый каталог был перемещен файл.

IN_CREATE

В наблюдаемом каталоге был создан файл.

IN_DELETE

В наблюдаемом каталоге был удален файл.

IN_DELETE_SELF

Сам наблюдаемый объект был удален.

IN_MOVE_SELF

Сам наблюдаемый объект был перемещен.

Также определены следующие события, группирующие в одно значение два или несколько событий:

IN_ALL_EVENTS

Все допустимые события.

IN_CLOSE

Все события, относящиеся к закрытию (в настоящее время это IN_CLOSE_WRITE и IN_CLOSE_NOWRITE).

IN_MOVE

Все события, связанные с перемещением (в настоящее время это IN_MOVED_FROM и IN_MOVED_TO).

Теперь можно посмотреть, как новый страж добавляется к существующему экземпляру inotify:

```
int wd;
wd = inotify_add_watch(fd, "/etc", IN_ACCESS | IN_MODIFY);
if (wd == -1) {
    perror("inotify_add_watch");
    exit(EXIT_FAILURE);
}
```

В этом примере добавляется страж, наблюдающий все операции считывания и записи в каталоге /etc. Если любой файл в каталоге /etc записывается или считывается, inotify отправляет событие дескриптору файла inotify, fd, указывая также дескриптор стража wd. Давайте взглянем, как inotify представляет эти события.

События inotify

События inotify представляют структура inotify_event, определенная в файле заголовка <inotify.h>:

```
#include <inotify.h>

struct inotify_event {
    int wd;           /* дескриптор стража */
    uint32_t mask;   /* маска событий */
    uint32_t cookie; /* уникальное значение cookie */
    uint32_t len;    /* размер поля name */
    char name[];    /* имя, завершающееся нулем */
};
```

Wd идентифицирует дескриптор стража, полученный при помощи inotify_add_watch(), а mask представляет события. Если wd представляет каталог и одно из наблюдаемых событий происходит для файла в этом каталоге, то в name сохраняется имя файла (относительное для данного каталога). Тогда значение len оказывается больше нуля. Обратите внимание, что значение len *не равно* длине строки name; в name может быть более одного завершающего нулевого символа, которые служат в качестве забивки, гарантирующей, что последующая структура inotify_event будет правильно выровнена. Следовательно, для вычисления смещения следующей структуры inotify_event в массиве нужно использовать len, а не strlen().

Например, если wd представляет /home/rlove, а значение mask равно IN_ACCESS, то, когда считывается файл /home/rlove/canon, поле name принимает значение canon, а len становится равным, как минимум, шести. Если бы мы наблюдали

файл /home/rlove/canon напрямую с той же маской, то значение len было бы равно 0, а name было бы нулевой длины — его трогать нельзя.

Поле cookie используется для связывания двух родственных, но непересекающихся событий. Мы поговорим об этом чуть далее.

Чтение событий inotify

Получать событий inotify легко: вы просто читаете данные из дескриптора файла, связанного с экземпляром inotify. С inotify связана возможность, называемая заглатыванием (slurping), которая позволяет считывать при помощи одного запроса на чтение сразу несколько событий — столько, сколько помещается в буфер, предоставляемый системному вызову read(). Благодаря тому что поле name обладает переменной длиной, это наиболее распространенный способ считывания событий inotify.

В предыдущем примере мы создали экземпляр inotify и добавили к нему страж. Теперь реализуем чтение ожидающих событий:

```
char buf[BUF_LEN]_attribute_(aligned(4));
ssize_t len, i = 0;

/* считать BUF_LEN байтов событий */
len = read (fd, buf, BUF_LEN);

/* пройти в цикле по всем считанным событиям */
while (i < len) {
    struct inotify_event *event =
        (struct inotify_event *) &buf[i];

    printf ("wd=%d mask=%d cookie=%d len=%d dir=%s\n".
            event->wd, event->mask,
            event->cookie, event->len,
            (event->mask & IN_ISDIR) ? "yes" : "no");

    /* если есть имя, то напечатать его */
    if (event->len)
        printf ("name=%s\n", event->name);

    /* обновить индекс, установив его на начало следующего события */
    i += sizeof (struct inotify_event) + event->len;
}
```

Так как дескриптор файла inotify работает как обычный файл, программы могут наблюдать за ним при помощи системных вызовов select(), poll() и epoll(). Это позволяет объединять события inotify с остальным файловым вводом-выводом из одного потока.

Расширенные события inotify

Помимо стандартных событий, inotify поддерживает генерирование прочих событий:

IN_IGNORED

Соответствующий дескриптору wd страж был удален. Это может происходить, когда пользователь вручную удаляет страж или когда наблюдаемый объект прекращает существовать. Мы обсудим это событие в следующем разделе.

IN_ISDIR

Затронутый объект является каталогом (если данный флаг не установлен, то затронутый объект является файлом).

IN_Q_OVERFLOW

Переполнение очереди `inotify`. Ядро ограничивает размер очереди событий, чтобы предотвращать неограниченное потребление памяти ядра. Как только число ожидающих событий становится на единицу меньше максимума, ядро генерирует данное событие и присоединяет его к хвосту очереди. Дальнейшие события не генерируются до тех пор, пока какое-либо событие не считывается из очереди, уменьшая таким образом, ее размер.

IN_UNMOUNT

Устройство, которому принадлежит наблюдаемый объект, было отсоединено. Следовательно, объект более недоступен; ядро удаляет страж и генерирует событие `IN_IGNORED`.

Эти события может генерировать любой страж; пользователю не нужно явно устанавливать их.

Программисты должны обрабатывать маску `mask` как битовую маску ожидающих событий. Следовательно, никогда *не* проверяйте события, применяя прямые проверки эквивалентности:

```
/* НЕ делайте так! */
```

```
if (event->mask == IN_MODIFY)
    printf ("В файл были записаны данные!\n");
else if (event->mask == IN_Q_OVERFLOW)
    printf ("Очередь переполнена!\n");
```

Вместо этого следует выполнять побитовые сравнения:

```
if (event->mask & IN_ACCESS)
    printf ("Из файла были считаны данные!\n");
if (event->mask & IN_UNMOUNTED)
    printf ("Отсоединенено устройство, которому принадлежит файл!\n");
if (event->mask & IN_ISDIR)
    printf ("Файл является каталогом!\n");
```

Объединение нескольких событий

Каждое из событий `IN_MOVED_FROM` и `IN_MOVED_TO` представляет только половину перемещения: первое – удаление из указанного местоположения, а второе – прибытие в новое местоположение. Таким образом, для того чтобы действительно иметь смысл в программах, старающихся разумно отслеживать перемещения файлов по файловой системе (возьмем, к примеру, индексатор, который не должен повторно индексировать перемещаемые файлы), процессы должны уметь связывать два или более событий в одно целое.

Здесь в игру вступает поле `cookie` структуры `inotify_event`.

Поле `cookie`, если оно не равно нулю, содержит уникальное значение, связывающее два события. Предположим, у нас есть процесс, наблюдающий за

каталогами `/bin` и `/sbin`. Предположим, что у `/bin` дескриптор стражи равен 7, а у `/sbin` дескриптор стражи равен 8. Если файл `/bin/compass` перемещается в `/sbin/compass`, то ядро генерирует два события `inotify`.

У первого события дескриптор `wd` равен 7, в поле `mask` содержится значение `IN_MOVED_FROM`, а в поле `name` — значение `compass`. У второго события поле `wd` равно 8, поле `mask` — `IN_MOVED_TO`, а поле `name` — также `compass`. В обоих событиях значение `cookie` одно и то же, скажем 12.

Если файл переименовывается, то ядро все равно генерирует два события. У этих событий одинаковые значения `wd`.

Обратите внимание, что если файл перемещается в каталог или из каталога, наблюдение за которым не настроено, то процесс не получает одно из соответствующих событий. В программе можно реализовать уведомление о том, что второе событие с таким же значением `cookie` не прибыло.

Расширенные параметры стражей

При создании нового стража можно добавлять к значению `mask` одно или несколько из следующих значений, чтобы управлять поведением стража:

IN_DONT_FOLLOW

Если это значение установлено и если цель пути `path` или любой из его компонентов является символьской ссылкой, то вызов не следует по ссылке, и `inotify_add_watch()` завершается ошибкой.

IN_MASK_ADD

Обычно, когда `inotify_add_watch()` вызывается на файле, для которого уже существует другой страж, маска стража обновляется, отражая переданное в параметрах последнего вызова значение. Если значение `mask` содержит данный флаг, то указанные события добавляются к существующей маске.

IN_ONESHOT

Если установлено данное значение, то ядро автоматически удаляет страж после генерирования первого события для указанного объекта. То есть страж, фактически, является «одноразовым».

IN_ONLYDIR

Если это значение установлено, то страж добавляется, только если указанный объект является каталогом. Если `path` представляет файл, а не каталог, то вызов `inotify_add_watch()` завершается ошибкой.

Например, в следующем фрагменте кода страж для `/etc/init.d` добавляется только в том случае, если `init.d` — это каталог и если ни `/etc`, ни `/etc/init.d` не являются символьскими ссылками:

`int wd;`

`/*`

`* Наблюдение за перемещением /etc/init.d, но только если это каталог
* и ни одна из составляющих его пути не является символьской ссылкой.`

```
/*
wd = inotify_add_watch (fd,
    "/etc/init.d",
    IN_MOVE_SELF |
    IN_ONLYDIR |
    IN_DONT_FOLLOW);
if (wd == -1)
    perror ("inotify_add_watch");
```

Удаление стражи inotify

Удалить страж из экземпляра inotify можно при помощи системного вызова `inotify_rm_watch()`:

```
#include <inotify.h>

int inotify_rm_watch (int fd, uint32_t wd);
```

Успешный вызов `inotify_rm_watch()` удаляет страж, представляемый дескриптором стража `wd`, из экземпляра `inotify`, указанного при помощи дескриптора файла `fd`, и возвращает значение 0. Например:

```
int ret.

ret = inotify_rm_watch (fd, wd);
if (ret)
    perror ("inotify_rm_watch");
```

В случае ошибки системный вызов возвращает -1 и присваивает переменной `errno` одно из следующих двух значений:

EBADF

Fd не является допустимым экземпляром `inotify`.

EINVAL

Wd не является допустимым дескриптором стража для указанного экземпляра `inotify`.

При удалении стражи ядро генерирует событие `IN_IGNORED`. Оно отправляет это событие не только во время удаления вручную, но и при разрушении стражи в качестве побочного эффекта другой операции. Например, когда наблюдаемый файл удаляется, также удаляются любые существующие для него стражи. Во всех подобных случаях ядро отправляет событие `IN_IGNORED`. Такое поведение позволяет приложениям объединять всю обработку удаления стражей в одном месте: в обработчике событий для `IN_IGNORED`. Это полезно для комплексных приложений, использующих `inotify`, которые управляют сложными структурами данных для каждого стража `inotify`, таких, как инфраструктура поиска Beagle от GNOME.

Проверка размера очереди событий

Размер очереди ожидающих событий можно проверить при помощи системного вызова `ioctl` на дескрипторе файла экземпляра `inotify`, добавив константу

FIONREAD. В первом аргументе запроса сохраняется размер очереди в байтах как целочисленное значение без знака:

```
unsigned int queue_len;
int ret;

ret = ioctl (fd, FIONREAD, &queue_len);
if (ret < 0)
    perror ("ioctl");
else
    printf ("Ожидает в очереди %u байтов\n", queue_len);
```

Обратите внимание, что этот запрос возвращает размер очереди в байтах, а не число событий в очереди. Программа может сделать оценку числа событий, исходя из количества байтов, используя известный размер структуры `inotify_event` (полученный при помощи `sizeof()`) и сделав предположение о среднем размере поля `name`. Зная количество ожидающих байтов, процесс может точно определять размер считывания.

Константа FIONREAD определяется в файле заголовка `<sys/ioctl.h>`.

Разрушение экземпляра inotify

Для того чтобы разрушить экземпляр `inotify` и все связанные с ним стражи, нужно просто закрыть дескриптор файла данного экземпляра:

```
int ret;

/* значение fd было получено при помощи inotify_init( ) */
ret = close (fd);
if (fd == -1)
    perror ("close");
```

Конечно же, как и для любого дескриптора файла, когда процесс завершается, ядро автоматически закрывает дескриптор файла и очищает все ресурсы.

8

Управление памятью

Память — это один из основных и в то же время самых важных ресурсов, используемых процессом. В этой главе мы рассмотрим управление данным ресурсом: выделение, манипулирование и последующее освобождение памяти.

Термин «выделять» (*allocate*), часто используемый для обозначения получения памяти, может вводить в заблуждение, так как он создает впечатление нормированного распределения скучного ресурса, спрос на который превышает предложение. Без сомнения, многие пользователи с удовольствием получили бы побольше памяти. Однако в современных системах проблема, в действительности, заключается не в том, чтобы поделить скромный ресурс на слишком большое количество пользователей, а в том, чтобы правильно использовать его и отслеживать расходование выделенного богатства.

В этой главе вы познакомитесь со всеми подходами к выделению памяти в различных областях программы, обращая особое внимание на преимущества и недостатки каждого метода. Мы также рассмотрим несколько способов установки и манипулирования содержимым произвольных областей памяти и познакомимся с тем, как блокировать память, чтобы данные оставались в ОЗУ, а ваша программа выполнялась оперативно, без затрат времени на ожидание, пока ядро скопирует данные из области подкачки.

Адресное пространство процесса

Linux, как и любая современная операционная система, виртуализирует свой физический ресурс памяти. Процессы не адресуют напрямую физическую память. Вместо этого ядро связывает с каждым процессом уникальное *виртуальное адресное пространство* (*virtual address space*). Это адресное пространство линейно, адреса в нем начинаются с нуля, и их значения увеличиваются до определенного максимума.

Страницы и подкачка страниц

Виртуальное адресное пространство состоит из *страниц* (page). Архитектура системы и тип машины определяют размер страницы. Это фиксированное значение, которое обычно равно 4 Кбайт (для 32-разрядных систем) или 8 Кбайт (для 64-разрядных систем)¹. Страницы бывают действительными и недействительными. *Действительная страница* (valid page) связана со страницей в физической памяти или на каком-то дополнительном запоминающем устройстве, это может быть раздел подкачки или файл на диске. Недействительная страница не связана ни с чем и представляет собой неиспользуемый, невыделенный фрагмент адресного пространства. Доступ к такой странице вызывает нарушение сегментации. Адресное пространство не обязательно бывает непрерывным. Хотя адресуется оно линейно, оно содержит множество неадресуемых промежутков.

Программа не может использовать страницу, существующую на дополнительном запоминающем устройстве, а не в физической памяти, до тех пор, пока эта страница не будет связана со страницей в физической памяти. Когда процесс пытается обратиться к адресу на такой странице, блок управления памятью (memory management unit, MMU) генерирует *ошибку обращения к несуществующей странице* (page fault). После этого вмешивается ядро, прозрачно подкачивая (page in) нужную страницу с дополнительного запоминающего устройства в физическую память. Так как виртуальной памяти намного больше, чем физической (во многих системах даже в одном виртуальном адресном пространстве!), ядро также постоянно *выгружает* (page out) физическую память на дополнительное устройство, освобождая место для подкачки новых данных. Ядро старается выгружать данные, которые с наименьшей вероятностью потребуются в ближайшем будущем, оптимизируя производительность.

Совместное использование и копирование при записи

Несколько страниц виртуальной памяти, даже в разных виртуальных адресных пространствах, принадлежащих разным процессам, могут соответствовать одной физической странице. Это позволяет разным виртуальным адресным пространствам *совместно использовать* (share) данные в физической памяти. Совместно используемая память может быть доступна только для чтения или для чтения и записи.

Когда процесс записывает данные в общую страницу, доступную для записи, реализуется одна из двух ситуаций. В простейшем случае ядро позволяет записи осуществляться, и тогда все процессы, совместно использующие страницу, могут увидеть результаты операции записи. Обычно для того чтобы несколько

¹ Некоторые системы поддерживают несколько размеров страницы. По этой причине размер страницы не является частью ABI. Приложения должны программно узнавать размер страницы во время выполнения. Мы говорили об этом в главе 4 и еще раз вернемся к этому вопросу в данной главе.

процессов могли читать или записывать данные в общую страницу, необходим какой-то уровень координирования и синхронизации.

В противном случае MMU может перехватывать операцию записи и создавать исключение; в ответ ядро прозрачно создает новую копию страницы для записывающего процесса и позволяет операции записи продолжиться на этой новой странице. Такой подход называется *копированием при записи* (copy-on-write, COW)¹. Фактически, процессам разрешается доступ на чтение к совместно используемым данным, что экономит пространство. Когда процесс желает записать данные в общую страницу, он получает уникальную копию этой страницы на лету, что позволяет ядру вести себя так, как будто у процесса всегда была собственная копия данной страницы. Если копирование при записи выполняется построчно, то огромный файл может совместно использоваться многими процессами и лишь отдельные процессы, которым необходимо записать что-то, получают уникальные физические страницы, соответствующие требуемым страницам файла, в которые и записывают данные.

Области памяти

Ядро организует страницы в блоки, обладающие определенными общими свойствами, такими, как разрешения на доступ. Эти блоки называются областями памяти, сегментами или отображениями. Некоторые типы областей памяти можно найти в каждом процессе:

- *текстовый сегмент* (text segment) содержит программный код процесса, строковые литералы, переменные с постоянными значениями и прочие доступные только для чтения данные. В Linux этот сегмент помечается как доступный только для чтения и отображается в память напрямую из объектного файла (исполняемого файла программы или библиотеки);
- в *стеке* (stack) содержится стек выполнения процесса, который растет и сжимается динамически, по мере того как глубина стека увеличивается и уменьшается. В стеке выполнения находятся локальные переменные и возвращаемые данные функций;
- в *сегменте данных* (data segment), или *куче* (heap), находится динамическая память процесса. Этот сегмент доступен для чтения и может расти и сжиматься. Именно эту память возвращает `malloc()` (подробнее об этом интерфейсе – в следующем разделе);
- сегмент *bss*² содержит не инициализированные глобальные переменные. В этих переменных находятся специальные значения (фактически, все нули) согласно стандарту C. Linux оптимизирует эти переменные двумя путями.

¹ В главе 5 говорилось, что системный вызов `fork()` применяет копирование при записи для дублирования и совместного использования родительского адресного пространства с потомком.

² Это историческое название, происходящее от термина *block started by symbol* (блок, начинаящийся с символа).

- Во-первых, поскольку сегмент `bss` предназначен для не инициализированных данных, компоновщик (`ld`) в действительности не хранит специальные значения в объектном файле. Это уменьшает размер двоичного файла. Во-вторых, когда этот сегмент загружается в память, ядро сопоставляет его путем копирования при записи со страницей, состоящей из одних нулей, фактически, устанавливая для переменных значения по умолчанию;
- большинство адресных пространств содержит несколько *отображенных файлов* (mapped file), таких, как сам исполняемый файл программы, библиотека С и прочие связанные библиотеки, а также файлы данных. Взгляните на `/proc/self/maps` или на вывод программы `ps`, и вы увидите отличный пример отображенных файлов в процессе.

В этой главе рассматриваются интерфейсы, предоставляемые в Linux для получения и возвращения памяти, создания и разрушения новых отображений и всего, что находится между ними.

Выделение динамической памяти

Память также может принимать форму автоматических и статических переменных, но основой любой системы управления памятью является выделение, использование и последующее возвращение *динамической памяти* (dynamic memory). Динамическая память выделяется во время выполнения, а не во время компилирования, и размер выделяемой памяти может быть неизвестен вплоть до момента выделения. Как у разработчика, у вас возникает необходимость в динамической памяти, когда объем требуемой памяти или длительность ее использования варьируется и точные значения неизвестны до тех пор, пока программа не начинает выполняться. Например, вы можете хранить в памяти содержимое файла или входные данные, считанные с клавиатуры. Так как размер файла неизвестен, а пользователь может ввести любое количество символов, размер буфера может быть разным, а вы должны динамически увеличивать его по мере считывания дополнительных данных.

В С переменная не может поддерживаться динамической памятью. Например, в языке С нет механизма для получения структуры `struct pirate_ship`, существующей в динамической памяти. Вместо этого С предоставляет механизм для выделения динамической памяти в объеме, достаточном для того, чтобы сохранить там структуру `pirate_ship`. После этого программист взаимодействует с памятью через указатель, в данном случае — `struct pirate_ship *`.

Классический интерфейс С для получения динамической памяти — это `malloc()`:

```
#include <stdlib.h>

void * malloc (size_t size);
```

Успешный вызов `malloc()` выделяет `size` байтов памяти и возвращает указатель на начало выделенной области. Содержимое памяти не определено; не

ожидайте, что она будет заполнена нулями. В случае сбоя `malloc()` возвращает `NULL` и присваивает переменной `errno` значение `ENOMEM`.

Использовать `malloc()` довольно просто; в данном примере мы применяем этот системный вызов для выделения фиксированного количества байтов:

```
char *p;
```

```
/* дайте мне 2 Кбайт! */
p = malloc (2048);
if (!p)
    perror ("malloc");
```

А в этом примере вызов используется для выделения структуры:
struct treasure_map *map;

```
/*
 * Выделить достаточно памяти, чтобы там поместилась структура treasure_map.
 * и сделать так, чтобы на нее указывала переменная map.
 */
map = malloc (sizeof (struct treasure_map));
if (!map)
    perror ("malloc");
```

Язык С автоматически повышает указатели на `void` при любых типах назначения. Таким образом, в этих примерах не требуется приводить тип возвращаемого значения `malloc()` к типу `lvalue`, используемому при назначении. Язык программирования C++, однако, не выполняет автоматическое повышение указателей `void`. Следовательно, пользователям C++ приходится приводить тип возвращаемого значения вызова `malloc()`, как показано далее:

```
char *name;
/* выделить 512 байт */
name = (char *) malloc (512);
if (!name)
    perror ("malloc");
```

Некоторым программистам на С нравится приводить тип результата любой функции, возвращающей указатель на `void`, в том числе `malloc()`. Я выступаю против этой практики, так как она скроет ошибку, если возвращаемое значение функции вдруг поменяется на что-то другое, отличное от указателя `void`. Помимо этого, подобное приведение типов также прячет ошибку, если функция неправильно объявлена¹. Хотя первый случай не является большим риском для `malloc()`, второй, определенно, скрывает опасность.

Так как `malloc()` может возвращать значение `NULL`, необходимо, чтобы разработчики всегда проверяли и обрабатывали условия ошибки. Во многих программах определяется и используется обертка для `malloc()`, выводящая сообщение

¹ Необъявленные функции по умолчанию возвращают значение типа `int`. Превращение целочисленного значения в указатель не может быть автоматическим, и при этом создается предупреждение. Приведение типа подавляет возвращение предупреждения.

об ошибке и прерывающей программу, если `malloc()` возвращает `NULL`. По традиции разработчики называют эту распространенную обертку `xmalloc()`:

```
/* аналогично malloc(), но завершается в случае ошибки */
void * xmalloc (size_t size)
{
    void *p;

    p = malloc (size);
    if (!p) {
        perror ("xmalloc");
        exit (EXIT_FAILURE);
    }

    return p;
}
```

Выделение массивов

Выделение динамической памяти может быть довольно сложным в том случае, если выделяемый размер сам является динамическим. Один из подобных примеров — динамическое выделение массивов, когда размер элемента массива может быть фиксированным, но число элементов переменное. Для упрощения такого случая в С предусмотрена функция `calloc()`:

```
#include <stdlib.h>

void * calloc (size_t nr, size_t size);
```

Успешный вызов `calloc()` возвращает указатель на блок памяти, вмещающий массив из `nr` элементов размером `size` каждый. Следовательно, объем запрашиваемой памяти в следующих двух вызовах идентичен (любой может вернуть больше памяти, чем запрошено, но никак не меньше):

```
int *x, *y;

x = malloc (50 * sizeof (int));
if (!x) {
    perror ("malloc");
    return -1;
}

y = calloc (50, sizeof (int));
if (!y) {
    perror ("calloc");
    return -1;
}
```

Но поведение, однако, не идентично. В отличие от вызова `malloc()`, который не дает никаких гарантий относительно содержимого выделенной памяти, `calloc()` обнуляет все байты в возвращаемом фрагменте памяти. Таким образом, все 50 элементов в целочисленном массиве `y` содержат значение 0, а содержимое элементов в `x` не определено. Если только программа не намеревается сразу же установить все 50 значений, программистам следует использовать

`calloc()`, чтобы гарантировать, что элементы массива не будут заполнены мусором. Обратите внимание, что двоичный ноль может отличаться от нуля с плавающей точкой!

Пользователям часто бывает нужно «обнулять» динамическую память, даже когда речь идет не о массивах. Далее в этой главе мы рассмотрим `memset()`, предоставляющий интерфейс для установки определенного значения для каждого байта во фрагменте памяти. Позволять `calloc()` обнулять память, однако, быстрее, потому что ядро в этом случае предоставляет память, которая уже содержит только нули.

В случае ошибки `calloc()`, как и `malloc()`, возвращает `NULL` и присваивает переменной `errno` код `ENOMEM`.

Почему в стандартах не определена функция «выделить и обнулить», отличная от `calloc()` для массивов, — загадка. Разработчики, однако, могут с легкостью определить собственный интерфейс:

```
/* работает идентично malloc( ), но память обнуляется */
void * malloc0 (size_t size)
{
    return calloc (1, size);
}
```

Удобно объединять интерфейс `malloc0()` с предыдущим `xmalloc()`:

```
/* аналогично malloc( ), но обнуляет память и завершается в случае ошибки */
void * xmalloc0 (size_t size)
{
    void *p.

    p = calloc (1, size);
    if (!p) {
        perror ("xmalloc0");
        exit (EXIT_FAILURE);
    }

    return p;
}
```

Изменение размера выделенной памяти

В языке С есть интерфейс для изменения размера (в большую или меньшую сторону) существующих областей выделенной памяти:

```
#include <stdlib.h>
```

```
void * realloc (void *ptr, size_t size);
```

Успешный вызов `realloc()` меняет размер области памяти, на которую указывает `ptr`, увеличивая или уменьшая ее до `size` байтов. Он возвращает указатель на область нового размера, причем указатель может совпадать или не совпадать с `ptr` — при увеличении области памяти, если `realloc()` не удается расширить существующий фрагмент памяти на том же месте, функция может выделить новую область памяти длиной в `size` байтов, скопировать старую

область в новую и освободить старую область. При любых операциях сохраняется либо все содержимое области памяти (увеличение), либо содержимое, помещающееся в новый размер (уменьшение). Из-за того что существует возможность копирования данных, увеличение области памяти при помощи `realloc()` потенциально может быть относительно дорогой операцией.

Если значение параметра `size` равно 0, то эффект оказывается таким же, как при вызове `free()` на указателе `ptr`.

Если значение `ptr` равно `NULL`, то результат операции аналогичен применению вызова `malloc()`. Если указатель `ptr` не равен `NULL`, то он должен содержать значение, ранее возвращенное вызовом `malloc()`, `calloc()` или `realloc()`.

В случае ошибки `realloc()` возвращает `NULL` и присваивает переменной `errno` код `ENOMEM`. Состояние памяти, на которую указывает `ptr`, не меняется.

Рассмотрим пример уменьшения области памяти. Сначала мы используем системный вызов `calloc()` для выделения достаточного количества памяти, чтобы сохранить массив структур `map` из двух элементов:

```
struct map *p;
```

```
/* выделение памяти для двух структур map */
p = calloc (2, sizeof (struct map));
if (!p) {
    perror ("calloc");
    return -1;
}
/* использование p[0] и p[1]... */
```

Теперь предположим, что мы нашли одно из двух сокровищ и нам больше не нужна вторая карта (речь идет о названии структур — `map` (карта)), поэтому мы решили изменить размер области памяти, вернув половину ее системе (это довольно бессмысленная операция, которая могла бы быть полезной, только если структура `map` была бы очень большой, а оставшуюся карту мы бы планировали хранить долгое время):

```
struct map *r;

/* теперь нам нужна память только для одной карты */
r = realloc (p, sizeof (struct map));
if (!r) {
    /* обратите внимание, что r до сих пор действует! */
    perror ("realloc");
    return -1;
}
/* использование r... */

free (r);
```

В этом примере `p[0]` сохраняется после вызова `realloc()`. Какие бы данные ни хранились в этом элементе, они все еще там. Если вызов завершается ошибкой, то значение `r` не меняется и, следовательно, все так же действует. Можно продолжать использовать его, и в конечном итоге нам понадобится освободить

этот массив. И наоборот, если вызов завершается успешно, то мы игнорируем *p* и вместо него используем массив *r* (который, вероятно, все равно совпадает с *p*, так как область практически точно сжималась на старом месте). Теперь наша задача — освободить *r*, когда мы закончим работу с этим массивом.

Освобождение динамической памяти

В отличие от автоматически выделяемых областей, которые автоматически освобождаются, когда стек возвращается в исходное состояние, динамически выделяемые области остаются постоянными частями адресного пространства процесса до тех пор, пока не освобождаются вручную. Программист, таким образом, несет ответственность за возвращение динамически выделенной памяти системе (и статически, и динамически выделенные области, конечно же, исчезают, когда процесс завершается).

Память, выделенную при помощи `malloc()`, `calloc()` или `realloc()`, после того как необходимость в ней отпадает, необходимо возвращать системе, используя для этого вызов `free()`:

```
#include <stdlib.h>
```

```
void free (void *ptr);
```

Вызов `free()` освобождает память, на которую указывает *ptr*. Значение параметра *ptr* должно быть значением, возвращенным вызовом `malloc()`, `calloc()` или `realloc()`. Это означает, что невозможно применить `free()` для освобождения части области памяти, например половины фрагмента, передав вызову указатель на середину выделенного блока.

Значение *ptr* может быть равно `NULL` — в этом случае `free()` молча возвраща-ет значение. Таким образом, распространенная практика, включающая проверку *ptr* на `NULL` перед вызовом `free()` в действительности не имеет смысла.

Рассмотрим пример:

```
void print_chars (int n, char c)
{
    int i.

    for (i = 0; i < n; i++) {
        char *s;
        int j;

        /*
         * Выделение и обнуление массива из I + 2 символьных
         * элементов. Обратите внимание, что sizeof (char)
         * всегда равно 1.
         */
        s = calloc (i + 2, 1);
        if (!s) {
            perror ("calloc");
            break;
    }
```

```
for (j = 0, j < i + 1, j++)
    s[j] = c;

printf ("%s\n", s);

/* Все закончилось. Можно возвращать память. */
free (s);
}

}
```

Код в этом примере выделяет n массивов значений типа `char`. Каждый последующий массив содержит на единицу больше элементов, чем предыдущий. В первом массиве два элемента (2 байт), в последнем $n + 1$ элементов ($n + 1$ байт). Затем в каждый массив символ `c` записывается в цикле в каждый байт, за исключением последнего (в последнем байте остается значение 0, которое там находится с самого начала). Массив выводится как строка, а затем динамически выделенная память освобождается.

Если вызвать `print_chars()`, когда n равно 5, а `c` содержит символ `X`, то мы получим

```
X
XX
XXX
XXXX
XXXXX
```

Конечно существуют значительно более эффективные способы реализации этой функции. Смысл заключается в том, что мы можем динамически выделять и освобождать память, даже если размер и число выделяемых областей становятся известны только во время выполнения.

ПРИМЕЧАНИЕ

Системы Unix, такие, как SunOS и SCO, предоставляют вариант вызова `free()` с именем `cfree()`, который, в зависимости от системы, работает так же, как `free()`, или принимает три параметра, отражая поведение `calloc()`. В Linux `free()` обрабатывает только память, полученную при помощи одного из механизмов выделения, о которых мы говорили выше. Не следует использовать `cfree()`, за исключением тех случаев, когда необходима обратная совместимость. Версия вызова в Linux совпадает с `free()`.

Обратите внимание на то, какими были бы последствия, если бы в данном примере не использовался вызов `free()`. Программа никогда бы не возвращала память системе, и, что еще хуже, она бы теряла единственную ссылку на память — указатель `s`, что делало бы совершенно невозможным дальнейший доступ к памяти. Мы называем этот тип программной ошибки *утечкой памяти* (*memory leak*). Утечка памяти и похожие ошибки использования динамической памяти наиболее распространены и, к сожалению, являются наиболее вредоносными в программировании на С. Так как язык С полностью перекладывает ответственность за управление памятью на программиста, программисты должны быть очень внимательными и осторожными, используя любые типы выделения памяти.

Еще одна часто встречающаяся ошибка в программировании на С – это использование после освобождения. Эта ситуация возникает, когда блок памяти освобождается, а затем к нему снова осуществляется доступ. После того как на блоке памяти выполняется вызов `free()`, программа больше никогда не должна снова обращаться к его содержимому. Программисты должны быть особенно аккуратны, следя за *повисшими указателями* (*dangling pointer*) и не равными `NULL` указателями, которые тем не менее указывают на недействительные блоки памяти. Два широко применяемых инструмента, которые помогают вам в этом, – это Electric Fence и valgrind¹.

Выравнивание

Выравнивание (*alignment*) данных относится к связи между адресом данных и фрагментами памяти, которыми оперирует аппаратное обеспечение. Переменная, находящаяся по адресу в памяти, кратному ее размеру, называется *естественно выровненной* (*naturally aligned*). Например, 32-разрядная переменная естественно выровнена, если ее адрес в памяти кратен четырем, то есть если два младших бита равны нулю. Таким образом, переменная, принадлежащая типу размером 2^n байт, должна обладать адресом с равными нулями наименее значимыми битами.

Правила, относящиеся к выравниванию, основаны на принципах работы аппаратного обеспечения. Некоторые машинные архитектуры накладывают очень строгие требования на выравнивание данных. В некоторых системах загрузка не выровненных данных приводит к процессорному прерыванию. В других системах можно безопасно обращаться к не выровненным данным, но следствием этого становится падение производительности. При написании переносимого кода следует избегать проблем выравнивания, и все типы должны быть естественно выровненными.

Выделение выровненной памяти

Большей частью компилятор и библиотека С прозрачно обрабатывают вопросы выравнивания. Стандарт POSIX указывает, что память, возвращаемая системными вызовами `malloc()`, `calloc()` и `realloc()`, должна быть правильно выровнена для использования с любыми стандартными типами С. В Linux эти функции всегда возвращают память, выровненную по 8-байтовой границе в 32-разрядных системах и по 16-байтовой границе в 64-разрядных системах.

Иногда программистам бывает нужно выравнивать динамическую память по более высокой границе, например по странице. Хотя причины этого могут быть разными, наиболее распространенная среди них – это необходимость правильно выравнивать буфера, применяемые для прямого блочного ввода-вывода или другого обмена данными между программным и аппаратным обеспечением.

¹ См. <http://perens.com/FreeSoftware/ElectricFence> и <http://valgrind.org> соответственно.

Для этой цели в POSIX 1003.1d предусматривается функция с именем `posix_memalign()`:

```
/* одна или другая – достаточно любой из них */
#define _XOPEN_SOURCE 600
#define _GNU_SOURCE

#include <stdlib.h>

int posix_memalign (void **memptr,
                    size_t alignment,
                    size_t size);
```

Успешный вызов `posix_memalign()` выделяет `size` байтов динамической памяти, гарантируя, что она выровнена по адресу памяти, кратному `alignment`. Значение параметра `alignment` должно быть степенью двойки и кратным размеру указателя `void`. Адрес выделенной памяти помещается в параметр `memptr`, а вызов возвращает значение 0.

В случае ошибки память не выделяется, значение `memptr` не определяется, а вызов возвращает один из следующих кодов ошибки:

`EINVAL`

Параметр `alignment` не является степенью двойки или не кратен размеру указателя `void`.

`ENOMEM`

Недостаточно памяти для выполнения запрошенного выделения памяти.

Обратите внимание, что переменная `errno` не устанавливается – функция напрямую возвращает эти ошибки.

Память, получаемая при помощи `posix_memalign()`, освобождается вызовом `free()`. Использовать эту функцию просто:

```
char *buf;
int ret;

/* выделение 1 Кбайт памяти по 256-байтовой границе */
ret = posix_memalign (&buf, 256, 1024);
if (ret) {
    fprintf (stderr, "posix_memalign: %s\n",
            strerror (ret));
    return -1;
}
/* использование buf... */

free (buf);
```

Старые интерфейсы

До того как в стандарте POSIX был определен вызов `posix_memalign()`, в BSD и SunOS использовались следующие интерфейсы соответственно:

```
#include <malloc.h>

void * valloc (size_t size);
void * memalign (size_t boundary, size_t size);
```

Функция `valloc()` работает аналогично `malloc()`, за исключением того, что выделяемая память выравнивается по границе страницы. В главе 4 говорилось, что размер страницы в конкретной системе легко узнать при помощи `getpagesize()`.

Функция `memalign()` похожа на нее, но выделяемая память выравнивается по границе в `boundary` байтов; значение этого параметра должно быть кратным двойке. В следующем примере обе функции возвращают блок памяти, достаточный для записи структуры `ship`, выровненный по границе страницы:

```
struct ship *pirate, *hms;
```

```
pirate = valloc (sizeof (struct ship));

if (!pirate) {
    perror ("valloc");
    return -1;
}

hms = memalign (getpagesize (), sizeof (struct ship));
if (!hms) {
    perror ("memalign");
    free (pirate);
    return -1;
}

/* использование pirate и hms... */

free (hms);
free (pirate);
```

В Linux память, полученная при помощи любой из этих функций, освобождается системным вызовом `free()`. В других системах Unix это может быть не так; некоторые системы не предлагают никакого механизма для безопасного освобождения памяти, выделенной этими функциями. В программах, в которых важна переносимость в другие системы, у программиста может не быть другого выбора, кроме как вообще не освобождать память, полученную через эти интерфейсы!

Программисты, пишущие для Linux, должны использовать эти две функции только в целях обеспечения переносимости на старые системы; `posix_memalign()` намного превосходит их. Все три перечисленных интерфейса нужны только в том случае, если требуется выравнивание по более высокой границе, чем дает `malloc()`.

Прочие вопросы выравнивания

Вопросы выравнивания не ограничиваются естественным выравниванием стандартных типов и выделением динамической памяти. Например, нестандартные и сложные типы предъявляют более запутанные требования, чем стандартные типы. Помимо этого, вопросы выравнивания становятся еще более важными, когда речь заходит о присвоении значений между указателями различных типов и использовании приведения типов.

Нестандартные типы

Нестандартные типы и комплексные типы данных накладывают на выравнивание дополнительные ограничения, помимо простого требования к естественному выравниванию. Существует четыре полезных правила:

- требование к выравниванию у структуры в целом совпадает с требованием к выравниванию, накладываемым наибольшим из ее составляющих типов. Например, если наибольший тип структуры — это 32-битное целое значение, выровненное по 4-байтовой границе, то структура также должна быть выровнена, по крайней мере, по 4-байтовой границе;
- со структурами также связана необходимость забивки, которая используется для того, чтобы гарантировать, что каждый из составляющих типов правильно выровнен согласно требованию самого этого типа. Таким образом, если за значением типа `char` (вероятно, характеризующимся однобайтовым выравниванием) следует значение типа `int` (которое, вероятно, выравнивается по четырехбайтовой границе), то компилятор вставляет три байта забивки между этими двумя типами, чтобы гарантировать выравнивание значения типа `int` по 4-байтовой границе. Иногда программисты выстраивают члены структуры в определенном порядке — например, по убыванию размера, — чтобы минимизировать пространство, «расходящееся попусту» из-за забивки. Параметр GCC `-Wpadded` помогает добиваться этого, а также генерирует предупреждение, когда компилятор вставляет неявную забивку;
- требование к выравниванию у объединения — это требование, накладываемое наибольшим из объединяемых типов;
- требование к выравниванию у массива совпадает с требованием к выравниванию его базового типа. Таким образом, массивы не накладывают никаких требований, превышающих размер одного экземпляра базового типа. Такое поведение обеспечивает естественное выравнивание всех членов массива.

Игра с указателями

Поскольку компилятор прозрачно обрабатывает большинство требований к выравниванию, он также прилагает определенные усилия, чтобы сообщать о возможных проблемах. Наиболее часто трудности возникают при работе с указателями и приведением типов.

Доступ к данным через указатель, приведенный от типа с выравниванием по меньшей границе к типу с выравниванием по более высокой границе, может вызывать загрузку процессором данных, не выровненных согласно требованию более крупного типа. Например, в следующем фрагменте кода присваивание значения с переменной `badnews` приводит к попытке считать с как значение типа `unsigned long`:

```
char greeting[] = "Ahoy Matey";
char *c = greeting[1];
unsigned long badnews = *(unsigned long *) c;
```

Значение типа `unsigned long`, вероятнее всего, выровнено по 4- или 8-байтовой границе; `c`, скорее всего, находится в одном байте от той же границе. Следо-

вательно, загрузка с при приведении типов вызывает нарушение выравнивания. В зависимости от архитектуры, это может приводить как к небольшому падению производительности, так и к аварийному сбою программы. В машинных архитектурах, которые распознают, но не могут правильно обрабатывать нарушения выравнивания, ядро отправляет процессу-виновнику сигнал SIGBUS, который прерывает процесс. Мы поговорим о сигналах в главе 9.

Примеры, подобные этому, встречаются довольно часто. Конечно же, в реальном мире они не выглядят так же глупо, но, с другой стороны, они и не настолько очевидны.

Управление сегментом данных

Системы Unix исторически предоставляют интерфейсы для управления сегментом данных напрямую. Однако большинство программ не применяют их, так как `malloc()` и прочие схемы выделения памяти гораздо проще в использовании и обладают более широкими возможностями. Я хочу упомянуть об этих интерфейсах здесь для того, чтобы удовлетворить любопытных читателей, а также для того, чтобы помочь редким программистам, желающим реализовать собственный механизм выделения на основе куч:

```
#include <unistd.h>

int brk (void *end);
void * sbrk (intptr_t increment);
```

Эти функции унаследовали свои имена от систем Unix старой школы, в которых куча и стек жили в одном сегменте. Выделение динамической памяти в куче выполнялось вверх, начиная со дна сегмента; стек рос вниз с верхушки сегмента, навстречу куче. Демаркационная линия, разграничающая их, называлась *остановом* (*break*) или *точкой останова* (*break point*). В современных системах, в которых сегмент данных живет в собственном отображении в памяти, мы продолжаем называть конечный адрес отображения точкой останова.

Вызов `brk()` устанавливает точку останова (конец сегмента данных) в адрес, указанный при помощи параметра `end`. В случае успеха он возвращает 0. В случае ошибки вызов возвращает значение -1 и присваивает переменной `errno` код ENOMEM.

Вызов `sbrk()` увеличивает адрес конца сегмента данных на значение `increment`, которое может быть положительным или отрицательным. `Sbrk()` возвращает новую точку останова. Таким образом, если передать вызову в качестве параметра `increment` значение 0, то он вернет текущую точку останова:

```
printf ("Текущая точка останова - %p\n". sbrk (0));
```

Ни в POSIX, ни в стандарте C намеренно не определяется ни одна из этих функций. Практически все системы Unix, однако, поддерживают одну или обе. В переносимых программах следует придерживаться только стандартизованных интерфейсов.

Анонимные отображения в памяти

Выделение памяти в glibc использует сегмент данных и отображения в памяти. Классический способ реализации системного вызова `malloc()` — поделить сегмент данных на последовательность разделов разных размеров, представляющих собой степени двойки, и удовлетворять запросы на выделение памяти, возвращая разделы, размеры которых находятся ближе всего к запрашиваемым. Для того чтобы освободить раздел, нужно всего лишь пометить его как «свободный». Если свободны соседние размеры, то их можно объединить в один более крупный раздел. Если верхушка кучи совершенно свободна, то система может применить `brk()` для того, чтобы опустить точку останова, сжимая кучу, и вернуть память ядру.

Этот алгоритм называется *схемой дружеского выделения памяти* (buddy memory allocation scheme). Он обладает такими преимуществами, как скорость и простота, но его недостаток — это два типа фрагментации. *Внутренняя фрагментация* (internal fragmentation) происходит в том случае, когда для удовлетворения запроса на выделение памяти используется больше памяти, чем фактически было запрошено. Результатом является неэффективное расходование доступной памяти. *Внешняя фрагментация* (external fragmentation) происходит в том случае, когда для удовлетворения запроса имеется достаточно свободной памяти, но она разбита на два или более несмежных фрагментов. Это может приводить к неэффективному расходованию памяти (из-за использования более крупных и менее подходящих по размеру блоков) или ошибкам выделения памяти (если альтернативного блока нет).

Помимо этого, данная схема позволяет одному фрагменту выделенной памяти «закреплять» другой, не давая glibc возвращать освобожденную память ядру. Представьте, что выделено два блока памяти, А и В. Блок А расположен прямо у точки останова, а блок В находится прямо под блоком А. Даже если программа освободит блок В, glibc не сможет изменить точку останова до тех пор, пока также не будет освобожден блок А. Таким образом, долгоживущие выделенные блоки могут закреплять все прочие выделенные фрагменты памяти.

Это не всегда становится проблемой, потому что glibc все равно в плановом порядке не возвращает память системе¹. В целом, куча не сжимается после каждого освобождения. Вместо этого glibc продолжает удерживать освобожденную память для последующего выделения. Только когда размер кучи значительно превышает объем выделенной памяти, glibc, наконец, уменьшает сегмент данных. Большой выделенный блок, однако, может мешать сжатию сегмента.

Следовательно, glibc не использует кучу для выделения больших блоков. Вместо этого glibc создает *анонимное отображение в памяти* (anonymous memory mapping), при помощи которого и удовлетворяет запрос на выделение памяти. Анонимные отображения в памяти похожи на файловые отображения,

¹ Glibc также использует значительно более сложный алгоритм выделения памяти, чем простая схема дружеского выделения, называемый алгоритмом арены (arena algorithm).

о которых мы говорили в главе 4, но они не поддерживаются каким-то файлом — отсюда прозвище «анонимные». В действительности, анонимное отображение в памяти — это просто большой заполненный нулями блок памяти, готовый к использованию. Его можно представлять себе как новую кучу, служащую только одному выделению памяти. Так как подобные отображения находятся за пределами кучи, они не увеличивают фрагментацию сегмента данных.

У выделения памяти при помощи анонимных отображений несколько преимуществ:

- никаких проблем с фрагментацией. Если программе больше не требуется анонимное отображение в памяти, оно просто отсоединяется, а память немедленно возвращается системе;
- размер анонимных отображений в памяти можно менять, а также модифицировать для них разрешения. Они могут получать советы, так же как и обычные отображения (см. главу 4);
- каждый выделенный блок памяти существует в отдельном отображении. Нет нужды поддерживать глобальную кучу.

У использования анонимных отображений в памяти по сравнению с кучей также есть два недостатка:

- размер каждого отображения в памяти является целочисленным кратным системного размера страницы. Следовательно, в выделенных блоках, которые не полностью заполняются страницами, есть «потерянное пространство». Это является проблемой для небольших блоков, когда потерянное пространство оказывается довольно большим по сравнению с размером самого выделения;
- создание нового отображения в памяти оказывает более сильную нагрузку на систему, чем возвращение памяти из кучи, которое может вообще не требовать никакого взаимодействия с ядром. Чем меньше выделение, тем правдивее это утверждение.

Исходя из этих за и против, вызов `malloc()` из библиотеки `glibc` использует сегмент данных для удовлетворения запросов на выделение небольших блоков памяти и анонимные отображения в памяти для выделения крупных блоков. Порог между ними можно настраивать (см. раздел «Расширенное выделение памяти» далее в этой главе), и он может быть разным в разных выпусках `glibc`. В настоящий момент порог равен 128 Кбайт: блоки, меньшие или равные 128 Кбайт, выделяются из кучи, тогда как блоки большего размера выделяются при помощи анонимных отображений в памяти.

Создание анонимных отображений в памяти

Например, вы желаете для определенного выделения принудительно включить использование отображений в памяти с применением кучи или пишете собственную систему выделения памяти. В обоих случаях в Linux очень просто вручную создать собственное анонимное отображение в памяти. В главе 4

рассказывалось, что системный вызов `mmap()` создает, а системный вызов `munmap()` разрушает отображение в памяти:

```
include <sys/mman.h>

void * mmap (void *start,
             size_t length,
             int prot,
             int flags,
             int fd,
             off_t offset);

int munmap (void *start, size_t length);
```

В действительности, создать анонимное отображение в памяти проще, чем создать отображение для файла, так как необходимости открывать файл и управлять им нет. Основное различие между двумя типами отображений заключается в присутствии специального флага, указывающего, что отображение является анонимным.

Рассмотрим пример:

```
void *p;

p = mmap (NULL, /* неважно где */
          512 * 1024, /* 512 Кбайт */
          PROT_READ | PROT_WRITE, /* чтение/запись */
          MAP_ANONYMOUS | MAP_PRIVATE, /* анонимное, закрытое */
          -1, /* дескриптор файла (игнорируется) */
          0); /* смещение (игнорируется) */

if (p == MAP_FAILED)
    perror ("mmap");
else
    /* p указывает на 512 Кбайт анонимной памяти... */
```

Для большинства анонимных отображений используются параметры `mmap()`, как в этом примере, за исключением того, что вы должны передавать другой размер отображения (в байтах), в зависимости от того, какой вам необходим. Прочие параметры перечислены далее:

- первому параметру, `start`, присваивается значение `NULL`, указывающее, что анонимное отображение может начинаться в любом месте памяти, как пожелает ядро. Можно также передавать и не равное `NULL` значение, если оно выровнено по странице, но это ограничивает переносимость программы. Очень редко программе действительно бывает важно, где в памяти находятся отображения!
- при помощи параметра `prot` обычно устанавливаются оба бита, `PROT_READ` и `PROT_WRITE`, что делает отображение доступным для чтения и записи. Пустое отображение не имеет никакого смысла, если вы не можете читать его и записывать в него данные. С другой стороны, исполнение кода из анонимного отображения — довольно опасная и редко необходимая задача, которая потенциально создает угрозу безопасности;

- при помощи параметра `flags` устанавливается бит `MAP_ANONYMOUS`, что делает данное отображение анонимным, а также бит `MAP_PRIVATE`, чтобы сделать отображение закрытым;
- параметры `fd` и `offset` игнорируются, когда установлен бит `MAP_ANONYMOUS`. Однако в некоторых старых системах ожидается, что значение `fd` будет равно `-1`, поэтому, если для вас важна переносимость программы, лучше передавать его.

Память, получаемая при помощи анонимного отображения, выглядит так же, как и память, получаемая при помощи кучи. Преимущество выделения памяти через анонимные отображения состоит в том, что страницы заранее заполняются нулями. Эта операция ничего не стоит, потому что ядро отображает анонимные страницы приложения на заполненную нулями страницу путем копирования при записи. Таким образом, к возвращаемой памяти не нужно применять вызов `memset()`. Действительно, это преимущество использования `calloc()` по сравнению с `malloc()`, за которым следует `memset()`: glibc знает, что анонимные отображения уже обнулены и что вызов `calloc()`, выполненный с использованием отображения, не требует явного обнуления.

Системный вызов `munmap()` освобождает анонимное отображение, возвращая выделенную память ядру:

```
int ret;  
  
/* мы закончили использовать р, поэтому можно вернуть  
 * отображение размером 512 Кбайт */  
ret = munmap (р, 512 * 1024);  
if (ret)  
    perror ("munmap");
```

ПРИМЕЧАНИЕ

Обсуждение вызовов `mmap()`, `munmap()` и отображений в целом вы найдете в главе 4.

Отображение `/dev/zero`

В прочих системах Unix, таких, как BSD, флаг `MAP_ANONYMOUS` отсутствует. Вместо этого аналогичное решение реализуется путем отображения специального файла устройства, `/dev/zero`. Этот файл устройства обеспечивает семантику, идентичную анонимной памяти. Отображение содержит страницы, копируемые при записи и состоящие из одних нулей; поведение аналогично принципу работы анонимной памяти.

В Linux всегда существовало устройство `/dev/zero` и предоставлялась возможность отображать этот файл и получать заполненную нулями память. И действительно, до появления флага `MAP_ANONYMOUS` программисты Linux применяли данный подход. Для обеспечения обратной совместимости с более старыми версиями Linux и переносимости на прочие системы Unix разработчики все так же могут продолжать отображать `/dev/zero` вместо создания анонимных отображений. Это совершенно не отличается от отображения любого другого файла:

```

void *p;
int fd;

/* открыть /dev/zero для чтения и записи */
fd = open ("/dev/zero", O_RDWR);
if (fd < 0) {
    perror ("open");
    return -1;
}

/* отобразить [0..размер страницы) файла /dev/zero */
p = mmap (NULL, /* неважно куда */
          getpagesize (),
          PROT_READ | PROT_WRITE,
          MAP_PRIVATE,
          fd,
          0); /* отобразить 1 страницу */
          /* отображение доступно
           * для чтения/записи */
          /* закрытое отображение */
          /* отобразить /dev/zero */
          /* без смещения */

if (p == MAP_FAILED) {
    perror ("mmap");
    if (close (fd))
        perror ("close");
    return -1;
}

/* закрыть /dev/zero, он нам больше не нужен */
if (close (fd))
    perror ("close");

/* p указывает на одну страницу памяти, используйте ее */

```

Отображение в памяти, выполненное таким способом, удаляется при помощи `munmap()`.

Данный подход включает нагрузку, создаваемую дополнительными системными вызовами для открытия и закрытия файла устройства. Таким образом, использование анонимной памяти — это более быстрое решение.

Расширенное выделение памяти

Многие операции выделения памяти, о которых шла речь в этой главе, ограничены и управляются параметрами ядра, которые позволено менять программисту. Для того чтобы настраивать такие параметры, используйте вызов `mallopt()`:

```
#include <malloc.h>
```

```
int mallopt (int param, int value);
```

Вызов `mallopt()` устанавливает для параметра, относящегося к управлению памятью и указанного при помощи аргумента `param`, значение, равное `value`. В случае успеха вызов возвращает ненулевое значение; в случае ошибки он возвращает 0. Обратите внимание, что `mallopt()` не устанавливает переменную `errno`. Обычно этот вызов всегда завершается успешно, поэтому не надейтесь,

что вам удастся извлечь какую бы то ни было полезную информацию из возвращаемого значения.

В настоящее время Linux поддерживает шесть значений `param`; все они определены в файле заголовка `<malloc.h>`:

`M_CHECK_ACTION`

Значение переменной среды `MALLOC_CHECK_` (подробнее о ней в следующем разделе).

`M_MMAP_MAX`

Максимальное число отображений, которое система создает для удовлетворения запросов на выделение динамической памяти. Когда достигается этот лимит, для всех выделений начинает использоваться только сегмент данных, пока не освободится хотя бы одно отображение. Значение 0 полностью отключает использование анонимных отображений в качестве базы для выделения динамической памяти.

`M_MMAP_THRESHOLD`

Пороговое значение размера памяти (измеряется в байтах), после которого запросы на выделение памяти начинают удовлетворяться при помощи анонимных отображений, а не с использованием сегмента данных. Обратите внимание, что по усмотрению системы запросы на выделение объема памяти, не достигающего порогового значения, также могут обрабатываться с помощью анонимных отображений. Значение 0 включает использование анонимных отображений для всех выделений памяти, фактически запрещая использование сегмента данных для выделения динамической памяти.

`M_MXFAST`

Максимальный размер (в байтах) быстрой корзины. *Быстрые корзины* (fast bin) – это специальные фрагменты памяти в куче, которые никогда не объединяются с соседними фрагментами и никогда не возвращаются системе, обеспечивая очень быстрое выделение памяти за счет увеличения фрагментации. Значение 0 отключает использование быстрых корзин.

`M_TOP_PAD`

Объем забивки (в байтах), применяемой при корректировке размера сегмента данных. Когда библиотека `glibc` выполняет вызов `brk()` для увеличения размера сегмента данных, она может запрашивать больше памяти, чем необходимо, в надежде снизить необходимость в дополнительном использовании `brk()` в ближайшем будущем. Аналогично, когда `glibc` сжимает сегмент данных, она удерживает дополнительную память, отдавая немного меньше, чем обычно. Эти дополнительные байты и представляют *забивку* (padding). Значение 0 полностью отключает использование забивки.

`M_TRIM_THRESHOLD`

Минимальный размер свободной памяти (в байтах), разрешенный наверху сегмента данных. Если объем памяти падает ниже данного порога, `glibc` вызывает `brk()`, чтобы отдать память обратно ядру.

Стандарт XPG, в котором в свободной форме определяется вызов `mallopt()`, включает еще три параметра: `M_GRAIN`, `M_KEEP` и `M_NLBLKS`. В Linux эти параметры определяются, но их установка ни на что не влияет. Полный список допустимых параметров, их значения по умолчанию и диапазоны допустимых значений рассматриваются в табл. 8.1.

Таблица 8.1. Параметры системного вызова `mallopt()`

Параметр	Происхождение	Значение по умолчанию	Допустимые значения	Специальные значения
<code>M_CHECK_ACTION</code>	Уникальный для Linux	0	0–2	
<code>M_GRAIN</code>	Стандарт XPG	Не поддерживается в Linux	≥ 0	
<code>M_KEEP</code>	Стандарт XPG	Не поддерживается в Linux	≥ 0	
<code>M_MMAP_MAX</code>	Уникальный для Linux	60×1024	≥ 0	0 отключает использование <code>mmap()</code>
<code>M_MMAP_THRESHOLD</code>	Уникальный для Linux	128×1024	≥ 0	0 отключает использование кучи
<code>M_MXFAST</code>	Стандарт XPG	64	0–80	0 отключает быстрые корзины
<code>M_NLBLKS</code>	Стандарт XPG	Не поддерживается в Linux	≥ 0	
<code>M_TOP_PAD</code>	Уникальный для Linux	0	≥ 0	0 отключает забивку

В программах любые вызовы `mallopt()` необходимо совершать до первого вызова `malloc()` или любого другого интерфейса выделения памяти. Использовать этот системный вызов просто:

```
int ret;
```

```
/* использовать mmap() для всех запросов на выделение памяти
/* объемом больше 64 Кбайт */
ret = mallopt(M_MMAP_THRESHOLD, 64 * 1024);
if (!ret)
    fprintf(stderr, "mallopt failed!\n");
```

Тонкая настройка при помощи `malloc_usable_size()` и `malloc_trim()`

В Linux есть пара функций, предлагающих низкоуровневый контроль над системой выделения памяти библиотеки glibc. Первая из этих функций позволяет программе спрашивать, сколько полезных байтов содержится в определенной выделенной области памяти:

```
#include <malloc.h>
```

```
size_t malloc_usable_size (void *ptr);
```

Успешный вызов `malloc_usable_size()` возвращает фактический размер выделенной памяти для фрагмента памяти, на который указывает `ptr`. Так как glibc иногда округляет выделенные области, чтобы они помещались в существующие фрагменты или анонимные отображения, полезное пространство в выделенной области может быть больше запрошенного. Выделенная область никогда не будет меньше запрошеннной. Вот пример использования функции:

```
size_t len = 21;  
size_t size;  
char *buf,
```

```
buf = malloc (len);  
if (!buf) {  
    perror ("malloc");  
    return -1;  
}
```

```
size = malloc_usable_size (buf);
```

```
/* фактически, можно использовать size байтов буфера buf... */
```

Вторая из этих функций позволяет программе заставлять glibc возвращать ядру всю память, для которой возможно немедленное освобождение:

```
#include <malloc.h>
```

```
int malloc_trim (size_t padding);
```

Успешный вызов `malloc_trim()` сжимает сегмент данных настолько, насколько это возможно, за исключением `padding` байтов, которые резервируются. После этого он возвращает значение 1. В случае ошибки вызов возвращает 0. Обычно glibc выполняет подобные сжатия сегмента данных автоматически, как только раз мер допускающей освобождение памяти достигает `M_TRIM_THRESHOLD` байтов. Размер забивки у нее равен `M_TOP_PAD`.

Вы практически никогда не будете применять эти две функции, разве что для отладки и в образовательных целях. Они не переносимы и предоставляют использующей их программе низкоуровневые подробности относительно системы выделения памяти библиотеки glibc.

Отладка выделения памяти

В программах можно устанавливать переменную среды `MALLOC_CHECK_`, чтобы использовать расширенные возможности отладки в подсистеме памяти. Когда применяются дополнительные отладочные проверки, выделение памяти выполняется менее эффективно, но лишняя нагрузка зачастую оправдывает себя на этапе отладки во время разработки приложения.

Так как переменная среды контролирует отладку, необходимости заново компилировать программу нет. Например, можно выполнить команду, аналогичную следующей:

```
$ MALLOC_CHECK_=1 ./rudder
```

Если значение `MALLOC_CHECK_` равно 0, то подсистема памяти бесшумно игнорирует любые ошибки. Если значение переменной равно 1, то на выход `stderr` выводится информационное сообщение. Если оно равно 2, то программа немедленно завершается при помощи вызова `abort()`. Так как `MALLOC_CHECK_` меняет поведение выполняющейся программы, программы `setuid` игнорируют данную переменную.

Получение статистики

В Linux присутствует функция `mallinfo()`, позволяющая получать статистику, касающуюся системы выделения памяти:

```
#include <malloc.h>
```

```
struct mallinfo mallinfo (void);
```

Вызов `mallinfo()` возвращает статистику в структуре `mallinfo`. Эта структура возвращается по значению, а не по указателю. Ее содержимое также определено в файле заголовка `<malloc.h>`:

```
/* all sizes in bytes */
struct mallinfo {
    int arena;      /* размер сегмента данных, используемого malloc */
    int ordblks;   /* число свободных фрагментов */
    int smblks;    /* число быстрых корзин */
    int hblkts;    /* число анонимных отображений */
    int hblkhd;    /* размер анонимных отображений */
    int usmblks;   /* максимальный общий размер выделенной памяти */
    int fsmblks;   /* размер доступных быстрых корзин */
    int wordblks;  /* размер общего выделенного пространства */
    int fordblks;  /* размер доступных фрагментов */
    int keepcost;  /* размер пространства, которое можно отрезать */
};
```

Использовать этот вызов просто:

```
struct mallinfo m;
m = mallinfo ();
printf ("свободные фрагменты %d\n", m.ordblks);
```

В Linux также существует функция `malloc_stats()`, которая выводит статистику, относящуюся к памяти, на выход `stderr`:

```
#include <malloc.h>
```

```
void malloc_stats (void);
```

Если применить `malloc_stats()` в программе, интенсивно использующей память, то можно получить довольно большие результаты:

```
Arena 0:  
system bytes      = 865939456  
in use bytes      = 851988200  
Total (incl. mmap):  
System bytes      = 3216519168  
in use bytes      = 3202567912  
max mmap regions = 65536  
max mmap bytes   = 2350579712
```

Выделение памяти на основе стека

Пока что все механизмы для выделения динамической памяти, которые мы изучили, применяли для получения динамической памяти кучу или отображения в памяти. Это логично, так как куча и отображения в памяти, несомненно, являются динамическими по своей природе. Другая распространенная конструкция в адресном пространстве программы, стек, — это место, где живут *автоматические переменные* (*automatic variables*) программы.

Нет никаких противопоказаний для использования стека для выделения динамической памяти. До тех пор пока выделенная область не переполняет стек, этот подход реализуется легко и работает весьма эффективно. Для выделения динамической памяти в стеке применяется системный вызов `alloca()`:

```
#include <alloca.h>
```

```
void * alloca (size_t size);
```

В случае успеха вызов `alloca()` возвращает указатель на `size` байтов памяти. Эта память живет в стеке и автоматически освобождается, когда вызывающая функция возвращает результат. Некоторые реализации в случае ошибки возвращают значение `NULL`, но большинство реализаций `alloca()` не могут создавать ошибки или не умеют сообщать о сбоях. Сбой проявляется себя в виде переполнения стека.

Использование данного вызова идентично `malloc()`, но вам не нужно (и, в действительности, вы не должны) освобождать выделяемую память. Далее приведен пример функции, которая открывает указанный файл в конфигурационном каталоге системы, вероятнее всего, `/etc`, но который ищется во время компиляции, что обеспечивает переносимость программы. Функция должна выделить новый буфер, скопировать имя системного конфигурационного каталога в буфер, а затем склеить этот буфер с указанным именем файла:

```
int open_sysconf (const char *file, int flags, int mode)
{
    const char *etc = SYSCONF_DIR; /* "/etc/" */
    char *name;

    name = alloca (strlen (etc) + strlen (file) + 1);
    strcpy (name, etc);
    strcat (name, file);

    return open (name, flags, mode);
}
```

После возвращения память, выделенная при помощи `alloca()`, автоматически освобождается, когда стек отматывается обратно до вызывающей функции. Это означает, что вы не можете использовать эту память после того, как функция, которая вызывает `alloca()`, возвращает результат! Однако так как «убирать за собой» при помощи вызова `free()` не требуется, результирующий код выглядит несколько чище. Далее показана та же функция, реализованная с помощью системного вызова `malloc()`:

```
int open_sysconf (const char *file, int flags, int mode)
{
    const char *etc = SYSCONF_DIR; /* "/etc/" */
    char *name;
    int fd;

    name = malloc (strlen (etc) + strlen (file) + 1);
    if (!name) {
        perror ("malloc");
        return -1;
    }

    strcpy (name, etc);
    strcat (name, file);
    fd = open (name, flags, mode);
    free (name);

    return fd;
}
```

Обратите внимание, что не следует использовать выделенную при помощи `alloca()` память в параметрах при вызове функций, так как выделенная память при этом оказывается в середине пространства стека, зарезервированного для параметров функции. Например, следующее запрещено:

```
/* НЕ ДЕЛАЙТЕ ТАК! */
ret = foo (x, alloca (10));
```

У интерфейса `alloca()` изменчивая история. Во многих системах он работал плохо или становился причиной неопределенного поведения. В системах с небольшим стеком фиксированного размера использование `alloca()` приводило к переполнению стека, что убивало программы. В еще каких-то системах интерфейс `alloca()` вообще не существовал. Со временем полные ошибок и противоречий реализаций создали `alloca()` плохую репутацию.

Таким образом, если ваша программа должна быть переносимой, не используйте `alloca()`. В Linux, однако, `alloca()` — это крайне полезное, но недостаточно широко используемое средство. Этот системный вызов работает превосходно — во многих архитектурах выделение памяти при помощи `alloca()` включает всего лишь увеличение значения указателя стека — и превышает по эффективности вызов `malloc()`. Если речь идет о выделении небольших объемов памяти в коде, предназначенном только для Linux, то `alloca()` может давать заметный выигрыш в производительности.

Дублирование строк в стеке

Широко распространенный вариант использования `alloca()` — для временного дублирования строки. Например:

```
/* нам нужно скопировать song */
char *dup;
dup = alloca (strlen (song) + 1);
strcpy (dup, song);

/* использование dup... */
return; /* dup автоматически освобождается */
```

Из-за частого возникновения подобной задачи, и так как `alloca()` обеспечивает большое преимущество в скорости исполнения, в системах Linux также предоставляются два варианта вызова `strdup()`, дублирующие указанную строку в стеке:

```
#define _GNU_SOURCE
#include <string.h>

char * strdupa (const char *s);
char * strndupa (const char *s, size_t n);
```

Вызов `strdupa()` возвращает дубликат строки `s`. Вызов `strndupa()` дублирует до `n` символов строки `s`. Если `s` больше `n`, то дублирование останавливается на `n`-м символе и функция добавляет нулевой байт. Эти функции обладают теми же преимуществами, что и `alloca()`. Дублированная строка автоматически освобождается, когда вызывающая функция возвращает результат.

В POSIX не определяются функции `alloca()`, `strdupa()` и `strndupa()`, и в других операционных системах они появляются крайне редко. Если для вас важна переносимость программы, то настоятельно не рекомендуется применять эти функции. В Linux, однако, `alloca()` и ее друзья работают довольно хорошо и могут давать значительный выигрыш в производительности, заменяя сложные ритуалы выделения динамической памяти простой корректировкой указателя кадра стека.

Массивы переменной длины

В стандарте C99 были впервые представлены *массивы переменной длины* (variable-length array, VLA), то есть массивы, геометрия которых задается во время выполнения, а не во время компиляции. Язык GNU C поддерживал массивы

переменной длины уже в течение какого-то времени, но стандартизация их в C99 стала дополнительной причиной для того, чтобы использовать эту отличную возможность. VLA позволяют избегать нагрузки, создаваемой выделением динамической памяти, практически так же, как это делает системный вызов `alloca()`.

Используются подобные массивы именно так, как и ожидается:

```
for (i = 0; i < n; ++i) {
    char foo[i + 1];
    /* использование foo... */
}
```

В этом фрагменте `foo` — это массив элементов типа `char` переменной длины $i + 1$. На каждой итерации цикла `foo` динамически создается и автоматически очищается, когда он выходит за пределы области определения. Если бы мы использовали `alloca()` вместо массива переменной длины, то память не освобождалась бы до тех пор, пока функция не возвращала результат. Использование VLA гарантирует освобождение памяти на каждой итерации цикла. Таким образом, при работе с VLA расходуется максимум n байт, тогда как вызов `alloca()` потребовал бы $n*(n+1)/2$ байт.

Применив массив переменной длины, мы можем переписать нашу функцию `open_sysconf()` следующим образом:

```
int open_sysconf (const char *file, int flags, int mode)
{
    const char *etc: = SYSCONF_DIR; /* "/etc/" */
    char name[strlen (etc) + strlen (file) + 1];

    strcpy (name, etc);
    strcat (name, file);

    return open (name, flags, mode);
}
```

Основное различие между `alloca()` и массивами переменной длины состоит в том, что память, полученная первым способом, существует на всем протяжении работы функции, тогда как выделенная вторым способом память существует до тех пор, пока удерживающая переменная не выходит за пределы области определения, что может происходить до того, как текущая функция возвращает результат. Такое поведение может быть в разных ситуациях желательным или нежелательным. В цикле `for`, который мы только что видели, возвращение памяти на каждой итерации цикла снижает общее потребление памяти без каких-либо побочных эффектов (нам не нужно было удерживать лишнюю память). Однако если по какой-то причине нам требуется, чтобы память существовала дольше одной итерации цикла, разумнее было бы использовать системный вызов `alloca()`.

ПРИМЕЧАНИЕ

Смешивание `alloca()` и массивов переменной длины в одной функции может приводить к несвойственному поведению. Избегайте риска и используйте в каждой конкретной функции либо одно либо другой.

Выбор механизма выделения памяти

Миллион параметров выделения памяти, о которых рассказывалось в этой главе, могут заставить программиста недоумевать, какое же решение лучше всего подходит для конкретной задачи. В большинстве ситуаций лучше всего полагаться на системный вызов `malloc()`. В некоторых случаях другой подход может предоставить более подходящий инструмент. В табл. 8.2 приводятся рекомендации по выбору механизма выделения памяти.

Таблица 8.2. Подходы к выделению памяти в Linux

Подход к выделению	Преимущества	Недостатки
<code>malloc()</code>	Простой, удобный, распространенный	Возвращаемая память не обязательно обнуляется
<code>calloc()</code>	Упрощает выделение массивов, обнуляет возвращаемую память	Слишком запутанный интерфейс для остальных задач, кроме выделения массивов
<code>realloc()</code>	Меняет размер существующих выделенных областей памяти	Полезен только для изменения размера существующих выделений
<code>brk()</code> и <code>sbrk()</code>	Обеспечивает глубокий контроль кучи	Слишком низкоуровневый для большинства пользователей
Анонимные отображения в памяти	Простые в использовании, допускают совместное использование, позволяют разработчику регулировать уровень защиты и давать советы; оптимальны для выделения больших областей	Не оптимальны для выделения маленьких областей; <code>malloc()</code> автоматически применяет анонимные отображения в памяти, когда это дает выгоду
<code>posix_memalign()</code>	Выделяет память, выровненную по любой разумной границе	Относительно новый и поэтому переносимость находится под вопросом; слишком серьезный способ для всех ситуаций, кроме тех, когда основным требованием является выравнивание памяти
<code>memalign()</code> и <code>valloc()</code>	Более распространены на других системах Unix, чем <code>posix_memalign()</code>	Не определены в стандарте POSIX, предоставляют меньше возможностей управления выравниванием, чем <code>posix_memalign()</code>

продолжение ↗

Таблица 8.2 (продолжение)

Подход к выделению	Преимущества	Недостатки
alloc()	Очень быстрое выделение, не нужно явно освобождать память; отлично подходит для выделения небольших областей	Не возвращает ошибки, не подходит для выделения больших областей, неправильно работает в некоторых системах Unix
Массивы переменной длины	То же, что и alloc(), но освобождает память, когда массив выходит за пределы области определения, а не когда функция возвращает результат	Полезны только для массивов; вариант освобождения памяти, используемый alloc(), в некоторых ситуациях может быть предпочтительнее; менее распространены в других системах Unix, чем alloc()

Помимо этого, не стоит забывать альтернативу всем этим вариантам: автоматическое и статическое выделение памяти. Выделение автоматических переменных в стеке или глобальных переменных в куче зачастую проще, не требует от программиста, чтобы он использовал указатели, и не заставляет беспокоиться об освобождении памяти.

Манипулирование памятью

Язык С предоставляет семейство функций для манипулирования сырьими байтами памяти. Работа этих функций во многих отношениях аналогична интерфейсам манипулирования строками, таким, как strcmp() и strcpy(), но они полагаются на предоставляемый пользователем размер буфера и не предполагают, что строки завершаются нулями. Обратите внимание, что ни одна из этих функций не возвращает ошибки. Предотвращение появления ошибок остается на совести программиста — передайте функции неправильную область памяти, и единственным вариантом развития событий будет нарушение сегментирования!

Установка байтов

В наборе функций для манипулирования памятью самая распространенная — это, конечно же, memset():

```
#include <string.h>
```

```
void * memset (void *s, int c, size_t n);
```

Вызов memset() устанавливает для *n* байт, начиная с *s*, значение *c* и возвращает *s*. Очень часто эту функцию применяют для обнуления блока памяти:

```
/* обнулить [s,s+256) */
memset (s, '\0', 256);
```

Bzero() — это устаревший и не рекомендуемый для использования интерфейс, который в BSD появился для того, чтобы разработчики могли выполнять ту же задачу. В новом коде следует применять memset(), но в Linux вызов bzero() предоставляется для обеспечения обратной совместимости и переносимости на старые системы:

```
#include <strings.h>

void bzero (void *s, size_t n);
```

Следующий вариант вызова идентичен предыдущему примеру с memset():
bzero (s, 256);

Обратите внимание, что bzero(), а также другие интерфейсы с именами, начинающимися на b, требуют использования файла заголовка `<strings.h>`, а не `<string.h>`.

ПРИМЕЧАНИЕ

Не используйте memset(), если можете использовать calloc()! Избегайте выделения памяти при помощи malloc() и немедленного последующего обнуления через memset(). Хотя результат может быть совершенно аналогичным, применять вместо двух функций один только вызов calloc(), который возвращает обнуленную память, намного лучше. Это не просто уменьшение на единицу количества вызовов — calloc() может получать от ядра уже обнуленную память, и в этом случае вы избегаете необходимости вручную устанавливать в ноль каждый байт, что повышает производительность.

Сравнение байтов

Аналогично strcmp(), memcmp() проверяет два фрагмента памяти на эквивалентность:

```
#include <string.h>
```

```
int memcmp (const void *s1, const void *s2, size_t n);
```

Вызов сравнивает первые n байт блоков s1 и s2 и возвращает 0, если блоки памяти эквивалентны, значение, меньшее нуля, если s1 меньше s2, и значение, большее нуля, если s1 больше s2.

И на этот случай в BSD есть теперь уже устаревший интерфейс, выполняющий ту же задачу:

```
#include <strings.h>
```

```
int bcmp (const void *s1, const void *s2, size_t n);
```

Вызов bcmp() сравнивает первые n байт фрагментов s1 и s2, возвращая 0, если блоки памяти эквивалентны, и ненулевое значение, если они отличаются.

Из-за того что структуры забиваются какими-то данными (см. раздел «Прочие вопросы выравнивания» ранее в этой главе), проверка двух структур на эквивалентность при помощи memcmp() или bcmp() ненадежна. В забивке может находиться не инициализированный мусор, отличающийся в двух экземплярах

структуры, нормальные данные в которых совершенно идентичны. Следовательно, подобный код нельзя назвать безопасным:

```
/* эти две шлюпки (dinghies) идентичны? (НЕПРАВИЛЬНЫЙ КОД) */
int compare_dinghies (struct dinghy *a, struct dinghy *b)
{
    return memcmp (a, b, sizeof (struct dinghy));
```

Вместо этого, если требуется сравнить структуры, нужно последовательно сравнить все элементы структуры, один за другим. Такой подход допускает определенную степень оптимизации, но требует больше усилий, чем ненадежный вариант с `memcmp()`. Вот эквивалентный код:

```
/* эти две шлюпки идентичны? */
int compare_dinghies (struct dinghy *a, struct dinghy *b)
{
    int ret;

    if (a->nr_oars < b->nr_oars)
        return -1;
    if (a->nr_oars > b->nr_oars)
        return 1;

    ret = strcmp (a->boat_name, b->boat_name);
    if (ret)
        return ret;

    /* и так далее для каждого элемента... */
}
```

Перемещение байтов

Вызов `memmove()` копирует первые *n* байт из *src* в *dst*, возвращая *dst*:

```
#include <string.h>
```

```
void * memmove (void *dst, const void *src, size_t n);
```

И снова, BSD предоставляет устаревший интерфейс для выполнения той же задачи:

```
#include <strings.h>
```

```
void bcopy (const void *src, void *dst, size_t n);
```

Обратите внимание, что, хотя обе функции принимают одни и те же параметры, порядок первых двух в `bcopy()` обратный.

И `bcopy()`, и `memmove()` умеют безопасно обрабатывать перекрывающиеся области памяти (скажем, если часть *dst* лежит внутри *src*). Это позволяет, например, перемещать байты памяти вверх или вниз в пределах одного региона. Поскольку такая ситуация довольно редка и так как программист точно знает, если она происходит, в стандарте C определяется вариант вызова `memmove()`, не поддерживающий перекрывающиеся области памяти. Этот вариант потенциально быстрее:

```
#include <string.h>

void * memcp (void *dst, const void *src, size_t n);
```

Данная функция работает аналогично `memmove()`, но области памяти `dst` и `src` не должны перекрываться. Если они перекрываются, то результат не определен.

Еще одна безопасная функция копирования — это `memccpy()`:

```
#include <string.h>

void * memccpy (void *dst, const void *src, int c, size_t n);
```

Функция `memccpy()` работает аналогично `memcp()`, но прекращает копирование, как только обнаруживает байт `c` в пределах первых `n` байт `src`. Вызов возвращает указатель на следующий байт в `dst` после `c` или значение `NULL`, если байт `c` обнаружен не был.

Наконец, можно применять `mempcpy()` для пошагового прохождения памяти:

```
#define __GNU_SOURCE
#include <string.h>

void * mempcpy (void *dst, const void *src, size_t n);
```

Функция `mempcpy()` работает так же, как и `memcp()`, но возвращает указатель на следующий байт после последнего скопированного байта. Это полезно в тех случаях, когда нужно скопировать набор данных в последовательные адреса в памяти, однако дает мало преимуществ, так как всего лишь возвращает значение `dst + n`. Данная функция уникальна для GNU.

Поиск байтов

Функции `memchr()` и `memrchr()` находят указанный байт в блоке памяти:

```
#include <string.h>

void * memchr (const void *s, int c, size_t n);
```

Функция `memchr()` сканирует `n` байт памяти, на которую указывает `s`, в поисках символа `c`, который интерпретируется как `unsigned char`:

```
#define __GNU_SOURCE
#include <string.h>

void * memrchr (const void *s, int c, size_t n);
```

Вызов возвращает указатель на первый байт, совпадающий с `c`, или же значение `NULL`, если байт `c` найден не был.

Функция `memrchr()` аналогична функции `memchr()`, за исключением того, что она выполняет поиск в обратном порядке, начиная с конца области из `n` байт, на которую указывает `s`, а не с начала. В отличие от `memchr()`, `memrchr()` — это расширение GNU, а не часть языка C.

Для более сложных заданий поиска используется функция с ужасным названием `memmem()`, которая ищет в блоке памяти произвольный массив байтов:

```
#define _GNU_SOURCE
#include <string.h>

void * memmem (const void *haystack,
               size_t haystacklen,
               const void *needle,
               size_t needlelen);
```

Функция `memmem()` возвращает указатель на первое вхождение подблока `needle` длиной `needlelen` байтов в блоке памяти `haystack` длиной `haystacklen` байтов. Если функция не находит `needle` в `haystack`, она возвращает значение `NULL`. Данная функция также является расширением GNU.

Прочее манипулирование байтами

Библиотека C в Linux предоставляет интерфейс для тривиального сворачивания байтов данных:

```
#define _GNU_SOURCE
#include <string.h>

void * memfrob (void *s, size_t n);
```

Вызов `memfrob()` затемняет первые `n` байт области памяти, начинающейся в `s`, применяя к каждому байту операцию исключающего ИЛИ (XOR) с числом 42. Вызов возвращает `s`.

Результат работы вызова `memfrob()` можно вернуть к первоначальному состоянию, вызвав `memfrob()` в той же области памяти еще раз. Таким образом, следующий фрагмент кода представляет собой холостую команду по отношению к `secret`:

```
memfrob (memfrob (secret, len), len);
```

Эта функция никоим образом не является правильной (или хотя бы приемлемой) заменой шифрования; ее использование ограничивается тривиальным затемнением строк. Она уникальна для GNU.

Блокирование памяти

В Linux реализуется вызов страниц по требованию (demand paging), это означает, что страницы подкачиваются с диска по необходимости и сбрасываются обратно на диск, когда необходимость в них отпадает. Это позволяет виртуальным адресным пространствам процессов в системе не иметь никаких прямых взаимоотношений с общим объемом физической памяти, так как область подкачки на диске может создавать иллюзию практически бесконечного запаса физической памяти.

Эта подкачка осуществляется прозрачно, и приложениям, в целом, не нужно заботиться (и даже знать) о работе системы постраничной подкачки ядра. Однако существует две ситуации, в которых приложениям может быть необходимо влиять на поведение этой системы:

Детерминизм

Приложениям, на которые накладываются временные ограничения, требуется детерминированное поведение. Если обращение к памяти приводит к ошибке доступа к несуществующей странице — и к дорогостоящим операциям дискового ввода-вывода, — приложение может нарушать поставленные для него временные ограничения. Гарантируя, что требуемые страницы всегда находятся в физической памяти и никогда не сбрасываются на диск, приложение получает возможность обращаться к памяти в уверенности, что ошибки доступа к несуществующей странице не будет, что обеспечивает постоянство, детерминизм и улучшенную производительность.

Безопасность

Если закрытые конфиденциальные сведения находятся в памяти, они могут быть сброшены на диск и сохранены в незашифрованном виде. Например, если закрытый ключ пользователя обычно хранится на диске в зашифрованной форме, незашифрованная копия ключа в памяти может попасть в файл подкачки. В средах с высокой безопасностью такое поведение может быть недопустимым. В приложениях, где это представляет проблему, можно защищивать, чтобы память, включающая ключ, всегда оставалась в физической памяти.

Конечно же, модификация поведения ядра может негативно влиять на общую производительность системы. Детерминизм или безопасность одного приложения улучшится, но, пока его страницы будут заблокированы в памяти, страницы другого приложения вместо этого будут сбрасываться на диск. Ядро, если мы доверяем его конструкции, всегда выбирает оптимальную страницу, чтобы сбросить ее на диск, то есть страницу, которая с наименьшей вероятностью потребуется в ближайшем будущем, так что, когда вы воздействуете на поведение ядра, ему приходится сбрасывать не оптимальные страницы.

Блокировка части адресного пространства

В POSIX 1003.1b-1993 определяются два интерфейса для «блокировки» одной или нескольких страниц в физической памяти, это гарантирует, что они никогда не будут сброшены на диск. Первый из них блокирует указанный интервал адресов:

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

```
int mlock (const void *addr, size_t len);
```

Вызов `mlock()` блокирует в физической памяти виртуальную память, начиная с адреса `addr` и на протяжении `len` байтов. В случае успеха вызов возвращает 0; в случае ошибки вызов возвращает -1 и присваивает переменной `errno` подходящее значение.

Успешный вызов может заблокировать все физические страницы, содержащие адреса `[addr,addr + len)` в памяти. Например, если в вызове указывается только один байт, то вся страница, которой принадлежит этот байт, блокируется в памяти. Стандарт POSIX диктует, что значение `addr` должно быть выровнено

по границе страницы. В Linux это требование не считается обязательным и система бесшумно округляет значение `addr` вниз до ближайшей страницы. Однако в программах, для которых обязательна переносимость на другие системы, необходимо гарантировать, что `addr` находится на границе страницы.

Допустимые коды `errno` включают:

EINVAL

Значение параметра `len` меньше нуля.

ENOMEM

Вызывающий попытался заблокировать больше страниц, чем допускает лимит ресурсов `RLIMIT_MEMLOCK` (см. раздел «Лимиты блокирования» далее).

EPERM

Значение лимита ресурсов `RLIMIT_MEMLOCK` было равно 0, но процесс не обладал характеристикой `CAP_IPC_LOCK` (см. раздел «Лимиты блокирования»).

ПРИМЕЧАНИЕ

Дочерний процесс не наследует заблокированную память после того, как срабатывает вызов `fork()`. Из-за поведения адресных пространств в Linux, основанного на копировании при записи, страницы дочернего процесса фактически блокируются в памяти до тех пор, пока потомок не перестает записывать в них данные.

В качестве примера предположим, что программа удерживает в памяти дешифрованную строку. Процесс может заблокировать страницу, содержащую строку, при помощи такого кода:

```
/* заблокировать secret в памяти */
ret = mlock (secret, strlen (secret));
if (ret)
    perror ("mlock");
```

Блокировка всего адресного пространства

Если процессу нужно заблокировать в физической памяти все свое адресное пространство, то использовать `mlock()` будет слишком неудобно. Для этой задачи, которая часто возникает в приложениях реального времени в POSIX, определяется системный вызов, блокирующий адресное пространство целиком:

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

int mlockall (int flags);
```

Вызов `mlockall()` блокирует все страницы из адресного пространства текущего процесса в физической памяти. Параметр `flags`, представляющий собой объединение следующих двух значений при помощи операции побитового ИЛИ, позволяет управлять поведением вызова:

MCL_CURRENT

Если оно установлено, это значение заставляет `mlockall()` блокировать все текущие отображенные страницы: стек, сегмент данных, отображенные файлы и т. д. — в адресном пространстве процесса.

MCL_FUTURE

Если оно установлено, это значение заставляет `mlockall()` гарантировать, что все страницы, которые будут отображены в адресное пространство в будущем, также будут блокироваться в памяти.

В большинстве приложений используются оба этих значения.

В случае успеха вызов возвращает 0; в случае ошибки он возвращает -1 и присваивает переменной `errno` один из следующих кодов ошибки:

EINVAL

Значение параметра `flags` меньше нуля.

ENOMEM

Вызывающий попытался заблокировать больше страниц, чем допускает лимит ресурсов `RLIMIT_MEMLOCK` (см. далее раздел «Лимиты блокировки»).

EPERM

Значение лимита ресурсов `RLIMIT_MEMLOCK` было равно 0, но процесс не обладал характеристикой `CAP_IPC_LOCK` (см. раздел «Лимиты блокировки»).

Разблокировка памяти

Чтобы разблокировать страницы в физической памяти, снова разрешив ядру сбрасывать их по необходимости на диск, в POSIX стандартизируется еще два интерфейса:

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

```
int munlock (const void *addr, size_t len);
int munlockall (void);
```

Системный вызов `munlock()` разблокирует страницы, начиная с адреса `addr` и на протяжении `len` байтов. Он отменяет эффект `mlock()`. Системный вызов `munlockall()` отменяет эффект `mlockall()`. Оба вызова возвращают 0 в случае успеха и -1 в случае ошибки, присваивая также переменной `errno` одно из следующих значений:

EINVAL

Недопустимое значение параметра `len` (только для `munlock()`).

ENOMEM

Одна из указанных страниц недопустима.

EPERM

Значение лимита ресурсов `RLIMIT_MEMLOCK` было равно 0, но процесс не обладал характеристикой `CAP_IPC_LOCK` (см. следующий раздел, «Лимиты блокирования»).

Блокировки памяти не вкладываются друг в друга. Таким образом, один вызов `munlock()` или `munlockall()` разблокирует заблокированную страницу независимо от того, сколько раз она была заблокирована при помощи `mlock()` или `mlockall()`.

Лимиты блокировки

Так как блокировка памяти может влиять на общую производительность системы, — действительно, если заблокировано слишком много страниц, выделение памяти может быть невозможно, — в Linux определяются пределы относительно того, как много страниц разрешается заблокировать процессу.

Процессы, обладающие характеристикой `SAP_IPC_LOCK`, могут блокировать в памяти любое число страниц. Процессы, не имеющие этой характеристики, могут заблокировать только `RLIMIT_MEMLOCK` байтов. По умолчанию значение данного лимита ресурсов равно 32 Кбайт — достаточно, чтобы зафиксировать в памяти секрет или два, но недостаточно много, чтобы можно было снизить производительность работы системы (лимиты ресурсов обсуждаются в главе 6; в ней рассказывается, как проверять и менять эти значения).

Находится ли страница в физической памяти?

Для целей отладки и диагностики в Linux предоставляется функция `mincore()`, которую можно использовать для того, чтобы проверять, находится данный диапазон памяти в физической памяти или он сброшен на диск:

```
#include <unistd.h>
#include <sys/mman.h>

int mincore (void *start,
             size_t length,
             unsigned char *vec);
```

Вызов `mincore()` предоставляет вектор, указывающий, какие страницы отображения находятся в физической памяти на момент выполнения системного вызова. Вызов `mincore()` возвращает этот вектор через аргумент `vec` и описывает страницы, начинающиеся с адреса `start` (который должен быть выровнен по странице) и продолжающиеся на `length` байтов (это значение не обязательно должно быть выровнено по странице). Каждый байт в векторе `vec` соответствует одной странице в указанном диапазоне. Первый байт описывает первую страницу и так далее. Следовательно, параметр `vec` должен быть, по крайней мере, настолько велик, чтобы вместить $(length - 1 + \text{размер_страницы}) / \text{размер_страницы}$ байтов. Бит низшего порядка в каждом байте равен 1, если соответствующая страница находится в физической памяти, и 0, если это не так. Прочие биты в настоящее время не определены и зарезервированы для использования в будущем.

В случае успеха вызов возвращает 0. В случае ошибки он возвращает `-1` и присваивает переменной `errno` одно из следующих значений:

`EAGAIN`

Недостаточно ресурсов ядра для выполнения запроса.

`EFAULT`

Параметр `vec` указывает на недопустимый адрес.

EINVAL

Значение параметра `start` не выровнено по границе страницы.

ENOMEM

[`address`,`address + 1`] содержит память, не являющуюся частью файлового отображения.

В настоящее время данный системный вызов правильно работает только для файловых отображений, созданных с использованием `MAP_SHARED`. Это значительно ограничивает варианты применения вызова.

Уступающее выделение памяти

В Linux применяется стратегия *уступающего выделения памяти* (*opportunistic allocation*). Когда процесс запрашивает у ядра дополнительную память, например, увеличивая свой сегмент данных или создавая новое отображение в памяти, ядро фиксирует *обязательство* (*commit*) выделить память, не представляя фактически никакое хранилище. Только когда процесс начинает записывать, ядро *удовлетворяет* (*satisfy*) запрос, преобразуя обязательство выделить память в физическое выделение памяти. Ядро делает это постранично, при необходимости выполняя отложенную подкачку страниц и копирование при записи.

У такого поведения несколько преимуществ. Во-первых, ленивое выделение памяти позволяет ядру откладывать большую часть работы до последнего момента — если вообще когда-нибудь действительно понадобится удовлетворять запросы на выделение памяти. Во-вторых, так как запросы удовлетворяются постранично и по требованию, только действительно потребляемая физическая память занимает физическое хранилище. Наконец, объем зафиксированной памяти может значительно превосходить объем физической памяти и даже доступного пространства подкачки. Эта последняя особенность называется принятием *чрезмерных обязательств* (*overcommitment*).

Принятие чрезмерных обязательств и ООМ

Принятие чрезмерных обязательств позволяет системе выполнять намного больше и намного более объемных приложений, чем она бы могла, если бы каждая запрошенная страница памяти должна была поддерживаться некоторым физическим хранилищем в точке выделения, а не в точке использования. Без принятия чрезмерных обязательств отображение 2-гигабайтного файла с копированием при записи потребовало бы у ядра выделить 2 Гбайт фактического хранилища. Однако с принятием чрезмерных обязательств для выполнения той же задачи используется только хранилище, по объему равное сумме страниц данных, которые процесс фактически записывает. Аналогичным образом без принятия чрезмерных обязательств каждый вызов `fork()` требовал бы достаточно

свободного пространства в хранилище, чтобы скопировать все адресное пространство, хотя подавляющее большинство страниц никогда не подпадает под копирование при записи.

Но что, если процессы пытаются удовлетворить больше незавершенных обязательств, чем может обеспечить система, учитывая имеющуюся физическую память и область подкачки? В таком случае одно или несколько обязательств должны завершиться ошибкой. Так как ядро уже зафиксировало обязательство выделить память — системный вызов, запрашивающий обязательство, вернул код успеха — и процесс пытается использовать эту обещанную память, единственным выходом для ядра становится убийство процесса, чтобы освободить доступную память.

Когда принятие чрезмерных обязательств приводит к тому, что памяти становится недостаточно для удовлетворения обещанных запросов, мы говорим, что происходит условие *нехватки памяти* (*out of memory*, ООМ). В ответ на условие ООМ ядро задействует убийцу ООМ (*OOM killer*), который выбирает процесс, «заслуживающий» завершения. Для этой цели ядро пытается найти наименее важный процесс, который потребляет больше всего памяти.

Условия ООМ редки — это основная причина, почему принятие чрезмерных обязательств используется настолько широко. Однако нужно отметить, что подобные условия крайне нежелательны и недетерминированное завершение процесса убийцей ООМ зачастую недопустимо.

В некоторых системах ядро позволяет отключать принятие чрезмерных обязательств через файл `/proc/sys/vm/overcommit_memory` и аналогичный параметр `sysctl vm.overcommit_memory`.

Значение по умолчанию этого параметра, равное 0, заставляет ядро применять эвристическую стратегию принятия чрезмерных обязательств, фиксируя обязательства на выделение памяти до определенных пределов, но не допуская принятия вовсе чрезмерных обязательств. Значение 1 разрешает фиксировать любые обязательства, вообще отключая какие бы то ни было предосторожности. Определенные приложения, обильно расходующие память, например научные, обычно запрашивают память в количестве, настолько превышающем их возможности когда-либо ее израсходовать, что этот параметр определенно имеет смысл использовать.

Значение 2 вообще отключает принятие чрезмерных обязательств и включает *строгий учет* (*strict accounting*). В этом режиме обязательства на выделение памяти ограничиваются размером области подкачки и конфигурируемым процентом физической памяти. Конфигурируемый процент устанавливается в файле `/proc/sys/vm/overcommit_ratio` или при помощи аналогичного параметра `sysctl vm.overcommit_ratio`. Значение по умолчанию равно 50, что ограничивает обязательства на выделение памяти размером области подкачки и половиной физической памяти. Так как физическая память содержит ядро, таблицы страниц, зарезервированные системой страницы, заблокированные страницы и так далее, только часть ее фактически допускает подкачуку и гарантированно может использоваться для удовлетворения обязательств.

Будьте аккуратны со строгим учетом! Многие разработчики систем, помня об убийце ООМ, полагают, что строгий учет – это панацея. Однако приложения зачастую делают множество ненужных запросов на выделение памяти, далеко заходящих на территорию чрезмерного принятия обязательств, и обеспечение возможности такого поведения было одной из основных мотиваций при разработке виртуальной памяти.

9 Сигналы

Сигналы — это программные прерывания, предоставляющие механизм для обработки асинхронных событий. Такие события могут происходить из-за пределов системы — например, из-за введения пользователем символа прерывания (обычно `Ctrl+C`) — или возникать вследствие действий в программе или ядре, например, когда процесс исполняет код, в котором выполняется деление на ноль. В виде примитивной формы *взаимодействия между процессами* (interprocess communication, IPC) один_процесс также может отправлять сигналы другому процессу.

Самое главное в сигналах — это не только то, что события происходят асинхронно, например пользователь может нажать `Ctrl+C` в любой момент выполнения программы, но и то, что программа асинхронно обрабатывает сигналы. Функции обработки сигналов регистрируются в ядре, которое асинхронно вызывает функции из оставшейся части программы, как только сигналы доставляются.

Сигналы существуют в Unix с самого рождения системы. Со временем, однако, они эволюционировали, наиболее значительно — в терминах надежности, поскольку когда-то сигналы могли теряться, и в терминах функциональности, поскольку теперь сигналы могут переносить данные, определяемые пользователем. Вначале в разных системах Unix вносились несовместимые друг с другом изменения в сигналы. К счастью, на помощь пришел стандарт POSIX, в котором обработка сигналов была стандартизирована. Именно этот стандарт обеспечивается в Linux и именно о нем мы будем говорить далее.

Эту главу мы начнем с обзора сигналов и обсуждения их правильного и неправильного использования. Затем мы рассмотрим различные интерфейсы Linux для управления и манипулирования сигналами.

Наиболее значимые приложения взаимодействуют с сигналами. Даже если вы намеренно разработаете свое приложение таким образом, чтобы оно не полагалось на сигналы в своих коммуникационных нуждах — а это часто оказывается хорошей идеей! — в определенных случаях вам все равно придется работать с сигналами, например, при обработке завершения программы.

Концепции сигналов

У сигналов очень точный жизненный цикл. Сначала сигнал *поднимается* (*raise*, также мы иногда говорим, что он *отправляется* (*send*) или *генерируется* (*generate*)). После этого ядро *хранит* (*store*) сигнал до тех пор, пока у него не появляется шанс доставить его. Наконец, как только такая возможность предоставляется, ядро соответствующим образом *обрабатывает* (*handle*) сигнал. Ядро может выполнить одно из трех действий, в зависимости от того, что у него запросил процесс:

Игнорировать сигнал

Некакое действие не предпринимается. Существуют два сигнала, которые не могут быть проигнорированы: SIGKILL и SIGSTOP. Причина этого заключается в том, что системным администраторам нужно иметь возможность убивать и останавливать процессы, и было бы нарушением такого права, если бы процесс мог просто игнорировать SIGKILL (что делало бы его неубиваемым) или SIGSTOP (что делало бы его неостанавливаемым).

Захватить и обработать сигнал

Ядро приостанавливает исполнение текущего пути кода процесса и переходит к ранее зарегистрированной функции. Затем процесс исполняет эту функцию. После того как процесс возвращается из этой функции, он перепрыгивает обратно в то место, где находился на тот момент, когда был захвачен сигнал.

SIGINT и SIGTERM — это два сигнала, которые захватываются чаще всего. Процессы захватывают SIGINT, чтобы обработать ситуацию, когда пользователь вводит символ прерывания, — например, терминал может захватить его и вернуться к главной строке приглашения. Процессы захватывают SIGTERM для выполнения необходимой уборки, например отсоединения от сети или удаления временных файлов перед завершением. Сигналы SIGKILL и SIGSTOP захватить невозможно.

Выполнить действие по умолчанию

Это действие зависит от того, какой сигнал отправляется. Действием по умолчанию часто бывает завершение процесса. Например, именно так обрабатывается сигнал SIGKILL. Однако многие сигналы предоставляются для специфических целей, которые беспокоят программистов лишь в определенных ситуациях, и эти сигналы по умолчанию игнорируются, так как многие программы просто в них не заинтересованы. Далее мы рассмотрим разнообразные сигналы и их действия по умолчанию.

Раньше, когда сигнал доставлялся, у функции, которая обрабатывала сигнал, не было никакой информации о том, что произошло, она знала только, что был создан определенный сигнал. Сегодня ядро умеет предоставлять программистам, заинтересованным в этом, обширный контекст, и сигналы даже могут передавать пользовательские данные — так работают современные более продвинутые механизмы взаимодействия между процессорами.

Идентификаторы сигналов

У каждого сигнала есть символическое имя, начинающееся с префикса SIG. Например, SIGINT — это сигнал, который отправляется, когда пользователь нажимает клавишное сочетание Ctrl+C, SIGABRT — это сигнал, отправляемый, когда процесс вызывает функцию `abort()`, а SIGKILL — сигнал, отправляемый, когда процесс принудительно завершается.

Все эти сигналы определены в файле заголовка, подключаемом через `<signal.h>`. Это просто определения препроцессора, представляющие положительные целые числа, то есть каждый сигнал также ассоциируется с целочисленным идентификатором. Отображение друг на друга имен и целых чисел зависит от реализации и варьируется в разных системах Unix, хотя первая дюжина сигналов или около того обычно совпадает во всех (например, SIGKILL — это широко известный *сигнал 9*). Хороший программист всегда использует имя сигнала, удобное для чтения, и никогда целочисленное значение.

Номера сигналов начинаются с единицы (обычно единице соответствует SIGHUP) и линейно продолжаются далее. Всего сигналов около 30, но в большинстве программ регулярно используется лишь несколько из них. Сигнала со значением 0 нет — это специальное значение, называемое *нулевым сигналом* (null signal). Он никакой особой роли не играет и не заслуживает отдельного имени, но некоторые системные вызовы (такие, как `kill()`) применяют значение 0 в специальных случаях.

Вывести список сигналов, поддерживаемых в вашей системе, можно при помощи команды `kill -l`.

Сигналы, поддерживаемые в Linux

В табл. 9.1 перечисляются сигналы, которые поддерживаются в Linux.

Таблица 9.1. Сигналы

Сигнал	Описание	Действие по умолчанию
SIGABRT	Отправляется функцией <code>abort()</code>	Завершиться с созданием дампа ядра
SIGNALRM	Отправляется функцией <code>alarm()</code>	Завершиться
SIGBUS	Аппаратная ошибка или ошибка выравнивания	Завершиться с созданием дампа ядра
SIGCHLD	Завершился дочерний процесс	Игнорировать
SIGCONT	Процесс продолжил выполняться после того, как был остановлен	Игнорировать
SIGFPE	Арифметическое исключение	Завершиться с созданием дампа ядра

Сигнал	Описание	Действие по умолчанию
SIGHUP	Управляющий терминал процесса был закрыт (чаще всего это связано с тем, что пользователь выходит из системы)	Завершиться
SIGILL	Процесс попытался выполнить недопустимую функцию	Завершиться с созданием дампа ядра
SIGINT	Пользователь ввел символ прерывания (Ctrl+C)	Завершиться
SIGIO	Асинхронное событие ввода-вывода	Завершиться ^a
SIGKILL	Завершение процесса, которое невозможно захватить	Завершиться
SIGPIPE	Процесс записал данные в конвейер, но читателей нет	Завершиться
SIGPROF	Истек таймер профилирования	Завершиться
SIGPWR	Сбой питания	Завершиться
SIGQUIT	Пользователь ввел символ выхода (Ctrl+\)	Завершиться с созданием дампа ядра
SIGSEGV	Нарушение доступа к памяти	Завершиться с созданием дампа ядра
SIGSTKFLT	Ошибка стека сопроцессора	Завершиться ^b
SIGSTOP	Приостановить выполнение процесса	Остановиться
SIGSYS	Процесс попытался выполнить недопустимый системный вызов	Завершиться с созданием дампа ядра
SIGTERM	Завершение процесса с возможностью захвата	Завершиться
SIGTRAP	Встретилась точка останова	Завершиться с созданием дампа ядра
SIGTSTP	Пользователь ввел символ приостановки (Ctrl+Z)	Остановиться
SIGTTIN	Фоновый процесс считал данные с управляющего терминала	Остановиться
SIGTTOU	Фоновый процесс записал данные в управляющий терминал	Остановиться
SIGURG	Срочное ожидание ввода-вывода	Игнорировать
SIGUSR1	Сигнал, определенный процессом	Завершиться

Таблица 9.1 (продолжение)

Сигнал	Описание	Действие по умолчанию
SIGUSR2	Сигнал, определенный процессом	Завершиться
SIGVTALRM	Генерируется функцией <code>setitimer()</code> , когда она вызывается с флагом <code>ITIMER_VIRTUAL</code>	Завершиться
SIGWINCH	Изменился размер окна управляющего терминала	Игнорировать
SIGXCPU	Превышены лимиты ресурсов процессора	Завершиться с созданием дампа ядра
SIGXFSZ	Превышены лимиты ресурсов файла	Завершиться с созданием дампа ядра

^a Поведение в других системах Unix, таких, как BSD, — игнорировать этот сигнал.

^b Ядро Linux более не генерирует этот сигнал; он остается только для обеспечения обратной совместимости.

Существует еще несколько значений сигналов, но в Linux они определяются как эквивалентные другим значениям: `SIGINFO` определяется как `SIGPWR1`, `SIGIOT` определяется как `SIGABRT`, а `SIGPOLL` и `SIGLOST` определяются как `SIGIO`.

Теперь, когда у нас есть справочная таблица, давайте подробнее рассмотрим каждый из сигналов:

SIGABRT

Функция `abort()` отправляет этот сигнал процессу, который вызывает ее. Затем процесс завершается и генерирует файл ядра. В Linux такие утверждения, как `assert()`, вызывают `abort()`, когда условие проваливается.

SIGALRM

Функции `alarm()` и `setitimer()` (с флагом `ITIMER_REAL`) отправляют этот сигнал процессу, который вызвал их, когда тревога оканчивается. Эти и похожие функции рассматриваются в главе 10.

SIGBUS

Ядро поднимает этот сигнал, когда в процессе происходит аппаратный сбой, отличный от защиты памяти, — в этом случае генерируется сигнал `SIGSEGV`. В традиционных системах Unix данный сигнал представлял различные неправильные ошибки, например невыровненный доступ к памяти. Ядро Linux автоматически исправляет большинство из этих ошибок, не генерируя сигнал. Ядро поднимает данный сигнал, когда процесс неправильно обращается

¹ Этот сигнал определяется только в архитектуре Alpha. Во всех остальных машинных архитектурах этот сигнал не существует.

к области памяти, созданной при помощи `mmap()` (отображения в памяти обсуждаются в главе 8). Если только этот сигнал не захватывается, ядро завершает процесс и генерирует дамп ядра.

SIGCHLD

Когда процесс завершается или останавливается, ядро отправляет этот сигнал предку процесса. Так как `SIGCHLD` по умолчанию игнорируется, процессы должны явно захватывать и обрабатывать данный сигнал, если они интересуются жизнями своих потомков. Обработчик этого сигнала обычно вызывает `wait()` (подробнее об этом системном вызове – в главе 5), чтобы узнать идентификатор `pid` потомка и код выхода.

SIGCONT

Ядро отправляет данный сигнал процессу, когда процесс возобновляется после остановки. По умолчанию этот сигнал игнорируется, но процессы могут захватывать его, если им нужно выполнять какое-то действие после того, как они возобновляются. Данный сигнал часто используется терминалами и редакторами, которые обновляют экран.

SIGFPE

Несмотря на свое имя, этот сигнал представляет любое арифметическое исключение, а не только те, которые относятся к операциям с плавающей точкой. Исключения включают переполнения, потерю значимости и деление на ноль. Действие по умолчанию – завершить процесс и сгенерировать файл ядра, но процессы по желанию могут захватывать и обрабатывать этот сигнал. Обратите внимание, что поведение процесса и результат операции, ставшей причиной сигнала, не определены, если процесс решает продолжать выполняться.

SIGHUP

Ядро отправляет этот сигнал лидеру сеанса, когда терминал сеанса отключается. Также ядро отправляет данный сигнал каждому процессу в приоритетной группе процессов, когда лидер сеанса завершается. Действие по умолчанию – завершиться, что логично, так как сигнал предполагает, что пользователь выходит из системы. Процессы-демоны «перегружают» этот сигнал механизмом, который заставляет их перезагружать свои конфигурационные файлы. Если отправить `SIGHUP`, например, Apache, то эта служба перезагрузит `httpd.conf`. Использование для этой цели сигнала `SIGHUP` очень распространено, но не обязательно. Такая практика безопасна, поскольку у демонов нет управляющих терминалов, и поэтому в обычных условиях они никогда не получают данный сигнал.

SIGILL

Ядро отправляет этот сигнал, когда процесс пытается выполнить недопустимую машинную инструкцию. Действие по умолчанию – завершить процесс и сгенерировать дамп ядра. Процессы могут захватывать и обрабатывать `SIGILL`, но их поведение после поднятия данного сигнала не определено.

SIGINT

Этот сигнал отправляется всем процессам в приоритетной группе процессов, когда пользователь вводит символ прерывания (обычно **Ctrl+C**). Поведение по умолчанию — завершиться, однако процессы могут захватывать и обрабатывать данный сигнал, и обычно они делают это в целях уборки перед завершением.

SIGIO

Этот сигнал отправляется, если генерируется асинхронное событие ввода-вывода в стиле **BSD**. Такой стиль ввода-вывода редко применяется в **Linux**. (Расширенные техники ввода-вывода, использующиеся в **Linux**, обсуждаются в главе 4.)

SIGKILL

Этот сигнал отправляется системным вызовом **kill()**; он существует для того, чтобы давать системным администраторам надежный способ безоговорочного убийства процесса. Этот сигнал невозможно захватить или игнорировать, и его результатом всегда является завершение процесса.

SIGPIPE

Если процесс записывает данные в конвейер, но читатель уже завершился, то ядро поднимает этот сигнал. Действие по умолчанию — завершить процесс, но данный сигнал можно захватывать и обрабатывать.

SIGPROF

Функция **setitimer()**, если она используется с флагом **ITIMER_PROF**, генерирует данный сигнал, когда таймер профилирования истекает. Действие по умолчанию — завершить процесс.

SIGPWR

Значение этого сигнала зависит от системы. В **Linux** он соответствует условию, когда заканчивается заряд аккумулятора (например, в источнике бесперебойного питания (**UPS**)). Демон мониторинга **UPS** отправляет этот сигнал процессу **init**, который отвечает уборкой и выключением системы — следует надеяться, что до того, как питание закончится!

SIGQUIT

Ядро поднимает данный сигнал для всех процессов в приоритетной группе процессов, когда пользователь вводит символ выхода из терминала (обычно **Ctrl+**). Действие по умолчанию — завершить процесс и создать дамп ядра.

SIGSEGV

Этот сигнал, название которого обозначает **segmentation violation** — нарушение сегментирования — отправляется процессу, когда тот делает попытку недопустимого доступа к памяти. Сюда входит доступ к не отображенной памяти, чтение из памяти, которая недоступна для чтения, выполнение кода в памяти, которая недоступна для выполнения, а также запись в память, недоступную для записи. Процессы могут захватывать и обрабатывать этот сигнал, но действие по умолчанию — завершить процесс и создать дамп ядра.

SIGSTOP

Этот сигнал отправляется только вызовом `kill()`. Он безусловно останавливает процесс, и его невозможно захватить или игнорировать.

SIGSYS

Ядро отправляет этот сигнал процессу, когда тот делает попытку выполнить недопустимый системный вызов. Это может происходить, когда двоичный файл собирается в более новой версии операционной системы (с новыми версиями системных вызовов), но затем выполняется в старой версии. Правильно собранные двоичные файлы, делающие свои системные вызовы через `glibc`, никогда не должны получать данный сигнал. Вместо этого недопустимые системные вызовы должны возвращать значение `-1` и присваивать переменной `errno` значение `ENOSYS`.

SIGTERM

Этот сигнал отправляется только вызовом `kill()`; он позволяет пользователю изящно завершить процесс (действие по умолчанию). Процессы могут захватывать этот сигнал и выполнять уборку перед завершением, но для процесса считается невежливым захватить его и не завершиться сразу же.

SIGTRAP

Ядро отправляет этот сигнал процессу, когда тот пересекает точку останова. В целом, отладчики захватывают данный сигнал, а остальные процессы игнорируют его.

SIGTSTP

Ядро отправляет этот сигнал всем процессам в приоритетной группе процессов, когда пользователь вводит символ приостановки (обычно `Ctrl+Z`).

SIGTTIN

Этот сигнал отправляется процессу, который находится в фоновом режиме, когда пытается считать данные из своего управляющего терминала. Действие по умолчанию — остановить процесс.

SIGTTOU

Этот сигнал отправляется процессу, который находится в фоновом режиме, когда пытается записать данные в свой управляющий терминал. Действие по умолчанию — остановить процесс.

SIGURG

Ядро отправляет этот сигнал процессу, когда в сокет прибывают внеполосные (out-of-band, OOB) данные. Внеполосные данные в этой книге не рассматриваются.

SIGUSR1 и SIGUSR2

Эти сигналы предназначены исключительно для пользователей; ядро никогда не поднимает их. Процессы могут использовать `SIGUSR1` и `SIGUSR2` в любых целях. Обычно они применяются для того, чтобы заставлять процесс-демон менять поведение. Действие по умолчанию — завершить процесс.

SIGVTALRM

Функция `setitimer()` отправляет этот сигнал, когда таймер, созданный с флагом `ITIMER_VIRTUAL`, истекает. Таймеры рассматриваются в главе 10.

SIGWINCH

Ядро поднимает этот сигнал для всех процессов в приоритетной группе процессов, когда размер их терминального окна меняется. По умолчанию процессы игнорируют этот сигнал, но они могут захватывать и обрабатывать его, если они осведомлены о размере своих терминальных окон. Хороший пример программы, захватывающей этот сигнал, — это `top`. Попробуйте изменить размер ее окна, пока она выполняется, и посмотрите, как она на это отреагирует.

SIGXCPU

Ядро поднимает этот сигнал, когда процесс превышает свой мягкий лимит процессора. Ядро продолжает поднимать его с периодичностью раз в секунду до тех пор, пока процесс не завершится или не превысит жесткий лимит процессора. Как только превышается жесткий лимит, ядро отправляет процессу сигнал `SIGKILL`.

SIGXFSZ

Ядро поднимает данный сигнал, когда процесс превышает свой лимит размера файла. Действие по умолчанию — завершить процесс, но если этот сигнал захватывается или игнорируется, системный вызов, который должен был изменить превышаемый лимит размера файла, возвращает `-1` и присваивает переменной `errno` значение `EFBIG`.

Простейшее управление сигналами

Простейший и самый старый интерфейс для управления сигналами — это функция `signal()`. Определенный стандартом ISO C89, который стандартизирует только наименьшее общее кратное поддержки сигналов, этот системный вызов элементарен. Linux может предоставить намного более тонкий контроль над сигналами через другие интерфейсы, о которых мы поговорим далее в этой главе. Так как системный вызов `signal()` самый простой и довольно распространенный благодаря определению в стандарте ISO C, мы рассмотрим его в первую очередь:

```
#include <signal.h>

typedef void (*sighandler_t)(int);

sighandler_t signal (int signo, sighandler_t handler);
```

Успешный вызов `signal()` удаляет текущее действие, предпринимаемое при получении сигнала `signo`, и вместо этого обрабатывает сигнал обработчиком, указанным при помощи аргумента `handler`. Значением аргумента `signo` может быть одно из названий сигналов, перечисленных в предыдущем разделе, напри-

мер SIGINT или SIGUSR1. Вспомните, что процесс не может захватывать ни SIGKILL, ни SIGSTOP, поэтому установка обработчика для любого из этих сигналов не имеет никакого смысла.

Функция `handler` должна возвращать значение `void`, что имеет смысл, так как (в отличие от обычных функций) в программе нет стандартного места, в котором данная функция должна была бы возвращать результат. Она принимает один аргумент — целочисленное значение, представляющее собой идентификатор обрабатываемого сигнала (например, `SIGUSR2`). Это позволяет одной функции обрабатывать несколько сигналов. Ее прототип имеет такую форму:

```
void my_handler (int signo);
```

В Linux для определения этого прототипа используется описание типа `sighandler_t`. В прочих системах Unix напрямую используются указатели функции; в некоторых системах есть собственные типы, которые не могут носить имя `sighandler_t`. В программах, нацеленных на переносимость, нельзя напрямую ссылаться на тип.

Когда ядро поднимает сигнал для процесса, который зарегистрировал обработчик сигнала, оно приостанавливает выполнение обычного потока инструкций программы и вызывает обработчик сигнала. Обработчику передается значение сигнала, содержащееся в аргументе `signo`, первоначально предоставленном вызову `signal()`.

Также можно использовать вызов `signal()` для того, чтобы заставить ядро игнорировать определенный сигнал для текущего процесса или восстановить поведение по умолчанию для сигнала. Это делается при помощи специальных значений параметра `handler`:

SIG_DFL

Восстановить поведение по умолчанию для сигнала, указанного при помощи аргумента `signo`. Например, в случае сигнала `SIGPIPE` процесс будет завершаться.

SIG_IGN

Игнорировать сигнал, указанный при помощи параметра `signo`.

Функция `signal()` возвращает предыдущее поведение сигнала, которое может принимать вид указателя на обработчик сигнала, `SIG_DFL` или `SIG_IGN`. В случае ошибки она возвращает значение `SIG_ERR`. Эта функция не устанавливает переменную `errno`.

Ожидание сигнала, любого сигнала

Удобный для отладки и для написания демонстрационных фрагментов кода, определенный в стандарте POSIX системный вызов `pause()` приостанавливает процесс до тех пор, пока тот не получит либо сигнал, допускающий обработку, либо сигнал, который заставит процесс завершиться:

```
#include <unistd.h>
```

```
int pause (void);
```

Вызов `pause()` возвращает результат только при получении сигнала, который процесс захватывает. В этом случае сигнал обрабатывается, а `pause()` возвращает значение `-1` и присваивает переменной `errno` код `EINTR`. Если ядро поднимает проигнорированный сигнал, то процесс не просыпается.

В ядре Linux системный вызов `pause()` является одним из самых простых. Он выполняет только два действия. Во-первых, он помещает процесс в прерываемое состояние ожидания. Во-вторых, он вызывает `schedule()`, чтобы запустить планировщик процессов Linux, который должен найти другой процесс и позволить ему выполниться. Так как процесс в действительности ничего не ожидает, ядро не будит его до тех пор, пока он не получает сигнал. Все эти страшания умещаются всего лишь в две строки кода на языке C¹.

Примеры

Давайте рассмотрим пару простых примеров. Код первого примера регистрирует для сигнала `SIGINT` обработчик, который просто выводит сообщение, а затем завершает программу (как `SIGINT` делает и сам по себе):

```
#include <stdlib.h>
#include <stdio.h>
#include <unistd.h>
#include <signal.h>

/* обработчик для SIGINT */
static void sigint_handler (int signo)
{
    /*
     * Технически, не следует использовать printf( )
     * в обработчике сигналов, но это не конец света.
     * Я расскажу, почему это так, в разделе
     * «Повторный вход».
     */
    printf ("Захвачен сигнал SIGINT!\n");
    exit (EXIT_SUCCESS);
}

int main (void)
{
    /*
     * Регистрируем sigint_handler как наш обработчик
     * сигналов для SIGINT.
     */
    if (signal (SIGINT, sigint_handler) == SIG_ERR) {
        fprintf (stderr, "Невозможно обработать SIGINT!\n");
        exit (EXIT_FAILURE);
    }

    for (;;)
```

¹ Таким образом, `pause()` по простоте занимает лишь второе место. Первое место делят между собой `getpid()` и `gettid()`, каждый из которых состоит лишь из одной строки.

```
    pause ( );
    return 0;
}
```

В следующем примере мы регистрируем тот же обработчик для сигналов SIGTERM и SIGINT. Также мы восстанавливаем поведение по умолчанию для сигнала SIGPROF (то есть прекращение процесса) и игнорируем SIGHUP (который в противном случае завершал бы процесс):

```
#include <stdlib.h>
#include <stdio.h>
#include <unistd.h>
#include <signal.h>

/* обработчик для SIGINT */
static void signal_handler (int signo)
{
    if (signo == SIGINT)
        printf ("Захвачен сигнал SIGINT!\n");
    else if (signo == SIGTERM)
        printf ("Захвачен сигнал SIGTERM!\n");
    else {
        /* это никогда не должно случаться */
        fprintf (stderr, "Неожиданный сигнал!\n");
        exit (EXIT_FAILURE);
    }
    exit (EXIT_SUCCESS);
}

int main (void)
{
    /*
     * Регистрируем signal_handler как наш обработчик сигнала
     * для SIGINT.
     */
    if (signal (SIGINT, signal_handler) == SIG_ERR) {
        fprintf (stderr, "Невозможно обработать SIGINT!\n");
        exit (EXIT_FAILURE);
    }

    /*
     * Регистрируем signal_handler как наш обработчик сигнала
     * для SIGTERM.
     */
    if (signal (SIGTERM, signal_handler) == SIG_ERR) {
        fprintf (stderr, "Невозможно обработать SIGTERM!\n");
        exit (EXIT_FAILURE);
    }

    /* Восстановление поведения по умолчанию для сигнала SIGPROF. */
    if (signal (SIGPROF, SIG_DFL) == SIG_ERR) {
        fprintf (stderr, "Невозможно сбросить SIGPROF!\n");
        exit (EXIT_FAILURE);
    }

    /* Игнорировать SIGHUP. */
}
```

```

    if (signal (SIGHUP, SIG_IGN) == SIG_ERR) {
        fprintf (stderr, "Невозможно игнорировать SIGHUP!\n");
        exit (EXIT_FAILURE);
    }

    for (;;) {
        pause ();
    }

    return 0;
}

```

Запуск и наследование

Когда процесс впервые запускается, для него устанавливается поведение всех сигналов по умолчанию, если только родительский процесс (тот, который запустил новый процесс) не игнорирует какие-то сигналы, — в таком случае только что созданный процесс также будет их игнорировать. Другими словами, любой сигнал, захватываемый предком, в новом процессе выполняет действие по умолчанию, а все прочие сигналы остаются такими, какими и были. Это имеет смысл, так как у нового процесса адресное пространство не совпадает с адресным пространством его предка, поэтому зарегистрированные обработчики сигналов могут не существовать.

У такого поведения при выполнении процессов есть удобное практическое преимущество: когда оболочка запускает процесс «в фоне» (или когда фоновый процесс запускает другой процесс), только что запущенный процесс должен игнорировать символы прерывания и выхода. Таким образом, перед тем как оболочка запустит фоновый процесс, она должна установить для сигналов SIGINT и SIGQUIT флаг SIG_IGN. Поэтому очень часто в программах, которые обрабатывают эти сигналы, сначала выполняется проверка, чтобы удостовериться, что сигналы не игнорируются. Например:

```

/* обрабатывать сигнал SIGINT, но только если он не игнорируется */
if (signal (SIGINT, SIG_IGN) != SIG_IGN) {
    if (signal (SIGINT, sigint_handler) == SIG_ERR)
        fprintf (stderr, "Невозможно обработать SIGINT!\n");
}

/* обрабатывать сигнал SIGQUIT, но только если он не игнорируется */
if (signal (SIGQUIT, SIG_IGN) != SIG_IGN) {
    if (signal (SIGQUIT, sigquit_handler) == SIG_ERR)
        fprintf (stderr, "Невозможно обработать SIGQUIT!\n");
}

```

Необходимость устанавливать поведение сигнала, для того чтобы иметь возможность проверять поведение сигнала, указывает на недостаток интерфейса `signal()`. Далее мы обсудим функцию, которая не обладает этим недостатком.

Как вы могли догадаться, поведение для `fork()` отличается. Когда процесс делает вызов `fork()`, потомок наследует ту же семантику сигналов, что и у его предка. Это также имеет смысл, так как у потомка и предка одно адресное пространство и обработчики сигналов родительского процесса продолжают существовать в дочернем.

Сопоставление номеров сигналов и строк

Пока что в наших примерах мы жестко кодировали имена сигналов. Но иногда удобнее иметь возможность преобразовывать номер сигнала в строковое представление его имени (это также бывает обязательным требованием). Есть несколько способов сделать это. Один из них — извлечь строку из статически определенного списка:

```
extern const char * const sys_siglist[];
```

Sys_siglist — это массив строк, содержащих имена сигналов, которые поддерживаются в системе, проиндексированный по номерам сигналов.

Альтернативой является определенный в BSD интерфейс psignal(), достаточно распространенный, чтобы он также поддерживался в Linux:

```
#include <signal.h>
```

```
void psignal (int signo, const char *msg);
```

Вызов psignal() печатает на выводе stderr строку, указанную при помощи аргумента msg, за которой следуют двоеточие, пробел и имя сигнала, соответствующего значению аргумента signo. Если значение signo недопустимо, то об этом будет говориться в сообщении.

Интерфейс strsignal() лучше. Он не стандартизирован, но Linux и многие отличные от Linux системы поддерживают его:

```
#define _GNU_SOURCE  
#include <string.h>
```

```
char *strsignal (int signo);
```

Вызов strsignal() возвращает указатель на описание сигнала, указанного при помощи signo. Если значение signo недопустимо, то об этом говорится в возвращаемом описании (в некоторых системах Unix, которые поддерживают данную функцию, возвращается вместо этого значение NULL). Возвращаемая строка имеет силу только до следующего вызова strsignal(), поэтому данную функцию небезопасно использовать с потоками.

Лучше всего придерживаться подхода с массивом sys_siglist. Применив его, можно переписать предыдущий обработчик сигналов следующим образом:

```
static void signal_handler (int signo)  
{  
    printf ("Захвачен сигнал %s\n", sys_siglist[signo]);  
}
```

Отправка сигнала

Системный вызов kill(), лежащий в основе распространенной утилиты kill, отправляет сигнал от одного процесса другому:

```
#include <sys/types.h>  
#include <signal.h>  
  
int kill (pid_t pid, int signo);
```

При нормальном использовании (то есть когда значение `pid` больше 0) `kill()` отправляет сигнал `signo` процессу, указанному при помощи идентификатора `pid`.

Если аргумент `pid` равен 0, то `signo` отправляется всем процессам в группе процессов вызывающего процесса.

Если аргумент `pid` равен -1, то `signo` отправляется всем процессам, для которых у вызывающего процесса есть разрешение на отправку сигнала, за исключением самого себя и процесса `init`. Мы обсудим разрешения, регулирующие доставку сигналов, в следующем подразделе.

Если аргумент `pid` меньше -1, то `signo` отправляется группе процессов `-pid`.

В случае успеха `kill()` возвращает 0. Вызов считается успешным, если хотя бы один сигнал был отправлен. В случае сбоя (ни одного сигнала не отправлено) вызов возвращает значение -1 и присваивает переменной `errno` один из следующих кодов:

EINVAL

Значение `signo` представляет недопустимый сигнал.

EPERM

У вызывающего процесса недостаточно полномочий для отправки сигнала какому-либо из запрошенных процессов.

ESRCH

Процесс или группа процессов, указанная при помощи `pid`, не существует или в случае процесса является зомби.

Разрешения

Для того чтобы отправить сигнал другому процессу, отправляющему процессу требуются определенные разрешения. Процесс, обладающий характеристикой `CAP_KILL` (обычно принадлежащий пользователю `root`), может отправлять сигналы любым процессам. У не имеющих этой характеристики отправляющих процессов действительный или реальный идентификатор пользователя должен быть равен реальному или сохраненному идентификатору пользователя получающего процесса. Проще говоря, пользователь может отправить сигнал только процессу, владельцем которого он является.

ПРИМЕЧАНИЕ

В системах Unix есть одно исключение, относящееся к `SIGCONT`: процесс может отправить данный сигнал любому другому процессу в том же сеансе. Идентификаторы пользователя не обязательно должны совпадать.

Если значение `signo` равно 0 (вышеупомянутый нулевой сигнал), то вызов не отправляет сигнал, но все же выполняет проверку ошибок. Это позволяет протестировать, обладает ли процесс необходимыми разрешениями для отправки сигнала определенному процессу или процессам.

Примеры

Далее показан пример отправки сигнала SIGHUP процессу с идентификатором процесса 1722:

```
int ret;  
  
ret = kill (1722, SIGHUP);  
if (ret)  
    perror ("kill");
```

Этот фрагмент кода делает то же самое, что и следующий вызов утилиты kill:

```
$ kill -HUP 1722
```

Чтобы проверить, есть ли у нас разрешение на отправку сигнала процессу с идентификатором 1722, не отправляя в действительности никакого сигнала, можно сделать следующее:

```
int ret;  
  
ret = kill (1722, 0);  
if (ret)  
    ; /* у нас нет разрешения */  
else  
    ; /* у нас есть разрешение */
```

Отправка сигнала себе

Функция `raise()` — это простой способ, при помощи которого процесс может отправить сигнал самому себе:

```
#include <signal.h>  
  
int raise (int signo);
```

Этот вызов:

```
raise (signo);  
эквивалентен такому вызову:
```

```
kill (getpid (), signo);
```

Вызов возвращает значение 0 в случае успеха и ненулевое значение в случае ошибки. Он не устанавливает переменную `errno`.

Отправка сигнала всей группе процессов

Еще одна удобная функция упрощает отправку сигнала всем процессам из одной группы, когда отрицание идентификатора группы процессов и использование вызова `kill()` кажется слишком обременительным:

```
#include <signal.h>
```

```
int killpg (int pgrp, int signo);
```

Этот вызов:

```
killpg (pgrp, signo);
```

эквивалентен следующему вызову:

```
kill (-pgrp, signo);
```

То же самое работает, даже если значение `pgrp` равно 0, — в этом случае `killpg()` отправляет сигнал `signo` каждому процессу в группе вызывающего процесса.

В случае успеха `killpg()` возвращает 0. В случае ошибки этот вызов возвращает -1 и присваивает переменной `errno` одно из следующих значений:

`EINVAL`

Значение `signo` представляет недопустимый сигнал.

`EPERM`

У вызывающего процесса отсутствуют разрешения на отправку сигнала какому-либо из запрошенных процессов.

`ESRCH`

Группа процессов, указанная при помощи аргумента `pgrp`, не существует.

Повторный вход

В момент, когда ядро поднимает сигнал, процесс может выполнять код в каком угодно месте. Например, он может находиться в середине важной операции, которая, если ее прервать, оставит процесс в противоречивом состоянии — например, с наполовину обновленной структурой данных или частично завершенным вычислением. Процесс может даже обрабатывать другой сигнал.

У обработчиков сигналов нет способа узнавать, какой код процесс выполняет в момент получения сигнала; обработчик может вклиниваться в любые операции. Поэтому очень важно, чтобы любой обработчик сигнала, устанавливаемый вашим процессом, соблюдал осторожность относительно действий, которые он выполняет, и данных, к которым он обращается. Обработчики сигналов должны заботиться о том, чтобы не делать никаких предположений касательно того, что процесс делал, когда он был прерван. В частности, они должны быть очень внимательны к модификации глобальных (то есть совместно используемых) данных. В целом, лучше, если обработчик сигнала никогда не будет менять глобальные данные; однако в следующем разделе мы рассмотрим способ временной блокировки доставки сигналов как способа обеспечить безопасное манипулирование данными, совместно используемыми обработчиком сигнала и оставшейся частью процесса.

А как обстоят дела с системными вызовами и другими библиотечными функциями? Что, если ваш процесс находится в середине операции записи файла или выделения памяти, а обработчик сигнала записывает данные в тот же файл или вызывает `malloc()`? Или если процесс не завершил вызов функ-

ции, использующей статический буфер, такой, как `strsignal()`, и в этот момент доставляется сигнал?

Некоторые функции, очевидно, не поддерживают повторный вход. Если программа находится в процессе выполнения такой функции, а в это время доставляется сигнал, после чего обработчик сигнала вызывает ту же не поддерживающую повторный вход функцию, результатом может стать полнейший хаос. *Функция, поддерживающая повторный вход* (reentrant function), — это функция, которую можно безопасно вызывать из самой себя (или параллельно из другого потока в том же процессе). Для того чтобы квалифицироваться как поддерживающая повторный вход, функция не должна манипулировать статическими данными, а должна манипулировать только данными, выделенными в стеке, или данными, предоставленными ей вызывающим, и не должна вызывать никакие другие функции, не поддерживающие повторный вход.

Функции, гарантированно поддерживающие повторный вход

При написании обработчика сигналов вы должны предполагать, что прерываемый процесс может находиться на стадии выполнения не поддерживающей повторный вход функции (или любой другой операции, если уж на то пошло). Таким образом, в обработчиках сигналов допускается использование только тех функций, которые поддерживают повторный вход.

В различных стандартах определяются списки функций, *безопасных для сигналов* (signal-safe), то есть поддерживающих повторный вход, которые можно применять внутри обработчиков сигналов. Например, POSIX.1-2003 и спецификация Single UNIX Specification диктуют список функций, которые гарантированно допускают повторный вход и безопасны для сигналов на всех совместимых платформах. Эти функции перечисляются в табл. 9.2.

Таблица 9.2. Функции, гарантированно безопасные для сигналов и поддерживающие повторный вход

<code>abort()</code>	<code>accept()</code>	<code>access()</code>
<code>aio_error()</code>	<code>aio_return()</code>	<code>aio_suspend()</code>
<code>alarm()</code>	<code>bind()</code>	<code>cgetispeed()</code>
<code>cgetospeed()</code>	<code>cfsetispeed()</code>	<code>cfsetospeed()</code>
<code>chdir()</code>	<code>chmod()</code>	<code>chown()</code>
<code>clock_gettime()</code>	<code>close()</code>	<code>connect()</code>
<code>creat()</code>	<code>dup()</code>	<code>dup2()</code>
<code>execle()</code>	<code>execve()</code>	<code>Exit()</code>
<code>_exit()</code>	<code>fchmod()</code>	<code>fchown()</code>

продолжение ↗

Таблица 9.2 (продолжение)

fcntl()	fdatasync()	fork()
fpathconf()	fstat()	fsync()
ftruncate()	getegid()	geteuid()
Getgid()	getgroups()	getpeername()
getpgrp()	getpid()	getppid()
getsockname()	getsockopt()	getuid()
kill()	link()	listen()
lseek()	lstat()	mkdir()
Mkfifo()	open()	pathconf()
Pause()	pipe()	poll()
posix_trace_event()	pselect()	raise()
read()	readlink()	recv()
recvfrom()	recvmsg()	rename()
rmdir()	select()	sem_post()
send()	sendmsg()	sendto()
Setgid()	setpgid()	setsid()
setsockopt()	setuid()	shutdown()
sigaction()	sigaddset()	sigdelset()
sigemptyset()	sigfillset()	sigismember()
Signal()	sigpause()	sigpending()
sigprocmask()	sigqueue()	sigset()
sigsuspend()	sleep()	socket()
socketpair()	stat()	symlink()
sysconf()	tcdrain()	tcflow()
tcflush()	tcgetattr()	tcgetpgrp()
tcsendbreak()	tcsetattr()	tcsetpgrp()
time()	timer_getoverrun()	timer_gettime()
timer_settime()	times()	umask()
uname()	unlink()	utime()
wait()	waitpid()	write()

Многие другие функции также безопасны, но Linux и прочие совместимые с POSIX системы гарантируют повторную входимость только для этих функций.

Наборы сигналов

Несколько функций, которые мы изучим далее в этой главе, манипулируют наборами сигналов, например набором сигналов, заблокированных процессом, или набором сигналов, ожидающих процесс. Такими наборами позволяют управлять *операции над наборами сигналов* (signal set operation):

```
#include <signal.h>

int sigemptyset (sigset_t *set);

int sigfillset (sigset_t *set);

int sigaddset (sigset_t *set, int signo);

int sigdelset (sigset_t *set, int signo);

int sigismember (const sigset_t *set, int signo);
```

Sigemptyset() инициализирует набор сигналов, указанный при помощи аргумента set, помечая его как пустой (из набора исключены все сигналы). Sigfillset() инициализирует набор сигналов, указанный при помощи аргумента set, помечая его как полный (все сигналы включены в набор). Обе функции возвращают значение 0. Перед дальнейшим использованием набора необходимо вызвать на этом наборе сигналов одну из этих двух функций.

Sigaddset() добавляет сигнал signo в набор сигналов, указанный при помощи аргумента set, тогда как sigdelset() удаляет signo из набора сигналов set. Обе функции возвращают 0 в случае успеха и -1 в случае ошибки. Также в случае ошибки переменной errno присваивается код ошибки EINVAL, указывающий, что signo содержит недопустимый идентификатор сигнала.

Sigmember() возвращает 1, если signo находится в наборе сигналов, указанном при помощи аргумента set, 0, если это не так, и -1 в случае ошибки. В последней ситуации переменной errno также присваивается код EINVAL, указывающий на недопустимое значение signo.

Больше функций для наборов сигналов

Все предыдущие функции стандартизированы POSIX, и их можно найти в любой современной системе Unix. В Linux также предоставляется несколько нестандартных функций:

```
#define _GNU_SOURCE
#define <signal.h>

int sigisemptyset (sigset_t *set);
```

```
int sigorset (sigset_t *dest, sigset_t *left, sigset_t *right);  
int sigandset (sigset_t *dest, sigset_t *left, sigset_t *right);
```

Sigisemptyset() возвращает значение 1, если набор сигналов, указанный при помощи параметра set, пуст, и 0 в противном случае.

Sigorset() помещает объединение (бинарное ИЛИ) наборов сигналов left и right в dest. Sigandset() помещает пересечение (бинарное И) наборов сигналов left и right в dest. Обе возвращают 0 в случае успеха и -1 в случае ошибки, присваивая переменной errno значение EINVAL.

Эти функции полезны, но в программах, где требуется полное соответствие POSIX, их следует избегать.

Блокировка сигналов

Ранее мы обсудили повторный вход и вопросы, связанные с тем, что обработчики сигналов выполняются асинхронно в любой момент времени. Мы поговорили о функциях, которые не следует вызывать из самих обработчиков сигналов, так как они не поддерживают повторный вход.

Но что делать, если программе нужно передать данные между обработчиком сигнала и любой другой частью программы? Что, если есть такие периоды выполнения программы, в которых крайне нежелательны любые перерывы, включая прерывания, создаваемые при обработке сигналов? Мы называем такие части программы *критическими регионами* (critical region) и защищаем их, временно приостанавливая доставку сигналов. Мы говорим, что такие сигналы заблокированы. Никакие сигналы, поднятые во время блокировки, не обрабатываются до тех пор, пока они не разблокируются. Процесс может заблокировать любое число сигналов; набор сигналов, заблокированных процессом, называется его *сигнальной маской* (signal mask).

POSIX определяет, а Linux реализует функцию для управления сигнальной маской процесса:

```
#include <signal.h>  
  
int sigprocmask (int how,  
                 const sigset_t *set,  
                 sigset_t *oldset);
```

Поведение sigprocmask() зависит от значения аргумента how, который может содержать один из следующих флагов:

SIG_SETMASK

Сигнальная маска для вызывающего процесса меняется на set.

SIG_BLOCK

Сигналы из set добавляются к сигнальной маске вызывающего процесса. Другими словами, сигнальная маска меняется на объединение (двоичное ИЛИ) текущей маски и set.

SIG_UNBLOCK

Сигналы из set удаляются из сигнальной маски вызывающего процесса. Другими словами, сигнальная маска меняется на пересечение (двоичное И) текущей маски и отрицания (двоичное НЕ) набора set. Невозможно разблокировать сигнал, который не был заблокирован.

Если значение oldset не равно NULL, то функция помещает предыдущий сигнальный набор в oldset.

Если значение set равно NULL, то функция игнорирует how и не меняет сигнальную маску, а помещает ее в аргумент oldset. Другими словами, передача нулевого значения в качестве set — это способ извлечь текущую сигнальную маску.

В случае успеха вызов возвращает значение 0. В случае сбоя он возвращает -1 и присваивает переменной errno либо код EINVAL, указывая на недопустимое значение how, либо код EFAULT, указывая, что set или oldset содержит недопустимый указатель.

Блокировка сигналов SIGKILL и SIGSTOP не допускается. sigprocmask() бесшумно игнорирует любые попытки добавить любой из этих сигналов в сигнальную маску.

Извлечение ожидающих сигналов

Если ядро поднимает заблокированный сигнал, он не доставляется. Мы называем такие сигналы ожидающими. Когда ожидающий сигнал разблокируется, ядро передает его процессу для обработки.

В POSIX определяется функция, позволяющая извлечь набор ожидающих сигналов:

```
#include <signal.h>

int sigpending(sigset_t *set);
```

Успешный вызов sigpending() помещает набор ожидающих сигналов в аргумент set и возвращает значение 0. В случае ошибки вызов возвращает -1 и присваивает переменной errno значение EFAULT, указывая, что set содержит недопустимый указатель.

Ожидание набора сигналов

Третья определенная в POSIX функция позволяет процессу временно менять свою сигнальную маску, а затем ожидать, пока не будет поднят сигнал, который либо завершит процесс, либо будет процессом обработан:

```
#include <signal.h>
```

```
int sigsuspend(const sigset_t *set);
```

Если сигнал завершает процесс, то sigsuspend() не возвращает никакой результат. Если сигнал поднимается и обрабатывается, то sigsuspend() возвращает -1 и, после того как обработчик сигнала возвращает результат, присваивает

переменной errno значение EINTR. Если set содержит недопустимый указатель, то переменной errno присваивается код EFAULT.

Распространенный сценарий использования `sigsuspend()` — для извлечения сигналов, которые могли прибыть и заблокироваться во время критического региона выполнения программы. Сначала процесс при помощи `sigprocmask()` блокирует набор сигналов, сохраняя старую маску в `oldset`. После выхода из критического региона процесс вызывает `sigsuspend()`, предоставляя `oldset` в качестве значения `set`.

Расширенное управление сигналами

Функция `signal()`, которую мы рассмотрели в начале этой главы, максимально проста. Так как она является частью стандартной библиотеки С и, следовательно, должна отражать минимальные предположения относительно возможностей операционной системы, на которой выполняется, данная функция может предложить лишь наименьшее общее кратное управления сигналами. В качестве альтернативы в POSIX стандартизируется системный вызов `sigaction()`, представляющий намного более широкие возможности управления сигналами. Среди прочего его можно применять для блокировки получения определенных сигналов во время выполнения обработчика и для извлечения широкого диапазона данных о системе и состоянии процесса на момент, когда был поднят сигнал:

```
#include <signal.h>

int sigaction (int signo,
               const struct sigaction *act,
               struct sigaction *oldact);
```

Вызов `sigaction()` меняет поведение сигнала, указанного при помощи аргумента `signo`, который может принимать любые значения, за исключением значений, ассоциирующихся с `SIGKILL` и `SIGSTOP`. Если значение `act` не равно `NULL`, то системный вызов меняет текущее поведение сигнала в соответствии со значением `act`. Если значение `oldact` не равно `NULL`, то вызов сохраняет предыдущее (или текущее, если аргумент `act` равен `NULL`) поведение указанного сигнала в структуре `oldact`.

Структура `sigaction` обеспечивает возможность тонкого управления сигналами. В заголовке `<sys/signal.h>`, который подключается через `<signal.h>`, эта структура определяется следующим образом:

```
struct sigaction {
    void (*sa_handler)(int); /* обработчик сигнала или действие */
    void (*sa_sigaction)(int, siginfo_t *, void *);
    sigset_t sa_mask;        /* блокируемые сигналы */
    int sa_flags;            /* флаги */
    void (*sa_restorer)(void); /* устаревшее поле. не входит в POSIX */
}
```

В поле `sa_handler` определяется действие, которое должно выполняться при получении сигнала. Как и для вызова `signal()`, данное поле может содержать

флаг `SIG_DFL`, обозначающий действие по умолчанию, флаг `SIG_IGN`, приказывающий ядру игнорировать сигнал для процесса, или указатель на функцию обработки сигнала. У функции тот же прототип, что и у обработчика сигналов, устанавливаемого `signal()`:

```
void my_handler (int signo);
```

Если флаг `SA_SIGINFO` установлен в поле `sa_flags`, то тогда поле `sa_sigaction`, но не `sa_handler` диктует используемую функцию обработки сигналов. Прототип этой функции немного отличается:

```
void my_handler (int signo, siginfo_t *si, void *ucontext);
```

Данная функция получает в качестве первого параметра номер сигнала, в качестве второго — структуру `siginfo_t` и в качестве третьего — структуру `ucontext_t` (приведенную к указателю `void`). Она не возвращает никакого значения. Структура `siginfo_t` предоставляет множество информации обработчику сигналов; мы рассмотрим ее чуть позже.

Обратите внимание, что на некоторых машинных архитектурах (и, возможно, в других системах Unix) `sa_handler` и `sa_sigaction` находятся в объединении, и не следует назначать значения обоим полям.

В поле `sa_mask` находится набор сигналов, которые система должна блокировать на протяжении выполнения обработчика сигналов. Это позволяет программистам внедрять должную защиту от повторной входимости среди множества обработчиков сигналов. Сигнал, который в данный момент обрабатывается, также блокируется, если только в `sa_flags` не содержится флаг `SA_NODEFER`. Невозможно заблокировать `SIGKILL` или `SIGSTOP`; вызов бесшумно игнорирует любое из этих значений, если оно добавляется в `sa_mask`.

Поле `sa_flags` — это битовая маска, включающая ноль, один или несколько флагов, меняющих способ обработки сигнала, указанного в `signo`. Я уже упоминал флаги `SA_SIGINFO` и `SA_NODEFER`; а далее перечисляются прочие значения поля `sa_flags`:

SA_NOCLDSTOP

Если значение `signo` равно `SIGCHLD`, то данный флаг заставляет систему не предоставлять уведомление, когда дочерний процесс останавливается или возобновляется.

SA_NOCLDWAIT

Если значение `signo` равно `SIGCHLD`, то данный флаг включает *автоматическое вычищение потомков* (automatic child reaping): потомки не превращаются в зомби при завершении, а предку не приходится (и он не может) делать для них системный вызов `wait()`. Подробное обсуждение потомков, зомби и вызова `wait()` приведено в главе 5.

SA_NOMASK

Данный флаг — это устаревший не входящий в POSIX эквивалент флага `SA_NODEFER` (о котором говорилось ранее в этом разделе). Используйте вместо него `SA_NODEFER`, но будьте готовы видеть это значение в старом коде.

SA_ONESHOT

Данный флаг — это устаревший не входящий в POSIX эквивалент флага `SA_RESETHAND` (который есть далее в этом списке). Используйте вместо него `SA_RESETHAND`, но будьте готовы видеть это значение в старом коде.

SA_ONSTACK

Данный флаг заставляет систему вызывать указанный обработчик сигналов на *альтернативном стеке сигналов* (*alternative signal stack*), указанном при помощи `sigaltstack()`. Если вы не предоставляете альтернативный стек, то используется стек по умолчанию, то есть система работает так, как если бы вы вообще не указывали этот флаг. Альтернативные стеки сигналов применяются редко, хотя они бывают полезны в некоторых приложениях `pthreads`, использующих небольшие стеки потоков, которые могут переполняться определенными действиями обработчика сигналов.

SA_RESTART

Данный флаг включает перезапуск системных вызовов, прерванных сигналами, в стиле BSD.

SA_RESETHAND

Данный флаг включает «одноразовый» режим. Для указанного сигнала восстанавливается поведение по умолчанию, как только обработчик сигнала возвращает результат.

Поле `sa_restorer` — это устаревшее поле, которое более не применяется в Linux. Оно не является частью POSIX. Считайте, что его вообще нет, и никогда не трогайте его.

Вызов `sigaction()` в случае успеха возвращает значение 0. В случае ошибки он возвращает -1 и присваивает переменной `errno` один из следующих кодов ошибки:

EFAULT

Аргумент `act` или `oldact` содержит недопустимый указатель.

EINVAL

Аргумент `signo` содержит недопустимый сигнал, `SIGKILL` или `SIGSTOP`.

Структура `siginfo_t`

Структура `siginfo_t` также определяется в `<sys/signal.h>`, как показано далее:

```
typedef struct siginfo_t {
    int si_signo;      /* номер сигнала */
    int si_errno;      /* значение errno */
    int si_code;       /* код сигнала */
    pid_t si_pid;      /* PID отправляющего процесса */
    uid_t si_uid;      /* действительный UID отправляющего процесса */
    int si_status;     /* значение выхода или сигнал */
    clock_t si_utime; /* потребленное пользовательское время */
    clock_t si_stime; /* потребленное системное время */
```

```
    si_val_t si_value; /* значение полезной нагрузки сигнала */
    int si_int;        /* сигнал POSIX.1b */
    void *si_ptr;     /* сигнал POSIX.1b */
    void *si_addr;    /* местоположение в памяти, вызвавшее сбой */
    int si_band;      /* событие полосы */
    int si_fd;        /* дескриптор файла */
};


```

Эта структура заполнена информацией, передаваемой обработчику сигналов (если вы используете `sa_sigaction` вместо `sa_sighandler`). В мире современной вычислительной техники многие считают сигнальную модель Unix ужасным способом взаимодействия между процессами. Вероятно, проблема заключается в том, что эти люди застряли на уровне `signal()`, тогда как им следовало бы использовать `sigaction()` с флагом `SA_SIGINFO`. Структура `sigaction_t` позволяет выжать из сигналов намного больше функциональности.

В этой структуре содержится множество интересных данных, включая информацию о процессе, отправившем сигнал, и о причине сигнала. Далее приведено детальное описание каждого из полей:

si_signo

Номер сигнала. В обработчике сигналов первый аргумент также предоставляет эту информацию (и избегает разыменования указателей).

si_errno

Если значение этого поля не равно нулю, то это код ошибки, связанный с сигналом. Данное поле применяется для всех сигналов.

si_code

Объяснение, почему и откуда процесс получил данный сигнал (например, от вызова `kill()`). Мы рассмотрим возможные значения в следующем разделе. Данное поле применяется для всех сигналов.

si_pid

Для сигнала `SIGCHLD` это идентификатор PID процесса, который завершился.

si_uid

Для сигнала `SIGCHLD` это UID владельца процесса, который завершился.

si_status

Для сигнала `SIGCHLD` это статус выхода процесса, который завершился.

si_utime

Для сигнала `SIGCHLD` это пользовательское время, потребленное процессом, который завершился.

si_stime

Для сигнала `SIGCHLD` это системное время, потребленное процессом, который завершился.

si_value

Объединение `si_int` и `si_ptr`.

si_int

Для сигналов, отправленных через `sigqueue()` (см. раздел «Отправка сигнала с полезной нагрузкой» далее в этой главе), это переданная полезная нагрузка, указанная в виде целочисленного значения.

si_ptr

Для сигналов, отправленных через `sigqueue()` (см. раздел «Отправка сигнала с полезной нагрузкой» далее в этой главе), это переданная полезная нагрузка, указанная в виде указателя `void`.

si_addr

Для сигналов `SIGBUS`, `SIGFPE`, `SIGILL`, `SIGSEGV` и `SIGTRAP` указатель `void` содержит адрес сбоя, ставшего причиной сигнала. Например, в случае `SIGSEGV` в данном поле находится адрес нарушения доступа к памяти (то есть очень часто – значение `NULL!`).

si_band

Для сигнала `SIGPOLL` это информация о внеполосных данных и приоритете для дескриптора файла, указанного в поле `si_fd`.

si_fd

Для сигнала `SIGPOLL` это дескриптор файла для файла, операция которого завершилась.

Поля `si_value`, `si_int` и `si_ptr` – это особенно сложные темы, так как процесс может применять эти поля для передачи произвольных данных другому процессу. Таким образом, их можно использовать для отправки либо простого целочисленного значения, либо указателя на структуру данных (обратите внимание, что указатель не имеет особого смысла, если процессы не используют совместно одно адресное пространство). Эти поля обсуждаются далее в разделе «Отправка сигнала с полезной нагрузкой».

Стандарт POSIX гарантирует, что только первые три поля применяются для всех сигналов. К остальным полям следует обращаться только при обработке соответствующего сигнала. Например, поле `si_fd` можно использовать, только если вы обрабатываете сигнал `SIGPOLL`.

Великолепный мир `si_code`

Поле `si_code` содержит причину сигнала. Для сигналов, отправленных пользователем, оно указывает на то, как сигнал был отправлен. Для сигналов, отправленных ядром, поле указывает, почему был отправлен сигнал.

Следующие значения `si_code` допустимы для любого сигнала. Они сообщают, как/почему сигнал был отправлен:

SI_ASYNCIO

Сигнал был отправлен по причине завершения асинхронного ввода-вывода (см. главу 5).

SI_KERNEL

Сигнал был поднят ядром.

SI_MESGQ

Сигнал был отправлен по причине изменения состояния очереди сообщений POSIX (она в этой книге не рассматривается).

SI_QUEUE

Сигнал был отправлен `sigqueue()` (см. следующий раздел).

SI_TIMER

Сигнал был отправлен по причине истечения таймера POSIX (см. главу 10).

SI_TKILL

Сигнал был отправлен `tkill()` или `tkill()`. Эти системные вызовы используются библиотеками обеспечения потоковой обработки и не рассматриваются в данной книге.

SI_SIGIO

Сигнал был отправлен вследствие постановки в очередь сигнала SIGIO.

SI_USER

Сигнал был отправлен вызовом `kill()` или `raise()`.

Следующие значения `si_code` допустимы только для сигнала SIGBUS. Они указывают тип произошедшей аппаратной ошибки:

BUS_ADRALN

В процессе произошла ошибка выравнивания (выравнивание обсуждалось в главе 8).

BUS_OBJERR

Процесс обратился к недопустимому физическому адресу.

BUS_OBJERR

Процесс вызывал какую-то другую форму аппаратной ошибки.

Для сигнала SIGCHLD следующие значения указывают, что сделал потомок, вследствие чего его предку был отправлен сигнал:

CLD_CONTINUED

Потомок был остановлен, но затем возобновлен.

CLD_DUMPED

Ненормальное завершение потомка.

CLD_EXITED

Нормальное завершение потомка через вызов `exit()`.

CLD_KILLED

Потомок был убит.

CLD_STOPPED

Потомок остановился.

CLD_TRAPPED

Потомок попал в ловушку.

Следующие значения допустимы только для сигнала SIGFPE. Они объясняют тип произошедшей арифметической ошибки:

FPE_FLTDIV

Процесс выполнил операцию с плавающей точкой, которая привела к делению на ноль.

FPE_FLTOVF

Процесс выполнил операцию с плавающей точкой, которая привела к переполнению.

FPE_FLTINV

Процесс выполнил недопустимую операцию с плавающей точкой.

FPE_FLTRES

Процесс выполнил операцию с плавающей точкой, которая дала неточный или недопустимый результат.

FPE_FLTSUB

Процесс выполнил операцию с плавающей точкой, которая привела к появлению выпадающего из диапазона нижнего индекса.

FPE_FLTUND

Процесс выполнил операцию с плавающей точкой, которая привела к потере значимости.

FPE_INTDIV

Процесс выполнил целочисленную операцию, которая привела к делению на ноль.

FPE_INTOVF

Процесс выполнил целочисленную операцию, которая привела к переполнению.

Следующие значения si_code допустимы только для сигнала SIGILL. Они объясняют природу выполнения недопустимой инструкции:

ILL_ILLADR

Процесс попытался войти в недопустимый режим адресации.

ILL_ILLOPC

Процесс попытался выполнить недопустимый код операции (opcode).

ILL_ILLOPN

Процесс попытался выполнить недопустимый operand.

ILL_PRVOPC

Процесс попытался выполнить привилегированный код операции (opcode).

ILL_PRVREG

Процесс попытался выполнить привилегированный регистр.

ILL_ILLTRP

Процесс попытался войти в недопустимую ловушку.

Для всех этих значений *si_addr* указывает на адрес сбоя.

Для сигнала SIGPOLL следующие значения идентифицируют событие ввода-вывода, которое привело к созданию сигнала:

POLL_ERR

Произошла ошибка ввода-вывода.

POLL_HUP

Устройство зависло или сокет отсоединился.

POLL_IN

В файле есть данные, доступные для чтения.

POLL_MSG

Доступно сообщение.

POLL_OUT

В файл можно записывать.

POLL_PRI

В файле есть доступные для чтения высокоприоритетные данные.

Следующие коды допустимы только для сигнала SIGSEGV и описывают два типа недопустимого доступа к памяти:

SEGV_ACCERR

Процесс обратился к допустимому региону памяти недопустимым способом, то есть процесс нарушил разрешения на доступ к памяти.

SEGV_MAPERR

Процесс обратился к недопустимому региону памяти.

Для обоих этих значений *si_addr* содержит адрес сбоя.

Для сигнала SIGTRAP следующие два значения *si_code* указывают на тип сработавшей ловушки:

TRAP_BRKPT

Процесс встретил точку останова.

TRAP_TRACE

Процесс встретил ловушку слежения.

Обратите внимание, что *si_code* – это поле значения, а не битовое поле.

Отправка сигнала с полезной нагрузкой

Как мы видели в предыдущем разделе, обработчикам сигналов, зарегистрированным с флагом SA_SIGINFO, отправляется параметр типа *siginfo_t*. Эта струк-

тура включает поле с именем `si_value`, которое может содержать необязательную полезную нагрузку, передаваемую от создателя сигнала его получателю.

Функция `sigqueue()`, определенная в POSIX, позволяет процессу отправлять сигналы с такой нагрузкой:

```
#include <signal.h>

int sigqueue (pid_t pid,
              int signo,
              const union sigval value);
```

Вызов `sigqueue()` работает аналогично `kill()`. В случае успеха сигнал, указанный при помощи аргумента `signo`, ставится в очередь к процессу или группе процессов, идентифицирующейся значением `pid`, а функция возвращает 0. Полезная нагрузка сигнала передается в параметре `value`, представляющем собой объединение целочисленного значения и указателя `void`:

```
union sigval {
    int sival_int;
    void *sival_ptr;
};
```

В случае ошибки вызов возвращает `-1` и присваивает переменной `errno` одно из следующих значений:

`EINVAL`

Значение аргумента `signo` соответствует недопустимому сигналу.

`EPERM`

У вызывающего процесса отсутствуют разрешения на отправку сигналов любым из запрошенных процессов. Разрешения, необходимые для отправки сигнала, аналогичны тем, которые требует `kill()` (см. раздел «Отправка сигнала» ранее в этой главе).

`ESRCH`

Процесс или группа процессов, указанная при помощи аргумента `pid`, не существует или, в случае процесса, является зомби.

Как и с вызовом `kill()`, для тестирования можно отправлять при помощи параметра `signo` нулевой сигнал (значение 0).

Пример

В этом примере процессу с `pid` 1722 отправляется сигнал `SIGUSR2` с полезной нагрузкой в виде целочисленного значения 404:

```
sigval value;
int ret;

value.sival_int = 404;

ret = sigqueue (1722, SIGUSR2, value);
if (ret)
    perror ("sigqueue").
```

Если процесс 1722 обрабатывает сигнал SIGUSR2 обработчиком SA_SIGINFO, то для него значение signo равно SIGUSR2, si->si_int – 404, а si->si_code – SI_QUEUE.

Заключение

У сигналов плохая репутация среди многих программистов в Unix. Это старый, вышедший из употребления механизм общения ядра и пользователя, который, в лучшем случае, может служить примитивной формой взаимодействия между процессами. В мире многопоточных программ и циклов событий сигналам места нет.

Тем не менее мы нуждаемся в сигналах. Сигналы представляют собой единственный способ получения многих уведомлений ядра (таких, как уведомление о выполнении недопустимого кода операции). Помимо этого, при помощи сигналов Unix (и, следовательно, Linux) завершает процессы и управляет взаимоотношениями между предками и потомками. Следовательно, нам никуда от сигналов не деться.

Одна из основных причин умаления достоинства сигналов заключается в том, что очень сложно написать правильный обработчик сигналов, который будет безопасен с точки зрения вопросов повторной входимости. Однако если вы будете делать свои обработчики простыми и использовать только функции, перечисленные в табл. 9.2 (если вообще собираетесь включать в них какие-то функции), то они будут достаточно безопасными.

Еще одна брешь в броне сигналов – то, что многие программисты до сих пор применяют для управления сигналами вызовы `signal()` и `kill()` вместо `sigaction()` и `sigqueue()`. Как показали последние два раздела, сигналы обладают значительно большей мощью и выразительной способностью, когда используются обработчики сигналов в стиле SA_SIGINFO. Хотя я сам не большой поклонник сигналов – по мне, так лучше бы на смену сигналам пришли основанные на дескрипторах файлов и поддерживающие опрашивание механизмы, которые в действительности уже рассматриваются как возможный вариант для будущих версий ядра Linux, – но все же обход их недостатков и применение расширенных интерфейсов управления сигналами Linux позволяет избавиться от большей части проблем (если уж не от нытья).

10 Время

Время служит множеству целей в современной операционной системе, и многим программам необходимо отслеживать его. Ядро измеряет течение времени тремя разными способами:

Стенное время (или реальное время) (wall time, real time)

Это фактическое время и дата в реальном мире — то есть то время, которое вы могли бы увидеть на часах, висящих на стене. Процессы используют стечное время при взаимодействии с пользователем и создании временных меток событий.

Время процесса (process time)

Это время, потребленное процессом либо напрямую в коде из пользовательского пространства, либо косвенно — когда ядро работало от имени процесса. Процессам требуется этот тип времени в основном для профилирования и статистики — например, чтобы измерить, как долго выполнялась определенная операция. Стенное время не подходит для оценки поведения процессов, так как, учитывая многозадачную природу Linux, для конкретной операции время процесса может оказаться намного меньше, чем стечное время. Процесс также может потратить достаточно много циклов, ожидая ввода-вывода (в частности, ввода с клавиатуры).

Монотонное время (monotonic time)

Этот источник времени строго линейно возрастает. Большинство операционных систем, включая Linux, используют продолжительность работы системы (время с момента загрузки). Стенное время может меняться — его может устанавливать пользователь или регулировать система, — а дополнительные неточности могут появляться из-за, скажем, корректировочных секунд. Продолжительность работы системы, с другой стороны, — это детерминированное и неизменяемое представление времени. Важным аспектом источника монотонного времени является не текущее значение, а гарантия, что источник времени строго линейно возрастает и, следовательно, позволяет вычислять разницу между двумя замерами.

Следовательно, монотонное время подходит для вычисления относительного времени, тогда как стечное время идеально для оценки абсолютного времени.

Эти три способа оценки времени могут представляться одним из двух форматов:

Относительное время

Это значение, которое является относительным по сравнению с каким-то контрольным значением, например, текущим моментом. Например, *через 5 секунд от настоящего момента или 10 минут назад*.

Абсолютное время

Время безотносительно контрольных значений, например 0 часов 25 марта 1968 года.

И у относительных, и у абсолютных форм времени много вариантов применения. Процессу может требоваться отменять запрос через 500 миллисекунд, обновлять экран 60 раз в секунду или отмечать, когда с момента начала операции проходит 7 секунд. Все это вычисления относительного времени. И наоборот, календарь может сохранять для пользователя дату вечеринки в мантиях как 8 февраля, файловая система может записывать полную дату и время момента создания файла (а не «пять секунд назад»), а часы пользователя могут отображать грекорианское время, а не число секунд, прошедших со времени загрузки системы.

Системы Unix представляют абсолютное время как число секунд, прошедшее с начала эпохи (*epoch*), которое определяется как 00:00:00 по универсальному глобальному времени утром 1 января 1970 года. Универсальное глобальное время (Universal Time, Coordinated) приблизительно равно среднему времени по Гринвичу (Greenwich Mean Time, GMT или Zulu time). Любопытно, что это означает, что в Unix даже абсолютное время на определенном уровне относительно. В Unix для хранения «секунд с начала эпохи» существует специальный тип данных, с которым мы познакомимся в следующем разделе.

Операционные системы отслеживают течение времени при помощи *программных часов* (software clock) — часов, поддерживаемых ядром в программном обеспечении. Ядро реализует периодический таймер, называемый *системным таймером* (system timer), который тикает с определенной частотой. Когда интервал таймера завершается, ядро увеличивает счетчик прошедшего времени на одну единицу, называемую *тиком* (tick) или *мигом* (jiffy). Счетчик тиков называется *счетчиком мигов* (jiffies counter). Бывший ранее 32-разрядным значением, начиная с версии ядра Linux 2.6, счетчик мигов теперь представляет собой 64-разрядный счетчик¹.

¹ В будущих версиях ядра Linux тики могут исчезнуть либо могут быть реализованы «динамические тики», чтобы ядро не отслеживало явное значение счетчика мигов. Вместо этого зависящие от времени операции ядра будут выполняться с использованием динамически реализуемых таймеров, а не системного таймера.

В Linux частота системного таймера называется HZ, так как его представляет определение препроцессора с тем же именем. Значение HZ зависит от архитектуры и не является частью ABI Linux, то есть программа не может зависеть или ожидать любое конкретное значение. Исторически в архитектуре x86 использовалось значение 100, то есть системный счетчик срабатывал 100 раз в секунду (и частота его составляла 100 герц). Это давало значение мига, равное 0,01 секунды — 1/HZ секунды. С выпуском ядра Linux 2.6 разработчики ядра подняли значение HZ до 1000, уменьшив каждый миг до 0,001 секунды. Однако в версии 2.6.13 и выше значение HZ равно 250, а значение мига — 0,004 секунды¹. Со значением HZ связан определенный компромисс: большие значения дают более высокое разрешение, но создают большую нагрузку таймера.

Хотя процессы не должны полагаться ни на какое фиксированное значение, HZ в POSIX определяет механизм для проверки частоты системного таймера во время выполнения:

```
long hz;  
  
hz = sysconf (_SC_CLK_TCK);  
if (hz == -1)  
    perror ("sysconf"); /* не должно никогда происходить */
```

Этот интерфейс полезен в том случае, когда программе нужно определить разрешение системного таймера, но нет необходимости преобразовывать значения системного времени в секунды, так как большинство интерфейсов POSIX экспортят уже преобразованные оценки времени или оценки, масштабированные до фиксированной частоты, независимой от HZ. В отличие от HZ, эта фиксированная частота является частью системного ABI; в архитектурах x86 ее значение равно 100. Функции POSIX, возвращающие время в терминах тиков часов, используют для представления фиксированной частоты значение CLOCKS_PER_SEC.

Иногда события выключают компьютер. В некоторых случаях компьютеры даже отключаются от сети, и все же после загрузки время на них остается правильным. Это происходит благодаря тому, что в большинстве компьютеров есть питающиеся от аккумулятора *аппаратные часы* (hardware clock), сохраняющие время и дату на тот момент, пока компьютер выключен. Когда ядро загружается, оно инициализирует свою концепцию текущего времени, броя его у аппаратных часов. Схожим образом, когда пользователь выключает систему, ядро записывает текущее время обратно в аппаратные часы. Администраторы системы могут синхронизировать время в другие моменты при помощи команды hwclock.

Управление течением времени в системе Unix включает несколько задач, из которых лишь некоторые беспокоят каждый конкретный процесс: они включают установку и извлечение текущего значения стенного времени, вычисление времени, прошедшего с некоторого момента, засыпание на определенный ин-

¹ HZ также является теперь параметром ядра времени компилирования, и в архитектуре x86 для него поддерживаются значения 100, 250 и 1000. Как бы то ни было, пользовательское пространство не должно зависеть ни от какого конкретного значения HZ.

тервал времени, выполнение высокоточных измерений времени и управление таймерами. В этой главе рассматривается полный набор задач, связанных с временем. Начнем мы со структур данных, при помощи которых в Linux представляется время.

Структуры данных для работы со временем

По мере того как системы Unix эволюционировали, реализуя собственные интерфейсы для управления временем, появлялось множество структур данных, представляющих кажущуюся простой концепцию времени. Диапазон таких структур данных включает как простые целочисленные значения, так и разнообразные многопольные структуры. Мы рассмотрим их здесь, перед тем как перейти к самим интерфейсам.

Исходное представление

Простейшая структура данных — это `time_t`, определенная в файле заголовка `<time.h>`. Первоначально тип `time_t` должен был быть непрозрачным типом. Однако в большинстве систем Unix, включая Linux, этот тип представляет собой простое приведение типа C `long`:

```
typedef long time_t;
```

`Time_t` представляет число секунд, прошедшее с начала эпохи. Обычно в ответ на это можно услышать заявление: «Очень скоро произойдет переполнение!» В действительности этого типа хватит очень надолго, дольше, чем вы можете ожидать, но переполнение все же наступит тогда, когда множество систем Unix будет продолжать активно использоваться. С учетом 32-разрядного типа `long` тип `time_t` может представлять до 2 147 483 647 секунд с начала эпохи. Это говорит о том, что проблема 2000 возникнет снова — теперь в 2038 году! Но мы все же надеемся, что к 22:14:07 понедельника 18 января 2038 года большинство систем и программного обеспечения уже будет 64-разрядными.

А теперь точность до микросекунд

Еще одна проблема, связанная с `time_t`, состоит в том, что за одну секунду может произойти очень многое. Структура `timeval` расширяет тип `time_t`, добавляя точность до микросекунд. В файле заголовка `<sys/time.h>` эта структура определяется следующим образом:

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

struct timeval {
    time_t tv_sec;          /* секунды */
    suseconds_t tv_usec;    /* микросекунды */
};
```

В поле `tv_sec` содержится число секунд, а в поле `tv_usec` — число микросекунд. Пугающий тип `suseconds_t` обычно бывает простым приведением целочисленного типа.

Еще лучше: точность до наносекунд

Неудовлетворенная разрешением с точностью до микросекунд, структура `timespec` повышает ставки до наносекунд. В файле заголовка `<time.h>` эта структура определяется следующим образом:

```
#include <time.h>

struct timespec {
    time_t tv_sec;      /* секунды */
    long tv_nsec;        /* наносекунды */
}:
```

При условии наличия выбора интерфейсы предпочитают точность до наносекунд¹, а не до микросекунд. Следовательно, с момента появления структуры `timespec` большинство связанных со временем интерфейсов переключились на нее, получив таким образом более высокую точность. Однако, как мы увидим далее, одна важная функция до сих пор использует `timeval`.

На практике ни одна из структур обычно не предоставляет объявленную точность, так как системный таймер не обеспечивает разрешение до наносекунд или даже микросекунд. Тем не менее лучше, чтобы у интерфейса была возможность использовать максимальное разрешение, чтобы он мог справляться с любой точностью времени, которую способна предоставлять система.

Разбор времени

Некоторые функции, которые мы рассмотрим, преобразуют время Unix в строковые значения и, наоборот, программно строят строки, представляющие определенную дату. Для упрощения этого процесса стандарт C предоставляет структуру `tm`, содержащую «разобранное на части» время в более удобном для восприятия человеком формате. Данная структура также определяется в файле заголовка `<time.h>`:

```
#include <time.h>

struct tm {
    int tm_sec;          /* секунды */
    int tm_min;          /* минуты */
    int tm_hour;         /* часы */
    int tm_mday;         /* день месяца */
    int tm_mon;          /* месяц */
    int tm_year;         /* год */
    int tm_wday;         /* день недели */
    int tm_yday;         /* день года */
```

¹ Помимо этого, в структуре `timespec` наконец избавились от бессмысленного типа `suseconds_t` в пользу простого и непритязательного типа `long`.

```
int tm_isdst;           /* летнее время? */
#endif /* _BSD_SOURCE
long tm_gmtoff;        /* смещение часового пояса от GMT */
const char *tm_zone;   /* аббревиатура часового пояса */
#endif /* _BSD_SOURCE */
};
```

Благодаря структуре `tm` очень просто сказать, относится значение `time_t`, скажем 314 159, к воскресенью или субботе (в действительности, к воскресенью). В терминах пространства очевидно, что такая структура — не очень хороший вариант для представления даты и времени, но ее удобно применять для преобразования пользовательских значений.

Структура включает следующие поля:

`tm_sec`

Число секунд после минуты. Обычно это значение в диапазоне от 0 до 59, но оно может достигать 61 и даже больше, указывая до двух корректировочных секунд.

`tm_min`

Число минут после часа. Значение в диапазоне от 0 до 59.

`tm_hour`

Число часов после полуночи. Значение в диапазоне от 0 до 23.

`tm_mday`

День месяца. Значение в диапазоне от 0 до 31. В POSIX значение 0 не определяется, однако в Linux оно используется для обозначения последнего дня предыдущего месяца.

`tm_mon`

Число месяцев с января. Значение в диапазоне от 0 до 11.

`tm_year`

Число лет с 1900 года.

`tm_wday`

Число дней с воскресенья. Значение в диапазоне от 0 до 6.

`tm_yday`

Число дней с 1 января. Значение в диапазоне от 0 до 365.

`tm_isdst`

Специальное значение, указывающее, действует ли летнее время на момент, определяемый прочими полями. Если это значение больше нуля, то летнее время действует. Если оно равно 0, то летнее время не действует. Если это значение меньше нуля, то состояние летнего времени неизвестно.

`tm_gmtoff`

Смещение в секундах текущего часового пояса от среднего времени по Гринвичу. Это поле присутствует только в том случае, если `_BSD_SOURCE` определяется до подключения `<time.h>`.

tm_zone

Аббревиатуры для текущего часового пояса, например EST. Это поле присутствует только в том случае, если `_BSD_SOURCE` определяется до подключения `<time.h>`.

Тип для времени процесса

Тип `clock_t` представляет тики часов. Это целочисленный тип, зачастую `long`. В зависимости от интерфейса, тики, которые представляет `clock_t`, обозначают либо фактическую частоту таймера системы (`HZ`), либо `CLOCKS_PER_SEC`.

Часы POSIX

Некоторые из системных вызовов, которые обсуждаются в этой главе, используют часы POSIX — стандарт для реализации и представления источников времени. Тип `clockid_t` представляет специальные часы POSIX, четыре типа которых поддерживает Linux:

CLOCK_MONOTONIC

Часы с монотонно увеличивающимся значением, которые не может устанавливать ни один процесс. Они представляют время, прошедшее с некоторого неопределенного начального момента, такого, как загрузка системы.

CLOCK_PROCESS_CPUTIME_ID

Часы с высоким разрешением, которые для каждого процесса предоставляет процессор. Например, в архитектуре i386 эти часы используют регистр *счетчика временной метки* (*timestamp counter*, TSC).

CLOCK_REALTIME

Часы, поддерживающие системное реальное (стенное) время. Для установки этих часов требуются специальные привилегии.

CLOCK_THREAD_CPUTIME_ID

Похожи на часы, предоставляемые каждому процессу, но уникальны для каждого потока в процессе.

Стандарт POSIX определяет все четыре источника времени, но требует наличия только `CLOCK_REALTIME`. Таким образом, хотя Linux надежно предоставляет все четыре вида часов, в переносимом коде следует полагаться только на `CLOCK_REALTIME`.

Разрешение источника времени

В стандарте POSIX определяется функция `clock_getres()` для проверки разрешения указанного источника времени:

```
#include <time.h>
```

```
int clock_getres (clockid_t clock_id,
                  struct timespec *res);
```

Успешный вызов `clock_getres()` сохраняет разрешение часов, указанных при помощи аргумента `clock_id`, в параметре `res`, если он не равен `NULL`, и возвращает значение `0`. В случае ошибки функция возвращает `-1` и присваивает переменной `errno` один из следующих двух кодов ошибки:

`EFAULT`

Аргумент `res` содержит недопустимый указатель.

`EINVAL`

`Clock_id` — это недопустимый источник времени в данной системе.

Код в следующем примере сообщает разрешение четырех источников времени, о которых я рассказал в предыдущем разделе:

```
clockid_t clocks[] = {
    CLOCK_REALTIME,
    CLOCK_MONOTONIC,
    CLOCK_PROCESS_CPUTIME_ID,
    CLOCK_THREAD_CPUTIME_ID,
    (clockid_t) -1 };
int i.

for (i = 0; clocks[i] != (clockid_t) -1; i++) {
    struct timespec res;
    int ret;

    ret = clock_getres (clocks[i], &res);
    if (ret)
        perror ("clock_getres");
    else
        printf ("clock=%d sec=%ld nsec=%ld\n",
               clocks[i], res.tv_sec, res.tv_nsec);
}
```

В современной системе x86 вывод получается приблизительно таким:

```
clock=0 sec=0 nsec=4000250
clock=1 sec=0 nsec=4000250
clock=2 sec=0 nsec=1
clock=3 sec=0 nsec=1
```

Обратите внимание, что `4 000 250` наносекунд — это `4` миллисекунды или `0,004` секунды. В свою очередь, `0,004` секунды — это разрешение системных часов x86 при условии, что значение `HZ` равно `250`, как говорилось в первом разделе этой главы. Таким образом, мы видим, что оба типа часов, `CLOCK_REALTIME` и `CLOCK_MONOTONIC`, привязаны к мигам и разрешению, обеспечиваемому системным таймером. И наоборот, `CLOCK_PROCESS_CPUTIME_ID` и `CLOCK_THREAD_CPUTIME_ID` используют источник времени с более высоким разрешением — на этой машине x86 это `TSC`, который, как мы видим, обеспечивает разрешение до наносекунд.

В Linux (и большинстве других систем Unix) все функции, использующие часы POSIX, требуют компоновки результирующего объектного файла с использованием `librt`. Например, для того чтобы скомпилировать предыдущий фрагмент в полный исполняемый файл, вам может потребоваться следующая команда:

```
$ gcc -Wall -W -O2 -lrt -g -o snippet snippet.c
```

Получение текущего времени дня

У приложений может быть несколько причин для того, чтобы пытаться узнать текущее время и дату: показать ее пользователю, вычислить относительное или затраченное время, установить временную метку события и так далее. Самый простой и исторически наиболее распространенный способ получения текущего времени включает применение функции `time()`:

```
#include <time.h>
time_t time (time_t *t);
```

Вызов `time()` возвращает текущее время в формате количества секунд, прошедшего с начала эпохи. Если значение параметра `t` не равно `NULL`, то функция также записывает текущее время в предоставленный ей указатель.

В случае ошибки функция возвращает значение `-1` (приведенное к типу `time_t`) и соответствующим образом устанавливает переменную `errno`. Единственная возможная ошибка — это `EFAULT`, указывающая, что `t` содержит недопустимый указатель.

Например:

```
time_t t;
printf ("текущее время: %ld\n", (long) time (&t));
printf ("то же значение: %ld\n", (long) t);
```

БЕЗЫСКУСНЫЙ ПОДХОД КО ВРЕМЕНИ

Представление в формате `time_t` — «количество секунд, прошедшее в начала эпохи» — в действительности не является фактическим числом секунд, которое прошло с того судьбоносного момента во времени. Вычисления в Unix предполагают, что високосные годы — это все годы, номера которых делятся на четыре, но полностью игнорируют корректировочные секунды. Смысл представления `time_t` не в том, что оно наиболее точное, а в том, что оно непротиворечивое, — и это самое главное.

Лучший интерфейс

Функция `gettimeofday()` расширяет возможности `time()`, обеспечивая разрешение до микросекунд:

```
#include <sys/time.h>
int gettimeofday (struct timeval *tv,
                  struct timezone *tz);
```

Успешный вызов `gettimeofday()` помещает текущее время в структуру `timeval`, на которую указывает `tv`, и возвращает значение `0`. Структура `timezone` и параметр `tz` устарели, и их не следует использовать в Linux. Всегда передавайте в качестве параметра `tz` значение `NULL`.

В случае ошибки вызов возвращает значение `-1` и присваивает переменной `errno` значение `EFAULT`; это единственная возможная ошибка, указывающая, что параметр `tv` или `tz` содержит недопустимый указатель. Например:

```
struct timeval tv;
int ret;
```

```

ret = gettimeofday (&tv, NULL);
if (ret)
    perror ("gettimeofday");
else
    printf ("секунды=%ld микросекунды=%ld\n",
           (long) tv.sec, (long) tv.usec);

```

Структура `timezone` не нужна, так как ядро не поддерживает часовые пояса, а `glibc` отказывается использовать поле `tz_dsttime` структуры `timezone`. Мы обсудим манипулирование часовыми поясами в следующем разделе.

Расширенный интерфейс

Стандарт POSIX предоставляет интерфейс `clock_gettime()` для получения времени от определенного источника времени. Еще одним преимуществом является то, что данная функция обеспечивает точность до наносекунд:

```

#include <time.h>

int clock_gettime (clockid_t clock_id,
                   struct timespec *ts);

```

В случае успеха вызов возвращает значение 0 и сохраняет текущее время указанного при помощи аргумента `clock_id` источника времени в аргументе `ts`. В случае ошибки вызов возвращает -1 и присваивает переменной `errno` одно из следующих значений:

EFAULT

Аргумент `ts` содержит недопустимый указатель.

EINVAL

`Clock_id` — это недопустимый источник времени в данной системе.

Код в следующем примере извлекает текущее время всех четырех стандартных источников времени:

```

clockid_t clocks[] = {
    CLOCK_REALTIME,
    CLOCK_MONOTONIC,
    CLOCK_PROCESS_CPUTIME_ID,
    CLOCK_THREAD_CPUTIME_ID,
    (clockid_t) -1 };
int i;
for (i = 0; clocks[i] != (clockid_t) -1; i++) {
    struct timespec ts;
    int ret;
    ret = clock_gettime (clocks[i], &ts);
    if (ret)
        perror ("clock_gettime");
    else
        printf ("часы=%d секунды=%ld наносекунды=%ld\n",
               clocks[i], ts.tv_sec, ts.tv_nsec);
}

```

Получение времени процесса

Системный вызов `times()` позволяет получить для выполняющегося процесса и его потомков время процесса в тиках часов:

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

struct tms {
    clock_t tms_utime: /* потребленное пользовательское время */
    clock_t tms_stime: /* потребленное системное время */
    clock_t tms_cutime: /* пользовательское время, */
                         /* потребленное потомками */
    clock_t tms_cstime: /* системное время, потребленное потомками */
};

clock_t times (struct tms *buf);
```

В случае успеха вызов заполняет предоставленную пользователем структуру `tms`, на которую указывает аргумент `buf`, данными о времени процесса, потребленномзывающим процессом и его потомками. Это время делится на пользовательское время и системное время. *Пользовательское время* (user time) – это время, затраченное на выполнение кода в пользовательском пространстве. *Системное время* (system time) – это время, затраченное на выполнение кода в пространстве ядра, например, во время системного вызова или обращения к несуществующей странице. Значения времени для каждого из потомков включаются в отчет только после того, как потомок завершается, а предок вызывает для этого процесса `waitpid()` (или аналогичную функцию). Вызов возвращает число тиков часов, монотонно увеличивающееся, начиная с произвольного момента в прошлом. Когда-то этой контрольной точкой была загрузка системы, – таким образом функция `times()` возвращала продолжительность работы системы в тиках, – но сейчас контрольная точка находится приблизительно в 429 миллионах секунд до загрузки системы. Разработчики ядра реализовали это изменение, чтобы захватить код ядра, который не мог обрабатывать обработка времени работы системы и попадание в ноль. Абсолютное значение данной функции не имеет никакой ценности; относительные изменения между двумя вызовами, однако, могут давать полезную информацию.

В случае ошибки вызов возвращает `-1` и присваивает переменной `errno` подходящее значение. В Linux единственным возможным кодом ошибки является `EFAULT`, указывающий, что аргумент `buf` содержит недопустимый указатель.

Установка текущего времени дня

В предыдущих разделах я описывал, как узнавать время, но приложениям периодически также требуется устанавливать текущее время и дату в соответствии с указанным значением. Практически всегда существует утилита, предназначенная исключительно для этой цели, такая, как `date`.

Эквивалентом `time()` для установки времени является системный вызов `stime()`:

```
#define _SVID_SOURCE  
#include <time.h>
```

```
int stime (time_t *t);
```

Успешный вызов `stime()` устанавливает системное время, присваивая ему значение, на которое указывает аргумент `t`, и возвращает значение 0. Вызов требует, чтобы у вызывающего пользователя была характеристика `CAP_SYS_TIME`. Как правило, она есть только у пользователя `root`.

В случае ошибки вызов возвращает -1 и присваивает переменной `errno` значение `EFAULT`, указывающее, что `t` содержит недопустимый указатель, или `EPERM`, указывающее, что вызывающий пользователь не обладает характеристикой `CAP_SYS_TIME`.

Применять этот системный вызов очень просто:

```
time_t t = 1;  
int ret;  
  
/* установить время через одну секунду после начала эпохи */  
ret = stime (&t);  
if (ret)  
    perror ("stime");
```

Мы рассмотрим функции, упрощающие преобразование удобных для человека форм времени в тип `time_t`, в следующем разделе.

Установка времени с точностью

Эквивалентом вызова `gettimeofday()` является вызов `settimeofday()`:

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

```
int gettimeofday (const struct timeval *tv ,  
                  const struct timezone *tz);
```

Успешный вызов `settimeofday()` устанавливает системное время, указанное при помощи аргумента `tv`, и возвращает значение 0. Как и для `gettimeofday()`, лучше всего передавать в качестве аргумента `tz` значение `NULL`. В случае ошибки вызов возвращает -1 и присваивает переменной `errno` одно из следующих значений:

`EFAULT`

`Tv` или `tz` указывает на недопустимую область память.

`EINVAL`

Поле в одной из указанных структур недопустимо.

`EPERM`

У вызывающего процесса отсутствует характеристика `CAP_SYS_TIME`.

Код в следующем примере устанавливает в качестве текущего времени субботу в середине декабря 1979 года:

```

struct timeval tv = { .tv_sec = 31415926,
                     .tv_usec = 27182818 };
int ret;

ret = settimeofday (&tv, NULL);
if (ret)
    perror ("settimeofday");

```

Расширенный интерфейс для установки времени

Точно так же, как `clock_gettime()` расширяет возможности `gettimeofday()`, вызов `clock_settime()` пришел на смену `settimeofday()`:

```

#include <time.h>

int clock_settime (clockid_t clock_id,
                   const struct timespec *ts);

```

В случае успеха вызов возвращает 0, а для источника времени, указанного при помощи аргумента `clock_id`, устанавливается время, переданное в параметре `ts`. В случае ошибки вызов возвращает -1 и присваивает переменной `errno` одно из следующих значений:

`EFAULT`

`Ts` содержит недопустимый указатель.

`EINVAL`

`Clock_id` представляет недопустимый источник времени в данной системе.

`EPERM`

У процесса отсутствуют необходимые разрешения для установки указанного источника времени или указанный источник времени невозможно установить.

В большинстве систем единственный источник времени, допускающий установку, — это `CLOCK_REALTIME`. Таким образом, единственным преимуществом данной функции по сравнению с `settimeofday()` является то, что она обеспечивает точность до наносекунд (а также то, что вам не приходится иметь дело с бесполезной структурой `timezone`).

Различные действия со временем

Системы Unix и язык С предоставляют семейство функций для преобразования разобранного на части времени (представления времени в строке ASCII) во время в формате `time_t` и наоборот. Системный вызов `asctime()` превращает структуру `tm` — разбитое на части время — в строку ASCII:

```

#include <time.h>

char * asctime (const struct tm *tm);
char * asctime_r (const struct tm *tm, char *buf);

```

Он возвращает указатель на статически выделенную строку. Последующий вызов любой функции времени может переопределить эту строку, поэтому вызов `asctime()` не безопасен при работе с потоками.

Таким образом, в многопоточных программах (а также если вы чувствуете отвращение к плохо структурированным интерфейсам) следует применять системный вызов `asctime_r()`. Вместо того чтобы возвращать указатель на статически выделенную строку, данная функция использует строку, предоставленную ей через аргумент `buf`, длина которой должна составлять, как минимум, 26 символов.

Обе функции возвращают `NULL` в случае ошибки.

`Mktime()` также преобразует структуру `tm`, но в тип `time_t`:

```
#include <time.h>
```

```
time_t mktime (struct tm *tm);
```

Помимо этого, вызов `mkttime()` устанавливает часовой пояс через `tzset()`, как указано в `tm`. В случае ошибки он возвращает значение `-1` (приведенное к типу `time_t`).

Вызов `ctime()` преобразует значение типа `time_t` в его представление в виде строки ASCII:

```
#include <time.h>
```

```
char * ctime (const time_t *timep);
char * ctime_r (const time_t *timep, char *buf);
```

В случае ошибки он возвращает `NULL`. Например:

```
time_t t = time (NULL);
```

```
printf ("время всего лишь одну строку назад: %s", ctime (&t));
```

Обратите внимание на отсутствие символа новой строки. Вероятно, это довольно непривычно, но `ctime()` добавляет символ новой строки к строке, которую он возвращает.

Как и `asctime()`, `ctime()` возвращает указатель на статическую строку. Так как это небезопасно при использовании потоков, в многопоточных программах следует вместо этого вызова применять вызов `ctime_r()`, который работает с буфером, указанным при помощи аргумента `buf`. Длина буфера должна составлять, как минимум, 26 символов.

`Gmtime()` преобразует переданное значение типа `time_t` в структуру `tm`, выраженную в терминах часового пояса UTC (Universal Time Coordinated, универсальное скординированное время):

```
#include <time.h>
```

```
struct tm * gmtime (const time_t *timep);
struct tm * gmtime_r (const time_t *timep, struct tm *result);
```

В случае ошибки этот вызов возвращает `NULL`.

Данная функция статически выделяет возвращаемую структуру, и, снова, это небезопасно для потоков. В многопоточных программах нужно использовать функцию `gmtime_r()`, работающую со структурой, на которую указывает `result`.

Вызовы `localtime()` и `localtime_r()` выполняют функции, аналогичные `gmtime()` и `gmtime_r()` соответственно, но они выражают указанное значение `time_t` в терминах часового пояса пользователя:

```
#include <time.h>
```

```
struct tm * localtime (const time_t *timep);
struct tm * localtime_r (const time_t *timep, struct tm *result);
```

Как и `mktime()`, вызов `localtime()` также вызывает `tzset()` и инициализирует часовой пояс. Выполняет ли это действие `localtime_r()` — неизвестно.

Вызов `difftime()` возвращает число секунд, прошедшее между двумя указанными значениями типа `time_t`, приведенное к типу `double`:

```
#include <time.h>
```

```
double difftime (time_t time1, time_t time0);
```

Во всех системах POSIX `time_t` — это арифметический тип, а вызов `difftime()` эквивалентен следующему, если игнорировать распознавание переполнения при вычитании:

```
(double) (time1 - time0)
```

В Linux, так как `time_t` является целочисленным типом, необходимости в приведении к типу `double` нет. Однако для обеспечения переносимости программ все же используйте `difftime()`.

Подстройка системных часов

Большие и неожиданные прыжки времени стенных часов могут вносить хаос в приложения, действия которых зависят от абсолютного времени. Рассмотрим пример с командой `make`, которая собирает программные проекты, как указано в файле `Makefile`. Команда не обязательно при каждом вызове полностью перестраивает деревья источников; если бы она делала это, то в больших программных проектах одно изменение в файле приводило бы к часам ожидания, пока завершится повторная сборка программы. Вместо этого `make` проверяет временные метки, указывающие последнюю модификацию исходного файла (например, `wolf.c`) и объектного файла (например, `wolf.o`). Если исходный файл или любой из его обязательных файлов, такой, как `wolf.h`, новее объектного файла, то `make` заново собирает исходный файл в обновленный объектный файл. Если исходный файл не новее объектного, то никакие действия не предпринимаются.

Помня об этом, представьте, что может произойти, если пользователь увидит, что его часы отстают или спешат на пару часов, и выполнит `date`, чтобы обновить системные часы. Если он после этого обновит и заново сохранит файл `wolf.c`, то это может привести к проблеме. Если пользователь переведет часы на-

зад, то `wolf.c` будет выглядеть старее, чем `wolf.o`, — хотя это не так! — и повторная сборка не будет проведена.

Для того чтобы предотвращать подобные трагедии, в Unix предусмотрена функция `adjtime()`, которая постепенно корректирует текущее время в направлении указанной дельты. Смысл заключается в том, чтобы фоновые процессы, такие, как демоны NTP (Network Time Protocol, сетевой протокол времени), которые постоянно регулируют время, корректируя последствия расфазировки тактовых сигналов, применяли `adjtime()` для минимизации своего влияния на систему:

```
#define _BSD_SOURCE
#include <sys/time.h>

int adjtime (const struct timeval *delta,
             struct timeval *olddelta);
```

Успешный вызов `adjtime()` заставляет ядро медленно начать корректировать время в соответствии со значением аргумента `delta`, а затем вернуть 0. Если время, указанное при помощи параметра `delta`, положительное, то ядро ускоряет системные часы на `delta` до тех пор, пока корректировка не будет завершена. Если время, указанное при помощи параметра `delta`, отрицательное, то ядро замедляет системные часы, пока корректировка не завершится. Ядро применяет все корректировки так, чтобы значение времени всегда монотонно увеличивалось и никогда не менялось слишком резко. Даже с отрицательным значением `delta` корректировка не приводит к движению часов назад — вместо этого часы просто замедляются до тех пор, пока системное время не сойдется с исправленным временем.

Если значение `delta` не равно NULL, ядро прекращает обработку любых ранее зарегистрированных корректировок. Однако та часть корректировки, которая уже выполнена (если таковая имеется), сохраняется. Если значение `olddelta` не равно NULL, то любые ранее зарегистрированные и еще не примененные корректировки записываются в указанную структуру типа `timeval`. Если передать значение `delta`, равное NULL, и допустимое значение `olddelta`, то это позволит извлечь все текущие корректировки.

Корректировки, применяемые с помощью `adjtime()`, должны быть небольшими; идеальный случай использования — это протокол NTP, который применяет минимальные корректировки (несколько секунд). Linux поддерживает минимальный и максимальный пороги корректировки, составляющие несколько тысяч секунд в каждом направлении.

В случае ошибки `adjtime()` возвращает -1 и присваивает переменной `errno` одно из следующих значений:

`EFAULT`

Параметр `delta` или `olddelta` содержит недопустимый указатель.

`EINVAL`

Корректировка, указанная в `delta`, слишком большая или слишком маленькая.

ЕPERM

Вызывающий пользователь не обладает характеристикой CAP_SYS_TIME.

В стандарте RFC 1305 определяется значительно более мощный и, соответственно, более сложный алгоритм корректировки часов, нежели постепенная коррекция, которую применяет `adjtime()`. Linux реализует этот алгоритм в системном вызове `adjtimex()`:

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

int adjtimex (struct timex *adj);
```

Вызов `adjtimex()` считывает параметры ядра, относящиеся к времени, в структуру `timex`, на которую указывает аргумент `adj`. В зависимости от значения поля `modes` этой структуры, системный вызов может также дополнительно устанавливать некоторые параметры.

В файле заголовка `<sys/timex.h>` структура `timex` определяется следующим образом:

```
struct timex {
    int modes;          /* селектор режима */
    long offset;        /* смещение времени (микросекунд) */
    long freq;          /* смещение частоты */
    /* (масштабированное значение ppm) */
    long maxerror;     /* максимальная ошибка (микросекунд) */
    long esterror;     /* расчетная ошибка (микросекунд) */
    int status;         /* статус часов */
    long constant;     /* константа времени PLL */
    long precision;    /* точность часов (микросекунд) */
    long tolerance;    /* допустимая погрешность частоты часов */
    /* (ppm) */
    struct timeval time; /* текущее время */
    long tick;          /* микросекунд между тиками часов */
};
```

Поле `modes` может содержать ноль или несколько следующих флагов, объединенных операцией побитового ИЛИ:

ADJ_OFFSET

Установить смещение времени при помощи поля `offset`.

ADJ_FREQUENCY

Установить смещение частоты при помощи поля `freq`.

ADJ_MAXERROR

Установить максимальную ошибку при помощи поля `maxerror`.

ADJ_ESTERROR

Установить расчетную ошибку при помощи поля `esterror`.

ADJ_STATUS

Установить статус часов при помощи поля `status`.

ADJ_TIMECONST

Установить константу времени PPL (phase-locked loop, цепь фазовой синхронизации) при помощи поля `constant`.

ADJ_TICK

Установить значение тика при помощи поля `tick`.

ADJ_OFFSET_SINGLESHOT

Установить смещение времени при помощи поля `offset` один раз, применив простой алгоритм, аналогично `adjtime()`.

Если значение поля `modes` равно 0, то никакие значения не устанавливаются. Только пользователь, обладающий характеристикой `CAP_SYS_TIME`, имеет право передавать ненулевые значения `modes`; любой пользователь может передать в поле `modes` значение 0, чтобы извлечь все параметры, но установить их он не вправе.

В случае успеха `adjtimex()` возвращает текущее состояние часов в одном из следующих кодов:

TIME_OK

Часы синхронизированы.

TIME_INS

Будет вставлена корректиrovочная секунда.

TIME_DEL

Будет удалена корректиrovочная секунда.

TIME_OOP

Корректиrovочная секунда в действии.

TIME_WAIT

Корректиrovочная секунда только что сработала.

TIME_BAD

Часы не синхронизированы.

В случае ошибки `adjtimex()` возвращает `-1` и присваивает переменной `errno` один из следующих кодов ошибки:

EFAULT

Аргумент `adj` содержит недопустимый указатель.

EINVAL

Одно из полей `modes`, `offset` или `tick` содержит недопустимое значение.

EPERM

Значение поля `modes` не равно нулю, но у вызывающего пользователя нет характеристики `CAP_SYS_TIME`.

Системный вызов `adjtimex()` уникален для Linux. В приложениях, нацеленных на переносимость, следует использовать вызов `adjtime()`.

Стандарт RFC 1305 определяет сложный алгоритм, поэтому полное обсуждение системного вызова `adjtimex()` лежит за пределами данной книги. Более подробную информацию вы найдете в описании алгоритма RFC.

Засыпание и ожидание

Разнообразные функции позволяют процессу засыпать (приостанавливать выполнение) на указанный промежуток времени. Первая из подобных функций, `sleep()`, помещает вызывающий процесс в состояние ожидания на количество секунд, указанное при помощи аргумента `seconds`:

```
#include <unistd.h>
```

```
unsigned int sleep (unsigned int seconds);
```

Вызов возвращает количество секунд, в течение которых процесс *не* спал. Таким образом, успешный вызов возвращает 0, но функция может возвращать и другие значения от 0 до `seconds` включительно (если, например, сигнал прерывает состояние ожидания). Функция не устанавливает переменную `errno`. Большинство пользователей `sleep()` не волнует, как долго процесс фактически будет спать, и, следовательно, они не проверяют возвращаемое значение:

```
sleep (7); /* спать семь секунд */
```

Если приостановка процесса на все указанное время действительно необходима, то можно продолжать вызывать `sleep()` с собственным возвращаемым значением до тех пор, пока функция не вернет 0:

```
unsigned int s = 5;
```

```
/* спать пять секунд: никаких если, и или но */
while ((s = sleep (s)))
```

```
;
```

Засыпание с точностью до микросекунд

Приостановка процесса с точностью до целых секунд довольно старомодна. Секунда — это целая вечность в современной системе, и программам очень часто бывает нужно засыпать с более высоким разрешением. Познакомьтесь с `usleep()`:

```
/* версия BSD */
#include <unistd.h>
```

```
void usleep (unsigned long usec);
```

```
/* версия SUSv2 */
#define _XOPEN_SOURCE 500
#include <unistd.h>
```

```
int usleep (useconds_t usec);
```

Успешный вызов `usleep()` помещает вызывающий процесс в состояние ожидания на `usec` микросекунд. К сожалению, стандарты BSD и Single UNIX Specification расходятся во мнениях относительно прототипа функции. Вариант BSD принимает значение типа `unsigned long` и не возвращает значение. Вариант SUS принимает тип `useconds_t` и возвращает значение типа `int`. Linux следует

определению для SUS, если константа `_XOPEN_SOURCE` определена со значением 500 или выше. Если `_XOPEN_SOURCE` не определена или ее значение меньше 500, то Linux следует стандарту BSD.

Версия SUS возвращает 0 в случае успеха и -1 в случае ошибки. Допустимые значения `errno` — это `EINTR`, если состояние ожидания было прервано сигналом, и `EINVAL`, если значение `usecs` было слишком велико (в Linux допустим полный диапазон типа, поэтому вторая ошибка никогда не возникает).

Согласно спецификации, тип `useconds_t` — это целое без знака, способное хранить значения до 1 000 000.

Из-за различий между конфликтующими прототипами и того факта, что некоторые системы Unix могут поддерживать одну версию, а другие — другую, но не обе, лучше никогда явно не включать тип `useconds_t` в свой код. Для максимальной переносимости программ предполагайте, что параметр относится к типу `unsigned int`, и не полагайтесь на возвращаемое значение `usleep()`:

```
void usleep (unsigned int usec);
```

Таким образом, используйте его так:

```
unsigned int usecs = 200;
```

```
usleep (usecs);
```

Это будет работать с любым вариантом функции, и проверка ошибок все также возможна:

```
errno = 0;
usleep (1000);
if (errno)
    perror ("usleep");
```

В большинстве программ, однако, ошибки `usleep()` не проверяются и никого не интересуют.

Засыпание с точностью до наносекунд

В Linux на смену функции `usleep()` приходит функция `nanosleep()`, обеспечивающая разрешение до наносекунд и более интеллектуальный интерфейс:

```
#define _POSIX_C_SOURCE 199309
#include <time.h>
```

```
int nanosleep (const struct timespec *req,
               struct timespec *rem);
```

Успешный вызов `nanosleep()` помещает вызывающий процесс в состояние ожидания на время, указанное при помощи аргумента `req`, и возвращает значение 0. В случае ошибки вызов возвращает -1 и соответствующим образом устанавливает переменную `errno`. Если состояние ожидания прерывается сигналом, вызов может возвращать значение до того, как истекает указанное время. В таком случае `nanosleep()` возвращает -1 и присваивает переменной `errno` значение `EINTR`. Если значение аргумента `rem` не равно `NULL`, то функция помещает оставшееся время ожидания (часть интервала `req`, которую процесс «не доспал»)

в `rem`. Программа может затем повторно выполнить вызов, передав `rem` в качестве `req` (как показано далее в этом разделе).

Другие возможные значения переменной `errno`:

`EFAULT`

Аргумент `req` или `rem` содержит недопустимый указатель.

`EINVAL`

Одно из полей структуры `req` недопустимо.

В простейшем случае вызов можно использовать следующим образом:

```
struct timespec req = { .tv_sec = 0,
                        .tv_nsec = 200 };
```

```
/* спать 200 наносекунд */
ret = nanosleep (&req, NULL);
if (ret)
    perror ("nanosleep");
```

В следующем примере второй параметр используется для того, чтобы продолжать находиться в состоянии ожидания в случае, если оно прерывается:

```
struct timespec req = { .tv_sec = 0,
                        .tv_nsec = 1369 };
struct timespec rem;
int ret;

/* спать 1369 наносекунд */
retry:
ret = nanosleep (&req, &rem);
if (ret) {
    if (errno == EINTR) {
        /* повторить с указанным оставшимся временем */
        req.tv_sec = rem.tv_sec;
        req.tv_nsec = rem.tv_nsec;
        goto retry;
    }
    perror ("nanosleep");
}
```

Наконец, далее приведен альтернативный подход (вероятно, наиболее эффективный, но менее удобный для чтения) к той же задаче:

```
struct timespec req = { .tv_sec = 1,
                        .tv_nsec = 0 };
struct timespec rem, *a = &req, *b = &rem;

/* спать 1 секунду */
while (nanosleep (a, b) && errno == EINTR) {
    struct timespec *tmp = a;
    a = b;
    b = tmp;
}
```

У системного вызова `nanosleep()` несколько преимуществ по сравнению со `sleep()` и `usleep()`:

- разрешение до наносекунд, в противоположность секундам и микросекундам;
- стандартизация в POSIX.1b;
- не реализуется с использованием сигналов (недостатки чего обсуждаются далее).

Несмотря на то, что он уже устарел, многие разработчики программ предпочитают использовать системный вызов `usleep()`, а не `nanosleep()`; но все меньше и меньше приложений, к счастью, включают сегодня вызов `sleep()`. Так как `nanosleep()` относится к стандарту POSIX и не использует сигналы, в новых программах следует предпочитать этот системный вызов (или интерфейс, о котором речь пойдет в следующем разделе) вызовам `sleep()` и `usleep()`.

Расширенный подход к засыпанию

Как и для всех классов функций времени, которые мы к этому времени изучили, семейство часов POSIX предоставляет наиболее расширенный интерфейс управления ожиданием:

```
#include <time.h>
```

```
int clock_nanosleep (clockid_t clock_id,
                     int flags,
                     const struct timespec *req,
                     struct timespec *rem);
```

`Clock_nanosleep()` работает аналогично `nanosleep()`. В действительности, этот вызов:

```
ret = nanosleep (&req, &rem);
```

делает то же самое, что и этот:

```
ret = clock_nanosleep (CLOCK_REALTIME, 0, &req, &rem);
```

Различие заключается в параметрах `clock_id` и `flags`. Первый позволяет указать источник времени, относительно которого требуется измерять период ожидания. Допустимы большинство источников времени, хотя нельзя указать часы процессора вызывающего процесса (то есть `CLOCK_PROCESS_CPUTIME_ID`); это не имело бы смысла, так как вызов приостанавливает выполнение процесса и, следовательно, время процесса прекращает увеличиваться.

Какой источник времени вы укажете — зависит от целей, с которыми ваша программа приостанавливается. Если нужно заставить ее спать до некоторого абсолютного значения времени, то `CLOCK_REALTIME` имеет наибольший смысл. Если же она приостанавливается на относительный интервал времени, то `CLOCK_MONOTONIC`, определенно, будет наилучшим выбором.

Параметр `flags` принимает значение либо `TIMER_ABSTIME`, либо 0. Если он равен `TIMER_ABSTIME`, но значение, указанное при помощи аргумента `req`, считается абсолютным, а не относительным. Это устраняет потенциальное условие состязания. Для того чтобы понять ценность этого параметра, предположим, что процесс желает приостановиться в момент времени T_0 до момента времени T_1 . В момент T_0 процесс вызывает `clock_gettime()`, чтобы получить текущее значение

времени T_0 . Затем он вычитает T_0 из T_1 , получая результат Y , который и передает вызову `clock_nanosleep()`. Однако какое-то количество времени все же проходит между моментом, в который процесс узнает текущее время, и моментом, когда процесс переходит в состояние ожидания. Еще хуже — что, если процессdezактивируется процессором, обращается к несуществующей странице или происходит нечто подобное? Всегда существует потенциальное условие состязания между получением текущего времени, вычислением разницы во времени и фактическим засыпанием.

Флаг `TIMER_ABSTIME` сводит на нет состязание, позволяя процессу напрямую указать время T_1 . Ядро приостанавливает процесс до тех пор, пока указанный источник времени не достигает момента T_1 . Если текущее время указанного источника времени уже превышает T_1 , то вызов возвращает результат немедленно.

Давайте рассмотрим и относительное, и абсолютное засыпание. В следующем примере процесс засыпает на 1,5 секунды:

```
struct timespec ts = { tv_sec = 1, tv_nsec = 500000000 };
int ret;

ret = clock_nanosleep (CLOCK_MONOTONIC, 0, &ts, NULL);
if (ret)
    perror ("clock_nanosleep").
```

И наоборот, в следующем примере процесс засыпает до тех пор, пока не достигается абсолютное значение времени — в частности одна секунда после момента, который вызов `clock_gettime()` возвращает для источника времени `CLOCK_MONOTONIC`:

```
struct timespec ts;
int ret;

/* нужно заснуть на одну секунду, начиная с ТЕКУЩЕГО МОМЕНТА */
ret = clock_gettime (CLOCK_MONOTONIC, &ts);
if (ret) {
    perror ("clock_gettime");
    return;
}

ts.tv_sec += 1;
printf ("Мы хотим заснуть до секунды=%ld наносекунды=%ld\n",
       ts.tv_sec, ts.tv_nsec);
ret = clock_nanosleep (CLOCK_MONOTONIC, TIMER_ABSTIME,
                      &ts, NULL);
if (ret)
    perror ("clock_nanosleep").
```

В большинстве программ требуется переход в состояние ожидания только на относительные промежутки времени, так как их требования к засыпанию не очень строгие. Однако у некоторых процессов реального времени могут быть очень точные требования к синхронизации, и они засыпают на абсолютные интервалы времени, чтобы избегать опасности потенциально разрушительных условий состязания.

Переносимый способ засыпания

Вспомните — в главе 2 мы обсуждали системный вызов `select()`:

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

```
int select (int n,
            fd_set *readfds,
            fd_set *writefds,
            fd_set *exceptfds,
            struct timeval *timeout);
```

Как упоминалось в той главе, `select()` обеспечивает переносимый способ перевода процесса в состояние ожидания с разрешением меньше секунды. Долгое время переносимые программы Unix были привязаны к вызову `sleep()` в своих нуждах, связанных с засыпанием: `usleep()` не был доступен достаточно широко, а `nanosleep()` был еще даже не написан. Разработчики обнаружили, что если передать вызову `select()` значение 0 в качестве аргумента `n`, `NULL` для всех трех указателей `fd_set` и желаемую продолжительность засыпания в качестве аргумента `timeout`, то этот вызов превращается в переносимый и эффективный способ перевода процессов в режим ожидания:

```
struct timeval tv = { tv_sec = 0,
                      tv_usec = 757 }.
```

```
/* заснуть на 757 микросекунд */
select (0, NULL, NULL, NULL, &tv).
```

Если переносимость на старые системы Unix для вас действительно важна, то лучшим выбором, вероятно, является системный вызов `select()`.

Превышение пределов

Все интерфейсы, о которых рассказывалось в этом разделе, гарантируют, что процесс будет спать, по крайней мере, в течение запрошенного интервала времени (если происходит обратное, они возвращают ошибку). Эти интерфейсы никогда не возвращаются успешно, пока запрошенный интервал не истекает. Процесс, однако, может проспать *дольше*, чем было запрошено.

Этот феномен может быть всего лишь следствием поведения планировщика: запрошенное время истекло, ядро разбудило процесс вовремя, но планировщик выбрал и позволил выполниться другой задаче.

Однако существует и другая причина: *превышение пределов таймера* (`timer overrun`). Это происходит, когда разрешение таймера грубее запрошенного интервала времени. Например, предположим, что системный таймер тикает с интервалами в 10 миллисекунд, а процесс запросил засыпание на 1 миллисекунду. Система способна измерять время и реагировать на связанные со временем события (такие, как необходимость разбудить заснувший процесс) только с интервалами в 10 миллисекунд. Если в момент, когда процесс делает запрос на засыпание, таймер находится в 1 миллисекунде от тика, то все будет в порядке — через 1 миллисекунду запрошенное время (1 миллисекунда) пройдет, и ядро

разбудит процесс. Однако если таймер срабатывает в точности в тот момент, когда процесс запрашивает засыпание, то очередного срабатывания не будет еще 10 миллисекунд. Следовательно, процесс проспит лишние 9 миллисекунд! Таким образом, случится девять превышений длиной в 1 миллисекунду. В среднем, таймер с периодом, равным X, характеризуется нормой превышения в X/2.

Использование высокоточных источников времени, таких, например, которые предоставляют часы POSIX, а также более высоких значений HZ минимизирует превышение пределов таймера.

Альтернативы засыпанию

По возможности следует избегать перевода процессов в режим ожидания. Зачастую это невозможно, но ничего страшного в этом нет — особенно если ваш код спит меньше секунды. Однако постоянное использование в коде засыпания в целях ожидания событий говорит о плохом дизайне. Код, блокирующийся на дескрипторе файла и позволяющий ядру обрабатывать засыпание и будить процесс, лучше. Вместо того чтобы крутить процесс в цикле, ожидая события, ядро может блокировать выполнение процесса и будить его только при необходимости.

Таймеры

Таймеры обеспечивают механизм, уведомляющий процесс о том, что определенный интервал времени истекает. Промежуток времени до того, как значение таймера *истекает* (expires), называется *задержкой* (delay) или *окончанием срока* (expiration). Как ядро уведомляет процесс о том, что таймер истек, — зависит от таймера. Ядро Linux предлагает несколько типов, и мы изучим их все.

Таймеры полезны по нескольким причинам. Примеры включают обновление экрана 60 раз в секунду или отмену ожидающей транзакции, если она все еще не завершилась в течение 500 миллисекунд.

Простые сигнализации

`Alarm()` — это простейший интерфейс таймера:

```
#include <unistd.h>
```

```
unsigned int alarm (unsigned int seconds);
```

Когда вы вызываете эту функцию, доставка сигнала `SIGALRM` вызывающему процессу планируется на момент, когда истечет `seconds` секунд реального времени. Если ранее запланированный сигнал уже ожидает доставки, вызов отменяет эту сигнализацию и заменяет ее только что запрошенней, возвращая число секунд, оставшихся до срабатывания предыдущей сигнализации. Если значение `seconds` равно 0, то предыдущая сигнализация, если таковая существует, отменяется, но новая не планируется.

Успешное использование данной функции таким образом также требует регистрации обработчика сигналов для сигнала SIGALRM. (Сигналы и обработчики сигналов рассматривались в предыдущей главе.) Далее показан фрагмент кода, который регистрирует обработчик SIGALRM, `alarm_handler()`, и устанавливает сигнализацию на пятисекундный интервал:

```
void alarm_handler (int signum)
{
    printf ("Пять секунд прошло!\n");
}

void func (void)
{
    signal (SIGALRM, alarm_handler);
    alarm (5);

    pause ();
}
```

Интервальные таймеры

Системные вызовы *интервальных таймеров* (interval timer), которые впервые появились в версии 4.2BSD, были с тех пор стандартизированы в POSIX и обеспечивают более высокую степень контроля, чем `alarm()`:

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

int getitimer (int which,
               struct itimerval *value);

int setitimer (int which,
               const struct itimerval *value,
               struct itimerval *ovalue);
```

Интервальные таймеры работают аналогично `alarm()`, но при необходимости могут автоматически заводиться повторно и функционируют в одном из трех режимов:

ITIMER_REAL

Измеряет реальное время. Когда указанный промежуток реального времени проходит, ядро отправляет процессу сигнал SIGALRM.

ITIMER_VIRTUAL

Значение таймера уменьшается только тогда, когда выполняется код из пользовательского пространства процесса. Как только указанный промежуток времени процесса истекает, ядро отправляет процессу сигнал SIGVTALRM.

ITIMER_PROF

Значение таймера уменьшается и когда выполняется процесс, и когда ядро работает от имени процесса (например, завершая системный вызов). Как только указанный промежуток времени истекает, ядро отправляет процессу сигнал SIGPROF. Этот режим обычно сочетается с режимом ITIMER_VIRTUAL, чтобы

программа могла измерять пользовательское время и время ядра, которое затрачивает процесс.

ITIMER_REAL измеряет то же время, что и `alarm()`; остальные два режима полезны для профилирования.

Структура `itimerval` позволяет пользователю указывать промежуток времени, после которого таймер должен завершаться, а также новый срок окончания (если необходимо), с которым таймер должен заново запускаться после первого завершения:

```
struct itimerval {
    struct timeval it_interval; /* следующее значение */
    struct timeval it_value;   /* текущее значение */
}:
```

Я уже упоминал, что структура `timeval` обеспечивает разрешение до микросекунд:

```
struct timeval {
    long tv_sec; /* секунды */
    long tv_usec; /* микросекунды */
}:
```

Системный вызов `setitimer()` устанавливает таймер типа `which` на срок завершения, указанный при помощи аргумента `it_value`. Как только время, соответствующее `it_value`, истекает, ядро заново заводит таймер, на этот раз используя время, указанное в аргументе `it_interval`. Таким образом, `it_value` — это время, оставшееся в текущем таймере. Как только `it_value` достигает нуля, таймеру присваивается значение `it_interval`. Если таймер истекает, а аргумент `it_interval` равен 0, то таймер заново не заводится. Схожим образом, если значение `it_value` активного таймера устанавливается на 0, то таймер останавливается, а не заводится заново.

Если значение `oalue` не равно `NULL`, то возвращаются предыдущие значения интервального таймера типа `which`.

Системный вызов `getitimer()` возвращает текущие значения интервального таймера типа `which`.

Обе функции возвращают 0 в случае успеха и -1 в случае ошибки. Также при ошибке переменной `errno` присваивается одно из следующих значений:

`EFAULT`

Аргумент `value` или `oalue` содержит недопустимый указатель.

`EINVAL`

`Which` не соответствует допустимому типу интервального таймера.

Код в следующем фрагменте создает обработчик сигнала `SIGALRM` (см. главу 9), а затем заводит интервальный таймер с первоначальным сроком истечения в пять секунд и последующими интервалами в одну секунду:

```
void alarm_handler (int signo)
{
    printf ("Таймер сработал!\n");
}
```

```
void foo (void) {
    struct itimerval delay;
    int ret;

    signal (SIGALRM, alarm_handler);

    delay.it_value.tv_sec = 5;
    delay.it_value.tv_usec = 0;
    delay.it_interval.tv_sec = 1;
    delay.it_interval.tv_usec = 0;
    ret = setitimer (ITIMER_REAL, &delay, NULL);
    if (ret) {
        perror ("setitimer");
        return;
    }

    pause ( );
}
```

В некоторых системах Unix системные вызовы sleep() и usleep() реализуются через сигнал SIGALRM и, что очевидно, alarm() и setitimer() также используют сигнал SIGALRM. Таким образом, программисты должны соблюдать осторожность, чтобы эти функции не перекрывались — результат такого сочетания не определен. Для небольших интервалов ожидания рекомендуется использовать вызов nanosleep(), который, согласно стандарту POSIX, не может применять сигналы. Для таймеров следует использовать setitimer() или alarm().

Расширенные таймеры

Самый мощный интерфейс таймера происходит из семейства часов POSIX.

При использовании таймеров, основанных на часах POSIX, акты реализации, инициализации и в конечном итоге удаления таймера разделяются на три разные функции: timer_create() создает таймер, timer_settime() инициализирует таймер, а timer_delete() разрушает его.

ПРИМЕЧАНИЕ

Интерфейсы таймеров для семейства часов POSIX, несомненно, обладают наибольшими возможностями, но они также самые новые (то есть наименее переносимые) и самые сложные в применении. Если основным требованием является простота или переносимость, то лучшим выбором будет setitimer().

Создание таймера

Для создания таймера используйте системный вызов timer_create():

```
#include <signal.h>
#include <time.h>

int timer_create (clockid_t clockid,
                  struct sigevent *evp,
                  timer_t *timerid);
```

Успешный вызов `timer_create()` создает новый таймер, связанный с часами POSIX `clockid`, сохраняет уникальный идентификатор таймера в аргументе `timerid` и возвращает значение 0. Этот вызов всего лишь определяет условия для работы таймера; фактически, ничего не происходит до тех пор, пока таймер не заводится, что демонстрируется в следующем разделе.

В этом примере создается новый таймер, связанный с часами POSIX `CLOCK_PROCESS_CPUTIME_ID`, а идентификатор таймера сохраняется в переменной `timer`:

```
int ret;

ret = timer_create (CLOCK_PROCESS_CPUTIME_ID,
                    NULL,
                    &timer);

if (ret)
    perror ("timer_create");
```

В случае ошибки вызов возвращает `-1`, значение `timerid` не определяется, а переменной `errno` присваивается одно из следующих значений:

`EAGAIN`

В системе недостаточно ресурсов для выполнения запроса.

`EINVAL`

Значению аргумента `clockid` соответствуют недопустимые часы POSIX.

`ENOTSUP`

Часы POSIX, указанные при помощи аргумента `clockid`, допустимы, но система не поддерживает использование часов для таймеров. Стандарт POSIX гарантирует, что все реализации поддерживают для таймеров часы `CLOCK_REALTIME`. Поддерживаются ли остальные часы – зависит от реализации.

Параметр `evp`, если его значение не равно `NULL`, содержит асинхронное уведомление, которое отправляется, когда таймер истекает. Структура определяется в файле заголовка `<signal.h>`. Ее содержимое должно быть непрозрачным для программиста, но, по крайней мере, она включает следующие поля:

```
#include <signal.h>

struct sigevent {
    union sigval sigev_value;
    int sigev_signo;
    int sigev_notify;
    void (*sigev_notify_function)(union sigval);
    pthread_attr_t *sigev_notify_attributes;
};

union sigval {
    int sival_int;
    void *sival_ptr;
};
```

Таймеры, основанные на часах POSIX, обеспечивают намного более высокую степень контроля над тем, как ядро уведомляет процесс об истечении тай-

мера, позволяя процессу точно указывать, какой сигнал ядро должно создавать, и даже позволяя ядру запускать потоки и выполнять функции в ответ на истечение таймера. Процесс определяет поведение при истечении таймера при помощи поля `sigev_notify`, которое может принимать одно из следующих трех значений:

SIGEV_NONE

«Нулевое» уведомление. При истечении таймера ничего не происходит.

SIGEV_SIGNAL

При истечении таймера ядро отправляет процессу сигнал, указанный в поле `sigev_signo`. В обработчике сигналов аргументу `si_value` присваивается значение `sigev_value`.

SIGEV_THREAD

При истечении таймера ядро запускает новый поток (в пределах того же процесса) и заставляет его выполнить `sigev_notify_function`, передавая в качестве единственного аргумента `sigev_value`. Поток завершается, как только он возвращается из этой функции. Если значение `sigev_notify_attributes` не равно `NULL`, то указанная структура `pthread_attr_t` определяет поведение нового потока.

Если значение аргумента `evp` равно `NULL`, как в нашем предыдущем примере, то уведомление об истечении таймера определяется так, как если бы значение `sigev_notify` было равно `SIGEV_SIGNAL`, `sigev_signo` — `SIGALRM`, а `sigev_value` — идентификатору таймера. Таким образом, по умолчанию эти таймеры отправляют уведомления так же, как интервальные таймеры POSIX. Однако благодаря пользовательской настройке они могут делать намного больше!

В следующем примере создается таймер, привязанный к часам `CLOCK_REALTIME`. Когда таймер истекает, ядро отправляет сигнал `SIGUSR1` и присваивает `si_value` значение, соответствующее адресу, где хранится идентификатор таймера:

```
struct sigevent evp;
timer_t timer;
int ret;

evp.sigev_value.sival_ptr = &timer;
evp.sigev_notify = SIGEV_SIGNAL;
evp.sigev_signo = SIGUSR1;
ret = timer_create(CLOCK_REALTIME,
                   &evp,
                   &timer);
if (ret)
    perror("timer_create");
```

Инициализация таймера

Таймер, созданный при помощи `timer_create()`, не заведен. Чтобы связать его с определенным моментом времени и запустить часы, используйте `timer_settime()`:

```
#include <time.h>

int timer_settime(timer_t timerid,
                  int flags,
                  const struct itimerspec *value,
                  struct itimerspec *ovalue);
```

Успешный вызов `timer_settime()` заводит таймер, указанный при помощи параметра `timerid`, на срок истечения `value`, который представляется структурой типа `itimerspec`:

```
struct itimerspec {
    struct timespec it_interval; /* следующее значение */
    struct timespec it_value;   /* текущее значение */
};
```

Как и для `setitimer()`, `it_value` указывает текущий срок истечения таймера. Когда таймер истекает, `it_value` обновляется и этому полю присваивается значение `it_interval`. Если в поле `it_interval` находится значение 0, то таймер не является интервальным таймером, поэтому отключается сразу же, как только истекает `it_value`.

Вспомним, что структура `timespec` обеспечивает разрешение до наносекунд:

```
struct timespec {
    time_t tv_sec; /* секунды */
    long tv_nsec;  /* наносекунды */
};
```

Если значение аргумента `flags` равно `TIMER_ABSTIME`, то время, указанное при помощи аргумента `value`, интерпретируется как абсолютное (в противоположность интерпретации по умолчанию, когда оно считается относительным в сравнении с текущим временем). Такое модифицированное поведение предотвращает условие состязания во время выполнения шагов по получению текущего времени, вычислению относительной разницы между этим временем и желаемым будущим временем и инициализации таймера. Подробное обсуждение — в разделе «Расширенный подход к засыпанию» выше.

Если значение `ovalue` не равно `NULL`, то предыдущее время истечения таймера сохраняется в предоставленной структуре `itimerspec`. Если таймер уже отключен, то всем полям структуры присваивается значение 0.

Используя значение `timer`, инициализированное ранее в `timer_create()`, код в следующем примере создает периодический таймер, истекающий каждую секунду:

```
struct itimerspec ts;
int ret;

ts.it_interval.tv_sec = 1;
ts.it_interval.tv_nsec = 0;
ts.it_value.tv_sec = 1;
ts.it_value.tv_nsec = 0;

ret = timer_settime(timer, 0, &ts, NULL);
if (ret)
    perror ("timer_settime");
```

Получение времени истечения таймера

В любой момент можно узнать время истечения таймера, не сбрасывая его, применив системный вызов `timer_gettime()`:

```
#include <time.h>

int timer_gettime(timer_t timerid,
    struct itimerspec *value);
```

Успешный вызов `timer_gettime()` сохраняет время истечения таймера, указанного при помощи аргумента `timerid`, в структуре, на которую указывает `value`, и возвращает 0. В случае ошибки вызов возвращает -1 и присваивает переменной `errno` одно из следующих значений:

EFAULT

Аргумент `value` содержит недопустимый указатель.

EINVAL

Аргумент `timerid` представляет недопустимый таймер.

Например:

```
struct itimerspec ts;
int ret;

ret = timer_gettime(timer, &ts);
if (ret)
    perror ("timer_gettime");
else {
    printf ("текущее время истечения: секунды=%ld наносекунды=%ld\n",
        ts.it_value.tv_sec, ts.it_value.tv_nsec);
    printf ("следующее время истечения: секунды=%ld наносекунды=%ld\n",
        ts.it_interval.tv_sec, ts.it_interval.tv_nsec);
}
```

Получение превышения таймера

В стандарте POSIX определяется интерфейс, позволяющий узнать, сколько превышений пределов таймера произошло для определенного таймера, если таймеры имели место:

```
#include <time.h>

int timer_getoverrun(timer_t timerid);
```

В случае успеха `timer_getoverrun()` возвращает число раз, когда таймер дополнительном истекал между первоначальным запланированным временем истечения и моментом отправки процессу уведомления об истечении таймера, — например, путем сигнала. Например, в нашем предыдущем примере, где таймер, запланированный на 1 мс, работал 10 мс, вызов вернул бы значение 9.

Если число превышений равно или больше константы `DELAYTIMER_MAX`, то вызов возвращает значение `DELAYTIMER_MAX`.

В случае ошибки функция возвращает `-1` и присваивает переменной `errno` значение `EINVAL` — единственное условие ошибки, указывающее, что аргументу `timerid` соответствует недопустимый таймер.

Например:

```
int ret;
```

```
ret = timer_getoverrun (timer);
if (ret == -1)
    perror ("timer_getoverrun");
else if (ret == 0)
    printf ("превышений не было\n");
else
    printf ("%d превышений\n", ret);
```

Удаление таймера

Удалить таймер просто:

```
#include <time.h>
```

```
int timer_delete (timer_t timerid);
```

Успешный вызов `timer_delete()` разрушает таймер, связанный с идентификатором `timerid`, и возвращает значение `0`. В случае ошибки вызов возвращает `-1` и присваивает переменной `errno` значение `EINVAL` — единственное условие ошибки, указывающее на то, что аргументу `timerid` соответствует недопустимый таймер.

Приложение. Расширения GCC языка C

Коллекция компиляторов GNU (GNU Compiler Collection, GCC) содержит множество расширений языка программирования C, некоторые из которых представляют особенную ценность для системных программистов. Большинство дополнений языка C, которые мы рассмотрим в этом приложении, предлагают программистам способы передачи компилятору дополнительной информации, описывающей поведение и назначение их кода. Компилятор, в свою очередь, использует эту информацию для генерирования более эффективного машинного кода. Прочие расширения заполняют бреши в языке программирования C, особенно на низких уровнях.

GCC предоставляет несколько расширений, которые теперь также доступны в новейшем стандарте C, ISO C99. Некоторые из этих расширений функционируют аналогично своим коллегам из C99, но другие расширения реализованы в ISO C99 совсем по-другому. В новом коде следует использовать варианты этих возможностей, определенные в ISO C99. Мы не будем рассматривать здесь подобные расширения — мы поговорим только о дополнениях, уникальных для GCC.

GNU C

Разновидность C, поддерживаемую GCC, зачастую называют GNU C. В 90-х годах XX в. GNU C заполнил несколько пробелов в языке C, предоставив такие возможности, как сложные переменные, массивы нулевой длины, подставляемые функции и именованные функции инициализации. Однако практически через десять лет язык C был наконец обновлен, а со стандартизацией ISO 99 расширения GNU C стали менее востребованы. Тем не менее GNU C продолжает предоставлять полезные возможности, и многие программисты, пишущие для Linux, все так же применяют поднабор функций GNU C — часто всего лишь одно-два расширения — в своем коде, совместимом с C90 и C99.

Одним из выдающихся примеров кода GCC является ядро Linux, написанное строго на GNU C. Недавно, однако, Intel инвестировала усилия по разработке в то, чтобы научить компилятор Intel C Compiler (ICC) понимать расширения

GNU C, используемые ядром. Следовательно, многие из этих расширений теперь становятся все менее привязанными к GCC.

Подставляемые функции

Компилятор копирует весь код «*подставляемой*» функции (inline function) в область, где эта функция вызывается. Вместо того чтобы хранить функцию внешне и переходить к ней, когда эта функция вызывается, компилятор выполняет содержимое функции напрямую. Такое поведение позволяет избегать нагрузки, создаваемой вызовом функции, и обеспечивает потенциальную оптимизацию в области вызова, так как компилятор может совместно оптимизировать вызывающего и вызываемого. Последнее особенно актуально, если параметры функции постоянны в области вызова. Конечно же, копирование функции в любой и каждый фрагмент кода, вызывающий ее, может оказывать негативное влияние на объем кода. Таким образом, функции следует подставлять, только если они небольшие и простые или же если они вызываются в малом количестве мест.

Многие годы GCC поддерживала ключевое слово `inline`, заставляющее компилятор подставлять указанную функцию. Это ключевое слово было formalизовано в C99:

```
static inline int foo (void) { /* ... */ }
```

Технически, однако, данное ключевое слово является всего лишь подсказкой — рекомендацией для компилятора, чтобы он рассмотрел возможность подстановки указанной функции. GCC также предоставляет расширение, позволяющее приказывать компилятору, чтобы он *всегда* подставлял определенную функцию:

```
static inline __attribute__ ((always_inline)) int foo (void) { /* ... */ }
```

Наиболее очевидный кандидат на роль подставляемой функции — это макрос препроцессора. Подставляемая функция в GCC выполняется так же, как и макрос, и дополнительно получает проверку типов. Например, вместо этого макроса:

```
#define max(a,b) ({ a > b ? a : b; })
```

можно было бы использовать соответствующую подставляемую функцию:

```
static inline max (int a, int b)
{
    if (a > b)
        return a;
    return b;
}
```

Среди программистов наблюдается тенденция к чрезмерному использованию подставляемых функций. Нагрузка, связанная с вызовом функции на большинстве современных архитектур — в частности, x86, — очень, очень мала. Только самые достойные среди функций должны рассматриваться в качестве кандидатов на подстановку!

Подавление подстановки

В своем самом агрессивном режиме оптимизации GCC автоматически выбирает функции, кажущиеся подходящими для подстановки, и подставляет их. Обычно это хорошая идея, но иногда программист знает, что функция будет работать неправильно, если ее подставить. Одним из возможных примеров этого является ситуация использования `_builtin_return_address` (речь об этом пойдет далее в приложении). Для подавления подстановки используйте ключевое слово `noinline`:

```
_attribute_(noinline) int foo (void) { /* ... */ }
```

Чистые функции

«Чистая» функция (pure function) — это функция без побочных эффектов, возвращаемое значение которой отражает только параметры функции и неразрушающиеся глобальные переменные. Доступ к любому параметру или глобальной переменной должен быть только на чтение. К таким функциям может применяться оптимизация циклов и исключение подвыражений. Функции помечаются как чистые при помощи ключевого слова `pure`:

```
_attribute_(pure) int foo (int val) { /* ... */ }
```

Распространенный пример чистой функции — это `strlen()`. При условии идентичного входа возвращаемое значение данной функции не меняется в зависимости от количества вызовов, и, таким образом, ее можно извлечь из цикла и вызвать только один раз. Например, рассмотрим следующий код:

```
/* посимвольно вывести каждую букву из p в верхнем регистре */
for (i = 0; i < strlen(p); i++)
    printf ("%c", toupper (p[i]));
```

Если бы компилятор не знал, что `strlen()` — чистая функция, то он мог бы вызвать ее на каждой итерации цикла!

Умные программисты — а также компиляторы, если `strlen()` помечается как чистая, — написали бы или сгенеририровали такой код:

```
size_t len;
len = strlen (p);
for (i = 0; i < len; i++)
    printf ("%c", toupper (p[i]));
```

Между прочим, еще более хитрые программисты (такие, как читатели этой книги) написали бы:

```
while (*p)
    printf ("%c", toupper (*p++));
```

Это недопустимо и действительно не имеет смысла, чтобы чистая функция возвращала значение типа `void`, так как возвращаемое значение является единственной сутью подобных функций.

Постоянные функции

«Постоянная» функция (constant function) — это более строгий вариант чистой функции. Подобные функции не могут обращаться к глобальным переменным и не могут принимать указатели как параметры. Таким образом, возвращаемое значение постоянной функции не отражает ничего, кроме параметров, переданных по значению. Для таких функций возможна дополнительная оптимизация, превышающая оптимизацию чистых функций. Примерами постоянных функций являются математические функции, такие, как `abs()` (в предположении, что они не сохраняют состояние или не выделяют другие фокусы во имя оптимизации). Для того чтобы пометить функцию как постоянную, нужно использовать ключевое слово `const`:

```
_attribute_(const) int foo (int val) { /* ... */ }
```

Как и с чистыми функциями, нет никакого смысла в том, чтобы постоянная функция возвращала значение типа `void`.

Функции, не возвращающие результат

Если функция не возвращает результат — например, потому, что она инвариантно вызывает `exit()`, — то программист может пометить функцию ключевым словом `noreturn`, подчеркнув этот факт для компилятора:

```
_attribute_(noreturn) void foo (int val) { /* ... */ }
```

В свою очередь, компилятор может выполнить дополнительную оптимизацию, понимая, что ни при каких обстоятельствах вызываемая функция никогда не вернет никакое значение. Для этой функции не имеет никакого смысла возвращать что-либо, отличное от значения типа `void`.

Функции, выделяющие память

Если функция возвращает указатель, который никогда не может ссылаться на уже существующую память¹, — практически всегда потому, что функция сама выделяет память и возвращает указатель на нее, — то программист может пометить функцию как таковую ключевым словом `malloc`, а компилятор, в свою очередь, проведет подходящую оптимизацию:

¹ Ссылка на существующую память (memory alias) возникает, когда две или больше переменные-указателя указывают на один и тот же адрес в памяти. Это может происходить в тривиальных ситуациях, когда указателю присваивается значение другого указателя, и в более сложных, менее очевидных случаях. Если функция возвращает адрес на только что выделенную память, то не могут существовать никакие другие указатели на тот же адрес.

```
_attribute_ _ ((malloc)) void * get_page (void)
{
    int page_size.

    page_size = getpagesize ( );
    if (page_size <= 0)
        return NULL.

    return malloc (page_size);
}
```

Как заставить вызывающего проверять возвращаемое значение

Не элемент оптимизации, а вспомогательное средство для программиста атрибут `warn_unused_result` заставляет компилятор генерировать предупреждение каждый раз, когда он замечает, что возвращаемое значение функции не сохраняется или не используется в условном операторе:

```
_attribute_ _ ((warn_unused_result)) int foo (void) { /* ... */ }
```

Это позволяет программисту удостоверяться, что все вызывающие проверяют и обрабатывают возвращаемые значения функций, когда эти значения имеют определенную важность. Функции с важными, но часто игнорируемыми возвращаемыми значениями, такие, как `read()`, являются превосходными кандидатами на добавление данного атрибута. Подобные функции не могут возвращать значение типа `void`.

Как пометить функцию как устаревшую

Атрибут `deprecated` заставляет компилятор генерировать предупреждение в области вызова каждый раз, когда вызывается данная функция:

```
_attribute_ _ ((deprecated)) void foo (void) { /* ... */ }
```

Это помогает программистам отвыкать от устаревших и неиспользуемых интерфейсов.

Как пометить функцию как используемую

Иногда не существует кода, видимого для компилятора, который бы вызывал определенную функцию. Если пометить функцию атрибутом `used`, то можно сообщить компилятору, что программа действительно использует данную функцию, несмотря на впечатление, что никаких ссылок на нее не существует.

```
static _attribute_ _ ((used)) void foo (void) { /* ... */ }
```

Таким образом, компилятор выводит результирующий код на языке ассемблера, не отображая предупреждение о неиспользуемой функции. Данный атрибут полезен, если статическая функция вызывается только из написанного

вручную кода ассемблера. Обычно, если компилятору неизвестно о вызовах функции, он создает предупреждение и, потенциально, выполняет оптимизацию без учета данной функции.

Как пометить функции или параметры как неиспользующиеся

Атрибут `unused` говорит компилятору, что данная функция или параметр функции не используется, и заставляет его не создавать соответствующие предупреждения:

```
int foo (long __attribute__((unused)) value) { /* ... */ }
```

Это полезно, если вы компилируете проект с флагом `-W` или `-Wunused` и желаете отловить неиспользующиеся параметры функций, но иногда у вас появляются функции, которые должны соответствовать предопределенной подписи (как часто бывает в основанном на событиях программировании графических интерфейсов пользователя или в обработчиках сигналов).

Упаковка структуры

Атрибут `packed` говорит компилятору, что тип или переменная должны быть упакованы в памяти с использованием минимального объема пространства, потенциально без учета требований к выравниванию. Если данный атрибут указан для структуры (`struct`) или объединения (`union`), то все входящие в них переменные также упаковываются. Если он указан только для одной переменной, то упаковывается только этот конкретный объект.

В следующем примере все переменные в структуре упаковываются в минимальное пространство:

```
struct __attribute__((packed)) foo { ... };
```

В качестве примера, структура, содержащая поле типа `char`, за которым следует поле типа `int`, вероятнее всего, будет размещена так, что целочисленное значение будет выровнено не по адресу в памяти, сразу же следующему за символьным значением, а, скажем, на три байта дальше. Компилятор выравнивает переменные, вставляя между ними байты неиспользуемой забивки. В упакованной структуре забивка отсутствует, что потенциально позволяет ей занимать меньше памяти, но при этом нарушаются архитектурные требования к выравниванию.

Увеличение границы выравнивания переменной

Помимо упаковки переменных, GCC также позволяет программистам указывать альтернативную минимальную границу выравнивания для определенной переменной. GCC выравнивает указанную переменную, как *минимум*, по этому

значению, в противоположность минимальной требуемой границе выравнивания, как ее диктует архитектура и ABI. Например, следующий оператор объявляет целочисленную переменную с именем `beard_length` с минимальным выравниванием в 32 байт (в противоположность типичному выравниванию в 4 байт на машинах с 32-битным целочисленными значениями):

```
int beard_length __attribute__ ((aligned (32))) = 0;
```

Принудительное определение границы выравнивания для типа, в целом, полезно только при работе с аппаратным обеспечением, которое может накладывать более высокие требования к выравниванию, чем сама архитектура, или когда вы вручную смешиваете код на языке C и код ассемблера и хотите использовать инструкции, требующие специально выровненных значений. Один из примеров, где применяется такая функциональность выравнивания, – это хранение часто используемых переменных на линиях процессорного кэша для оптимизации поведения кэша. Ядро Linux задействует такую технику.

В качестве альтернативы указанию определенной минимальной границы выравнивания можно заставить GCC выравнивать определенный тип по самой большой из минимальных границ выравнивания среди всех типов данных. Например, следующий оператор говорит GCC, что переменную `parrot_height` необходимо выравнивать по самой большой границе, вероятно, соответствующей границе выравнивания типа `double`:

```
short parrot_height __attribute__ ((aligned)) = 5;
```

Данное решение обычно включает компромисс между использованием пространства и времени: переменные, выровненные таким способом, занимают больше места, но их копирование в обе стороны (а также другие сложные манипуляции) может осуществляться быстрее, так как компилятору удается выдавать машинные инструкции, действующие более крупные объемы памяти.

Максимальные пределы выравнивания переменных могут определяться различными аспектами архитектуры и цепочки инструментов системы. Например, в некоторых архитектурах Linux компоновщик не способен распознавать границы выравнивания за пределами относительно небольшого значения по умолчанию. В этом случае выравнивание, указанное при помощи данного ключевого слова, округляется вниз до наименьшей допустимой границы. Например, если вы запрашиваете выравнивание по 32 байт, но компоновщик системы не умеет выравнивать по границе, превышающей 8 байт, то переменная выравнивается по 8-байтовой границе.

Помещение глобальных переменных в регистр

GCC позволяет программистам помещать глобальные переменные в специальный машинный регистр, где переменные остаются на всем протяжении выполнения программы. В GCC подобные переменные называются глобальными регистровыми переменными.

Синтаксис требует, чтобы программист указывал машинный регистр. В следующем примере используется регистр ebx:

```
register int *foo asm ("ebx");
```

Программист должен выбрать переменную, которая не затирается функцией: то есть выбранная переменная должна быть доступна для использования локальными функциями, сохраняться и восстанавливаться при вызове функций и не определяться для какой бы то ни было специальной цели архитектурой или ABI операционной системы. Компилятор генерирует предупреждение, если выбранный регистр не подходит. Если регистр приемлем — ebx, упомянутый в данном примере, отлично подходит для архитектуры x86, — то компилятор, в свою очередь, сам прекращает использовать регистр.

Подобная оптимизация может обеспечивать большой выигрыш в производительности, если переменная часто используется. Хороший пример — виртуальная машина. Помещение переменной, в которой хранится, предположим, указатель кадра виртуального стека, в регистр может привести к значительному выигрышу. С другой стороны, если с самого начала в архитектуре недостаточно регистров (как в архитектуре x86), то такая оптимизация имеет мало смысла.

Глобальные регистровые переменные нельзя использовать в обработчиках сигналов или одновременно в нескольких потоках исполнения. Также они не могут иметь первоначальные значения, так как не существует механизма, позволяющего исполняемым файлам передавать содержимое по умолчанию для регистров. Объявления глобальных регистровых переменных должны предшествовать любым определениям функций.

Комментирование ветвей

GCC позволяет программистам комментировать ожидаемые значения выражений — например, говорить компилятору, должно быть условное выражение истинным или ложным. В свою очередь, GCC может менять порядок блоков и выполнять другие типы оптимизации, повышая производительность условных ветвей.

Синтаксис GCC для комментирования ветвей ужасающе уродлив. Для того чтобы сделать комментирование ветвей более приятным для глаза, мы используем макрос препроцессора:

```
#define likely(x) __builtin_expect (!!x, 1)
#define unlikely(x) __builtin_expect (!!x, 0)
```

Программисты могут помечать выражения как вероятно правдивые или вряд ли правдивые, заворачивая их в likely() и unlikely() соответственно.

Код в следующем примере помечает ветвь как вряд ли правдивую (то есть вероятно ложную):

```
int ret;
ret = close (fd);
if (unlikely (ret))
    perror ("close");
```

И наоборот, в следующем примере ветвь помечается как вероятно правдивая:

```
const char *home.  
  
home = getenv ("HOME").  
if (likely (home))  
    printf ("Ваш домашний каталог - %s\n", home).  
else  
    fprintf (stderr, "Переменная среды HOME не установлена!\n").
```

Как и с подставляемыми функциями, у программистов есть тенденция излишне усердно применять комментирование ветвей. Как только вы начинаете комментировать выражения, вам хочется сделать это для *всех* выражений. Однако будьте осторожны — помечать ветви как вероятно правдивые или ложные следует только тогда, когда вы знаете *заранее* и практически несомненно, что выражение будет истинным или ложным *практически всегда* (скажем, с уверенностью в 99 %). Редко происходящие ошибки — это хорошие кандидаты для *unlikely()*. Помните, однако, что ложное предсказание хуже, чем отсутствие любого предсказания.

Получение типа выражения

В GCC есть ключевое слово `typeof()`, предназначенное для проверки типа указанного выражения. Семантически данное слово работает так же, как и `sizeof()`. Например, это выражение возвращает тип того, на что указывает `x`:

```
typeof (*x)
```

Его можно использовать для объявления массива у элементов этого типа:

```
typeof (*x) y[42];
```

Популярный вариант использования `typeof()` — написание «безопасных» макросов, которые могут функционировать на любом арифметическом значении и оценивают свои параметры только один раз:

```
#define max(a,b) ({  
    typeof (a) _a = (a); \  
    typeof (b) _b = (b); \  
    _a > _b ? _a : _b; \  
})
```

Получение границы выравнивания типа

В GCC предусмотрено ключевое слово `_alignof_`, предназначенное для проверки выравнивания указанного объекта. Это значение зависит от архитектуры и интерфейса ABI. Если текущая архитектура не предлагает необходимое выравнивание, то данное ключевое слово возвращает рекомендуемое выравнивание

интерфейса ABI. В противном случае ключевое слово возвращает минимальное требуемое выравнивание.

Синтаксис идентичен синтаксису `sizeof()`:

```
_alignof_(int)
```

В зависимости от архитектуры, это выражение, вероятно, вернет 4, так как 32-битные целочисленные значения обычно выравниваются по 4-байтовой границе.

Данное ключевое слово также работает на значениях `lvalue`. В этом случае возвращаемая граница выравнивания — это минимальное выравнивание для соответствующего типа, а не фактическая граница для указанного значения `lvalue`. Если минимальное выравнивание было изменено атрибутом `aligned` (о котором было рассказано выше в разделе «Увеличение границы выравнивания переменной»), то ключевое слово `_alignof_` отражает это изменение.

Например, рассмотрим такую структуру:

```
struct ship {
    int year_built;
    char canons;
    int mast_height;
};
```

и такой фрагмент кода:

```
struct ship my_ship;

printf ("%d\n", _alignof_(my_ship.canons));
```

`_alignof_` в этом фрагменте вернет 1, несмотря на то, что из-за забивки структуры `canons`, вероятно, занимает четыре байта.

Смещение члена структуры

В GCC есть встроенное ключевое слово, позволяющее узнавать смещение члена структуры в этой структуре. Макрос `offsetof()`, определенный в файле заголовка `<stddef.h>`, является частью стандарта ISO C. Большинство определений ужа-сающими включают непристойную арифметику с указателями и код, неподходящий для просмотра несовершеннолетними лицами. Расширение GCC проще и потенциально быстрее:

```
#define offsetof(type, member) __builtin_offsetof (type, member)
```

Вызов возвращает смещение члена `member` в типе `type` — то есть число байтов, начиная с нуля, от начала структуры до этого члена. Например, возьмем следующую структуру:

```
struct rowboat {
    char *boat_name;
    unsigned int nr_oars;
    short length;
};
```

Фактическое смещение зависит от размера переменных, требований к выравниванию в данной архитектуре и способа забивки, но на 32-разрядной машине можно ожидать, что вызов `offsetof()` на структуре `struct rowboat` из `boat_name`, `nr_oars` и `length` вернет 0, 4 и 8 соответственно.

В системе Linux макрос `offsetof()` должен определяться с использованием ключевого слова `GCC` и не должен переопределяться.

Получение адреса возврата функции

В `GCC` предоставляется ключевое слово, позволяющее узнавать адрес возврата текущей функции или одного из вызывающих текущей функции:

```
void * __builtin_return_address (unsigned int level)
```

Параметр `level` указывает функцию в цепочке вызовов, адрес которой необходимо вернуть. Значение 0 позволяет запросить адрес возврата текущей функции, значение 1 — адрес возврата функции, вызывающей текущую функцию, значение 2 — адрес возврата функции, вызывающей вызывающую функцию и так далее.

Если текущая функция — это подставляемая функция, то возвращаемый адрес принадлежит вызывающей функции. Если это недопустимо, то используйте ключевое слово `noinline` (описанное выше в разделе «Подавление подстановки») для того, чтобы заставить компилятор не подставлять функцию.

Существует несколько вариантов использования ключевого слова `__builtin_return_address`, например, в целях отладки или информирования. Еще один вариант — применять его для раскручивания цепочки вызовов, для реализации самодиагностики, утилиты создания аварийного дампа, отладчика и т. д.

Обратите внимание, что некоторые архитектуры поддерживают возвращение только адреса вызывающей функции. В таких архитектурах ненулевое значение параметра может приводить к возврату случайного значения. Таким образом, любые параметры, отличные от нуля, исключают переносимость и должны применяться только в целях отладки.

Диапазоны оператора case

`GCC` позволяет в метках оператора `case` указывать диапазоны значений для отдельных блоков. В целом, синтаксис таков:

```
case low ... high:
```

Например:

```
switch (val) {
case 1 ... 10:
    /* ... */
    break;
case 11 ... 20:
    /* ... */
```

```

break;
default:
/* ... */
}

```

Данная функциональность также довольно полезна для диапазонов значений ASCII:

```
case 'A' ... 'Z':
```

Обратите внимание, что до и после многоточия должны стоять пробелы. В противном случае компилятор может не понять написанное, особенно когда речь идет о диапазонах целочисленных значений. Всегда делайте так:

```
case 4 ... 8:
```

и никогда так:

```
case 4...8:
```

Арифметика указателей типа `void` и указателей на функции

В GCC операции сложения и вычитания допустимы на указателях типа `void` и указателях на функции. Обычно ISO C не разрешает выполнять арифметические операции на подобных указателях, потому что размер значения типа `void` — это глупая концепция, а размер зависит от того, на что указатель в действительности указывает. Для упрощения такой арифметики GCC считает размер ссылаемого объекта равным одному байту. Таким образом, следующий код увеличивает значение `a` на единицу:

```
a++; /* a — это указатель типа void */
```

Параметр `-Wpointer-arith` заставляет GCC генерировать предупреждение, когда используются эти расширения.

Более переносимо и красиво

Давайте признаем это — синтаксис `_attribute_` далек от совершенства. Некоторые из расширений, с которыми мы познакомились, действительно требуют макросов препроцессора, чтобы их было приятно использовать, но любые только выигрывают, если привести их внешний вид в порядок.

Это совсем несложно, если применить немного препроцессорного волшебства. Помимо этого, одновременно мы можем сделать расширения GCC переносимыми, определив их на случай не поддерживающего GCC компилятора (что бы это ни было).

Для этого нужно вставить в заголовок следующий код и включить его в исходные файлы:

```

#if __GNUC__ >= 3
#undef inline
#define inline inline __attribute__((always_inline))
#define __noinline __attribute__((noinline))
#define __pure __attribute__((pure))
#define __const __attribute__((const))
#define __noreturn __attribute__((noreturn))
#define __malloc __attribute__((malloc))
#define __must_check __attribute__((warn_unused_result))
#define __deprecated __attribute__((deprecated))
#define __used __attribute__((used))
#define __unused __attribute__((unused))
#define __packed __attribute__((packed))
#define __align(x) __attribute__((aligned (x)))
#define __align_max __attribute__((aligned))
#define Likely(x) __builtin_expect (!!x, 1)
#define unlikely(x) __builtin_expect (!!x, 0)
#else
#define __noinline /* не noinline */
#define __pure /* не pure */
#define __const /* не const */
#define __noreturn /* не noreturn */
#define __malloc /* не malloc */
#define __must_check /* не warn_unused_result */
#define __deprecated /* не deprecated */
#define __used /* не used */
#define __unused /* не unused */
#define __packed /* не packed */
#define __align(x) /* не aligned */
#define __align_max /* не align_max */
#define Likely(x) (x)
#define unlikely(x) (x)
#endif

```

Например, следующий код помечает функцию как чистую, используя наше сокращение:

```
_pure int foo (void) { /* ... */
```

Если GCC используется, то функция помечается атрибутом `pure`. Если применяется другой компилятор, отличный от GCC, то препроцессор заменяет маркер `_pure` пустой операцией. Обратите внимание, что в одно определение можно поместить несколько атрибутов, поэтому в одном определении без проблем можно использовать несколько перечисленных конструкций.

Проще, удобнее и переносимее!

Библиография

В этом списке литературы перечисляются рекомендуемые книги, связанные с системным программированием, которые я разбил на четыре подкатегории. Ни одна из этих работ не относится к обязательному чтению. Они всего лишь представляют мою выборку лучших книг по данной теме. Если вы обнаружите себя в поисках более подробной информации относительно того, что обсуждалось в этой книге, то советую обратиться к моим любимым источникам.

В некоторых из перечисленных книг обсуждается материал, с которым читатель данной книги уже должен быть знаком, например язык программирования C. Прочие тексты дополняют эту книгу, например работы, посвященные gdb, Subversion (svn) и дизайну операционной системы. Помимо этого, обсуждаются темы, выходящие за пределы круга интересов данной книги, такие, как многопоточная работа сокетов. Как бы то ни было, я рекомендую все эти источники. Конечно же, следующий список далеко не исчерпывающий — не стесняйтесь исследовать и другие ресурсы.

Книги по языку программирования С

Эти книги документируют язык программирования C, *lingua franca* системного программирования. Если вы не кодируете на C с такой же ревностью, как говорите на родном языке, то одна или несколько из следующих работ (в сочетании с большим количеством 'тренировок!') должны помочь вам устранить этот недостаток. Если вы не найдете ничего другого, то первая книга — широко известная под псевдонимом K&R, — это настоящее сокровище для интересующегося читателя. Ее краткость объясняет простоту C.

Kernighan Brian W., Ritchie Dennis M. *The C Programming Language*. 2-nd edition. Prentice Hall, 1988.

Эта книга, написанная автором языка программирования C и его тогдашним коллегой, представляет собой библию программирования на C.

Prinz Peter, Crawford Tony. *C in a Nutshell*. O'Reilly Media, 2005.

Отличная книга, в которой раскрывается и язык C, и стандартная библиотека C.

Prinz Peter, Kirch-Prinz Ulla. *C Pocket Reference*. O'Reilly Media, 2002.

Краткое справочное руководство по языку C, умело обновленное в соответствии со стандартом ANSI C99.

Van der Linden Peter. *Expert C Programming*. Prentice Hall, 1994.

Великолепное обсуждение малоизвестных аспектов языка программирования C, сопровождающееся удивительным остроумием и чувством юмора. Эта книга изобилует неочевидными шутками, я просто обожаю ее.

Summit Steve. *C Programming FAQs: Frequently Asked Questions*. Addison-Wesley, 1995.

Эта потрясающая книга содержит более 400 часто задаваемых вопросов (с ответами) по языку программирования C. Ответы на многие вопросы кажутся очевидными в глазах мастеров C, но некоторые весомые вопросы и ответы впечатляют даже самых эрудированных программистов на C. Вы настоящий ниндзя C, если можете ответить на все эти каверзные вопросы! Единственный недостаток этой книги — то, что она не обновлялась в соответствии со стандартом ANSI C99, а он, определенно, привнес некоторые изменения (я вношу исправления в свою копию вручную). Обратите внимание, что существует и интерактивная версия книги, которая, вероятно, включает более свежие обновления.

Книги по программированию в Linux

Следующие книги посвящены программированию в Linux, в том числе темам, которые не были рассмотрены в данной книге (сокеты, взаимодействие между процессами и потоки pthread), а также инструментам программирования в Linux (CVS, GNU Make и Subversion).

Стивенс У., Феннер Б., Рудофф Э. Unix: Разработка сетевых приложений. 3-е изд. СПб.: Питер, 2004.

Определяющая работа по интерфейсу прикладного программирования сокетов; к сожалению, она не сосредоточена на особенностях Linux, но, к счастью, была недавно обновлена для IPv6.

Стивенс У. Unix: взаимодействие процессов. СПб.: Питер, 2001.

Отличное обсуждение взаимодействия между процессами (interprocess communication, IPC).

Nichols Bradford. *PThreads Programming: A POSIX Standard for Better Multiprocessing*. O'Reilly Media, 1996.

Обзор интерфейса прикладного программирования потоков в POSIX, pthreads.

Mecklenburg Robert. *Managing Projects with GNU Make*. O'Reilly Media, 2004.

Отличное обсуждение GNU Make – классического инструмента для сборки программных проектов в Linux.

Versperman Jennifer. *Essential CVS*. O'Reilly Media, 2006.

Превосходное обсуждение CVS – классического инструмента для управления версиями и исходным кодом в системах Unix.

Collins-Sussman Ben. *Version Control with Subversion*. O'Reilly Media, 2004.

Феноменальный подход к Subversion –циальному инструменту для управления версиями и управления исходным кодом в системах Unix, предпринятый тремя авторами Subversion.

Robbins Arnold. *GDB Pocket Reference*. O'Reilly Media, 2005.

Удобное карманное руководство по gdb – отладчику для Linux.

Siever Ellen. *Linux in a Nutshell*. O'Reilly Media, 2005.

Ураганный справочник по всем внутренностям Linux, включая многие инструменты, составляющие среду разработки Linux.

Книги, посвященные ядру Linux

Две книги, указанные здесь, рассказывают о ядре Linux. Есть три причины, по которым программист может интересоваться этой темой. Во-первых, ядро предоставляет для пользовательского пространства интерфейс системных вызовов и, таким образом, является основой системного программирования. Во-вторых, поведение и характерные черты ядра проливают свет на взаимодействие между ним и приложениями, которые ядро выполняет. Наконец, ядро Linux – это великолепный пример кода, а эти книги действительно интересные.

Love Robert. *Linux Kernel Development*. Novell Press, 2005.

Эта книга идеально подходит для системных программистов, желающих узнать о конструкции и реализации ядра Linux (и, конечно, было бы непростительным упущением с моей стороны не упомянуть о собственной работе на эту тему!). Это не справочник по API, здесь вы найдете отличное обсуждение используемых алгоритмов и решений, принимаемых ядром Linux.

Corbet Jonathan. *Linux Device Drivers*. O'Reilly Media, 2005.

Это великолепное руководство по написанию драйверов устройств для ядра Linux с отличным справочником по API. Сфокусированная в основном на драйверах устройств, эта книга будет отличным помощником программистам любых вероисповеданий, включая системных программистов, просто желающих проникнуть в махинации, производимые ядром Linux. Отличное дополнение к моей книге о ядре Linux.

Книги по дизайну операционной системы

Эти две книги, хотя и не относятся конкретно к Linux, рассматривают вопросы дизайна операционной системы в целом. Как я подчеркивал ранее, хорошее понимание системы, в которой вы пишете код, только улучшает ваши результаты.

Deitel Harvey. *Operating Systems*. Prentice Hall, 2003.

Tour de force по теории дизайна операционных систем в сочетании с перво-сортными практическими примерами, переводящими теорию в практику. Из всех учебников по дизайну операционных систем этот — мой любимый: он современный, полный и прекрасно читается.

Schimmel Curt. *UNIX Systems for Modern Architectures: Symmetric Multiprocessing and Caching for Kernel Programming*. Addison-Wesley, 1994.

Несмотря на то, что она лишь косвенно относится к системному программированию, эта книга предлагает такой великолепный подход к опасностям параллелизма и современного кэширования, что я порекомендовал бы ее даже стоматологам.

Лаэ Роберт
Linux. Системное программирование
Перевела с английского Е. Шикарева

Заведующий редакцией
Руководитель проекта
Ведущий редактор
Художественный редактор
Корректоры
Верстка

*А. Сандрыйкин
П. Манишин
О. Некруткина
Л. Адуевская
Н. Филатова, И. Тимофеева
Л. Егорова*

Подписано в печать 18.06.08. Формат 70x100/16. Усл. п. л. 33,54 Тираж 2500. Заказ № 540.
ООО «Питер Пресс», 198206, Санкт-Петербург, Петергофское шоссе, д. 73, лит. А29.
Налоговая льгота — общероссийский классификатор продукции ОК 005-93, том 2; 95 3005 — литература учебная
Отпечатано по технологии СИР в Отпечатано с готовых диапозитивов в ГП ПО «Псковская областная типография».
180004, г. Псков, ул. Ротная, 34.

LINUX. СИСТЕМНОЕ ПРОГРАММИРОВАНИЕ



Эта книга о системном программировании — искусстве написания системного программного обеспечения для Linux. Системный код находится на низком уровне, взаимодействуя напрямую с ядром и системными библиотеками. В книге приводится описание функций и интерфейсов, включая стандартные и специальные интерфейсы Linux. Это издание — своеобразное руководство по созданию качественного, быстрого кода. Автор, хакер ядра Linux — Роберт Лав, объясняет не только то, как системные интерфейсы должны работать, но и то, как они реально работают и как эффективно и безопасно их использовать.

Темы, рассмотренные в книге:

- что такое интерфейс системного уровня и как писать приложения системного уровня в Linux;
- системные вызовы для процессов управления, включая процессы реального времени;
- каталоги и файлы: создание, копирование, перемещение, удаление и управление ими;
- управление памятью: интерфейсы для распределенной памяти, оптимизация доступа к памяти.

Роберт Лав (Robert Love) — пользователь Linux и хакер с юных лет.

Он интересуется разработкой ядра Linux и среды рабочего стола GNOME, а также принимает активное участие в жизни Linux-сообщества. Его последним вкладом в ядро Linux стала работа над событиями уровня ядра и inotify. А вклад в GNOME включает Beagle, GNOME Volume Manager, Network Manager и Project Utopia. В настоящее время Роберт работает в отделе Open Source Program Office в Google.

ISBN: 978-5-388-00014-9



9 785388 000149

 **ПИТЕР**®

Заказ книг:

197198, Санкт-Петербург, а/я 619
тел.: (812) 703-73-74, postbook@piter.com

61093, Харьков-93, а/я 9130
тел.: (057) 758-41-45, 751-10-02, piter@kharkov.piter.com

www.piter.com — вся информация о книгах и веб-магазин

O'REILLY®