

Брайан Гоетс

при участии Тима Перлса, Джошуа Блоха, Джозефа Боубира,
Дэвида Холмса и Дага Ли

Параллельное программирование в JAVA на практике

Перевод: Сорокин А.В.



О книге

Java Concurrency на практике

Мне посчастливилось работать с фантастической командой над дизайном и реализацией функций параллелизма, добавленных к платформе Java в Java 5.0 и Java 6. Теперь эта же команда даёт лучшее объяснение этих новых функций и параллелизма в целом. Параллелизм больше не является темой только для опытных пользователей. Каждый Разработчик Java должен прочитать эту книгу.

—Martin Buchholz
JDK Concurrency Czar, Sun Microsystems

В течении прошедших 30 лет рост производительности компьютеров обусловливался законом Мура; отныне он будет зависеть от закона Амдала. Написание кода, эффективно использующего несколько процессоров, может быть очень сложным. Книга *Java Concurrency in Practice* знакомит вас с концепциями и техниками, необходимыми для написания безопасных и масштабируемых Java-программ для сегодняшних - и завтраших - систем.

—Doron Rajwan
Research Scientist, Intel Corp

Эта книга вам необходима, если вы пишите – или разрабатываете, или отлаживаете, или развёртываете, или обдумываете – многопоточные Java-программы. Если вам когда-либо приходилось «синхронизировать» метод, и вы небыли уверены в том, для чего это делаете, вы задолжали самому себе и вашим пользователям обязанность прочитать эту книгу, от корки до корки.

—Ted Neward
Author of Effective Enterprise Java

Брайан рассматривает фундаментальные проблемы и сложности параллелизма с необычайной ясностью. Эту книгу должен прочесть каждый из тех, кто использует потоки и заботится о производительности.

—Kirk Pepperdine
CTO, JavaPerformanceTuning.com

Эта книга охватывает очень глубокую и тонкую тему в очень ясной и краткой манере, что делает её превосходным справочным руководством по параллелизму Java. Каждая страница заполнена проблемами (и решениями!) с которыми программисты борются каждый день. Эффективное использование параллелизма становится все более и более важным, ведь теперь Закон Мура, предоставляет большее количество ядер, но они не становятся быстрее, и эта книга покажет вам, как это сделать.

—Dr. Cliff Click
Senior Software Engineer, Azul Systems

Я очень сильно интересуюсь параллелизмом, и, вероятно, написал большее количество взаимоблокировок потоков и допустил большее количество ошибок синхронизации, чем большинство других программистов. Книга Брайана наиболее читаема по теме параллельности и многопоточности в Java, и имеет дело с этой сложной темой с прекрасным практическим подходом. Это книга, которую я рекомендую всем моим читателям в The Java Specialists' Newsletter, потому что она интересна, полезна и релевантна к проблемам, с которыми разработчики Java сталкиваются лицом к лицу каждый день.

—Dr. Heinz Kabutz
The Java Specialists' Newsletter

Я сосредоточил свою карьеру на упрощении простых проблем, но эта книга честолюбиво и эффективно работает над упрощением сложной, но важной темы: параллелизма. Книга *Java Concurrency in Practice* является революционной в своем подходе, гладкой и с легким стилем, и очень своевременной в появлении — ей суждено стать очень важной книгой.

—Bruce Tate
Author of Beyond Java

Java Concurrency in Practice является неоценимой компиляцией ноу-хау о потоках для разработчиков Java. Я обнаружил, что чтение этой книги интеллектуально захватывающее, отчасти потому, что это превосходное введение в API параллелизма Java, но в основном потому, что она полностью и доступно охватывает экспертные знания в области потоков, которые нелегко найти в другом месте.

—Bill Venners
Author of Inside the Java Virtual Machine

Параллельное программирование в Java на практике

Java Concurrency in Practice

Brian Goetz

with

Tim Peierls

Joshua Bloch

Joseph Bowbeer

David Holmes

and Doug Lea

◆ Addison-Wesley

Upper Saddle River, NJ • Boston • Indianapolis • San Francisco New York • Toronto •
Montreal • London • Munich • Paris • Madrid Capetown • Sydney • Tokyo • Singapore •
Mexico City

Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trade-marks. Where those designations appear in this book, and the publisher was aware of a trademark claim, the designations have been printed with initial capital letters or in all capitals.

The authors and publisher have taken care in the preparation of this book, but make no expressed or implied warranty of any kind and assume no responsibility for errors or omissions. No liability is assumed for incidental or consequential damages in connection with or arising out of the use of the information or programs contained herein.

The publisher offers excellent discounts on this book when ordered in quantity for bulk purchases or special sales, which may include electronic versions and/or custom covers and content particular to your business, training goals, marketing focus, and branding interests. For more information, please contact:

U.S. Corporate and Government Sales
(800) 382-3419
corpsales@pearsontec.com
hgroup.com

For sales outside the United States, please contact:

International Sales
international@pearsoned.com

Visit us on the Web: www.awprofessional.com



This Book Is Safari Enabled

The Safari® Enabled icon on the cover of your favorite technology book means the book is available through Safari Bookshelf. When you buy this book, you get free access to the online edition for 45 days.

Safari Bookshelf is an electronic reference library that lets you easily search thousands of technical books, find code samples, download chapters, and access technical information whenever and wherever you need it.

To gain 45-day Safari Enabled access to this book:

- Go to <http://www.awprofessional.com/safarienabled>
- Complete the brief registration form
- Enter the coupon code UUIR-XRJG-JWWF-AHGM-137Z

If you have difficulty registering on Safari Bookshelf or accessing the online edition, please e-mail customer-service@safaribooksonline.com.

Library of Congress Cataloging-in-Publication Data

Goetz, Brian.

Java Concurrency in Practice / Brian Goetz, with Tim Peierls.

. . [et al.] p. cm.

Includes bibliographical references and index.

ISBN 0-321-34960-1 (pbk. : alk. paper)

1. Java (Computer program language) 2. Parallel programming (Computer science) 3. Threads (Computer programs) I. Title.

QA76.73.J38G588 2006

005.13'3--dc22

2006012205

Copyright © 2006 Pearson Education, Inc.
ISBN 0-321-34960-1

Джессике

От переводчика

Как родился этот перевод? Да, в общем-то, достаточно просто. Всегда хотелось прочитать эту книгу на русском языке, но годы шли, а издательства книгу всё не переводили. В какой-то момент появилось понимание, что книгу переводить не будут, т.к. она «устарела», ей уже 13 лет. И постепенно стала зреТЬ мысль, что можно её перевести самому. Вот так и вызрел план перевода. Когда я только планировал перевод, всё казалось легко и просто, думал, кавалерийским насоком, быстренько, за месяц всё сделаю. А потом началась работа над переводом, и бросать это дело было уже поздно. Перевод – довольно сложный процесс.

Приходится держать в голове множество понятий. Множество раз приходилось пробегаться по тексту книги, когда подбирал более удачный перевод для какого-либо термина. От этого в переводе могут быть небольшие несогласованности и опечатки. Но, несмотря на сложность, было очень интересно. Переведя эту книгу практически от корки до корки (не считая какой-то мишуры от издательства) могу сказать, что это фундаментальный труд, и он ни сколько не устарел. Считаю, что перевод этой книги является ценным вкладом в копилку знаний сообщества Java-разработчиков. Может быть, кто-то последует моему примеру, и также возьмётся перевести те книги, которые уже никогда не увидят свет на русском языке по коммерческим соображениям.

Репозиторий с проектом перевода книги:

[git@github.com:Lulay82/JavaConcurrencyInPracticeTranslation.git](https://github.com/Lulay82/JavaConcurrencyInPracticeTranslation.git)

Сорокин А.В.

Екатеринбург, июнь 2018 - декабрь 2018.

Оглавление

От переводчика	8
Оглавление	9
Предисловие.....	15
Как читать эту книгу	16
Примеры кода	18
Часть I Основы.....	19
Глава 1 Введение	20
1.1 Очень краткая история параллелизма.....	20
1.2 Преимущества потоков	22
1.2.1 Использование нескольких процессоров.....	22
1.2.2 Упрощение моделирования	22
1.2.3 Упрощённая обработка асинхронных событий	23
1.2.4 Более отзывчивый пользовательский интерфейс	24
1.3 Риски, которые несут потоки.....	24
1.3.1 Угрозы безопасности.....	25
1.3.2 Угрозы живучести потока	27
1.3.3 Угрозы производительности.....	28
1.4 Потоки есть везде	28
Глава 2 Потокобезопасность	31
2.1 Что такое потокобезопасность?.....	33
2.1.1 Пример: сервлет без сохранения состояния (stateless)	34
2.2 Атомарность.....	35
2.2.1 Условия гонок	36
2.2.2 Пример: состояние гонки в отложенной инициализации	37
2.2.3 Составные действия	38
2.3 Блокировка	40
2.3.1 Внутренние блокировки.....	41
2.3.2 Повторная входимость	42
2.4 Защита состояния с помощью блокировок	43
2.5 Живучесть и производительность	46
Глава 3 Совместно используемые объекты.....	50
3.1 Видимость	50
3.1.1 Устаревшие данные	52
3.1.2 Неатомарные 64 битные операции	53
3.1.3 Блокировки и видимость	53
3.1.4 Volatile переменные.....	54
3.2 Публикация и побег	56
3.2.1 Методы безопасного построения	58
3.3 Ограничение потока	59
3.3.1 Специальные ограничения потока	60
3.3.2 Ограничение стека.....	60
3.3.3 Класс ThreadLocal.....	62
3.4 Неизменяемость	63
3.4.1 Поля типа final	65
3.4.2 Использование <i>volatile</i> для публикации неизменяемых объектов	65
3.5 Безопасная публикация	66
3.5.1 Некорректная публикация: когда хорошие объекты ведут себя плохо	67
3.5.2 Неизменяемые объекты и безопасность инициализации	68
3.5.3 Идиомы безопасной публикации	69

3.5.4 Фактически неизменяемые объекты	70
3.5.5 Изменяемые объекты	70
3.5.6 Безопасное совместное использование объектов	71
Глава 4 Компоновка объектов	72
4.1 Проектирование потокобезопасных классов	72
4.1.1 Сбор требований к синхронизации	73
4.1.2 Операции, зависящие от состояния	74
4.1.3 Владелец состояния	74
4.2 Ограничение экземпляра	76
4.2.1 Шаблон Java “Монитор”	78
4.2.2 Пример: отслеживание парка автомобилей.....	78
4.3 Делегирование потокобезопасности	81
4.3.1 Пример: транспортный трекер с использованием делегирования	81
4.3.2 Независимые переменные состояния.....	83
4.3.3 Когда делегирование не работает	84
4.3.4 Публикация базовых переменных состояния	86
4.3.5 Пример: трекер транспортных средств публикующий своё состояние....	86
4.4 Добавление функциональности в существующие потокобезопасные классы	88
4.4.1 Блокировка на стороне клиента.....	89
4.4.2 Композиция	91
4.5 Документирование политики синхронизации	92
4.5.1 Интерпретация расплывчатой документации	93
Глава 5 Строительные блоки	95
5.1 Синхронизированные коллекции	95
5.1.1 Проблемы с синхронизированными коллекциями	95
5.1.2 Итераторы и ConcurrentModificationException.....	97
5.1.3 Скрытые итераторы.....	99
5.2 Параллельные коллекции.....	100
5.2.1 Класс ConcurrentHashMap.....	101
5.2.2 Дополнительные атомарные операции интерфейса Map	102
5.2.3 Класс CopyOnWriteArrayList	103
5.3 Блокирующие очереди и шаблон производитель-потребитель	104
5.3.1 Пример: поиск в компьютере	106
5.3.2 Последовательное ограничение потока	108
5.3.3 Двусторонние очереди и кражा работы	108
5.4 Методы блокирования и прерывания	109
5.5 Синхронизаторы	111
5.5.1 Защёлки	111
5.5.2 Класс FutureTask	113
5.5.3 Семафоры	115
5.5.4 Барьеры.....	117
5.6 Создание эффективного и масштабируемого кэша результатов.....	119
Итоги части I	126
Часть II Структурирование параллельных приложений	127
Глава 6 Выполнение задач	128
6.1 Выполнение задач в потоках	128
6.1.1 Последовательное выполнение задач	128
6.1.2 Явное создание потоков для задач	129
6.1.3 Недостатки неограниченного создания потоков	130
6.2 Фреймворк Executor	132
6.2.1 Пример: веб-сервер использующий фреймворк Executor.....	132
6.2.2 Политики выполнения.....	134

6.2.3 Пулы потоков	134
6.2.4 Жизненный цикл экземпляра Executor	136
6.2.5 Отложенные и периодические задачи	138
6.3 Поиск уязвимостей параллелизма	139
6.3.1 Пример: последовательный генератор страниц	139
6.3.2 Задачи возвращающие результат: Callable и Future	140
6.3.3 Пример: отрисовка страниц с помощью Future	142
6.3.4 Ограничения распараллеливания разнородных задач	143
6.3.5 CompletionService: Executor пересекается с BlockingQueue	144
6.3.6 Пример: рендер страниц с использованием CompletionService	145
6.3.7 Ограничение времени выполнения задач	146
6.3.8 Пример: портал бронирования путешествий	147
6.4 Итоги	149
Глава 7 Отмена и завершение	150
7.1 Отмена задачи	150
7.1.1 Прерывание	152
7.1.2 Политики прерывания	156
7.1.3 Ответ на прерывание	157
7.1.4 Пример: запуск по времени	159
7.1.5 Выполнение отмены с помощью интерфейса Future	161
7.1.6 Работа с непрерываемыми блокировками	162
7.1.7 Инкапсуляция нестандартной отмены с помощью метода newTaskFor	164
7.2 Остановка служб основанных на потоках	166
7.2.1 Пример: Сервис логирования	166
7.2.2 Завершение работы ExecutorService	169
7.2.3 Ядовитые пилюли	170
7.2.4 Пример: одноразовая служба выполнения	172
7.2.5 Ограничения метода shutdownNow	173
7.3 Обработка аварийного завершения потока	175
7.3.1 Обработчики неперехваченных исключений	177
7.4 Завершение работы JVM	178
7.4.1 Завершающие хуки	179
7.4.2 Потоки демоны	180
7.4.3 Финализаторы	181
7.5 Итоги	181
Глава 8 Применение пулов потоков	182
8.1 Неявные связи между задачами и политиками выполнения	182
8.1.1 Взаимоблокировка потоков, вызванная голоданием	183
8.1.2 Долговременные задачи	185
8.2 Размеры пулов потоков	185
8.3 Конфигурирование класса ThreadPoolExecutor	186
8.3.1 Создание и удаление потока	187
8.3.2 Управление задачами в очереди	188
8.3.3 Политики насыщения	190
8.3.4 Фабрики задач	192
8.3.5 Настройка экземпляра ThreadPoolExecutor после построения	194
8.4 Расширение класса ThreadPoolExecutor	194
8.4.1 Пример: добавление статистики в пул потоков	195
8.5 Распараллеливание рекурсивных алгоритмов	196
8.5.1 Пример: фреймворк для решения головоломок	198
8.6 Итоги	203
Глава 9 Приложения GUI	204

9.1 Почему фреймворки GUI однопоточны?	204
9.1.1 Последовательная обработка событий	206
9.1.2 Ограничение потока в Swing	206
9.2 Кратковременные задачи GUI	208
9.3 Долговременные задачи GUI.....	210
9.3.1 Отмена	211
9.3.2 Индикатор прогресса и завершения.....	212
9.3.3 Класс SwingWorker	214
9.4 Совместно используемые модели данных	215
9.4.1 Потокобезопасные модели данных	215
9.4.2 Разделение моделей данных	216
9.5 Другие формы однопоточных подсистем.....	217
9.6 Итоги.....	217
Часть 3 III Живучесть, производительность и тестирование	218
Глава 10 Предотвращение возникновения угроз живучести.....	219
10.1 Взаимоблокировки	219
10.1.1 Взаимоблокировки, вызванные порядком наложения блокировок.....	220
10.1.2 Взаимоблокировки, вызванные динамическим порядком блокировок	221
10.1.3 Взаимоблокировки между взаимодействующими объектами	225
10.1.4 Открытые вызовы	226
10.1.5 Взаимоблокировки ресурсов	229
10.2 Предотвращение и диагностика взаимоблокировок.....	229
10.2.1 Блокировки, ограниченные по времени.....	230
10.2.2 Анализ взаимоблокировок с использованием дампов потоков	230
10.3 Прочие угрозы живучести	232
10.3.1 Голодание	232
10.3.2 Плохая отзывчивость	233
10.3.3 Динамическая взаимоблокировка	234
10.4 Итоги.....	234
Глава 11 Производительность и масштабируемость.....	236
11.1 Размышления о производительности.....	236
11.1.1 Производительность и масштабируемость	237
11.1.2 Оценка компромиссов производительности	238
11.2 Закон Амдала	240
11.2.1 Пример: скрытое последовательное выполнение в фреймворках	243
11.2.2 Качественное применение закона Амдала	244
11.3 Затраты, вводимые потоками	245
11.3.1 Переключение контекста	245
11.3.2 Синхронизация памяти	246
11.3.3 Блокирование	247
11.4 Уменьшение конкуренции за блокировку	248
11.4.1 Сужение области действия блокировки (“Get in, get out”)	249
11.4.2 Уменьшение детализации блокировки	251
11.4.3 Чередование блокировок	253
11.4.4 Как избежать использования “горячих полей”	254
11.4.5 Альтернативы монопольным блокировкам.....	255
11.4.6 Мониторинг использования CPU	256
11.4.7 Просто скажите “нет” помещению объектов в пул	257
11.5 Пример: Сравнение производительности реализаций Map	258
11.6 Сокращение накладных расходов на переключение контекста	260
11.7 Итоги.....	261
Глава 12 Тестирование параллельных программ	263

12.1 Тестирование на корректность	264
12.1.1 Простые модульные тесты.....	265
12.1.2 Тестирование блокирующих операций	266
12.1.3 Тестирование безопасности	268
12.1.4 Тестирование управления ресурсами	272
12.1.5 Использование обратных вызовов	274
12.1.6 Увеличение степени чередования	275
12.2 Тестирование производительности	276
12.2.1 Расширение класса PutTakeTest с добавлением учёта времени	276
12.2.2 Сравнение нескольких алгоритмов	279
12.2.3 Измерение отзывчивости	280
12.3 Как избежать ошибок тестирования производительности	282
12.3.1 Сборка мусора.....	282
12.3.2 Динамическая компиляция	283
12.3.3 Нереалистичная выборка веток выполнения кода	284
12.3.4 Нереалистичные предположения о степенях конкуренции.....	285
12.3.5 Устранение мёртвого кода.....	285
12.4 Комплементарные подходы к тестированию	287
12.4.1 Ревю кода	287
12.4.2 Инструменты для проведения статического анализа	288
12.4.3 Аспектно-ориентированные подходы к тестированию.....	290
12.4.4 Профилировщики и инструменты мониторинга.....	290
12.5 Итоги.....	290
Часть IV Дополнительные темы	292
Глава 13 Явные блокировки	293
13.1 Интерфейсы Lock и ReentrantLock.....	293
13.1.1 Опрашиваемый и ограниченный по времени захват блокировки	294
13.1.2 Прерываемый захват блокировки	296
13.1.3 Блокировка с не-блочной структурой.....	297
13.2 Вопросы производительности	298
13.3 Справедливость.....	299
13.4 Выбор между synchronized и ReentrantLock.....	301
13.5 Блокировки на чтение-запись	302
13.6 Итоги.....	305
Глава 14 Разработка собственных синхронизаторов.....	306
14.1 Управление зависимостью от состояния	306
14.1.1 Пример: распространение сбоев предусловий на вызывающий код	308
14.1.2 Пример: грубая блокировка путем опроса и сна	310
14.1.3 Очереди условий на освобождение.....	311
14.2 Использование очередей условий	313
14.2.1 Предикат условия	313
14.2.2 Раннее пробуждение.....	314
14.2.3 Пропущенные сигналы.....	316
14.2.4 Уведомление	316
14.2.5 Пример: класс затвора.....	318
14.2.6 Вопросы безопасности подклассов	319
14.2.7 Инкапсулирование очередей условий.....	320
14.2.8 Протоколы входа и выхода	321
14.3 Явные объекты условия	321
14.4 Анатомия синхронизатора	324
14.5 Класс AbstractQueuedSynchronizer	326
14.5.1 Простая защёлка	328

14.6 AQS в классах-синхронизаторах из пакета <code>java.util.concurrent</code>	329
14.6.1 Класс <code>ReentrantLock</code>	329
14.6.2 Классы <code>Semaphore</code> и <code>CountDownLatch</code>	330
14.6.3 Класс <code>FutureTask</code>	331
14.6.4 Класс <code>ReentrantReadWriteLock</code>	332
14.7 Итоги.....	332
Chapter 15 Атомарные переменные и неблокирующая синхронизация	333
15.1 Недостатки блокировки	333
15.2 Аппаратная поддержка параллелизма	335
15.2.1 Сравнить и обменять	335
15.2.2 Неблокирующий счётчик.....	337
15.2.3 Поддержка операций CAS в среде JVM	338
15.3 Классы атомарных переменных	339
15.3.1 <code>Atomics</code> as “better volatiles”	340
15.3.2 Производительность: блокировки против атомарных переменных	341
15.4 Неблокирующие алгоритмы	344
15.4.1 Неблокирующий стек	344
15.4.2 Неблокирующий связанный список.....	346
15.4.3 Обновления атомарного поля	349
15.4.4 Проблема ABA	350
15.5 Итого.....	351
Глава 16 Модель памяти Java	352
16.1 Что такое модель памяти и зачем она мне нужна?	352
16.1.1 Основа моделей памяти	353
16.1.2 Переупорядочивание	354
16.1.3 О модели памяти Java менее чем в 500 словах	355
16.1.4 Комбинирование в синхронизации	357
16.2 Публикация	360
16.2.1 Небезопасная публикация.....	360
16.2.2 Безопасная публикация	361
16.2.3 Идиомы безопасной инициализации	362
16.2.4 Блокировка с двойной проверкой	363
16.3 Безопасность инициализации	364
16.3 Итоги.....	366
Приложение А Описание аннотаций	367
A.1 Аннотации уровня классов	367
A.2 Аннотации уровня полей и методов	367
Библиография	369

Предисловие

Только сейчас, когда пишется книга, многоядерные процессоры стали достаточно доступными, для использования в настольных системах среднего уровня. Не случайно, что множество команд разработчиков в своих проектах получают всё больше и больше отчётов об ошибках, связанных с потоками. В недавнем посте, на сайте разработчиков NetBeans, один из ключевых мэнейнеров (*maintainers*) отметил, что единственный класс был поправлен более 14 раз, для исправления проблем, связанных с потоками. Дин Алмайер, бывший редактор TheServerSide, недавно написал в своём блоге (после мучительной сессии отладки, в конечном итоге выведшей к ошибке с потоками), что в большинстве Java-программ так распространены ошибки многопоточности, что они работают только “случайно”.

В самом деле, разработка, тестирование и отладка многопоточных программ может быть исключительно трудной, потому что ошибки параллельности проявляют себя непредсказуемо. И когда они всплывают на поверхность, это часто происходит в наихудший из возможных моментов времени – в продуктиве, под большой нагрузкой.

Одним из вызовов в разработке параллельных программ на Java, является несоответствие между возможностями параллелизма, предлагаемыми платформой и тем, как разработчики вынуждены думать о параллельности в своих программах. Язык предоставляет низкоуровневые механизмы, такие как синхронизация (*synchronization*) и условие ожидания (*condition waits*), но эти механизмы должны использоваться согласованно, для реализации протоколов или политик уровня приложения. Без таких политик очень просто создавать программы, которые компилируются и представляются работающими, но, тем не менее, содержат ошибки. Множество других превосходных книг о параллельности не достигают своей цели, фокусируясь на низкоуровневых механизмах и API, а не на политиках и шаблонах уровня проектирования.

В Java 5.0 был сделан гигантский шаг вперёд в разработке параллельных приложений Java, предоставляя новые, высокоуровневые компоненты, и дополнительные, низкоуровневые механизмы, что одинаково упростило построение параллельных приложений, как для новичков, так и для экспертов. Авторы являются основными членами экспертной группы JCP, создавшей эти средства; в дополнение, к описанию их поведения и возможностей, мы представляем лежащие в основе шаблоны проектирования и ожидаемые сценарии использования, что мотивировало их включение в библиотеки платформы.

Мы ставили себе цель дать пользователю набор правил проектирования и мыслительных моделей, которые делают проще, – и веселее – построение корректных, производительных параллельных классов и приложений Java.

Мы надеемся, что вам понравится книга *Java Concurrency in Practice*.

Brian Goetz
Williston, VT
March 2006

Как читать эту книгу

Для устранения абстрактного несоответствия между низкоуровневыми механизмами Java и необходимыми политиками уровня проектирования, мы представили *упрощённый* набор правил для написания параллельных программ. Эксперты могут посмотреть на эти правила и сказать, - “Хмм, это не совсем так: класс С потокобезопасен, даже если он нарушает правило R”. Хотя возможно писать корректные программы, нарушающие наши правила, это требует глубокого понимания низкоуровневых деталей модели памяти Java (*Java Memory Model*), и мы хотели, чтобы разработчики имели возможность писать корректные параллельные программы, *без необходимости освоения* этих деталей. Последовательное следование нашим упрощенным правилам позволит создавать правильные и сопровождаемые параллельные программы.

Мы предполагаем, что читатель уже в некоторой степени знаком с базовыми механизмами параллелизма в Java. *Java Concurrency in Practice* не является введением в тему параллелизма – для этого, смотрите главу о потоках любого достойного талмуда, такого как *The Java Programming Language* (Arnold et al., 2005). И это не энциклопедический справочник для всех вещей параллелизма – за этим обращайтесь к *Concurrent Programming in Java* (Lea, 2000). Напротив, книга предлагает практические правила проектирования, помогающие разработчикам в таком сложном процессе, как создание безопасных и производительных параллельных классов. При необходимости, мы ссылаемся на соответствующие разделы из книг *The Java Programming Language*, *Concurrent Programming in Java*, *The Java Language Specification* (Gosling et al., 2005), and *Effective Java* (Bloch, 2001) используя сокращения [JPL n.m], [CPJ n.m], [JLS n.m], and [EJ Item n].

После введения (Глава 1), книга делится на четыре части:

Основы. Часть I (Главы 2 - 5) фокусируется на базовых концепциях параллелизма, потокобезопасности и на том, как составлять потокобезопасные классы из параллельных строительных блоков, предоставляемых библиотекой классов. “Шпаргалка” в разделе “Итоги части I” суммирует наиболее важные из правил, представленных в части I.

Главы 2 (потокобезопасность) и 3 (совместно используемые объекты) формируют основу книги. Почти все правила позволяют избегать рисков параллелизма, конструировать потокобезопасные классы и проверять безопасность потоков. Читатели, предпочитающие “практику” вместо “теории”, могут податься соблазну перейти к части II, но должны обязательно вернуться к главам 2 и 3 до написания любого параллельного кода!

Глава 4 (Компоновка объектов) покрывает техники для компоновки потокобезопасных классов в большие потокобезопасные классы. Глава 5 (Строительные блоки) покрывает тему параллельных строительных блоков – потокобезопасных коллекций и синхронизации – предоставляемых библиотекой платформы.

Структурирование параллельных программ. В части II (Главы 6-9) описывается, как эксплуатировать потоки для увеличения пропускной способности и отзывчивости параллельных приложений Глава 6 (Выполнение задач) затрагивает тему идентификации параллелизуемых задач и выполнения их внутри выполняющих задачи фреймворков (*task-execution framework*). Глава 7 (Отмена и завершение) имеет дело с техниками убеждения задач и потоков завершаться до того, как они смогут это сделать нормально; то, как программы имеют дело с

отменой и завершением, часто является одним из факторов, который отделяет действительно надёжные параллельные приложения от тех, что просто работают. В главе 8 (Применение пула потоков) рассматриваются некоторые более продвинутые возможности фреймворков выполнения задач. Глава 9 (Приложения GUI) фокусируется на техниках по повышению отзывчивости в однопоточных подсистемах.

Жизнеспособность, производительность и тестирование. Часть III (Главы 10-12) заботятся о том, чтобы параллельные программы действительно делали то, что вы хотите, и делали это с приемлемой производительностью. Глава 10 (Как избежать угроз живучести) описывает, как предотвращать проблемы живучести, которые могут помешать программам, продвигаться вперёд. Глава 11 (Производительность и масштабируемость) описывает техники, для увеличения производительности и масштабируемости в параллельном коде. Глава 12 (Тестирование параллельных программ) рассматривает техники тестирования параллельного кода на предмет корректности и производительности.

Дополнительные темы. Часть IV (Главы 13-16) затрагивает темы, которые понравятся и будут интересны только опытным разработчикам: явные блокировки, атомарные переменные, неблокирующие алгоритмы, и разработка собственных синхронизаторов (*custom synchronizers*).

Примеры кода

В то время как множество концепций в этой книге применимы к версиям Java, предшествующим Java 5.0 и даже к не Java окружению, большинство примеров кода (и все утверждения о модели памяти Java) предполагают Java 5.0 или старше. Некоторые примеры кода могут использовать функции библиотеки, добавленные в Java 6.

Примеры кода были сжаты для уменьшения размера и для выделения соответствующих частей. Полные версии примеров кода, а также дополнительные примеры и список опечаток, доступны на сайте книги <http://www.javaconcurrencyinpractice.com>.

Примеры кода делятся на три категории: “хорошие” примеры, “не очень хорошие” примеры и “плохие” примеры. Хорошие примеры иллюстрируют техники, которым необходимо следовать. Плохие примеры иллюстрируют техники, которые, определённо, не должны воспроизводиться, и они, чтобы дать ясно понять, что это токсичный код, представляются с иконкой “**Mr. Yuk**”¹ (см. листинг 1). Не очень хорошие примеры иллюстрируют техники, которые необязательно ошибочные, но являются хрупкими, рискованными или имеют плохую производительность, они помечаются иконкой “**Mr. Could Be Happier**”, как показано в листинге 2.

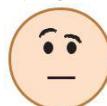
```
public <T extends Comparable<? super T>> void sort(List<T> list) {  
    // Never returns the wrong answer!  
    System.exit(0);  
}
```



Листинг 1. Плохой способ сортировки списка. Не делайте так!

Некоторые читатели могут задать вопрос о роли *плохих* примеров в этой книге; в конце концов, книга должна показать, как делать все правильно, а не ошибаться. Плохие примеры преследуют две цели. Они иллюстрируют наиболее часто встречающие подводные камни, но, что более важно, они демонстрируют, как анализировать программы на потокобезопасность, и лучший путь сделать это, это увидеть пути, которые приводят к тому, что потокобезопасность компрометируется.

```
public <T extends Comparable<? super T>> void sort(List<T> list) {  
    for (int i=0; i<1000000; i++)   
        doNothing();  
    Collections.sort(list);  
}
```



Листинг 2. Не самый оптимальный способ сортировки списка

¹ **Mr. Yuk** является зарегистрированной торговой маркой детского госпиталя в Питтсбурге и используется с разрешения.

Часть I Основы

Глава 1 Введение

Написать корректную программу трудно. Написать многопоточную корректную программу ещё труднее. В параллельной программе может произойти огромное количество вещей, которые выполняются не так, как ожидалось, хотя, вполне корректно, сработают в последовательной программе.

Так почему мы связываемся с многопоточностью? Потоки являются неизбежной частью языка Java, и они могут существенно упростить разработку комплексных систем путём превращения сложного асинхронного кода в более простой прямолинейный код. В дополнение, потоки являются простейшим способом использовать вычислительную мощь многопроцессорных систем. По мере увеличения количества процессоров, важность эффективного использования параллелизма будет только возрастать.

1.1 Очень краткая история параллелизма

В давнем прошлом, компьютеры не имели операционных систем; они выполняли единственную программу от начала и до конца, и программа имела прямой доступ ко всем ресурсам машины. Писать программы, которые работали на голом железе, было не единственной сложностью, другой была возможность запуска за раз только одной программы, что вызывало неэффективное использование дорогостоящих и ограниченных ресурсов компьютера.

Эволюционировав, операционные системы позволили запускать более одной программы за раз, запуская программы в индивидуальных процессах: изолированных, независимо выполняемых программах, для которых операционная система выделяет ресурсы, такие как память, указатели на файлы и ограничения системы безопасности. Также, если необходимо, процессы могут взаимодействовать с другими процессами через различные механизмы коммуникации: сокеты, обработчики сигналов, разделяемую память, семафоры и файлы.

Несколько мотивирующих факторов привели к разработке операционных систем, которые позволяли одновременно выполнять несколько программ:

Использование ресурсов. Программы иногда вынуждены ожидать завершения внешних операций, таких как ввод или вывод, и пока длится ожидание, не могут выполнять полезную работу. Гораздо эффективнее использовать время ожидания, позволив выполняться другой программе.

Справедливость. Множество ресурсов и программ могут иметь эквивалентные запросы к ресурсам машины. Предпочтительно позволить им разделять ресурсы компьютера на равно-порционные куски времени, чем ожидать завершения одной программы и запускать другую

Удобство. Часто проще и предпочтительнее написать несколько программ, каждая из которых выполняет одну задачу и имеет возможность координироваться с другими (если это необходимо), чем написать одну программу, которая выполняет сразу все задачи.

В ранних системах с разделением времени, каждый процесс являлся виртуальным Фон-Неймановским компьютером, он имел пространство в памяти

для совместного хранения инструкций и данных, выполнялись инструкции последовательно, в соответствии с семантикой языка машины, и взаимодействовали с внешним миром через набор примитивов ввода/вывода операционной системы. Для каждой выполняемой инструкции была четко определена «следующая инструкция» и контроль потока выполнения программы осуществлялся строго в соответствии с правилами набора инструкций. Почти все широко используемые сегодня языки программирования следуют этой последовательной программной модели, где спецификация языка программирования чётко определяет, «что следует далее» после выполнения некоторой команды.

Модель последовательного программирования интуитивна и естественна, и, фактически повторяет модель работы человека: делать по одной вещи в один момент времени, в основном – последовательно. Встать с постели, облачиться в халат, спустится вниз и приготовить чай. Как и в языках программирования, каждое из действий реального мира является абстракцией для последовательности небольших действий - открыть шкаф, выбрать сорт чая, отмерить некоторое количество в чайник, проверить, достаточно ли в чайнике воды, если нет, то долить немного воды в него, поставить чайник на плиту, ожидать пока вода закипит, и так далее. Этот последний шаг – ожидание, пока закипит вода – предполагает некоторую степень *асинхронности* (*asynchrony*). Пока вода нагревается, вы можете выбрать, что делать – просто подождать, или выполнить другую задачу в это время, такую как начать жарить тосты (другая асинхронная задача) или почитать новости, учитывая, что в скором времени вашего внимания потребует чайник. Производители чайников и тостеров знают, что их продукты часто используются в асинхронной манере, поэтому устройства подают звуковой сигнал, когда завершают свою работу. Поиск правильной балансировки последовательных и асинхронных действий является характеристикой эффективного человека – это также справедливо для программ.

То же самое (утилизация ресурсов, справедливое распределение и удобство) что мотивировал разработку процессов, также мотивировало разработку потоков (*threads*). Потоки позволяют нескольким потокам управления программным потоком (*control flow*) сосуществовать в процессе. Они разделяют общие ресурсы процесса, такие как память и указатели на файлы, о каждый поток имеет свой собственный программный счётчик, стек, и локальные переменные. Потоки также предоставляют естественный механизм декомпозиции для использования аппаратного параллелизма в многопроцессорных системах; множество потоков внутри программы могут быть запущены одновременно на множество ЦПУ.

Потоки иногда называют легковесными процессами, и большинство современных операционных систем рассматривают потоки, не процессы, как простейшие единицы планирования. В отсутствие явной координации, потоки выполняются одновременно и асинхронно относительно друг друга. С момента обращения потоков к общей разделяемой памяти собственного процесса, все потоки внутри процесса имеют доступ к переменным и выделяют объекты из кучи (*heap*), которая позволяет организовать более изящное совместное использование данных (*data sharing*), чем механизмы межпроцессного взаимодействия. Но без явной синхронизации для координации доступа к совместно используемым данным, поток может изменить переменную, пока другой поток находится в процессе использования переменной (*in the middle of using*), что приведёт к непредсказуемым результатам.

1.2 Преимущества потоков

Когда потоки используются должным образом, они могут снизить сложность разработки и стоимость сопровождения кода и повысить общую производительность приложения. Потоки упрощают формирование модели поведения и взаимодействия человека, путём превращения асинхронных рабочих процессов в большей части последовательные. В другом случае они могут превратить запутанный код в прямолинейный, который проще писать, читать и сопровождать.

Потоки полезны в приложениях GUI для повышения отзывчивости пользовательского интерфейса, и в серверных приложениях для улучшения использования ресурсов и повышения пропускной способности. Они также упрощают реализация JVM – сборщик мусора обычно запускается в одном или более выделенных потоках. Наиболее нетривиальные приложения на Java при создании в большей степени полагаются на потоки.

1.2.1 Использование нескольких процессоров

Ранее многопроцессорные системы использовались редко и стоили дорого, как правило, они находились только в больших ЦОД-ах (*data centers*) или применялись в научном оборудовании. Сегодня они дешевы и многочисленны, даже малые сервера или среднего уровня рабочие станции часто имеют множество процессоров. Эта тенденция будет только возрастать; так как всё сложнее становится увеличивать тактовую частоту, производители процессоров будут вместо этого располагать большее количество ядер на одном чипе. Все ведущие производители чипов начали этот переход, и мы уже видим машины с очень большим количеством процессоров.

Поскольку основной единицей планирования является поток, программа только с одним потоком может запускаться не более чем на одном процессоре за раз. На двух процессорной системе, однопоточная программа получит доступ к половине доступных ресурсов ЦПУ; на 100-процессорной системе она откажется от 99% доступных ресурсов. С другой стороны, программы с множеством активных потоков могут одновременно выполняться на множестве процессоров. Будучи должным образом спроектирована, многопоточная программа может увеличить пропускную способность путём более эффективного использования доступных ресурсов процессора.

Использование многопоточности может также помочь достигнуть большей пропускной способности на однопроцессорной машине. Если программа однопоточная, процессор просто ожидает (*idle*) завершения синхронных операций ввода/вывода. Если программа многопоточная, другой поток может все ещё выполняться, пока первый поток ожидает завершения операций ввода/вывода, позволяя приложению продолжать работу во время блокировки ввода/вывода. (Это похоже на чтение новостей, во время ожидания закипания воды в чайнике, вместо того, чтобы ждать пока вода закипит, прежде чем начать читать.)

1.2.2 Упрощение моделирования

Часто вам проще управлять временем, когда вы выполняете задачу только одного вида (правите эти 12 ошибок(*bugs*)), чем, когда вы имеете несколько их (исправляете эти ошибки, проводите интервью с кандидатами на замену системного администратора, оценивание производительность вашей команды, и

создаёте слайды для презентации на следующей неделе). Когда вы заняты задачей только одного вида, вы можете начать сверху вороха и взять её в работу, и так до тех пор, пока вся куча не будет разобрана; (или вы); вам нет необходимости затрачивать мыслительную энергию для понимания, какая работа будет следующей. С другой стороны, управление множеством приоритетов и сроков и переключение от задачи к задаче обычно вносит накладные расходы.

Также это справедливо и для программного обеспечения: программа, которая последовательно выполняет только один тип задач, проще в написании, менее подвержена ошибкам, и проще в тестировании, чем одна, управляющая множеством различных типов задач за раз. Назначение потоков каждому типу задач, или даже элементам, в симуляции создаёт иллюзию последовательности и изолированности логики домена (*domain logic*) от деталей планирования, чередования операций, асинхронного ввода/вывода, и ожидания ресурсов. Сложные, асинхронные рабочие процессы (*workflows*) могут быть разбиты на конечное количество простых, синхронных рабочих процессов, каждый из которых запущен в отдельном потоке, взаимодействующих друг с другом только в определённых точках синхронизации (*synchronization points*).

Эти преимущества часто эксплуатируются фреймворками, такими как сервлеты или RMI (удалённый вызов процедур). Фреймворки берут на себя детали управления запросами, создания потоков, балансировкой нагрузки, диспетчериизуют порции запросов обрабатываемых соответствующими компонентами приложения в соответствующих точках рабочего процесса. Тем, кто пишет сервлеты, нет необходимости беспокоиться о том, как много запросов будет обработано в один момент времени или о блоке ввода и вывода сокета; когда сервисный метод сервлета (*service method*) вызывается в ответ на веб-запрос, он может обрабатывать запрос асинхронно, как если бы это была однопоточная программа. Это может упростить разработку компонентов и снизить кривую обучения для использования фреймворков.

1.2.3 Упрощённая обработка асинхронных событий

Серверное приложение, принимающее подключения сокетов от множества удалённых клиентов может быть проще в разработке, когда под каждое подключение выделяется собственный поток и в нём позволяет использовать синхронный ввод/вывод.

Если приложение приступает к чтению из сокета, когда данные недоступны, блоки будут считываться до тех пор, пока данные не станут доступны. В однопоточном приложении это значит, что приостановится не только обработка соответствующего запроса, но будет приостановлена обработка всех запросов, пока один поток заблокирован. Для обхода этой проблемы, однопоточные серверные приложения вынуждены использовать неблокирующий ввод/вывод, который намного сложнее и более подвержен ошибкам, чем синхронные операции ввода/вывода. Однако, если каждый запрос имеет свой собственный поток, тогда блокировка не оказывает влияние на обработку других запросов.

Исторически так сложилось, что операционные системы располагали относительно низким лимитом на количество потоков, которые мог создать процесс, всего несколько сотен (или даже меньше). В результате операционные системы разработали эффективные средства для мультиплексированного ввода/вывода, такие как системные вызовы UNIX *select* и *poll*, и для доступа к этим средствам библиотекой классов Java был добавлен набор пакетов (*java.nio*) для неблокирующего ввода/вывода. Тем не менее, операционные системы стали

поддерживать значительно большее количество потоков, таким образом, предоставив возможность практической реализации модели *отдельный-поток-для-каждого-клиента*, даже для большого количества клиентов для некоторых платформ²

1.2.4 Более отзывчивый пользовательский интерфейс

Приложения GUI ранее были однопоточными, что приводило к тому, что-либо приходилось часто опрашивать код на возникновение входящих событий (которые беспорядочны и назойливы), либо выполнять весь код приложения косвенно, через “основной цикл событий (*main event loop*)”. Если код, вызванный из основного цикла событий, выполнялся слишком долго, то пользовательский интерфейс “замерзал” в ожидании завершения его работы, потому что последующие события GUI не могут быть обработаны до тех пор, пока не будет возвращено управление основному циклу событий.

Современные фреймворки GUI³, такие как AWT или Swing toolkits, заменили основной цикл событий механизмом рассылки событий (*event dispatch thread, edt*). Когда возникает событие пользователяского интерфейса, например, такое как нажатие кнопки, в потоке событий будут вызваны объявленные в приложении обработчики. Большинство фреймворков GUI имеют однопоточные подсистемы, так что основной цикл событий всё ещё присутствует, но он запускается в собственном потоке под управлением инструментария GUI, а не приложения.

Если в потоке событий выполняются кратковременные задачи, интерфейс остаётся отзывчивым, поскольку поток событий всегда способен обрабатывать действия пользователя достаточно быстро. Однако, обработка долго выполняющихся задач в потоке событий, таких как проверка орфографии, в большом документе или получение ресурсов по сети, ухудшает отзывчивость. Если пользователь выполняет какие-то действия, пока выполняется такая задача, появится большая задержка до того момента, как поток событий сможет обработать их или даже просто узнать о них. Сгущая краски, отметим, что пользовательский интерфейс не только становится недоступен, но также невозможно нажать кнопку отмены задачи, даже если UI имеет таковую, потому что поток событий занят и не может обработать событие нажатия кнопки до завершения длинной задачи! Однако, если вместо этого, долго работающая задача выполняется в отдельном потоке, поток событий остаётся свободным для обработки событий UI, делая UI более отзывчивым.

1.3 Риски, которые несут потоки

Встроенная в Java поддержка потоков — это обоюдоострый меч. Хотя это упрощает разработку многопоточных приложений, предоставляя поддержку со стороны языка и библиотек и формальную кроссплатформенную модель памяти (именно эта формальная кроссплатформенная модель памяти Java делает возможной разработку однажды написанных, запускающихся где угодно многопоточных приложений), это также поднимает планку уровня разработчиков, потому что всё больше количества программ будет использовать потоки. Когда

² Пакет потоков NPTL, сейчас являющийся частью большинства дистрибутивов Linux, был спроектирован для поддержки сотен тысяч потоков. Неблокирующий ввод/вывод имеет свои преимущества, но улучшенная поддержка ОС для потоков означает, что существует немного ситуаций, для которых необходимо его использование.

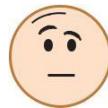
³ Для тех лет это было круто, по современным меркам – ад.

потоки были вещью скорее эзотерической, параллельность была “подвинутой” темой; сейчас общий тренд таков, что разработчики должны быть осведомлены о безопасном использовании потоков.

1.3.1 Угрозы безопасности

Потокобезопасность может быть неожиданно тонким понятием: из-за отсутствия достаточной синхронизации, порядок операций во множестве потоков непредсказуем и иногда случаются сюрпризы. Класс `UnsafeSequence` в листинге 1.1, который должен генерировать последовательность уникальных целых чисел, предлагает простую иллюстрацию того, как чередование действий в множестве потоков может привести к нежелательным результатам. Он ведёт себя корректно в пределах однопоточного окружения, но в многопоточном окружении всё иначе.

```
@NotThreadSafe
public class UnsafeSequence {
    private int value;
    /* Returns a unique value. */
    public int getNext() {
        return value++;
    }
}
```



Листинг 1.1 Непотокобезопасный генератор последовательностей.

Проблема с `UnsafeSequence` возникает в неудачный момент времени, два потока могут вызвать `getNext` и получить *некоторое значение*. На рис. 1.1 показано, как это может произойти. Нотация инкремента, `someVariable++`, может казаться единой операцией, но это фактически *три* раздельных операции: чтения значения, добавления единички к нему, и запись нового значения в переменную. Так как операции в множестве потоков могут произвольно чередоваться во время выполнения, вполне возможно для двух потоков прочитать значения в один момент времени, и затем добавить по единичке к ним. В результате один и тот же порядковый номер возвращается из нескольких вызовов в разных потоках.

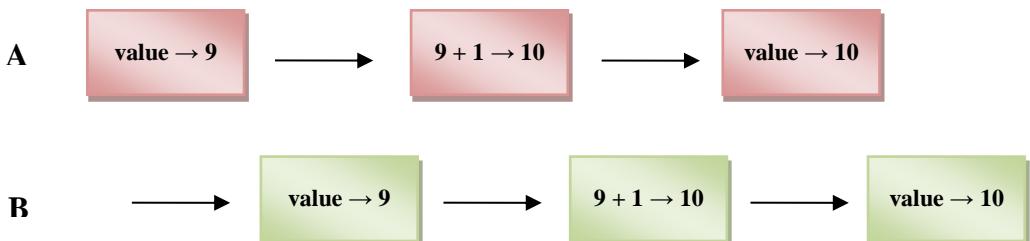


Рис. 1.1 Неудачное выполнение `UnsafeSequence.getNext`.

Диаграммы на рис. 1.1 показывают возможное чередование операций в различных потоках. На этих диаграммах, движется слева направо, и каждая линия представляет активность различных потоков. Эти чередующиеся

диаграммы обычно показывают наихудший случай⁴ и предназначены для того, чтобы показать опасность неверного предположения, о том, что все произойдет в строго определенном порядке.

Класс `UnsafeSequence` использует нестандартную аннотацию: `@NotThreadSafe`. Это одна из нескольких пользовательских аннотаций, используемых на протяжении всей книги, для документирования параллелизма свойств и членов классов. (Аналогичным образом используются другие аннотации уровня класса - `@ThreadSafe` и `@Immutable`; см. Приложение А) Аннотации, документирующие потокобезопасность, полезны для нескольких аудиторий. Если классы аннотированы с помощью `@ThreadSafe`, пользователи могут с уверенностью использовать их в многопоточном окружении, со своей стороны, разработчик ставит в известность, что потокобезопасность гарантирована и должна быть сохранена, и инструменты для анализа кода могут определить возможные ошибки кодирования.

Класс `UnsafeSequence` показывает распространённую проблему параллелизма, называемую *условием гонок* (*race condition*). Возвращает ли `getNext` уникальное значение при вызове из нескольких потоков, как того требует спецификация, или нет, зависит от того, как среда выполнения (*runtime*) чередует операции - это является нежелательным состоянием дел.

Из-за того, что потоки разделяют общее адресное пространство и запускаются параллельно, они могут получать доступ и модифицировать те переменные, которые могут использовать другие потоки. Это потрясающе удобно, потому что позволяет совместно использовать данные проще, чем в любых других межпоточных механизмах коммуникации. Но, такой механизм коммуникации, также несёт в себе значительные риски: потоки могут быть дискредитированы (*confused*) при неожиданном изменении данных. Возможность множества потоков получить доступ и модифицировать некоторые переменные вводит элемент непоследовательности в *последовательную* модель программирования, что может сбивать с толку и быть сложным для понимания. Для многопоточных программ поведение должно быть предсказуемым, доступ к совместно используемым (*shared*) переменным должен быть правильно скоординирован так, чтобы потоки не вмешивались в работу друг друга. К счастью, Java предоставляет механизмы синхронизации для координации такого доступа.

Класс `UnsafeSequence` может быть исправлен путём превращения метода `getNext` в синхронизированный (*synchronized*), как показано в листинге 1.2⁵, тем самым предотвращая неудачное взаимодействие как на рис. 1.1. (Подробно о том, как это работает, является темой глав 2 и 3)

```
@ThreadSafe
public class Sequence {
    @GuardedBy("this") private int value;

    public synchronized int getNext() {
        return value++;
    }
}
```

⁴ Фактически, как мы увидим в Главе 3, наихудший случай может быть ещё сложнее, чем обычно показывается на этих диаграммах, из-за возможности переупорядочивания.

⁵ Аннотация `@GuardedBy` описывается в секции 2.4; Она документирует политику синхронизации (*synchronization policy*) для класса `Sequence`.

```
}
```

Листинг 1.2. Потокобезопасный генератор последовательности

В отсутствие синхронизации, компилятору, железу и среде выполнения позволяются существенные вольности в обращении с таймингами и порядком выполнения действий, такие как кэширование переменных в регистрах или в локальных кэшах процессора, где они временно (или даже постоянно) недоступны для других потоков. Эти ухищрения помогают улучшить производительность, и в общем случае желательны, но они возлагают бремя ответственности за чёткое обозначение того, где данные должны совместно использоваться между потоками, на разработчика, который должен сделать так, чтобы эти оптимизации не подрывали безопасность. В главе 16 приводятся подробные сведения о том, как именно JVM гарантирует порядок и какое влияние оказывает синхронизация на эти гарантии, но если вы последуете правилам в главах 2 и 3, то можете свободно опустить эти низкоуровневые подробности.

1.3.2 Угрозы живучести потока

Критически важно обратить внимание на вопросы потокобезопасности, когда разрабатывается параллельный код: безопасность не может быть скомпрометирована. Важность безопасности не является уникальной чертой многопоточных программ – однопоточные программы также должны заботиться о сохранении безопасности и корректности – но использование потоков представляет собой дополнительные вызовы безопасности, не представленные в однопоточных программах. Подобным образом, использование потоков порождает дополнительные формы сбоев живучести (*liveness failure*), которые не происходят в однопоточных программах.

Пока *безопасность* означает “*ничего плохого никогда не случается*”, *живучесть* заботится о дополнительной цели, которая гласит “*что-то хорошее когда-нибудь случается*”. Сбои в работе возникают, когда активность попадает в такое состояние, что постоянно невозможен прогресс. Одной из форм сбоя в работе, который может происходить в последовательной программе, является непреднамеренное создание бесконечного цикла, из-за которого последующий код никогда не получит возможность выполнится. Использование потоков вводит дополнительные риски для живучести. Например, если поток A ожидает ресурс, который поток B удерживает эксклюзивно, и поток B никогда не освободит его, A будет ожидать вечно. Глава 10 описывает различные формы сбоев живучести потока и способы их избежать, включая взаимоблокировки (*deadlock*)⁶, голодание (*starvation*)⁷ и динамическая взаимоблокировка (*livelock*)⁸.

Подобно большинству ошибок параллельного выполнения, ошибки вызывающие сбои живучести потока могут быть неуловимы, потому что зависят от времени происхождения событий (*timing of events*) в различных потоках, и, поэтому, не всегда проявляют себя в разработке или в тестировании.

⁶ **Взаимная блокировка (deadlock**, секция 10.1) описывает ситуацию, когда два или более потока блокируются навсегда, каждый ожидая другого.

⁷ **Голодание (starvation**, секция 10.3.1) описывает ситуацию, когда поток не может получить доступ к совместно используемым ресурсам и не может продвинуться в своём выполнении дальше.

⁸ Поток часто реагирует на события из другого потока. Если действие другого потока тоже является ответом на событие из другого потока, то может произойти **динамическая взаимоблокировка (livelock**, секция 10.3.3). Подробнее в <https://urvanov.ru/2016/05/27/java-8-многопоточность/#liveness>.

1.3.3 Угрозы производительности

С живучестью потока связана производительность. В то время как живучесть означает, что что-то хорошее непременно произойдёт, это в конечном итоге может быть недостаточно хорошо - мы часто хотим, чтобы хорошие вещи происходили быстро. Вопросы производительности относятся к широчайшему спектру проблем, включая плохое время обслуживания, отзывчивость, пропускную способность, потребление ресурсов или масштабируемость. Как и в случае безопасности и живучести, многопоточные программы подвержены всем проблемам производительности однопоточных программ, а также прочим, добавляющимся при использовании потоков.

В хорошо спроектированных параллельных приложениях использование потоков приводит к чистому приросту производительности, тем не менее, потоки приводят к большим накладным расходам во время выполнения. *Переключение контекста (Context switches)* – это ситуация, когда планировщик временно приостанавливает активный поток, так что другой поток может запуститься – что наиболее часто происходит в приложениях с множеством потоков и стоит довольно дорого: сохранение и восстановление контекста выполнения, потеря локальности и время ЦП затраченное на планирование потоков, вместо их запуска. Когда потоки совместно используют данные, они должны использовать механизмы синхронизации, которые могут препятствовать компилятору в проведении оптимизаций, путём сброса кэшей памяти или пометки их устаревшими, и, таким образом, создавая траффик синхронизации в шине разделяемой памяти. Все эти факторы вводят дополнительные расходы производительности; В Главе 11 описывается техника анализа и снижения стоимости.

1.4 Потоки есть везде

Даже если ваша программа никогда явно не создаёт потоки, фреймворки могут создавать потоки от вашего имени, и код, вызываемый из этих потоков, должен быть потокобезопасным. Это оказывает серьёзное влияние на проектирование и реализацию приложения со стороны разработчика, так как разработка потокобезопасных классов требует больше внимания и анализа, чем разработка не потокобезопасных классов.

Все приложения Java используют потоки. Когда JVM запускается, она создаёт потоки для своих внутренних нужд (сборка мусора, завершение) и главных поток для запуска метода `main`. Фреймворки пользовательского интерфейса AWT (*Abstract Window Toolkit*) и Swing создают потоки для управления событиями пользовательского интерфейса. Класс `Timer` создаёт потоки для выполнения отложенных задач. Компонентные фреймворки, такие как сервлеты и RMI, создают пулы потоков и вызывают методы компонентов в этих потоках.

Если вы пользуетесь этими возможностями – как и многие другие разработчики – вы знакомитесь с параллельностью и безопасностью потоков, потому что эти фреймворки создают потоки и вызывают ваши компоненты из них. Хотелось бы верить, что параллельность это некоторая опциональная или продвинутая возможность языка, но реальность такова, что почти все приложения Java многопоточны и эти фреймворки не изолируют вас от необходимости правильной координации доступа к состоянию приложения (*application state*)⁹.

⁹ Имеются в виду члены экземпляров классов, например – поля.

Когда фреймворком в приложение вводится параллельность, обычно невозможно ограничить осведомлённость о параллелизме кодом фреймворка, потому что фреймворк, по своей природе, создаёт обратные вызовы к компонентам приложения, которые, в свою очередь, получают доступ к состоянию приложения. Подобным образом, необходимость в потокобезопасности не заканчивается на компонентах, вызываемых фреймворком – она распространяется на все пути выполнения кода, по которым происходит обращение к состоянию программы, доступные из этих компонентов. Так, потребность в потокобезопасности, становится становиться заразительной.

Фреймворки вводят многопоточность в приложение путём вызова компонентов приложения из потоков. Компоненты неизменно получают доступ к состоянию приложения, тем самым требуя, чтобы все пути выполнения кода, обращающиеся к этому состоянию, были потокобезопасными.

Описанные ниже средства приводят к тому, что код приложения, вызывается из потоков, не управляемых приложением. При необходимости, потокобезопасность может начаться с этих средств, хотя редко ограничивается ими; вместо этого она струится через всё приложение.

Таймер. Класс `Timer` является удобным механизмом для планирования отложенного запуска задач, единовременно или периодически. Введение класса `Timer` может дополнить последовательную программу, потому что класс `TimerTask` выполняется в потоке, управляемом классом `Timer`, а не приложением. Если `TimerTask` получает доступ к данным, к которым уже получен доступ из других потоков приложения, тогда не только `TimerTask` должен делать это в потокобезопасной манере, но также должны поступать и любые другие классы, получающие доступ к данным. Часто, простейшим способ достичь этого, является гарантия того, что объекты, к которым получает доступ класс `TimerTask`, сами по себе потокобезопасны, таким образом, инкапсулируя потокобезопасность внутри совместно-используемых объектов.

Сервлеты и JavaServer Pages (JSPs). Фреймворк сервлетов предназначается для обработки задач всей инфраструктуры развертывания веб приложения и диспетчеризации запросов от удалённых HTTP клиентов. Запрос, поступающий на сервер, отправляется, возможно, через цепочку фильтров, соответствующей странице JSP или сервлету. Каждый сервле́т представляет собой компонент логики приложения, и для крупных вебсайтов является нормой, что множество клиентов могут обращаться к одному и тому же сервлету одновременно. Спецификация сервлета требует, чтобы сервле́т был готов к тому, что может быть вызван одновременно во множестве потоков. Другими словами, сервле́т должен быть покобезопасен.

Даже если вы сможете гарантировать, что сервле́т будет вызываться одновременно только из одного потока, вам всё равно придётся обратить внимание на безопасность потоков при создании веб приложения. Сервлеты часто получают доступ к информации о состоянии, разделяемой (*shared*) с другими сервлетами, такой, как объекты с областью видимостью уровня приложения¹⁰ (те, что хранятся

¹⁰ Область видимости **application-scoped**

в `ServletContext`), или объекты с областью видимостью уровня сессии¹¹ (те, что хранятся для каждого клиента в `HttpSession`). Когда сервлет получает доступ к объектам, разделяемым между сервлетами и запросами, он обязаненным образом координировать доступ, чтобы множество запросов могло обращаться к ним одновременно из разных потоков. Сервлеты и JSP, так же как фильтры сервлетов и объекты, хранимые в контейнерах с областью видимости, подобные `ServletContext` и `HttpSession`, просто должны быть потокобезопасны.

Удалённый вызов методов (RMI) RMI позволяет вам вызывать методы объектов, запущенных в других JVM. Когда вы вызываете удалённый метод с RMI, аргументы метода упаковываются (маршализуются) в поток байтов и отправляются по сети к удалённой (*remote*) JVM, где они распаковываются (демаршализуются) и передаются удалённому (*remote*) методу.

Когда код RMI вызывает удалённый объект, в каком потоке происходит этот вызов? Вы не знаете, но определённо не в том потоке, который создали вы – ваш объект будет вызван в потоке, управляемом RMI. Как много потоков может создать RMI? Может ли быть вызван некоторый удалённый метод в некотором удалённом объекте одновременно во множестве потоков RMI?¹²

Удаленный объект должен иметь защиту от двух угроз безопасности потоку: должным образом, координированный доступ к состоянию, которое может совместно использоваться с другими объектами, и должным образом скоординированный доступ к состоянию самого удаленного объекта (так как тот же объект может быть вызван в нескольких потоках одновременно). Подобно сервлетам, объекты RMI должны быть готовы к множеству одновременных вызовов и должны обеспечивать собственную потокобезопасность.

Swing и AWT. Приложения GUI по своей природе асинхронны. Пользователи могут выбрать элемент меню или нажать на кнопку в любое время, и они ожидают, что приложение будет реагировать быстро, даже если оно находится в процессе выполнения чего-нибудь другого. Swing и AWT для решения этой проблемы создают отдельный поток, который занимается обработкой инициированных пользователем событий и обновлением графического представления, отображаемого пользователю.

Компоненты Swing, такие как `JTable` не потокобезопасны. Вместо этого, программы на Swing достигают собственной потокобезопасности, ограничивая все обращения к компонентам GUI только потоком событий. Если приложение хочет манипулировать GUI извне потока событий, оно должно вызывать код¹³, который будет управлять GUI, для запуска в потоке событий.

Когда пользователь выполняет UI¹⁴ действие, в потоке событий вызывается обработчик событий для выполнения любой операции, запрошенной пользователем. Если обработчику требуется доступ к состоянию приложения,

доступ к которому также происходит из других потоков (например, редактируемому документу), тогда обработчик событий, наряду с другим кодом, обращающимся к этому состоянию, должен выполнять его потокобезопасным способом.

¹¹ Область видимости **session-scoped**

¹² Ответ: да, но это не все, что понятно из Javadoc — для большей информации прочитайте спецификацию RMI.

¹³ Можно из кода отправить компоненту сообщение о событии.

¹⁴ UI – user interface/пользовательский интерфейс

Глава 2 Потокобезопасность

Возможно, для вас будет сюрпризом, что параллельное программирование — это не столько о потоках или блокировках, скорее, если провести аналогию, не более чем гражданское строительство, о заклёпках или двутавровых балках. Конечно, строительство мостов, которые затем не рухнут вниз, требует правильного использования множества заклёпок и двутавровых балок, также, как и параллельные программы требуют корректного использования потоков и блокировок. Но всё это просто механизмы — средства для достижения цели. Написание потокобезопасного кода является, в своей основе, только управлением доступом к состоянию, и, в частности, к *разделяемому (shared), изменяемому состоянию (mutable state)*.

Неофициально, состояние объекта — это его данные, хранимые в переменных состояния, таких как экземпляры или статические поля. Состояние объекта может включать поля от различных, зависимых объектов; Состояние объекта `HashMap` частично хранится в самом объекте `HashMap`, но также во множестве объектов `Map.Entry`. Состояние объекта включает в себя любые данные, которые могут оказывать влияние на его внешнее видимое поведение.

Под разделяемым состоянием мы понимаем то, что переменная может быть доступна из множества потоков; под изменяемым состоянием мы понимаем то, что значение переменной может изменяться в течение всей её жизни. Мы можем говорить о потокобезопасности, как коде, но на самом деле, мы пытаемся защитить *данные от неконтролируемого параллельного доступа*.

Необходимость объекта быть потокобезопасным зависит от того, будет ли он доступен из нескольких потоков. Это свойство того, как объект *используется* в программе, а не того, что он выполняет. Для создания потокобезопасного объекта требуется использовать синхронизацию, обеспечивающую координацию доступа к его изменяемому состоянию; невыполнение этого требования может привести к повреждению данных и другим нежелательным последствиям.

Всякий раз, когда более чем один поток обращается к данной переменной состояния, и один из них может записать что-то в неё, все они должны координировать свой доступ к ней с помощью синхронизации. Главным механизмом синхронизации в Java является ключевое слово `synchronized`, которое обеспечивает эксклюзивную блокировку, но термин “синхронизация” также включает в себя использование переменных `volatile`, явных блокировок и атомарных переменных.

Следует избегать соблазна думать, что существуют “особые” ситуации, в которых это правило не применяется. Программа, которая опускает необходимую синхронизацию, может казаться работающей, проходить свои тесты и хорошо работать в течение многих лет, но она все еще неисправна и может потерпеть крах в любой момент.

Если несколько потоков могут получить доступ к изменяемому состоянию переменной без соответствующей синхронизации, ваша программа имеет заложенные ошибки. Существует три пути для её исправления:

- Не разделяйте переменные между потоками;
- Сделайте переменные состояния неизменяемыми;

- Используйте синхронизацию, когда бы ни обращались к состоянию переменной.

Если вы не задумывались при проектировании класса о параллельном доступе к нему, некоторые из этих подходов могут потребовать внесения значительных изменений в его дизайн, так что исправление проблемы может быть не таким тривиальным, как озвучивание этого совета. *Гораздо проще сразу проектировать класс, который будет потокобезопасным, чем позже его модифицировать для обеспечения потокобезопасности.*

В большой программе довольно сложно выяснить, могут ли несколько потоков получить доступ к данной переменной. К счастью, те же объектно-ориентированные методы, которые помогают создавать хорошо организованные, поддерживаемые классы, — такие как инкапсуляция и скрытие данных — также могут помочь в создании потокобезопасных классов. Чем меньше кода имеет доступ к определенной переменной, тем проще обеспечить, чтобы все это использовало правильную синхронизацию, и тем легче рассуждать об условиях, при которых к данной переменной можно получить доступ. Язык Java не принуждает вас к инкапсуляции состояния — вполне допустимо хранить состояние в публичных полях (даже в статических публичных полях) или публиковать ссылку на другой внутренний объект — но, чем лучше инкапсулировано состояние программы, тем проще сделать ее потокобезопасной и помочь сопровождающим¹⁵ впоследствии сохранить ее такой.

При проектировании потокобезопасных классов хорошие объектно-ориентированные методы, — инкапсуляция, неизменяемость и четкая спецификация инвариантов — ваши лучшие друзья.

Настанут времена, когда хорошие объектно-ориентированные методы проектирования станут противоречить реальным требованиям; в этих случаях может потребоваться поступиться правилами хорошего дизайна ради производительности или ради совместимости с устаревшим кодом (*legacy code*). Иногда абстракция и инкапсуляция расходятся с производительностью, хотя и не так часто, как полагают многие разработчики — но всегда рекомендуется сначала сделать код правильным, а затем уже заняться его ускорением. Даже в этом случае, оптимизацию следует проводить только тогда, когда это продиктовано требованиями и показателями производительности, и только в том случае, если измерения, проведённые в реалистичных условиях¹⁶, говорят вам, что ваша оптимизация действительно изменила ситуацию.

Если вы решите, что просто должны разрушить инкапсуляцию, ещё не всё потеряно. Вашу программу всё еще возможно сделать потокобезопасной, это просто будет намного сложнее. Кроме того, потокобезопасность вашей программы станет более хрупкой, увеличится не только стоимость и риски разработки, но также возрастут расходы и риски на сопровождение. В главе 4 характеризуется

¹⁵ Давно известно, что львиную долю жизненного цикла программы, занимает сопровождение — исправление ошибок, различные доработки, добавление «бантиков и рюшечек».

¹⁶ В параллельном коде этой практики следует придерживаться даже больше, чем обычно. Поскольку ошибки параллелизма очень сложно воспроизвести и отладить, преимущество получения небольшого прироста производительности в некоторых, редко используемых, ветвях кода, вполне может быть карликовым, из-за риска того, что программа потерпит крах в этой области.

условия, при которых можно безопасно ослабить инкапсуляцию переменных состояния.

До сих пор мы использовали термины “потокобезопасный класс” и “потокобезопасная программа” почти взаимозаменяемо. Является ли потокобезопасной программа, полностью построенная из потокобезопасных классов? Не обязательно - программа, полностью состоящая из потокобезопасных классов, может быть не потокобезопасной, а потокобезопасная программа может содержать классы, не являющиеся потокобезопасными. Вопросы, связанные с составом потокобезопасных классов, также рассматриваются в главе 4. В любом случае концепция потокобезопасного класса имеет смысл только в том случае, если класс инкапсулирует свое собственное состояние. “Потокобезопасность” может быть термином, применяемым к коду, но речь идет о состоянии, поэтому термин может быть применен только ко всему телу кода, инкапсулирующему состояние, которое может быть объектом или всей программой.

2.1 Что такое потокобезопасность?

Определение потокобезопасности удивительно запутанно. Более формальные попытки настолько сложны, что предлагают небольшое практическое руководство или требуют интуитивного понимания, другие же, неофициальные описания, могут показаться “вещью в себе”. Быстрый поиск в Google возвращает множество “определений”, похожих на эти:

... могут быть вызваны из нескольких потоков программы, без нежелательных взаимодействий между потоками.

... могут вызываться более чем одним потоком за раз, не требуя каких-либо других действий со стороны вызывающих.

Учитывая такие определения, неудивительно, что мы находим потокобезопасность запутанной! Они звучат подозрительно похоже на “класс потокобезопасен, если он может быть безопасно использован из множества потоков”. Фактически, вы не можете оспорить это утверждение, но оно также почти не предлагает практической помощи. Как отличить потокобезопасный класс от непотокобезопасного? Что мы вообще подразумеваем под термином “безопасный”?

Ядром любого разумного объяснения потокобезопасности является концепция *корректности*. Если наше определение потокобезопасности нечеткое, это только потому, что нам не хватает определения термина “*корректность*”.

Корректность означает, что класс *соответствует своей спецификации*. Хорошая Спецификация определяет инварианты, ограничивающие состояние объекта и постусловия, описывающие эффекты от его операций. Поскольку мы часто не пишем адекватных спецификаций для наших классов, как мы можем узнать, что они корректны? Мы не можем, но это в любом случае не мешает нам их использовать, так как однажды мы убедили самих себя, что “код работает”. Это “доверие коду” относительно близко к тому, как многие из нас добиваются корректности, так позволим себе предположить, что корректность однопоточных программ, это что-то вроде “мы это знаем, когда мы видим это”. Примем оптимистичное определение “*корректности*” как чего-то, что можно распознать, теперь мы можем определить потокобезопасность несколько менее запутанным способом: класс потокобезопасен, когда он продолжает вести себя корректно в моменты одновременного обращения из нескольких потоков.

Класс является *потокобезопасным*, если он ведет себя правильно при доступе из нескольких потоков, независимо от планирования или чередования выполнения этих потоков средой выполнения, и без дополнительной синхронизации или другой координации со стороны вызывающего кода.

Поскольку любая однопоточная программа также является допустимой (*valid*) многопоточной программой, она не может быть потокобезопасной, если не является корректной в однопоточной среде¹⁷. Если объект реализован корректно, никакая последовательность операций - вызовов публичных (*public*) методов и операций чтения или записи открытых полей - не должна нарушать ни один из его инвариантов или постусловий. *Никакой набор операций, выполняемых последовательно или параллельно на экземплярах потокобезопасного класса, не может привести к тому, что экземпляр окажется в недопустимом состоянии (invalid).*

Потокобезопасные классы инкапсулируют любые необходимые синхронизации, так что клиенты не должны заботиться об этом.

2.1.1 Пример: сервлет без сохранения состояния (stateless)

В главе 1 мы перечислили ряд фреймворков, которые создают потоки и вызывают ваши компоненты из этих потоков, оставляя вас ответственными за то, чтобы сделать ваши компоненты потокобезопасными. Очень часто, требования к потокобезопасности основаны не на необходимости напрямую использовать потоки, а вследствие решения использовать средства, подобные фреймворку сервлетов. Мы собираемся разработать простой пример — основанную на сервлете службу факторизации — и медленно расширять его, добавляя функции, пока сохраняется его потокобезопасность.

В листинге 2.1 показан наш простой сервлет факторизации. Он распаковывает факторизуемое число из запроса сервлета, факторизует его, и упаковывает результаты в ответ сервлета

```
@ThreadSafe
public class StatelessFactorizer implements Servlet {
    public void service(ServletRequest req, ServletResponse resp) {
        BigInteger i = extractFromRequest(req);
        BigInteger[] factors = factor(i);
        encodeIntoResponse(resp, factors);
    }
}
```

Листинг 2.1 Сервлет без сохранения состояния

Класс `StatelessFactorizer`, подобно большинству сервлетов, не имеет состояния (*stateless*): он не имеет полей и не ссылается на поля в других классах. Переходное состояние для конкретного вычисления существует только в локальных

¹⁷ Если вам не нравится свободное использование термина “корректность”, вы можете предпочесть думать о потокобезопасном классе, как о том, который в многопоточной среде подвержен ошибкам не более, чем в однопоточной.

переменных, которые хранятся в стеке потока и доступны только исполняющему потоку. Один поток, обращающийся к StatelessFactorizer, не может оказать влияние на результат другого потока, обращающегося к тому же StatelessFactorizer; поскольку два потока не разделяют состояния, это выглядит, как если бы они обращались к различным экземплярам. Поскольку действия потока, обращающегося к объекту без состояния, не могут повлиять на корректность операций в других потоках, объекты без состояния являются потокобезопасными.

Объекты без состояния всегда потокобезопасны.

Тот факт, что большинство сервлетов может быть реализовано без какого-либо состояния, позволяет не учитывать требование потокобезопасности, что значительно снижает нагрузку по их созданию. Только когда сервлеты хотят переписывать значения от одного запроса к другому, только тогда требование потокобезопасности становится проблемой.

2.2 Атомарность

Что происходит, когда мы добавляем один элемент состояния к тому, что было объектом без состояния? Предположим, мы хотим добавить "счетчик посещений", который измеряет количество обработанных запросов. Очевидным подходом является добавление поля типа long к сервлету и увеличение его значения при каждом запросе, как показано в классе UnsafeCountingFactorizer в листинге 2.2.

```
@NotThreadSafe
public class UnsafeCountingFactorizer implements Servlet {
    private long count = 0;

    public long getCount() { return count; }

    public void service(ServletRequest req, ServletResponse resp) {
        BigInteger i = extractFromRequest(req);
        BigInteger[] factors = factor(i);
        ++count;
        encodeIntoResponse(resp, factors);
    }
}
```



Листинг 2.2 Сервлет подсчитывающий запросы без необходимой синхронизации. Не делайте так.

К сожалению, класс UnsafeCountingFactorizer не является потокобезопасным, хотя он будет прекрасно работать в однопоточной среде. Так же, как класс UnsafeSequence на странице 22, он чувствителен к *потерянным обновлениям*. Оператор инкремента, `++count`, может выглядеть как одно действие из-за его компактного синтаксиса, но он не является атомарным, что означает, что он не выполняется как одна, неделимая операция. Вместо этого он является сокращением для последовательности из трех дискретных операций: извлечения текущего значения, добавление к нему единицы и запись нового значение обратно. Это пример операции **read-modify-write**, в которой результирующее состояние является производным от предыдущего состояния.

Нас рисунке 1.1 на стр. 22 показано, что может произойти, если два потока попытаются увеличить счетчик одновременно и без синхронизации. Если счетчик первоначально имел значение 9, в некоторый *неудачный* времени каждый из потоков мог прочитать значение, увидеть, что оно равно 9, добавить к нему единичку, и установить счетчику значение 10. Очевидно, произошло не то, что предполагалось; по ходу дела, приращение было потеряно, и счетчик посещений теперь постоянно меньше на единицу.

Вы можете подумать, что наличие немного неточного количества обращений к веб-службе является приемлемой потерей точности, и иногда допустимо. Но, если счетчик используется для создания последовательностей (sequences) или уникальных идентификаторов объектов, возврат одного и того же значения из нескольких вызовов может вызвать серьезные проблемы с целостностью данных¹⁸. Ситуация, при которой существует возможность получения некорректных результатов в неудачный момент времени, настолько важна в параллельном программировании, что получила название: *состояние гонки (race condition)*.

2.2.1 Условия гонок

Класс `UnsafeCountingFactorizer` имеет несколько условий гонок, которые делают его результаты ненадежными. Условие гонки возникает, когда правильность вычисления зависит от относительного момента времени или от чередования нескольких потоков во время выполнения; другими словами, получение правильного ответа зависит от удачного момента времени¹⁹. Наиболее распространенным типом состояния гонки является `check-then-act`, где потенциально устаревшее наблюдение используется для принятия решения о том, что делать дальше.

Мы часто сталкиваемся с условиями гонок в реальной жизни. Допустим, вы планируете встречу с другом в полдень, в Starbucks на Университетской авеню. Но когда вы доберетесь туда, вы поймете, что на Университетской авеню расположены два Starbucks, и вы не уверены, в котором из них договорились встретиться. В 12:10 вы не видите своего друга в Starbucks A, поэтому идёте в Starbucks B, чтобы посмотреть, не пришёл ли он туда, но его там тоже нет. Есть несколько возможностей: ваш друг опаздывает, и он не в Starbucks; или ваш друг прибыл в Starbucks после того, как вы покинули его; или ваш друг был в Starbucks B, но пошел искать вас, и теперь находится на пути к Starbucks A. Предположим худшее и скажем, что это была последняя возможность. Сейчас 12:15, вы оба были в Starbucks, и вы оба задаетесь вопросом, встали ли вы. Чем вы сейчас занимаетесь? Вернитесь к другому Starbucks? Сколько раз вы собираетесь ходить туда и обратно? Если Вы не договорились о протоколе, вы оба можете провести день, ходя вверх и вниз по Университетской Авеню, разочарованные и унылые.

¹⁸ Подход, принятый в классах `UnsafeSequence` и `UnsafeCountingFactorizer`, имеет другие серьезные проблемы, включая возможность устаревших данных (раздел 3.1.1).

¹⁹ Термин "**состояние гонки (race condition)**" часто путают со связанным термином "**гонка данных (data race)**", который возникает, когда синхронизация не используется для координации всего доступа к разделяемому не финальному полю (*shared nonfinal field*). Вы рискуете **получить гонку** данных всякий раз, когда поток пишет значение в переменную, которая затем может быть прочитана другим потоком, или читает переменную, которая могла быть в последний раз записана другим потоком, если оба потока не используют синхронизацию; код с **гонками данных** не имеет никакой определенной полезной семантики в модели памяти Java. Не все **условия гонки** - это **гонки данных**, и не все **гонки данных** - это **условия гонки**, но оба условия могут привести к непредсказуемому сбою параллельных программ. Класс `UnsafeCountingFactorizer` имеет как **условия гонки**, так и **гонки данных**. Подробнее о гонках данных см. главу 16

Проблема с подходом “я просто выйду на улицу и посмотрю, находится ли он на другом её конце” заключается в том, что пока вы идете по улице, ваш друг, возможно, уже ушёл. Вы смотрите вокруг Starbucks A, видите, что “его здесь нет” и идете искать его. И Вы можете сделать то же самое для Starbucks B, но *в другой момент времени*. Чтобы пройти по улице, вам потребуется несколько минут, и в течение этих нескольких минут, *состояние системы могло измениться*.

Пример Starbucks иллюстрирует состояние гонки, потому что достижение желаемого результата (встреча с вашим другом) зависит от относительного времени событий (когда каждый из вас прибывает в тот Starbucks или другой, как долго вы ждете там перед переключением и т. д.). Наблюдение, что он не в Starbucks A, становится потенциально недействительным, как только вы выходите из входной двери; он мог бы войти через заднюю дверь, и Вы об этом не узнаете. Именно это аннулирование (*invalidation*) наблюдений характеризует большинство условий гонки - **использование потенциально устаревшего** наблюдения для принятия решения или выполнения вычислений. Этот тип условия гонок называется ***check-then-act***: вы наблюдаете что-то, что является истинным (файл X не существует), а затем принимаете решение на основе этого наблюдения (создание X); но, фактически, это наблюдение может стать недействительным в момент времени между тем, когда вы наблюдали *это* и моментом времени, когда вы воздействовали на *это* (в промежутке между этими моментами времени, кто-то другой создал файл X), что вызовет проблемы (неожиданные исключения (*unexpected exception*), перезапись данных, повреждения файлов).

2.2.2 Пример: состояние гонки в отложенной инициализации

Распространенной идиомой, использующей *check-then-act*, является **ленивая инициализация**. Целью ленивой инициализации является отсрочка инициализации объекта до тех пор, пока он не будет фактически необходим, в то же время, гарантируя, что он будет инициализирован только один раз. Класс LazyInitRace в листинге 2.3 иллюстрирует идиому ленивой инициализации. Метод `getInstance` сначала проверяет, был ли объект `ExpensiveObject` уже инициализирован, и в случае положительного ответа, он возвращает существующий экземпляр; в противном случае, он создает новый экземпляр и возвращает его после сохранения ссылки на него, чтобы последующие вызовы могли сэкономить на формировании объекта.

```
@NotThreadSafe
public class LazyInitRace {
    private ExpensiveObject instance = null;

    public ExpensiveObject getInstance() {
        if (instance == null)
            instance = new ExpensiveObject();
        return instance;
    }
}
```



Листинг 2.3 Условие гонок при ленивой инициализации. Не делайте так.

Класс `LazyInitRace` имеет условия гонки, которые могут подорвать его корректность. Скажем, что потоки *A* и *B* одновременно выполняют метод

`getInstance`. A видит, что экземпляр имеет значение `null` и создает новый объект `ExpensiveObject`. B также проверяет значение экземпляра на `null`. Имеет ли экземпляр значение `null` в этой точке, непредсказуемо зависит от момента времени, включая капризы планировщика и того времени, которое требуется A для создания экземпляра `ExpensiveObject` и установки значения поля `instance`. Если экземпляр имеет значение `null`, когда B проверяет его, оба потока, вызывающих метод `getInstance`, могут получать два разных результата, несмотря на то, что метод `getInstance` всегда должен возвращать один и тот же экземпляр.

Операция подсчета очков в `UnsafeCountingFactorizer` имеет состояние гонки другого вида. Операция *read-modify-write*, такая как приращение счетчика, определяет преобразование состояния объекта с точки зрения его предыдущего состояния. Чтобы увеличить счетчик, вы должны знать его предыдущее значение и убедиться, что никто больше не изменяет или не использует это значение, пока вы находитесь в процессе его обновления.

Подобно большинству ошибок параллелизма, условия гонки не всегда приводят к сбою: также требуется некоторый неудачный момент времени. Но условия гонки могут вызвать серьезные проблемы. Если `LazyInitRace` используется для создания экземпляра реестра, применяемого в масштабах всего приложения, то возврат различных экземпляров из нескольких вызовов может привести к потере регистраций или к тому, что несколько действий будут иметь несогласованные представления наборов зарегистрированных объектов. Если класс `UnsafeSequence` используется для генерации идентификаторов сущностей в фреймворке хранения данных (*persistence framework*), два различных объекта могут в конечном итоге получить один и тот же идентификатор, при этом вызвав нарушение ограничения целостности идентификаторов²⁰.

2.2.3 Составные действия

Как класс `LazyInitRace`, так и класс `UnsafeCountingFactorizer` содержат в себе последовательность операций, которые должны быть атомарными или неделимыми по отношению к другим операциям в том же состоянии. Чтобы избежать условий гонки, должен быть способ предотвратить использование переменной другими потоками, пока мы находимся в процессе ее изменения, так мы сможем гарантировать, что другие потоки смогут читать или изменять состояние только до начала или после завершения изменения, но не в процессе.

Операции A и B являются атомарными по отношению друг к другу, если с позиции потока, выполняющего операцию A, когда другой поток выполняет операцию B, либо вся операция B выполняется, либо никакие действия из неё. Атомарная операция - это операция единая по отношению ко всем операциям, работающим в том же состоянии, включая саму себя²¹.

Если бы операция инкремента в классе `UnsafeSequence` была атомарной, условие гонки, показанное на рис. 1.1, не могло бы произойти, и каждое выполнение операции инкремента всегда бы увеличивало значение счетчика ровно на единицу. Для обеспечения потокобезопасности, операции *check-then-act*

²⁰ Как правило, каждая запись в хранилище имеет уникальный идентификатор, и нарушение этого ограничения приводит к ошибке “unique key violation...” или аналогичной.

²¹ Операция не сможет оказать воздействие на саму себя, либо она будет выполнена, либо нет.

(схожая с отложенной инициализацией) и ***read-modify-write*** (схожая с операцией инкремента) всегда должны быть атомарными. Мы обобщённо относим последовательности действий ***check-then-act*** и ***read-modify-write*** к типу «составных» (*compound actions*). Составные действия - это последовательности операций, которые должны быть выполнены атомарно, чтобы оставаться потокобезопасными. В следующем разделе мы рассмотрим блокировку (*locking*) - встроенный механизм Java для обеспечения атомарности. На данный момент мы исправим проблему другим способом, используя существующий потокобезопасный класс, как показано в классе `CountingFactorizer`, в листинге 2.4.

```
@ThreadSafe
public class CountingFactorizer implements Servlet {
    private final AtomicLong count = new AtomicLong(0);

    public long getCount() { return count.get(); }

    public void service(ServletRequest req, ServletResponse resp) {
        BigInteger i = extractFromRequest(req);
        BigInteger[] factors = factor(i);
        count.incrementAndGet();
        encodeIntoResponse(resp, factors);
    }
}
```

Листинг 2.4 Сервлет, подсчитывающий количество запросов, используя класс `AtomicLong`

Пакет `java.util.concurrent.atomic` содержит классы *атомарных переменных* (*atomic variable*), позволяющие атомарно выполнять переходы состояний в операциях с числами и объектными ссылками. Заменяя счетчик типа `long` на счётчик типа `AtomicLong`, мы гарантируем, что все действия (*actions*), которые обращаются к состоянию счетчика, являются атомарными²². Поскольку состояние сервлета определяется состоянием счетчика, а счетчик потокобезопасен, наш сервлет также потокобезопасен.

Мы могли добавить счетчик запросов к сервлету факторинга и поддерживать потокобезопасность, используя существующий потокобезопасный класс для управления состоянием счетчика - класс `AtomicLong`. При добавлении единственного элемента состояния в класс без сохранения состояния (*stateless*), результирующий класс будет потокобезопасным, если состояние полностью управляется потокобезопасным объектом. Но, как мы увидим в следующем разделе, переход от одной переменной состояния к нескольким, не обязательно также прост, как переход от нуля к единице.

Там, где это целесообразно, используйте для управления состоянием вашего класса существующие потокобезопасные объекты, такие как `AtomicLong`.

Гораздо проще рассуждать о возможных состояниях и переходах состояний для существующих потокобезопасных объектов, чем для произвольных

²² Класс `CountingFactorizer` вызывает метод `incrementAndGet` для увеличения значения счётчика, который вернёт инкрементированное значение; в этом случае результат игнорируется.

переменных состояния, и это упрощает сопровождение и проверку безопасности потоков.

2.3 Блокировка

Мы могли добавить одну переменную состояния к своему сервлету, для сохранения потокобезопасности используя потокобезопасный объект, управляющий всем состоянием сервлета. Но если мы хотим добавить к своему сервлету больше состояний, можем ли мы просто добавить потокобезопасных переменных состояния?

Представьте, что мы хотим улучшить производительность нашего сервлета, кэшируя последний вычисленный результат, на случай, если два последовательных клиента запросят факторизацию одного и того же числа. (Это не похоже на эффективную стратегию кэширования, мы предлагаем вариант лучше в разделе 5.6) Для реализации этой стратегии необходимо запомнить две вещи: последнее факторизуемое число и его факторы.

Мы использовали класс `AtomicLong` для управления состоянием счетчика потокобезопасным способом; можем ли мы использовать его двоюродного брата, класс `AtomicReference`²³, для управления последним факторизуемым числом и его факторами? Пример попытки показан в классе `UnsafeCachingFactorizer` в листинге 2.5.

К сожалению, такой подход не работает. Несмотря на то, что атомизированные ссылки индивидуально потокобезопасны, класс `UnsafeCachingFactorizer` имеет условия гонки, которые могут привести к неправильному ответу.

```
@NotThreadSafe
public class UnsafeCachingFactorizer implements Servlet {
    private final AtomicReference<BigInteger> lastNumber
        = new AtomicReference<BigInteger>();
    private final AtomicReference<BigInteger[]> lastFactors
        = new AtomicReference<BigInteger[]>();

    public void service(ServletRequest req, ServletResponse resp) {
        BigInteger i = extractFromRequest(req);
        if (i.equals(lastNumber.get()))
            encodeIntoResponse(resp, lastFactors.get());
        else {
            BigInteger[] factors = factor(i);
            lastNumber.set(i);
            lastFactors.set(factors);
            encodeIntoResponse(resp, factors);
        }
    }
}
```



Листинг 2.5 Сервлет, пытается кэшировать свой последний результат без адекватного использования атомарности. Не делайте так.

²³ Так же, как `AtomicLong` является потокобезопасным классом-держателем (holder class) для длинного целого числа, `AtomicReference` является потокобезопасным классом-держателем для ссылки на объект. Атомарные переменные и их преимущества описаны в главе 15.

Определение потокобезопасности требует сохранения инвариантов независимо от моментов времени или чередования операций в нескольких потоках. Одним из инвариантов класса `safeCachingFactorizer` является то, что произведение факторов, кэшированных в переменной `lastFactors`, равно значению, кэшированному в `lastNumber`; наш сервлет корректен, только если этот инвариант всегда выполняется. Когда несколько переменных участвуют в инварианте, они не являются независимыми: значение одной ограничивает допустимое значение(я) других. Таким образом, при обновлении значения одной переменной, вы должны обновить значения других переменных в *той же атомарной операции*.

В некоторые неудачные моменты времени, класс `UnsafeCachingFactorizer` может нарушать эти инварианты. Используя атомарные ссылки, мы не можем обновлять как `lastNumber`, так и `lastFactors` одновременно, даже если каждый вызов `set` атомарен; все еще есть окно уязвимости, когда один был изменен, а другой нет, и в течение этого времени другие потоки могли видеть, что инвариант не удерживается. Точно так же, два значения не могут быть выбраны одновременно: между моментом времени, когда поток *A* получает два значения, поток *B* мог бы изменить их, и снова *A* может заметить, что инвариант не удерживается.

Чтобы сохранить состояние согласованным, обновляйте связанные переменные состояний в одной атомарной операции.

2.3.1 Внутренние блокировки

Java предоставляет встроенный механизм блокировки для обеспечения атомарности: блок `synchronized`. (Также существует еще один важный аспект блокировки и других механизмов синхронизации – видимость (*visibility*) - который рассматривается в главе 3). Блок `synchronized` состоит из двух частей: ссылки на объект, который будет служить блокировкой, и блока кода, который должен быть защищен этой блокировкой. Метод `synchronized` является сокращением для блока `synchronized`, охватывающего всё тело метода, и чья блокировка является объектом, в котором вызывается метод. (Статический метод `synchronized` использует объект `Class` для блокировки).

```
synchronized (lock) {  
    // Access or modify shared state guarded by lock  
}
```

Каждый объект может неявно выступать в качестве блокировки в целях синхронизации; эти встроенные блокировки называются *внутренними блокировками* (*intrinsic locks*) или *мониторами* (*monitor locks*). Блокировка автоматически приобретается (*acquired*) исполняющим потоком перед входом в блок `synchronized` и автоматически освобождается, когда поток выходит из блока `synchronized`, как при нормальном пути выполнения, так и при возбуждении исключения в блоке. Единственный способ получить встроенную блокировку – войти в блок `synchronized` или в метод, защищенный этой блокировкой.

Встроенные блокировки Java выступают в роли мьютексов (*mutex*²⁴), это означает, что блокировка может принадлежать не более чем одному потоку. Когда

²⁴ Mutual exclusion locks – взаимоисключающая блокировка.

поток *A* пытается приобрести блокировку удерживаемую потоком *B*, поток *A* должен подождать, или заблокироваться, до тех пор, пока *B* не освободит её. Если поток *B* никогда не освободит блокировку, поток *A* будет ждать вечно.

Поскольку, за раз, только один поток может выполнить блок кода, защищенный данной блокировкой, блоки `synchronized`, защищенные одной и той же блокировкой, выполняются атомарно относительно друг друга. В контексте параллелизма термин “*атомарность*” означает то же самое, что и в транзакционных приложениях - это группа операторов, выполняемая как единая, неделимая единица. Ни один поток, выполняющий блок `synchronized`, не может наблюдать за другим потоком, находящимся в процессе выполнения блока `synchronized`, защищенного одной и той же блокировкой.

Механизм синхронизации позволяет легко восстановить потокобезопасность в сервлете факторинга. В листинге 2.6 выполняется синхронизация²⁵ метода `service`, вследствие чего, в метод `service`, в конкретный момент времени, может входить только один поток. Класс `SynchronizedFactorizer` сейчас потокобезопасен; однако, этот подход является довольно экстремальным, так как препятствует возможности использовать сервlet факторинга нескольким клиентам одновременно - это приводит к недопустимо низкой отзывчивости. Эта проблема - проблема производительности, а не проблемы безопасности потоков - рассматривается в разделе 2.5.

```
@ThreadSafe
public class SynchronizedFactorizer implements Servlet {
    @GuardedBy("this")
    private BigInteger lastNumber;
    @GuardedBy("this")
    private BigInteger[] lastFactors;

    public synchronized void service(ServletRequest req,
                                    ServletResponse resp) {
        BigInteger i = extractFromRequest(req);
        if (i.equals(lastNumber))
            encodeIntoResponse(resp, lastFactors);
        else {
            BigInteger[] factors = factor(i);
            lastNumber = i;
            lastFactors = factors;
            encodeIntoResponse(resp, factors);
        }
    }
}
```



Листинг 2.6 Сервlet, кэширующий последний результат запроса, имеет неприемлемо низкий уровень параллелизма. Не делайте так.

2.3.2 Повторная входимость

Когда поток запрашивает блокировку, которая уже удерживается другим потоком, запрашивающий поток блокируется. Но, поскольку внутренние блокировки

²⁵ Добавление оператора `synchronized`.

являются реентерабельными (*reentrant*), если поток пытается получить блокировку, которую ранее он уже захватил, запрос выполняется успешно. Реентерабельность (*reentrancy*) означает, что блокировки приобретаются для каждого потока, а не для каждого вызова²⁶. Реентерабельность реализуется путем связывания с каждой блокировкой счетчика полученных захватов и владеющего потока. Если число захватов равно нулю, блокировка считается отпущенной. Когда поток получает ранее отпущенную блокировку, JVM записывает владельца и устанавливает счетчику захватов значение единицы. Если тот же самый поток снова получает блокировку, счетчик увеличивается, и когда владеющий блокировкой поток выходит из блока `synchronized`, счетчик уменьшается. Когда значение счетчика достигает нуля, блокировка снимается.

Реентерабельность способствует инкапсуляции поведения блокировок и, таким образом, упрощает разработку объектно-ориентированного параллельного кода. Без реентерабельных блокировок (*reentrant locks*), очень естественно смотрящийся код в листинге 2.7, в котором подкласс переопределяет `synchronized` метод, а затем вызывает метод суперкласса, вызовет взаимоблокировку (*deadlock*). Поскольку методы `doSomething` в классах `Widget` и `LoggingWidget` синхронизированы (`synchronized`), каждый пытается получить блокировку класса `Widget` перед продолжением. Но если бы встроенные блокировки не были реентерабельными, вызов метода `super.doSomething` никогда не смог бы получить блокировку, потому что блокировка считалась бы уже захваченной, и поток был бы вынужден постоянно ожидать блокировку, которую никогда не смог захватить. Реентерабельность спасает нас от взаимоблокировок в ситуациях, подобных этой.

```
public class Widget {  
    public synchronized void doSomething() {  
        ...  
    }  
}  
  
public class LoggingWidget extends Widget {  
    public synchronized void doSomething() {  
        System.out.println(toString() + ": calling doSomething");  
        super.doSomething();  
    }  
}
```

Листинг 2.7. Код, который без реентерабельности будет вызывать взаимоблокировку.

2.4 Защита состояния с помощью блокировок

Поскольку блокировки обеспечивают сериализованный²⁷ доступ (*serialized access*) к веткам кода, которые они защищают, их можно использовать для создания протоколов, гарантирующих монопольный доступ к разделяемому состоянию. Следование этим протоколам может гарантировать согласованность состояний.

²⁶ Это отличается от поведения блокировок по умолчанию для мьюнексов pthreads (потоки POSIX), при котором блокировки предоставляются при каждом вызове.

²⁷ Сериализация доступа к объекту не имеет ничего общего с сериализацией объекта (превращением объекта в поток байтов); сериализация доступа означает, что потоки обращаются к объекту *по очереди*, а не одновременно.

Составные действия (*compound action*) над разделяемым состоянием, такие как увеличение счетчика попаданий (*read-modify-write*) или отложенная инициализация (*check-then-act*), должны быть атомарными, чтобы избежать условий гонки. Удержание блокировки на протяжении всего времени выполнения составного действия, может сделать его атомарным. Однако недостаточно просто обернуть составное действие блоком `synchronized`; если синхронизация используется для координации доступа к переменной, она необходима *во всех местах*, где *происходит обращение к этой переменной*. Более того, при использовании блокировок для координации доступа к переменной, одна и та же блокировка должна использоваться *везде*, где доступна эта переменная.

Распространенной ошибкой является предположение, что синхронизацию необходимо использовать только при записи в общие переменные; *это просто не соответствует действительности*. (Причины для этого замечания станут яснее в разделе 3.1)

К каждой изменяемой переменной состояния, к которой может обращаться более одного потока, *все* обращения должны выполняться с одной и той же блокировкой. В этом случае мы можем сказать, что *переменная защищена* этой блокировкой.

В классе `SynchronizedFactorizer`, в листинге 2.6, переменные `lastNumber` и `lastFactors` защищены внутренней блокировкой объекта сервлета; это отмечено аннотацией `@GuardedBy`.

Между внутренней блокировкой объекта и его состоянием нет связи; поля объекта не обязательно должны быть защищены встроенной блокировкой, хотя это вполне допустимое соглашение о блокировке, используемое многими классами. Захват блокировки, связанной с объектом, не препятствует другим потокам в получении доступа к этому объекту – единственno, что, захват блокировки предотвращает захват этой же блокировки любыми другими потоками. Тот факт, что каждый объект имеет встроенную блокировку, является просто удобством, поэтому вам не нужно явно создавать объекты блокировки²⁸. Создание *протоколов блокировки* (*locking protocols*) или *политик синхронизации*, обеспечивающих безопасный доступ к разделяемому состоянию, и их согласованное использование в рамках всего кода программы, зависит от вас.

Каждая разделяемая, изменяемая переменная должна быть защищена только одной блокировкой. Пишите код прозрачно для тех, кто будет его сопровождать: блокировка должна явно соотноситься с переменной.

Общее соглашение о блокировках заключается в инкапсуляции изменяемого состояния в объекте и защите его от параллельного доступа - путем синхронизации любой ветки кода, обращающейся к изменяемому состоянию - с помощью встроенной блокировки объекта. Такой подход используется во множестве потокобезопасных классов, таких как `Vector` и других синхронизированных коллекциях. В таких случаях, все переменные состояния объекта, защищаются встроенной блокировкой объекта. Однако в этом шаблоне нет ничего особенного, и ни компилятор, ни среда выполнения не применяют этот (или любой другой)

²⁸ Оглядываясь назад можно сказать, что такое проектное решение, вероятно, было плохим: оно не только может ввести в заблуждение, но и заставляет разработчиков JVM идти на компромиссы между размером объекта и производительностью блокировки.

шаблон блокировки²⁹. Протокол блокировки можно с легкостью случайно разрушить, добавив новый метод или ветку кода и забыв использовать синхронизацию.

Не все данные должны быть защищены блокировками - только изменяемые данные, которые будут доступны из нескольких потоков. В главе 1 мы описывали, как добавление простого асинхронного события, такого как `TimerTask`, может создать требования к потокобезопасности, которые скажутся на всём коде программы, особенно, если состояние программы плохо инкапсулировано. Рассмотрим однопоточную программу, обрабатывающую большой объём данных. Однопоточные программы не требуют синхронизации, потому что нет потоков, обращающихся к разделяемым данным. Теперь представьте, что вы хотите добавить функцию для создания периодических снимков прогресса программы, таким образом, чтобы ей не пришлось начинать заново, в случае если она упадёт с ошибкой (*crashes*) или должна будет быть остановлена. Вы можете сделать это с помощью класса `TimerTask`, который отключается каждые десять минут, сохраняя состояние программы в файле.

Поскольку класс `TimerTask` будет вызван из другого потока (один из которых управляет классом `Timer`), любые данные, участвующие в снимке, становятся доступны из двух потоков: основного потока программы и потока класса `Timer`. Это означает, что не только код класса `TimerTask` должен использовать синхронизацию при доступе к состоянию программы, но и любая другая ветка кода остальной части программы, касающаяся тех же данных. То, что раньше не требовало синхронизации, теперь требует синхронизации во всей программе.

Когда переменная защищена блокировкой - это означает, что доступ к этой переменной каждый раз выполняется с этой блокировкой - вы гарантируете, что только один поток, в конкретный момент времени, может получить доступ к этой переменной. Когда класс имеет инварианты, которые включают более одной переменной состояния, существует дополнительное требование: каждая переменная, участвующая в инварианте, должна быть защищена *одной и той же* блокировкой. Это позволяет вам получить доступ к ним или обновить их в одной атомной операции, сохранив инвариант³⁰. Это правило демонстрируется в классе `SynchronizedFactorizer`: и кэшированное число, и кэшированный фактор защищены внутренней блокировкой объекта сервлета.

Для каждого инварианта, включающего более одной переменной, все переменные, участвующие в этом инварианте, должны быть защищены *одной и той же* блокировкой.

Если синхронизация является “лекарством” от условий гонок, почему бы просто не объявить каждый метод как `synchronized`? Оказывается, такое беспорядочное применение оператора `synchronized`, может привести либо к слишком большой, либо к слишком малой синхронизации приложения. Простой синхронизации каждого метода, как в классе `Vector`, недостаточно для отображения составных действий на атомарность класса `Vector`:

```
if (!vector.contains(element))
    vector.add(element);
```

²⁹ Инструменты аудита кода, такие как FindBugs, могут определять ситуации, когда переменная часто, но не всегда, доступна с блокировкой, что может указывать на ошибку.

³⁰ В данном случае, под *инвариантом* понимается согласованность состояния объекта.

Попытка провести операцию *put-if-missing* содержит условие гонки, хотя оба метода, *contains* и *add*, атомарны. Хотя синхронизированные методы могут сделать отдельные операции атомарными, при объединении нескольких операций в составное действие требуется дополнительная блокировка. (См. раздел 4.4, в нём рассказывается о некоторых методах для безопасного добавления дополнительных атомарных операций к потокобезопасным объектам.) В то же время, синхронизация каждого метода может привести к проблемам с живучестью или производительностью, как мы видели в классе *SynchronizedFactorizer*.

2.5 Живучесть и производительность

В классе *UnsafeCachingFactorizer*, мы ввели в сервлет факторинга механизмы кэширования, в надежде, что это улучшит производительность. Кэширование требует некоторого разделяемого состояния, которое, в свою очередь, требует синхронизации для обеспечения целостности этого состояния. Но то, каким образом мы использовали синхронизацию в классе *SynchronizedFactorizer*, только ухудшило его производительность. Политика синхронизации для класса *SynchronizedFactorizer* заключается в защите каждой переменной состояния с помощью внутренней блокировки объекта сервлета, и эта политика была реализована путем синхронизации всего метода *service*. Этот простой, грубый подход восстановил безопасность, но обошёлся дорого.

Поскольку метод *service* помечен как *synchronized*, только один поток может выполнить его в один момент времени. Это подрывает предназначение фреймворка сервлетов – заключающееся в том, что сервлеты должны быть в состоянии обрабатывать несколько запросов одновременно – и может привести к разочарованию пользователей, если нагрузка достаточно высока. Если сервлет занят факторингом большого числа, другие клиенты будут вынуждены ожидать, пока текущий запрос не будет завершен, прежде чем сервлет сможет запуститься с новым числом. Если в системе имеется несколько CPUs, процессоры могут оставаться в состоянии простоя (*idle*) даже при высокой нагрузке. В любом случае, даже кратковременные запросы, например запросы, для которых кэшируется значение, могут занять неожиданно много времени, потому что вынуждены ожидать завершения выполнения предыдущих длительных запросов.

На рисунке 2.1 показано, что происходит при поступлении нескольких запросов к синхронизированному сервлету факторинга: они помещаются в очередь и обрабатываются последовательно.

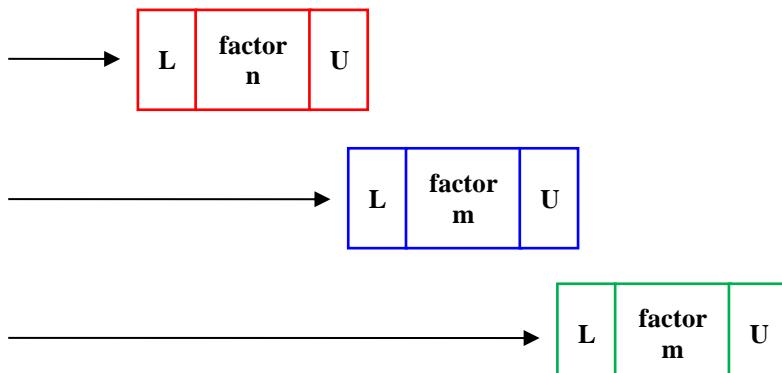


Рисунок 2.1 Плохой уровень параллелизма в классе *SynchronizedFactorizer*.

Мы бы описали это веб-приложение, как демонстрирующее низкий уровень параллелизма: количество одновременных вызовов ограничено не доступностью ресурсов обработки, а структурой самого приложения. К счастью, достаточно легко улучшить параллелизм сервлета, при этом сохранив его потокобезопасность, применив сужение области действия блока `synchronized`. Вы должны позаботиться о том, чтобы область действия блока `synchronized` была не слишком ограничена; вы бы не хотели разделить операцию, которая должна быть атомарной, на выполнение более чем в одном блоке `synchronized`. Разумно попытаться исключить из блоков `synchronized` длительные операции, которые не влияют на разделяемое состояние объекта, чтобы другие потоки не ограничивались в доступе к разделяемому состоянию, пока происходит выполнение длительной операции.

Класс `CachedFactorizer`, представленный в листинге 2.8, реструктурировал сервлет для использования двух отдельных блоков `synchronized`, каждый из которых охватывает небольшой кусок кода. Один из них защищает последовательность действий *check-then-act*, которая проверяет, можем ли мы просто вернуть кэшированный результат, а другой защищает обновление переменных состояния - как кэшированного числа, так и кэшированного фактора. В качестве бонуса, мы повторно ввели счетчик посещений и добавили счетчик "попаданий в кэш", обновляя их в начальном блоке `synchronized`. Поскольку эти счетчики составляют разделяемое изменяемое состояние, мы должны использовать синхронизацию везде, где к ним происходит обращение. Части кода, которые находятся за пределами блоков `synchronized`, работают исключительно с локальными (находящимися в стеке) переменными, которые не разделяются между потоками и поэтому не требуют синхронизации.

```
@ThreadSafe
public class CachedFactorizer implements Servlet {
    @GuardedBy("this") private BigInteger lastNumber;
    @GuardedBy("this") private BigInteger[] lastFactors;
    @GuardedBy("this") private long hits;
    @GuardedBy("this") private long cacheHits;

    public synchronized long getHits() { return hits; }
    public synchronized double getCacheHitRatio() {
        return (double) cacheHits / (double) hits;
    }

    public void service(ServletRequest req, ServletResponse resp) {
        BigInteger i = extractFromRequest(req);
        BigInteger[] factors = null;
        synchronized (this) {
            ++hits;
            if (i.equals(lastNumber)) {
                ++cacheHits;
                factors = lastFactors.clone();
            }
        }
        if (factors == null) {
```

```
    factors = factor(i);
    synchronized (this) {
        lastNumber = i;
        lastFactors = factors.clone();
    }
}
encodeIntoResponse(resp, factors);
}
```

Листинг 2.8 Сервлет кэширующий последний запрос и результат его выполнения.

Класс `CachedFactorizer` более не использует тип `AtomicLong` для подсчёта посещений, вместо этого вернувшись к использованию поля типа `long`. Было бы безопаснее использовать тип `AtomicLong`, но это принесёт меньше пользы, чем было в случае класса `CountingFactorizer`. Атомарные переменные полезны для выполнения атомарных операций над одной переменной, но, поскольку мы уже используем блоки `synchronized` для построения атомарных операций, использование двух различных механизмов синхронизации может вызвать путаницу и не приведёт к приросту производительности или повышению уровня безопасности.

Реструктуризация класса `CachedFactorizer` обеспечивает баланс между простотой (синхронизация всего метода) и параллелизмом (синхронизация возможно кратчайших веток выполнения кода). Захват и снятие блокировки несёт в себе накладные расходы, поэтому нежелательно слишком часто помещать код в блоки `synchronized` (например, факторинг `++hits` в своём собственном блоке `synchronized`), даже если это не подвергает атомарность угрозе. Класс `CachedFactorizer` удерживает блокировку в моменты доступа к переменным состояния и на время выполнения составных действий, но освобождает ее перед выполнением потенциально длительной операции факторинга. Это позволяет сохранять потокобезопасность, без чрезмерного влияния на параллелизм; ветки кода в каждом из блоков `synchronized` являются "достаточно короткими".

Решение о том, насколько большими или малыми должны быть блоки `synchronized`, может потребовать компромиссов между конкурирующими силами, оказывающими влияние на дизайн, включая безопасность (которая не должна быть скомпрометирована), простоту и производительность. Иногда простота и производительность противоречат друг другу, хотя, как показывает класс `CachedFactorizer`, разумный баланс обычно может быть найден.

Часто существует противоречие между *простотой и производительностью*. При реализации политики синхронизации, не поддавайтесь искушению преждевременно пожертвовать простотой (потенциально снижающей безопасность) ради производительности.

Всякий раз, когда вы используете блокировку, вы должны знать, что в блоке делает код и насколько вероятно, что потребуется много времени на его выполнение. Длительное удержание блокировки в связи с интенсивными вычислениями, либо из-за выполнения потенциально блокирующей операции, создает риск возникновения проблем с производительностью или живучестью.

Избегайте удержания блокировок во время длительных вычислений или операций, которые могут не завершиться быстро, например, сетевой или консольный ввод/вывод.

Глава 3 Совместно используемые объекты

В начале главы 2 мы заявили, что написание корректных параллельных программ в первую очередь касается управления доступом к разделяемому, изменяемому состоянию. В прошлой главе речь шла об использовании синхронизации для предотвращения одновременного доступа нескольких потоков к одним и тем же данным; в этой главе рассматриваются методы совместного использования (*sharing*)³¹ и публикации объектов, чтобы они могли быть безопасно доступны из нескольких потоков. Вместе они закладывают основу для создания потокобезопасных классов, и для безопасного структурирования параллельных приложений с использованием классов библиотеки `java.util.concurrent`.

Мы видели, как синхронизированные блоки и методы могут гарантировать, что операции выполняются атомарно, но существует распространенное заблуждение о том, что оператор `synchronized` касается *только* описания атомарности или разграничения «критических секций». Синхронизация также имеет другой важный и тонкий аспект: *видимость памяти*. Мы хотим не только запретить одному потоку изменять состояние объекта, когда он используется другим потоком, но и гарантировать, что когда поток изменит состояние объекта, другие потоки смогут *видеть* фактически внесенные изменения. Но без синхронизации, этого может не произойти. Вы можете обеспечить безопасную публикацию объектов либо с помощью явной синхронизации, либо путем использования синхронизации, встроенной в классы библиотек.

3.1 Видимость

Видимость – понятие тонкое, потому что вещи, которые могут пойти не так, противоречат здравому смыслу. В однопоточной среде, если вы пишете значение в переменную и позже читаете значение этой переменной без промежуточных записей, вы можете ожидать, что получите то же самое значение обратно. Это кажется естественным. Сначала может быть сложно это принять, но когда операции чтения и записи происходят в разных потоках, *это не так*. В общем, нет никакой гарантии, что читающий поток увидит значение, написанное другим потоком, своевременно или даже вообще. Чтобы обеспечить видимость записи в память между потоками, вы должны использовать синхронизацию.

В классе `NoVisibility` в листинге 3.1 показано, что может пойти не так, когда потоки совместно используют разделяемые данные без синхронизации. Два потока, главный поток и поток читатель, получают доступ к разделяемым переменным `ready` и `number`. Главный поток запускает поток-читатель и затем устанавливает переменной `number` значение 42, а переменной `ready` значение `true`. Поток читатель будет прокручиваться до тех пор, пока не увидит, что переменная `ready` имеет значение `true`, и затем напечатает значение переменной `number`. Хотя может показаться очевидным, что класс `NoVisibility` будет печатать 42, на самом деле, возможно, что он будет печатать ноль или вообще никогда не завершит свою работу! Поскольку он не использует соответствующую синхронизацию, нет

³¹ Разделения доступа к объекту между несколькими потоками.

никакой гарантии, что значения переменных `ready` и `number`, записанные основным потоком, будут видны потоку-читателю.

```
public class NoVisibility {  
    private static boolean ready;  
    private static int number;  
  
    private static class ReaderThread extends Thread {  
        public void run() {  
            while (!ready)  
                Thread.yield();  
            System.out.println(number);  
        }  
    }  
  
    public static void main(String[] args) {  
        new ReaderThread().start();  
        number = 42;  
        ready = true;  
    }  
}
```



Листинг 3.1 Совместное использование переменных без синхронизации. Не делайте так.

Класс `NoVisibility` может зацикливаться навсегда, потому что значение переменной `ready` может никогда не стать видимым для потока читателя. Еще более странно, что класс `NoVisibility` может напечатать ноль, потому что запись в переменную `ready` может быть сделана видимой для потока читателя до записи в переменную `number`, это явление известно как *переупорядочение(reordering)*. Нет никакой гарантии, что операции в одном потоке будут выполняться в порядке заданном программой, при условии, что переупорядочение не будет обнаружено изнутри этого потока - *даже если переупорядочение будет очевидным для других потоков*³². Когда главный поток, без использования синхронизации, сперва пишет значение в переменную `number`, а затем в переменную `ready`, поток читатель может увидеть, что эти записи происходят в обратном порядке - или нет.

В отсутствие синхронизации компилятор, процессор и среда выполнения могут делать некоторые совершенно странные вещи с порядком выполнения операций. Попытки рассуждать о порядке, в котором “должны” происходить действия с памятью, в недостаточно синхронизированных многопоточных программах почти наверняка будут неправильными.

Класс `NoVisibility` примерно так же просто, как и параллельная программа, может получить два потока и две общие переменные - и все же все еще слишком легко прийти к неправильным выводам о том, что он делает или даже будет ли

³² Такое решение может показаться признаком плохого дизайна, но это позволяет JVM полностью использовать преимущества современных многопроцессорных систем. Например, при отсутствии синхронизации модель памяти Java позволяет компилятору переупорядочить операции и значения кэша в регистрах и позволяет процессорам переупорядочивать операции и значения кэша, специфичных для процессора. Подробнее см. главу 16.

завершён. Рассуждения о недостаточно синхронизированных параллельных программах непомерно сложны.

Все это может показаться немного пугающим, и так и должно быть. К счастью, есть простой способ избежать этих сложных проблем: всегда используйте правильную синхронизацию *при совместном использовании данных в потоках*.

3.1.1 Устаревшие данные

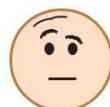
Класс `NoVisibility` продемонстрировал один из путей, следуя которым недостаточно синхронизированные программы могут получать неожиданные результаты: *устаревшие данные (stale data)*. Когда поток-читатель проверяет переменную `ready`, он может увидеть устаревшее значение. Если синхронизация не используется *каждый раз при доступе к переменной*, с некоторой долей вероятности можно увидеть устаревшее значение этой переменной. Хуже того, устаревание это не “всё или ничего”: поток может видеть обновленное значение одной переменной, но также и устаревшее значение другой переменной, которая была записана первой.

Когда еда несвежая, она обычно все ещё съедобна - просто менее приятна. Но устаревшие данные могут быть более опасными. В то время как устаревший счетчик посещений в веб-приложении может быть не так уж плох³³, устаревшие значения могут вызвать серьезные сбои безопасности или живучести. В классе `NoVisibility` устаревшие значения могут привести к тому, что будет напечатано неверное значение или выполнение программы не сможет быть прервано. Еще сложнее обстоят дела с устаревшими значениями ссылок на объекты, такими как указатели ссылок в реализации связанного списка. *Устаревшие данные могут вызывать серьезные и запутанные сбои, такие как непредвиденные исключения, повреждение структур данных, неточные вычисления и бесконечные циклы.*

Класс `MutableInteger` в листинге 3.2 не потокобезопасен, потому что значение поля `value` доступно из обоих методов `get` и `set` без синхронизации. Среди прочих опасностей, он подвержен устареванию значений: если один поток вызывает `set`, другие потоки, вызывающие `get`, могут видеть или не видеть это обновление.

```
@NotThreadSafe
public class MutableInteger {
    private int value;

    public int get() { return value; }
    public void set(int value) { this.value = value; }
}
```



Листинг 3.2 Непотокобезопасный изменяемый холдер типа `Integer`.

Мы можем сделать класс `MutableInteger` потокобезопасным, синхронизировав геттер и сеттер, как показано в классе `SynchronizedInteger` в листинге 3.3. Синхронизация только сеттера будет недостаточной: потоки, вызывающие метод `get`, все равно смогут видеть устаревшие значения.

³³ Чтение данных без синхронизации аналогично использованию уровня изоляции `READ_UNCOMMITTED` в базе данных, где вы готовы разменивать точность на производительность. Однако, в случае использования несинхронизированных операций чтения, вы теряете большую степень точности, так как видимое значение разделяемой переменной может устареть в любой момент времени.

```
@ThreadSafe
public class SynchronizedInteger {
    @GuardedBy("this") private int value;

    public synchronized int get() { return value; }
    public synchronized void set(int value) { this.value = value; }
}
```

Листинг 3.3 Потокобезопасный изменяемый холдер типа Integer.

3.1.2 Неатомарные 64 битные операции

Когда поток читает переменную без синхронизации, он может увидеть устаревшее значение, но, по крайней мере, он видит значение, которое было фактически помещено туда каким-то потоком, а не какое-то случайное значение. Эта гарантия безопасности называется “*безопасностью из неоткуда*” (*out-of-thin-air safety*).

Безопасность “из ниоткуда” применяется ко всем переменным, с одним исключением: 64 битные числовые переменные (`double` и `long`), не объявленные с ключевым словом `volatile` (см. раздел [3.1.4](#)). Модель памяти Java требует, чтобы операции выборки и хранения были атомарными, но для **не** `volatile` переменных типа `long` и `double` JVM разрешено рассматривать 64 битное чтение или запись как две отдельные 32 битные операции. Если операции чтения и записи происходят в разных потоках, существует вероятность прочитать **не** `volatile` переменную типа `long` и в результате получить старшие 32 бита от одного значения и младшие 32 бита от другого³⁴. Таким образом, даже если вы не заботитесь об устаревании значений переменных, небезопасно использовать разделяемые переменные типа `long` или `double` в многопоточных программах, если они не объявлены с ключевым словом `volatile` или не защищены блокировкой.

3.1.3 Блокировки и видимость

Внутреннюю блокировку можно использовать для того, чтобы гарантировать, что один поток видит изменения внесённые другим потоком в предсказуемом виде, как показано на рисунке 3.1. Когда поток *A* выполняет синхронизированный блок, и впоследствии поток *B* входит в блок `synchronized`, защищаемый той же блокировкой, значения переменных, которые были видны потоку *A* до освобождения блокировки, гарантированно будут видны потоку *B* после захвата блокировки. Другими словами, всё в потоке *A*, выполненное внутри или до блока `synchronized`, видимо для потока *B*, когда он выполняет блок `synchronized`, защищенный одной и той же блокировкой. *Без синхронизации таких гарантий нет.*

Теперь мы можем привести другую причину следовать правилу, требующему, чтобы все потоки синхронизировались на одной и той же блокировке при доступе к разделяемой изменяемой переменной - чтобы гарантировать, что значения, записанные одним потоком, станут видимыми для других потоков. В противном случае, если поток читает переменную, не удерживая соответствующую блокировку, он может увидеть устаревшее значение.

³⁴ В то время, когда была написана спецификация виртуальной машины Java, множество широко используемых процессорных архитектур не могли эффективно выполнять атомарные 64-разрядные арифметические операции.

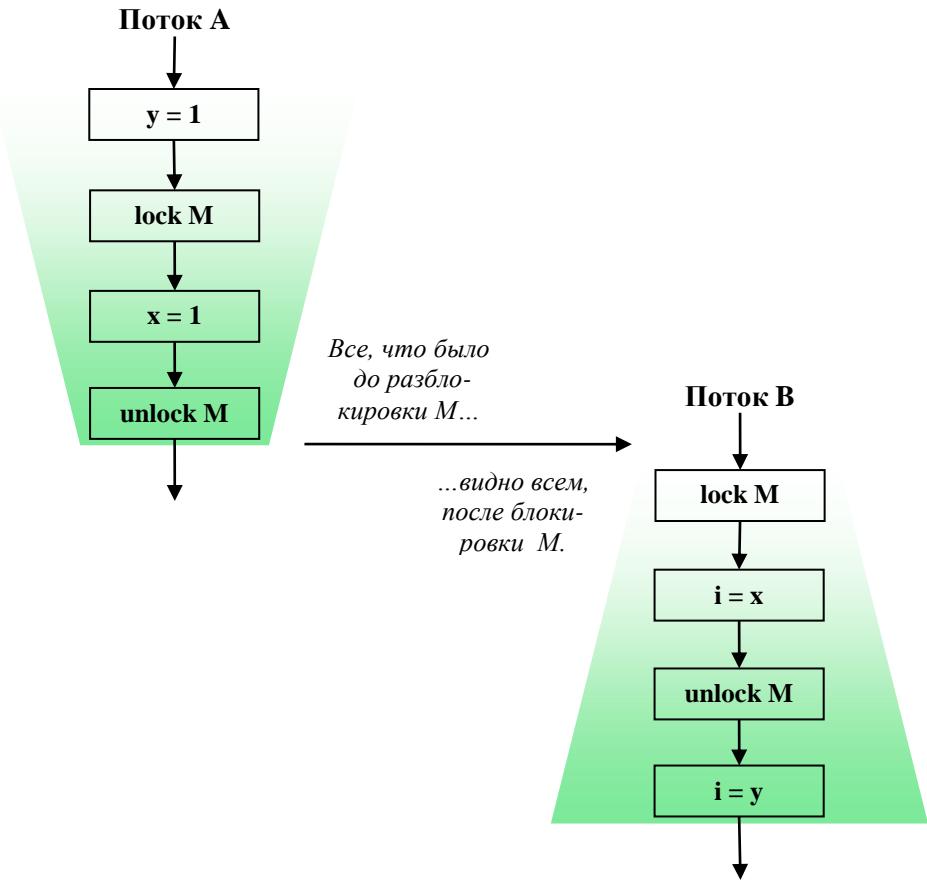


Рисунок 3.1 Видимость гарантированная синхронизацией.

Блокировка - это не просто взаимное исключение, но также и видимость памяти. Чтобы убедиться, что все потоки видят самые последние значения разделяемых изменяемых переменных, потоки чтения и записи должны синхронизироваться на общей блокировке

3.1.4 Volatile переменные

Язык Java также предоставляет альтернативную, более слабую форму синхронизации, *изменчивые переменные (volatile variables)*, чтобы гарантировать, что обновления переменной предсказуемо распространяются к другим потокам. Когда поле объявлено с ключевым словом `volatile`, компилятор и среда выполнения уведомляются о том, что эта переменная является разделяемой и что операции над ней не должны переупорядочиваться с другими операциями памяти. Изменчивые переменные не кэшируются в регистрах или в кэшах, где они скрыты от других процессоров, поэтому чтение изменчивой переменной всегда возвращает самую последнюю запись, сделанную любым потоком.

Хороший способ восприятия `volatile` переменных заключается в том, чтобы представить, будто бы они, в грубом приближении, похожи на класс

`SynchronizedInteger` из листинга 3.3, в котором операции чтения и записи `volatile` переменной заменены вызовами `get` и `set`³⁵. Тем не менее, доступ к `volatile` переменной не требует захвата блокировки и поэтому не может привести к блокировке выполняющегося потока, что делает переменную, объявленную как `volatile`, более легковесным механизмом синхронизации, по сравнению с блоком `synchronized`³⁶.

Эффекты видимости `volatile` переменной выходят за пределы значения самой `volatile` переменной. Когда поток *A* записывает значение в `volatile` переменную, а затем поток *B* считывает значение этой же переменной, значения всех переменных, которые были видны потоку *A* перед записью в `volatile` переменную, становятся видимыми потоку *B*, после прочтения `volatile` переменной. Таким образом, с точки зрения видимости памяти, запись `volatile` переменной похожа на выход из блока `synchronized`, а чтение `volatile` переменной похоже на вход в блок `synchronized`. Однако мы не рекомендуем слишком сильно полагаться на `volatile` переменные для обеспечения видимости; код, который опирается на `volatile` переменные для обеспечения видимости произвольного состояния, является более хрупким и его сложнее понять, чем код, который использует блокировки.

Используйте переменные `volatile` только в том случае, если они упрощают реализацию и проверку политики синхронизации; избегайте использования переменных `volatile`, когда корректность проверки потребует тонких рассуждений о видимости. Хороший стиль использования переменных `volatile` включает обеспечение видимости их собственного состояния, состояния объекта, на который они ссылаются, или указание на то, что произошло важное событие жизненного цикла (например, инициализация или завершение работы).

В листинге 3.4 приведён типичный пример использования `volatile` переменных: проверка флага состояния для определения момента выхода из цикла. В этом примере наш “очеловеченный” поток пытается уснуть с помощью проверенного временем метода подсчета овец. Для того, чтобы это пример работал, переменная-флаг `asleep` должна быть `volatile`. В ином случае, поток может не заметить установку значения переменной `asleep` другим потоком³⁷. Вместо этого, мы могли бы использовать блокировку, чтобы обеспечить видимость изменений в переменной `asleep`, но это сделало бы код более громоздким.

```
volatile boolean asleep;
...
while (!asleep)
```

³⁵ Эта аналогия не точна; эффекты видимости памяти немного сильнее, чем у изменчивых (`volatile`) переменных. См. главу 16.

³⁶ Чтение `volatile` переменных лишь немного затратнее чтения не `volatile` переменных, в большинстве современных процессорных архитектур.

³⁷ **Совет по отладке:** для серверных приложений всегда указывайте параметр командной строки JVM для переключения в серверный режим, даже для разработки и тестирования. Серверная JVM выполняет больше оптимизаций, чем клиентская JVM, такие как подъем переменных из цикла, которые не изменяются в цикле; код, который может работать в среде разработки (клиентская JVM) может сломаться в среде развертывания (серверная JVM). Например, если мы “забыли” объявить переменную в листинге 3.4, серверная JVM может поднять тест из цикла (превратив его в бесконечный цикл), но клиентская JVM этого не сделает. Бесконечный цикл, который проявляется в разработке, обходится намного дешевле того, который появляется только в продуктиве.

```
countSomeSheep();
```

Листинг 3.4 Класс CountingSheep

Переменные `volatile` удобны, но у них есть ограничения. Наиболее часто переменные `volatile` используются в качестве флагов завершения, прерывания или состояния, как переменная-флаг `asleep` в листинге 3.4. Переменные `volatile` можно также использовать для получения других типов сведений о состоянии, но при этом требуется уделять больше внимания. Например, семантика `volatile` недостаточно сильна, чтобы сделать операцию инкремента (`count++`) атомарной, если только вы не можете гарантировать, что переменная записывается только из одного потока. (Атомарные переменные обеспечивают атомарную поддержку чтения-изменения-записи и часто могут использоваться как "лучшие `volatile` переменные"; см. главу [15.](#).)

Блокировки могут гарантировать и видимость, и атомарность; `volatile` переменные только видимость.

Переменные `volatile` можно использовать только при соблюдении следующих условий:

- Запись в переменную не зависит от ее текущего значения, или вы можете гарантировать, что только один поток обновляет значение;
- Переменная не участвует в инвариантах с другими переменными состояния;
- Блокировка не требуется по какой-либо другой причине во время обращения к переменной.

3.2 Публикация и побег

Публикация (*publishing*) объекта означает, что он доступен для кода вне его текущей области видимости, например, путем сохранения ссылки на него там, где другой код может найти её, или путём возвращения его из неприватного (*nonprivate*) метода, или путём передачи его методу в другом классе. Во многих ситуациях мы хотим иметь гарантию того, что объекты и их внутренние компоненты не будут публиковаться. В других ситуациях мы хотим опубликовать объект для общего использования, но это может потребовать использования синхронизации. Публикация внутренних переменных состояния может нарушить инкапсуляцию и затруднить сохранение инвариантов; публикация объектов до их полного построения может поставить под угрозу потокобезопасность. Объект, который был опубликован в тот момент, когда этого не должно было произойти, называют **сбежавшим** (*escaped*). В разделе 3.5 рассматриваются идиомы безопасной публикации; прямо сейчас, мы посмотрим, как объект может сбежать.

Наиболее вопиющей формой публикации является сохранение ссылки в открытом статическом поле, где любой класс и поток могут видеть его, как в листинге 3.5. Метод `initialize` создает новый объект `HashSet` и публикует его, сохраняя ссылку на него в переменной `knownSecrets`.

```
public static Set<Secret> knownSecrets;

public void initialize() {
    knownSecrets = new HashSet<Secret>();
```

}

Листинг 3.5 Публикация объекта. Класс `Secrets`.

Публикация одного объекта может вызвать косвенную публикацию других. Если вы добавите объект класса `Secret` к опубликованному набору объектов `knownSecrets`, вы также опубликуете и объект класса `Secret`, потому что любой код сможет перебрать значения объекта класса `Set` и получить ссылку на новый объект класса `Secret`. Аналогичным образом, возврат ссылки из не приватного метода также публикует возвращенный объект. Класс `UnsafeStates` в листинге 3.6 публикует якобы приватный массив аббревиатур названий штатов.

```
class UnsafeStates {  
    private String[] states = new String[] {  
        "AK", "AL" ...  
    };  
    public String[] getStates() { return states; }  
}
```



Листинг 3.6 Предоставление возможности «сбежать» внутреннему изменяемому состоянию. Не делайте так.

Публикация объекта `states`, таким образом, является проблематичной, потому что любой вызывающий код может изменить его содержимое. В этом случае массив `states` “сбегает” из предполагаемой области видимости, потому что то, чему предполагалось быть приватным, фактически было сделано публичным.

Публикация объекта также публикует все объекты, на которые ссылаются его публичные поля. Обобщая сказанное, любой объект, который доступен из опубликованного объекта, следя некоторой цепочке публичных ссылок на поля и вызовы методов, также был опубликован.

С точки зрения класса `C`, чужой метод - один из тех, чье поведение не полностью определяется классом `C`. Это понятие включает в себя методы в других классах, а также переопределенные методы (ни `private`, ни `final`) в самом классе `C`. Передача объекта чужому методу, также должна рассматриваться как публикация этого объекта. Поскольку вы не можете предполагать, какой код фактически будет вызван, вы не можете знать, захочет ли чужой метод опубликовать объект или сохранить ссылку на него, которая впоследствии может быть использована из другого потока.

Действительно ли другой поток сделает что-либо с опубликованной ссылкой, на самом деле не имеет значения, потому что риск неправильного использования существует в любом случае³⁸. Как только объект убегает, вы должны предположить, что другой класс или поток могут, злонамеренно или по неосторожности, использовать его некорректно. Это является веской причиной для использования инкапсуляции: она позволяет выполнять практический анализ программ на корректность и усложняет случайное нарушение ограничений, заложенных на стадии проектирования.

Последним механизмом, с помощью которого можно опубликовать объект или его внутреннее состояние, является публикация внутреннего экземпляра класса, как показано в классе `ThisEscape`, в листинге 3.7. Когда класс `ThisEscape` публикует `EventListener`, он также неявно публикует вложенный класс

³⁸ Если кто-то украл ваш пароль и поместил его в группе новостей `alt.free-passwords`, эта информация “убежала”: независимо от того, использовал ли кто-либо эти учетные данные для причинения вреда, ваша учетная запись все еще скомпрометирована. Публикация ссылки представляет собой такой же риск.

`ThisEscape`, потому что экземпляры внутреннего класса содержат скрытую ссылку на вложенный экземпляр.

```
public class ThisEscape {  
    public ThisEscape(EventSource source) {  
        source.registerListener(  
            new EventListener() {  
                public void onEvent(Event e) {  
                    doSomething(e);  
                }  
            }  
        );  
    }  
}
```



Листинг 3.7 Неявное позволяние "сбежать" ссылке на `this`. Не делайте так.

3.2.1 Методы безопасного построения

Класс `ThisEscape` иллюстрирует важный частный случай побега - когда ссылка на `this` ссылка "сбегает" во время построения. Когда публикуется внутренний экземпляр `EventListener`, то же самое относится и к вложенному экземпляру `ThisEscape`. Но объект находится в предсказуемом и согласованном состоянии только после *возврата* из конструктора, поэтому публикация объекта из конструктора может привести к публикации не полностью построенного объекта. Это справедливо, даже если публикация выполняется последним выражением в конструкторе. Если ссылка на `this` сбегает в процессе построения, объект считается *построенным неправильно*³⁹. This is true even if the publication is the last statement in the constructor.

Не позволяйте ссылке на `this` "сбегать" в процессе построения.

Распространённой ошибкой, позволяющей сбегать ссылке на `this` во время построения, является запуск потока из конструктора. Когда объект создает поток из своего конструктора, он почти всегда делится своей ссылкой с новым потоком либо явно (передавая его конструктору), либо неявно (поскольку `Thread` или `Runnable` является внутренним классом объекта-владельца). Новый поток сможет увидеть объект-владелец, прежде чем он будет полностью построен. Нет ничего плохого в создании потока в конструкторе, но лучше не запускать поток сразу. Вместо этого предоставьте методы `start` или `initialize`, которые запустят собственный поток. (Дополнительные сведения о проблемах жизненного цикла службы см. в главе 7.) Вызов переопределяемого метода экземпляра (который не является ни `private`, ни `final`) из конструктора также может позволить сбежать ссылке на `this`.

Если у вас возникнет соблазн зарегистрировать слушателя событий или запустить поток из конструктора, вы можете избежать неправильного построения с

³⁹ Если говорить более конкретно, ссылка на `this` не должна покидать *поток* до тех пор, пока конструктор не вернёт управление. Ссылка на `this` может быть где-нибудь сохранена конструктором, но до завершения построения, она не может быть использована другими потоками. Класс `SafeListener` в листинге 3.8 использует эту технику.

помощью приватного конструктора и открытого фабричного метода, как показано в классе `SafeListener` в листинге 3.8.

```
public class SafeListener {
    private final EventListener listener;

    private SafeListener() {
        listener = new EventListener() {
            public void onEvent(Event e) {
                doSomething(e);
            }
        };
    }

    public static SafeListener newInstance(EventSource source) {
        SafeListener safe = new SafeListener();
        source.registerListener(safe.listener);
        return safe;
    }
}
```

Листинг 3.8 Использование фабричного метода для предотвращения побега ссылки на `this` в процессе построения.

3.3 Ограничение потока

Доступ к общим изменяемым данным требует использования синхронизации; один из способов избежать этого требования - *не предоставлять общий доступ*. Если доступ к данным осуществляется только из одного потока, синхронизация не требуется. Этот метод, *ограничение потока*, является одним из самых простых способов обеспечения безопасности потоков. Когда объект ограничен потоком, такое использование автоматически потокобезопасно, даже если ограниченный объект сам по себе не потокобезопасен [CPJ 2.3.2]⁴⁰.

Swing использует метод ограничения потока в широких пределах. Визуальные компоненты и модели данных Swing не потокобезопасны; напротив, потокобезопасность достигается путём ограничения доступа к ним только потоком обработки событий Swing. Чтобы правильно использовать Swing, код, выполняющийся в потоке, отличном от потока событий, не должен обращаться к этим объектам. (Чтобы упростить работу, Swing предоставляет механизм `invokeLater` для назначения `Runnable` на выполнение в потоке событий.) Множество ошибок параллелизма в приложениях Swing происходит из-за неправильного использования ограниченных объектов из других потоков.

Другим распространенным применением ограничения потоков является использование объединенных в пул объектов соединения JDBC (Java Database Connectivity) `Connection`. Спецификация JDBC не требует, чтобы объекты `Connection` были потокобезопасны.⁴¹ В типичных серверных приложениях поток получает соединение из пула, использует его для обработки одного запроса и

⁴⁰ Doug Lea. Concurrent Programming in Java, Second Edition. Addison–Wesley, 2000.

⁴¹ Реализации *пулов* соединений, предоставляемые серверами приложений, потокобезопасны; как правило, доступ к пулам соединений осуществляется из нескольких потоков, поэтому реализации, не являющиеся потокобезопасными, не имеют смысла.

возвращает его. Поскольку большинство запросов, таких как запросы сервлета или вызовы EJB (Enterprise JavaBeans), обрабатываются синхронно одним потоком, и пул не будет передавать одно и то же соединение другим потокам до тех пор, пока оно не будет возвращено, этот шаблон управления соединением неявно ограничивает объект `Connection` этим потоком на время выполнения запроса.

Так же, как язык не имеет механизмов для принудительного обеспечения того, чтобы переменная была защищена блокировкой, у него нет средств ограничения объекта потоком. Ограничение потоков - это элемент дизайна вашей программы, который должен обеспечиваться его реализацией. Язык и основные библиотеки предоставляют механизмы, которые могут помочь в поддержании ограничения потока - локальные переменные и класс `ThreadLocal` - но даже с ними, только программист несёт ответственность за гарантию того, что ограниченные потоком объекты не покинут предназначенного им потока.

3.3.1 Специальные ограничения потока

Специальное ограничение потока (Ad-hoc thread confinement) описывает, когда ответственность за поддержание ограничения потока полностью ложится на реализацию. Ограничение потока может быть хрупким, поскольку ни одна из возможностей языка, такая как модификаторы видимости или локальные переменные, не помогает ограничить объект целевым потоком. На самом деле ссылки на объекты, ограниченные потоками, такие как визуальные компоненты или модели данных в приложениях с графическим пользовательским интерфейсом, часто хранятся в открытых полях.

Решение использовать ограничение потока часто является следствием решения реализовать определенную подсистему, такую как GUI, как однопоточную подсистему. Однопоточные подсистемы могут иногда предложить такое преимущество, как простота, которое перевешивает хрупкость специального ограничения потока⁴².

Частный случай ограничения потока применяется и к `volatile` переменным. Безопасно выполнять операции *read-modify-write* с использованием общих переменных `volatile`, пока вы гарантируете, что `volatile` переменная будет записана только из одного потока. В этом случае вы ограничиваете модификацию только одним потоком, чтобы предотвратить условия гонки, а гарантии видимости `volatile` переменных устанавливают, что другие потоки будут видеть последнее изменённое значение.

В связки с хрупкостью, ограничение потока следует использовать крайне скрупульно; если возможно, используйте одну из более сильных форм ограничения потоков (ограничение стека или `ThreadLocal`).

3.3.2 Ограничение стека

Ограничение стека (Stack confinement) является частным случаем ограничения потока, при котором доступ к объекту может быть получен только через локальные переменные. Подобно тому, как инкапсуляция может упростить сохранение инвариантов, локальные переменные могут упростить ограничение объектов потоком. Локальные переменные внутренне ограничены исполняемым потоком; они существуют в стеке исполняемого потока, который недоступен для других

⁴² Еще одна причина сделать подсистему однопоточной - исключение возможности взаимоблокировок; это одна из основных причин того, почему большинство платформ GUI однопоточны. Однопоточные подсистемы описаны в главе 9.

потоков. Ограничение стека (также называемое внутри-поточным (*within-thread*) или локально-поточным (*thread-local*) использованием потока, не путать с библиотечным классом `ThreadLocal`!) проще поддерживать и оно менее хрупко, чем специальное ограничение потока.

В случае локальных переменных примитивных типов, таких как `numPairs` в методе `loadTheArk` в листинге 3.9, вы не сможете нарушить ограничение стека, даже если попробуете. Невозможно получить ссылку на примитивную переменную, поэтому семантика языка гарантирует, что примитивные локальные переменные всегда будут ограничены стеком.

```
public int loadTheArk(Collection<Animal> candidates) {
    SortedSet<Animal> animals;
    int numPairs = 0;
    Animal candidate = null;

    // animals confined to method, don't let them escape!
    animals = new TreeSet<Animal>(new SpeciesGenderComparator());
    animals.addAll(candidates);
    for (Animal a : animals) {
        if (candidate == null || !candidate.isPotentialMate(a))
            candidate = a;
        else {
            ark.load(new AnimalPair(candidate, a));
            ++numPairs;
            candidate = null;
        }
    }
    return numPairs;
}
```

Листинг 3.9 Ограничение потока локальными примитивами и ссылочными переменными в классе `Animals`.

Поддержание ограничения стека для ссылок на объекты требует немного больше помощи от программиста, чтобы гарантировать, что ссылающийся объект не сбежит. В классе `loadTheArk` мы создаем объект `TreeSet` и храним ссылку на него в переменной `animals`. На данный момент существует только одна ссылка на объект `Set`, содержащаяся в локальной переменной и, следовательно, ограничивающая исполняющим потоком. Однако если бы мы опубликовали ссылку на объект `Set` (или любую его внутреннюю часть), то ограничение было бы нарушено, и “животные” убежали бы.

Использование непотокобезопасного объекта в контексте “внутри потока” по-прежнему является потокобезопасным. Однако будьте осторожны: требование к дизайну, чтобы объект ограничивался исполняемым потоком, или осознание того, что закрытый объект не является потокобезопасным, часто существует только в голове разработчика в момент написания кода. Если допущение о внутрипоточном использовании объекта явно не задокументировано, те, кто будут в будущем сопровождать код, могут по ошибке позволить этому объекту сбежать.

3.3.3 Класс ThreadLocal

Более формальным средством поддержания ограничения потока является класс `ThreadLocal`, который позволяет связать значение каждого потока с объектом хранения значения. Класс `ThreadLocal` предоставляет методы доступа `get` и `set`, которые поддерживают отдельную копию значения для каждого потока, который использует его, поэтому `get` возвращает самое последнее значение, переданное `set` из текущего выполняемого потока.

Локально-поточные переменные часто используются для предотвращения совместного использования в проектах на основе изменяемых синглтонов или глобальных переменных. Например, однопоточное приложение может поддерживать подключение к глобальной базе данных, инициализируемое при запуске, чтобы избежать передачи объекта `Connection` каждому методу. Так как соединения JDBC не могут быть потокобезопасными, многопоточное приложение, которое использует глобальное соединение без дополнительной координации доступа, также не является потокобезопасным. Используя класс `ThreadLocal` для хранения соединения JDBC, как в методе `connectionHolder` в листинге 3.10, каждый поток будет иметь свое собственное соединение

```
private static ThreadLocal<Connection> connectionHolder
    = new ThreadLocal<Connection>() {
        public Connection initialValue() {
            return DriverManager.getConnection(DB_URL);
        }
    };

public static Connection getConnection() {
    return connectionHolder.get();
}
```

Листинг 3.10 Использование `ThreadLocal` для обеспечения ограничения потока.

Этот метод может также использоваться, когда часто используемая операция требует временного объекта, такого как буфер, и хочет избежать перераспределения временного объекта при каждом вызове. Например, до Java 5.0 метод `Integer.toString` использовал класс `ThreadLocal` для хранения 12-байтового буфера, используемого для форматирования его результата, вместо использования общего статического буфера (который потребовал бы блокировки) или выделения нового буфера для каждого вызова⁴³.

Когда поток вызывает метод `ThreadLocal.get` в первый раз, метод `initialValue` используется для обеспечения потока начальным значением. В концептуальном плане вы можете думать о классе `ThreadLocal<T>` как о холдере `Map<Thread, T>`, который сохраняет значения, зависящие от потока, хотя, на самом деле, это реализовано не так. Значения, специфичные для потока, хранятся в самом объекте `Thread`; когда поток завершается, значения, зависящие от потока, будут убраны “сборщиком мусора”.

⁴³ Этот метод вряд ли даст выигрыши производительности, если операции выполняются очень часто или выделение памяти необычайно дорого. В Java 5.0 он был заменен более прямолинейным подходом, заключающимся в выделении нового буфера для каждого вызова, в связи с предположением, что для чего-то столь же обыденного, как временный буфер, он не даст выигрыша производительности.

Если вы переносите однопоточное приложение в многопоточную среду, вы можете сохранить безопасность потоков путем преобразования общих глобальных переменных в объекты `ThreadLocal`, если это позволяет семантика общих глобальных переменных; кэш уровня приложения не был бы столь же полезным, если бы он был размещён в нескольких локально-поточных кэшах..

Класс `ThreadLocal` широко используется в реализации фреймворков. Например, контейнеры J2EE связывают контекст транзакции (*transaction context*) с исполняемым потоком на всём протяжении вызова EJB. Это легко реализуется с использованием статического объекта `ThreadLocal`, содержащего контекст транзакции: когда фреймворк необходимо определить, какая транзакция выполняется в настоящий момент, он извлекает контекст транзакции из объекта `ThreadLocal`. Это удобно в том смысле, что снижает необходимость передавать информацию о контексте выполнения в каждый метод, но связывает любой код, который использует этот механизм в фреймворке.

Легко злоупотребить классом `ThreadLocal`, трактуя его свойство ограничения потока как лицензию на использование глобальных переменных или как средство создания “скрытых” аргументов методов. Как и глобальные переменные, поточно-локальные переменные могут отвлекать от повторного использования и вводить скрытые связи между классами, и поэтому их следует использовать с большой осторожностью.

3.4 Неизменяемость

Другой конечной целью синхронизации является использование неизменяемых объектов [EJ Item 13]. The other end-run around the need to synchronize is to use *immutable* objects [EJ Item 13]. Почти все опасности атомарности и видимости, которые мы описывали до сих пор, такие как просмотр устаревших значений, потеря обновлений или наблюдение за объектом, находящимся в несогласованном состоянии, связаны с превратностями доступа из нескольких потоков, пытающихся получить доступ к одному и тому же изменяемому состоянию в одно и то же время. Если состояние объекта не может быть изменено, эти риски и сложности просто исчезают.

Неизменяемым (*immutable*) объектом является объект, состояние которого не может быть изменено после построения. Неизменяемые объекты по своей сути являются потокобезопасными; их инварианты устанавливаются конструктором, и если их состояние не может быть изменено, то эти инварианты всегда удерживаются.

Неизменяемые объекты всегда потокобезопасны.

Неизменяемые объекты *просты*. Они могут находиться только в одном состоянии, которое тщательно контролируется конструктором. Один из самых сложных элементов дизайна программы - рассуждение о возможных состояниях сложных объектов. С другой стороны, рассуждение о состоянии неизмененных объектов тривиально.

Неизменяемые объекты также *безопаснее*. Передача изменяемого объекта в ненадежный код или публикация в то место, где ненадежный код может найти его – опасна. Ненадежный код может изменить состояние объекта или, что еще хуже, сохранить ссылку на него и изменить его состояние позже из другого потока. С другой стороны, неизменяемые объекты не могут быть скомпрометированы с

помощью вредоносного или ошибочного кода, поэтому они могут безопасно совместно использоваться (*share*) и публиковаться без необходимости создавать защищённые копии [EJ Item 24].

Ни спецификация языка Java, ни модель памяти Java формально не определяют неизменность, но неизменяемость *не* эквивалентна простому объявлению всех полей объекта как `final`. Объект, все поля которого являются `final`, может по-прежнему быть изменяемым, поскольку `final` поля могут содержать ссылки на изменяемые объекты.

Объект *неизменяемый* если:

- Его состояние не может быть изменено после построения;
- Все поля объявлены как `final`⁴⁴;
- Он правильно построен (ссылка на `this` не может “сбежать” в процессе построения)).

Неизменяемые объекты могут по-прежнему использовать изменяемые объекты для управления своим состоянием, как показано в классе `ThreeStooges` в листинге 3.11. Хотя объект `Set`, в котором хранятся имена, изменяемый, дизайн класса `ThreeStooges` не позволяет изменить объект `Set` после построения. Ссылка на переменную `stooges` помечена как `final`, поэтому всё состояние объекта достигается через поле типа `final`. Последнее требование - правильное построение - легко выполняется, поскольку конструктор ничего не делает, чтобы эта ссылка стала доступной для кода, отличного от конструктора и кода, вызвавшего его.

```
@Immutable
public final class ThreeStooges {
    private final Set<String> stooges = new HashSet<String>();

    public ThreeStooges() {
        stooges.add("Moe");
        stooges.add("Larry");
        stooges.add("Curly");
    }

    public boolean isStooge(String name) {
        return stooges.contains(name);
    }
}
```

Листинг 3.11 Неизменяемый класс построенный с использованием изменяемых объектов.

Поскольку состояние программы изменяется все время, вы можете испытывать соблазн подумать, что неизменяемые объекты имеют ограниченное использование, но это не так. Существует разница между неизменяемым *объектом* и

⁴⁴ Технически возможно иметь неизменяемый объект, если не все поля объявлены как `final`. `String` - пример такого класса, но это всё опирается на тонкие рассуждения о “доброточныхенных” гонках данных, что требует глубокого понимания модели памяти Java. (Для любопытных: `String` вычисляет хэш-код в ленивой манере при первом вызове метода `hashCode` и кширует его в не финальном поле, но это работает только потому, что это поле может принимать только одно значение `nondefault`, которое является одинаковым каждый раз, когда оно вычисляется, потому что оно детерминировано выводится из неизменяемого состояния. Не пытайтесь повторить это дома.)

неизменяемой ссылкой на него. Состояние программы, сохраненное в неизменяемых объектах, все еще может быть обновлено путем «замены» неизменяемых объектов новыми экземплярами, содержащими новое состояние; в следующем разделе приведен пример использования этой методики⁴⁵.

3.4.1 Поля типа final

Ключевое слово `final`, более ограниченная версия механизма `const` из C++, поддерживает построение неизменяемых объектов. Поля типа `final` не могут быть изменены (хотя объекты, на которые они ссылаются, могут быть изменены, если они изменяемы), но они также имеют специальную семантику в модели памяти Java. Именно использование полей типа `final` делает возможной гарантию *безопасной инициализации* (см. раздел [3.5.2](#)), которая позволяет получать свободный доступ к неизменяемым объектам без синхронизации.

Даже если объект является изменяемым, пометка некоторых полей как `final` может упростить рассуждения о его состоянии, поскольку ограничение изменяемости объекта ограничивает набор возможных состояний объекта. Объект, который является “в основном неизменяемым”, но имеет одну или две изменяемые переменные состояния, по-прежнему проще в понимании, чем тот, который имеет множество изменяемых переменных. Объявление полей `final` также документирует для сопровождающих, что эти поля не должны изменяться.

Хорошая практика заключается в том, чтобы делать все поля `private`, если они не нуждаются в большей видимости [EJ Item 12], также хорошей практикой является сделать все поля `final`, если они не должны быть изменяемыми.

3.4.2 Использование `volatile` для публикации неизменяемых объектов

В классе `UnsafeCachingFactorizer` мы попытались использовать два объекта `AtomicReferences` для хранения последнего числа и последнего фактора, но это было не потокобезопасно, потому что мы не могли получить или обновить оба связанных значения атомарно. Использование `volatile` переменных для этих значений не будет потокобезопасным по той же причине. Однако неизменяемые объекты могут иногда обеспечивать слабую форму атомарности.

Сервлет факторинга выполняет две операции, которые должны быть атомарными: обновление кэшированного результата и условная выборка кэшированных факторов, если кэшированный номер соответствует запрошенному номеру. Всякий раз, когда группа связанных элементов данных должна действовать атомарно, рассмотрите возможность создания для них неизменяемого класса-холдера, такого как `OneValueCache`⁴⁶ в листинге 3.12.

Условия гонки при доступе или обновлении нескольких связанных переменных можно устранить, используя неизменяемый объект для хранения всех переменных.

⁴⁵ Многие разработчики опасаются, что такой подход вызовет проблемы с производительностью, но эти опасения обычно необоснованы. Выделение памяти дешевле, чем вы могли бы подумать, а неизменяемые объекты предлагают дополнительные преимущества производительности, такие как сокращение потребности в блокировках или защитных копиях и снижение нагрузки на сборщик мусора.

⁴⁶ Класс `OneValueCache` не был бы неизменяемым без вызова метода `copyOf` в конструкторе и геттере. Метод `Arrays.copyOf` был добавлен для удобства в Java 6; метод `clone` также будет работать.

```
@Immutable
class OneValueCache {

    private final BigInteger lastNumber;
    private final BigInteger[] lastFactors;

    public OneValueCache(BigInteger i,
                         BigInteger[] factors) {
        lastNumber = i;
        lastFactors = Arrays.copyOf(factors, factors.length);
    }

    public BigInteger[] getFactors(BigInteger i) {
        if (lastNumber == null || !lastNumber.equals(i))
            return null;
        else
            return Arrays.copyOf(lastFactors, lastFactors.length);
    }
}
```

Листинг 3.12 Неизменяемый холдер, кэширующий число и его фактор.

С изменяемым объектом `holder` вам придется использовать блокировку для обеспечения атомарности; с неизменяемым, как только поток получает ссылку на него, ему не нужно беспокоиться о другом потоке, изменяющем его состояние. Если переменные должны быть обновлены, создается новый объект `holder`, но все потоки, работающие с предыдущим объектом `holder`, по-прежнему видят его в согласованном состоянии.

Класс `VolatileCachedFactorizer` в листинге 3.13 использует класс `OneValueCache` для хранения кэшированного числа и факторов. Когда поток устанавливает `volatile` полю `cache` ссылку на новый экземпляр класса `OneValueCache`, новые кэшированные данные сразу становятся видимыми другим потокам.

Операции, связанные с кэшем, не могут оказывать влияние друг на друга, поскольку класс `OneValueCache` неизменяемый, и поле `cache` доступно только единожды в каждой из соответствующих веток кода. Эта комбинация неизменяемого объекта-холдера для нескольких переменных состояния, связанных с инвариантом, и ссылки на поле типа `volatile`, используемой для гарантии своевременной видимости, позволяет классу `VolatileCachedFactorizer` быть потокобезопасным, даже если он не содержит явных блокировок.

3.5 Безопасная публикация

До сих пор мы были сосредоточены на обеспечении того, чтобы объект *не* публиковался, предполагая, что он должен ограничиваться потоком или внутренним состоянием другого объекта. Конечно, иногда мы хотим совместно использовать объекты в разных потоках, и в этом случае мы должны делать это безопасно. К сожалению, простое сохранение ссылки на объект в поле типа `public`, как в листинге 3.14, является *не* достаточным условием для безопасной публикации объекта.

```
@ThreadSafe
public class VolatileCachedFactorizer implements Servlet {
    private volatile OneValueCache cache = new
        OneValueCache(null, null);

    public void service(ServletRequest req, ServletResponse resp) {
        BigInteger i = extractFromRequest(req);
        BigInteger[] factors = cache.getFactors(i);
        if (factors == null) {
            factors = factor(i);
            cache = new OneValueCache(i, factors);
        }
        encodeIntoResponse(resp, factors);
    }
}
```

Листинг 3.13 Кэширование результата выполнения последней операции с использованием `volatile` ссылки на неизменяемый объект холдера.

```
// Unsafe publication
public Holder holder;

public void initialize() {
    holder = new Holder(42);
}
```



Листинг 3.14 Публикация объекта без адекватной синхронизации. Не делайте так.

Вы можете удивиться тому, к насколько серьёзным последствиям может привести крах этого безобидно выглядящего примера. Из-за проблем с видимостью класс `Holder` может оказаться в другом потоке в несогласованном состоянии, несмотря на то, что его инварианты были правильно установлены его конструктором! Некорректная публикация может привести к тому, что другие потоки смогут наблюдать *частично сконструированный объект*.

3.5.1 Некорректная публикация: когда хорошие объекты ведут себя плохо

Вы не можете полагаться на целостность частично построенных объектов. Наблюдающий поток мог видеть объект в несогласованном состоянии, а затем увидеть, что его состояние внезапно изменилось, хотя оно не изменялось после публикации. Фактически, если класс `Holder` в листинге 3.15 опубликован с использованием небезопасной идиомы публикации из листинга 3.14, а поток, отличный от публикующего потока, должен был вызвать метод `assertSanity`, могло быть порождено исключение `AssertionError`⁴⁷

⁴⁷ Проблема здесь не в классе `Holder` как таковом, а в том, что класс `Holder` не опубликован должным образом. Однако класс `Holder` может быть защищен от некорректной публикации, путём объявления поля `n` как `final`, что сделает владельца неизменяемым; см. раздел 3.5.2.

```
public class Holder {  
    private int n;  
  
    public Holder(int n) { this.n = n; }  
  
    public void assertSanity() {  
        if (n != n)  
            throw new AssertionError("This statement is false.");  
    }  
}
```



Листинг 3.15 Класс с заложенной возможностью краха из-за некорректной публикации.

Поскольку синхронизация не использовалась, чтобы сделать класс `Holder` видимым для других потоков, мы говорим, что класс `Holder` был *опубликован некорректно*. Две вещи могут пойти не так, как ожидалось, с некорректно опубликованными объектами. Другие потоки могли видеть устаревшее значение поля типа `holder` и, таким же образом, видеть нулевую ссылку или другое более старое значение, даже если значение было помещено в `holder`. Но, что гораздо хуже, другие потоки могли увидеть текущую ссылку на класс `holder`, но устаревшее значение состояния класса `holder`⁴⁸. Чтобы сделать вещи еще менее предсказуемыми, поток может увидеть устаревшее значение при первом чтении поля, а затем более актуальное значение в следующий раз, поэтому метод `assertSanity` может бросить исключение `AssertionError`.

Рискуя повторить ранее сказанное, отметим, что могут происходить очень странные вещи, когда данные разделяются между потоками без достаточной синхронизации.

3.5.2 Неизменяемые объекты и безопасность инициализации

Поскольку неизменяемые объекты настолько важны, модель памяти Java предлагает специальную гарантию *безопасности инициализации* для совместного использования неизменяемых объектов. Как мы видели ранее, то, что ссылка на объект становится видимой для другого потока, не обязательно означает, что состояние этого объекта видно потребляющему потоку. Чтобы гарантировать согласованное представление состояния объекта, необходима синхронизация.

С другой стороны, с неизменяемыми объектами можно безопасно обращаться, *даже если синхронизация не используется для публикации ссылки на объект*. Для обеспечения гарантии безопасности инициализации должны соблюдаться все требования к неизменяемости: не модифицируемое состояние, все поля объявлены как `final` и корректное построение. (Если бы класс `Holder` в листинге 3.15 был неизменяемым, метод `assertSanity` не смог бы возбудить исключение `AssertionError`, даже если класс `Holder` не был опубликован должным образом.)

Неизменяемые объекты могут безопасно использоваться любым потоком без дополнительной синхронизации, даже если синхронизация не используется для их публикации.

⁴⁸ Хотя может показаться, что значения полей, заданные в конструкторе, являются первыми значениями, записываемыми в эти поля, и поэтому нет «более старых» значений, называемых устаревшими, конструктор класса `Object`, перед запуском конструкторов подкласса, сперва записывает для всех полей значения по умолчанию. Следовательно, существует вероятность увидеть значение по умолчанию установленное полю, как устаревшее значение.

Эта гарантия распространяется на значения всех конечных полей правильно построенных объектов; конечные поля могут быть безопасно доступны без дополнительной синхронизации. Однако, если поля типа `final` ссылаются на изменяемые объекты, синхронизация по-прежнему требуется для доступа к состоянию объектов, на которые они ссылаются.

3.5.3 Идиомы безопасной публикации

Объекты, которые не являются неизменяемыми, должны быть безопасно опубликованы, что обычно влечет за собой синхронизацию как публикующего, так потребляющего потоков. На данный момент давайте сосредоточимся на том, что потребляющий поток может видеть объект в его опубликованном состоянии; мы будем иметь дело с видимостью изменений, сделанных в ближайшее, после публикации, время.

Для безопасной публикации объекта, ссылка на объект и его состояние должны быть одновременно видны другим потокам. Правильно построенный объект может быть безопасно опубликован с помощью:

- Инициализации ссылки на объект из статического инициализатора;
- Хранения ссылки на него в поле типа `volatile` или `AtomicReference`;
- Сохранения ссылки на него в поле типа `final` корректно построенного объекта;
- Хранения ссылки на него в поле, должным образом защищенным блокировкой.

Внутренняя синхронизация в потокобезопасных коллекциях означает, что размещение объекта в потокобезопасной коллекции, такой как `Vector` или `synchronizedList`, удовлетворяет последнему из этих требований. Если поток *A* помещает объект *X* в потокобезопасную коллекцию и поток *B* впоследствии извлекает его, поток *B* гарантированно видит состояние объекта *X* таким, каким поток *A* оставил его, даже если код приложения, который передает объект *X* подобным образом, не имеет явной синхронизации. Потокобезопасные библиотеки коллекций предлагают следующие гарантии безопасной публикации, даже если в Javadoc об этом почти ничего не сказано:

- Размещение ключа или значения в классах `Hashtable`, `synchronizedMap` или `ConcurrentMap` безопасно публикует его в любом потоке, который извлекает его из класса `Map` (напрямую или через итератор);
- Размещение элемента в классах `Vector`, `CopyOnWriteArrayList`, `CopyOnWriteArraySet`, `synchronizedList`, или `synchronizedSet` безопасно публикует его в любой поток, который извлекает его из коллекции;
- Размещение элемента в `BlockingQueue` или `ConcurrentLinkedQueue` безопасно публикует его в любом потоке, который извлекает его из очереди.

Другие механизмы передачи состояния в библиотеке классов (такие как `Future` и `Exchanger`) также представляют собой безопасную публикацию; мы определим их как обеспечивающие безопасную публикацию по мере их представления.

Использование статического инициализатора часто является самым простым и безопасным способом публикации объектов, которые могут быть построены статически:

```
public static Holder holder = new Holder(42);
```

Статические инициализаторы выполняются JVM во время инициализации класса; из-за внутренней синхронизации в JVM, этот механизм гарантирует безопасную публикацию любых объектов, инициализированных таким образом [JLS 12.4.2].

3.5.4 Фактически неизменяемые объекты

Безопасная публикация является достаточным условием для безопасного доступа других потоков к объектам, которые не будут изменяться после публикации, без использования дополнительной синхронизации. Механизмы безопасной публикации гарантируют, что опубликованное состояние объекта станет видимым для всех потоков, получающих доступ к нему, как только станет видна ссылка на него, и если состояние объекта впредь не будет изменяться, это будет достаточным условием для гарантии того, что доступ к нему будет безопасен.

Объекты, технически не являющиеся неизменяемыми, состояние которых не будет изменяться после публикации, называются *фактически неизменяемыми*. Они не должны следовать строгим определениям неизменяемости из раздела 3.4; они просто должны рассматриваться программой так, как если бы они были неизменяемыми после их публикации. Использование *фактически неизменяемых объектов* может упростить разработку и повысить производительность за счет снижения расходов на синхронизацию.

Безопасно опубликованные фактически неизменяемые объекты могут безопасно использоваться любым потоком без дополнительной синхронизации.

Например, объект Date изменяемый⁴⁹, но если вы используете его так, как если бы он был неизменяемым, вы можете исключить блокировку, которая в противном случае потребовалась бы при совместном использовании несколькими потоками объекта Date. Предположим, вы хотите сохранить объект Map, в котором содержится время последнего входа каждого пользователя:

```
public Map<String, Date> lastLogin =
    Collections.synchronizedMap(new HashMap<String, Date>());
```

Если значения типа Date не изменяются после их размещения в объекте класса Map, то синхронизации реализованной в классе synchronizedMap достаточно для безопасной публикации значений типа Date, и при доступе к ним дополнительная синхронизация не требуется.

3.5.5 Изменяемые объекты

Если объект может быть изменен после построения, безопасная публикация обеспечивает только видимость опубликованного состояния. Синхронизация должна использоваться не только для публикации изменяемого объекта, но и при

⁴⁹ Вероятно, была допущена ошибка в дизайне библиотеки классов.

каждом обращении к объекту, чтобы обеспечить видимость последующих изменений. Для обеспечения безопасного совместного использования изменяемых объектов, они должны быть безопасно опубликованы *и* быть потокобезопасными или защищёнными блокировкой.

Требования к публикации объекта зависят от его изменяемости::

- *Неизменяемые объекты* могут быть опубликованы с использованием любых механизмов;
- *Фактически неизменяемые* объекты должны быть безопасно опубликованы;
- *Изменяемые объекты должны быть безопасно опубликованы и, также, должны быть потокобезопасными или защищёнными блокировкой.*

3.5.6 Безопасное совместное использование объектов

Всякий раз, когда вы приобретаете ссылку на объект, вы должны знать, что вам разрешено делать с ним. Вам необходимо захватить блокировку перед его использованием? Можете ли вы изменить его состояние или только прочитать его? Множество ошибок параллелизма берёт своё начала с неспособности понять «правила взаимодействия» с совместно используемыми объектами. Когда вы публикуете объект, вы должны документировать, как можно получить к нему доступ

Наиболее полезными политиками для использования и предоставления совместного доступа к объектам в параллельной программе являются:

Ограничение потока. Объект, ограниченный потоком, принадлежит исключительно одному потоку и ограничен им, и может быть изменен только собственным потоком.

Совместный доступ “только на чтение”. Объект, с правом совместного доступа “только на чтение”, может быть доступен одновременно нескольким потокам без дополнительной синхронизации, но не может быть изменен ни одним потоком. Совместно используемые объекты с правом “только на чтение” включают неизменяемые и фактически неизменяемые объекты.

Совместное использование потокобезопасных объектов.
Потокобезопасный объект выполняет синхронизацию внутри себя, поэтому несколько потоков могут свободно обращаться к нему через его открытые интерфейсы, без необходимости, в дальнейшем, применять синхронизацию.

Защищённые объекты. К защищённому объекту можно получить доступ только с определенной блокировкой. К защищённым объектам относятся объекты, инкапсульрованные в другие потокобезопасные объекты и опубликованные объекты, которые, как известно, защищены определенной блокировкой.

Глава 4 Компоновка объектов

До сих пор мы рассмотрели низкоуровневые основы потокобезопасности и синхронизации. Но мы не хотим анализировать каждое обращение к памяти, чтобы гарантировать, что наша программа является потокобезопасной; мы хотим иметь возможность создавать потокобезопасные компоненты и безопасно составлять из них более крупные компоненты или программы. В этой главе рассматриваются шаблоны для структурирования классов, которые могут облегчить обеспечение их потокобезопасности и поддержки, без риска случайного нарушения гарантий безопасности.

4.1 Проектирование потокобезопасных классов

Хотя можно написать потокобезопасную программу, которая хранит все свое состояние в общедоступных статических полях, гораздо сложнее проверить ее на потокобезопасность или изменить ее так, чтобы она оставалась потокобезопасной, в противоположность той, что использует инкапсуляцию соответствующим образом. Инкапсуляция позволяет прийти к выводу о потокобезопасности класса без необходимости изучения всей программы.

Процесс проектирования потокобезопасного класса должен включать три базовых элемента:

- Определение переменных, формирующих состояние объекта;
- Определение инвариантов, ограничивающих переменные состояния;
- Установление политики для управления параллельным доступом к состоянию объекта.

Определение состояния объекта начинается с его полей. Если все они примитивных типов, поля содержат в себе всё состояние. Класс `Counter` в листинге 4.1 имеет только одно поле, поэтому в поле `value` содержится всё его состояние. Состояние объекта с n примитивными полями - это всего лишь n -кортеж⁵⁰ значений его полей; состояние класса `2D Point` определяется значением (x, y) . Если у объекта имеются поля, содержащие ссылки на другие объекты, его состояние также будет включать поля из объектов, на которые имеются ссылки. Например, состояние класса `LinkedList` включает в себя состояние всех ссылок на объекты узлов, принадлежащих списку.

```
@ThreadSafe
public final class Counter {
    @GuardedBy("this")
    private long value = 0;

    public synchronized long getValue() {
        return value;
    }
}
```

⁵⁰ Кортеж представляет собой вектор или одномерный массив значений.

```
public synchronized long increment() {
    if (value == Long.MAX_VALUE)
        throw new IllegalStateException("counter overflow");
    return ++value;
}
```

Листинг 4.1 Простой потокобезопасный счётчик, использующий Java шаблон “монитор”

Политика синхронизации определяет, каким образом объект координирует доступ к своему состоянию, для того, чтобы избежать нарушения своих инвариантов или постусловий. Она определяет, какая комбинация из неизменяемости, ограничения потока и блокировки используется для обеспечения потокобезопасности и какие переменные защищены какими блокировками. Чтобы убедиться, что класс можно анализировать и поддерживать, документируйте политику синхронизации.

4.1.1 Сбор требований к синхронизации

Сделать класс потокобезопасным – это значит гарантировать, что его инварианты будут “удержаны”⁵¹ при параллельном доступе; это требует рассуждений о его состоянии. Объекты и переменные имеют *пространство состояний*: диапазон возможных состояний, которые они могут принимать. Чем меньше это пространство состояний, тем легче о нём рассуждать. Используя поля типа `final` везде, где это практически возможно, вы упрощаете анализ возможных состояний объекта. (В крайнем случае, неизменяемые объекты могут находиться только в одном состоянии.)

Многие классы имеют инварианты, которые определяют некоторые состояния как *допустимые* (*valid*) или *недопустимые* (*invalid*). Поле `value` в классе `Counter` типа `long`. Пространство состояний типа `long` составляет диапазон от `Long.MIN_VALUE` до `Long.MAX_VALUE`, но класс `Counter` накладывает собственное ограничение на поле `value`; отрицательные значения не допустимы.

Аналогичным образом операции могут иметь постусловия, которые идентифицируют определенные *переходы состояний* как недопустимые. Если текущее состояние объекта `Counter` равно 17, *единственным допустимым* следующим значением состояния является 18. Когда следующее состояние выводится из текущего состояния, операция обязательно является составным действием. Не все операции налагают ограничения на переходы состояний; при обновлении переменной, которая содержит текущую температуру, ее предыдущее состояние не оказывает влияние на вычисления.

Ограничения, накладываемые на состояния или переходы состояний инвариантов и постусловий, создают дополнительные требования к использованию синхронизации или инкапсуляции. Если некоторые состояния недопустимы, то базовые переменные состояния должны быть инкапсулированы, иначе клиентский код может перевести объект в недопустимое состояние. Если операция имеет недопустимые переходы состояний, она должна быть атомарной. С другой стороны, если класс не накладывает таких ограничений, мы можем ослабить требования к инкапсуляции или сериализации, чтобы получить большую гибкость или лучшую производительность.

⁵¹ Состояние инвариантов не будет изменяться при параллельном доступе.

Класс также может иметь инварианты, ограничивающие несколько переменных состояния. Класс с числовым диапазоном, например, такой как `NumberRange` в листинге 4.10, обычно поддерживает сохранение нижней и верхней границ диапазона в переменных состояния. Эти переменные должны подчиняться ограничению, заключающемуся в том, что нижняя граница диапазона была меньше или равна верхней границе. Многомерные инварианты, подобные этому, создают требования к атомарности: связанные переменные должны извлекаться или обновляться в одной атомарной операции. Вы не можете обновить переменную, снять и повторно захватить блокировку, а затем обновить другие переменные, так как это может повлечь за собой оставление объекта в недопустимом состоянии, когда блокировка будет снята. Когда в инварианте участвует несколько переменных, блокировка, которая их защищает, должна удерживаться в процессе выполнения любой операции, обращающейся к связанным переменным.

Невозможно обеспечить потокобезопасность без понимания инвариантов и постусловий объекта. Ограничение допустимых значений или переходов состояний для переменных состояний может потребовать применения атомарности и инкапсуляции.

4.1.2 Операции, зависящие от состояния

Инварианты классов и постусловия методов ограничивают допустимые состояния и переходы состояний объекта. Некоторые объекты также имеют методы с предварительными условиями на основе состояния. Например, нельзя удалить элемент из пустой очереди; очередь должна находиться в состоянии “не пусто”, прежде чем можно будет удалить элемент. Операции с предварительными условиями на основе состояния называются *зависимыми от состояния* [CPJ 3].

В однопоточной программе, если предварительное условие не выполняется, операция не имеет выбора, кроме как завершиться с ошибкой. Но в параллельной программе предварительное условие может быть выполнено позже из-за действия другого потока. Параллельные программы добавляют возможность ожидания до тех пор, пока условие не станет истинным, а затем продолжат операцию.

Встроенные механизмы эффективного ожидания выполнения условия - `wait` и `notify` - тесно связаны со встроенными блокировками и могут быть сложны в правильном использовании. Для создания операций, ожидающих выполнения предварительных условий перед выполнением, часто проще использовать существующие классы библиотек, такие как блокирующие очереди или семафоры, чтобы обеспечить требуемое поведение, зависящее от состояния. Блокирующие библиотечные классы, такие как `BlockingQueue`, `Semaphore` и другие *синхронизаторы*, рассматриваются в главе 5; создание зависимых от состояния классов с использованием низкоуровневых механизмов, предоставляемых платформой и библиотекой классов, рассматривается в главе 14.

4.1.3 Владелец состояния

В разделе 4.1 мы подразумевали, что состояние объекта может быть подмножеством полей в графе объектов, внедренных в этот объект. Почему это может быть подмножество? При каких условиях поля, достижимые из данного объекта, не являются частью состояния этого объекта?

Определяя, какие переменные формируют состояние объекта, мы хотим рассматривать только те данные, которыми объект *владеет*. Владелец явно в языке не реализован, но является элементом дизайна класса. Если вы размещаете и заполняете объект `HashMap`, вы создаете несколько объектов: объект `HashMap`, несколько объектов `Map.Entry`, используемых реализацией `HashMap`, и, возможно, другие внутренние объекты. Логическое состояние объекта `HashMap` включает состояние всех объектов `Map.Entry` и внутренних объектов, даже если они реализованы как отдельные объекты.

К лучшему или к худшему, сборка мусора позволяет нам не думать о владельце. При передаче объекта методу в C++, необходимо достаточно тщательно обдумывать вопрос о том, передаете ли вы право собственности, занимаетесь ли вы краткосрочным кредитом или предполагаете долгосрочное совместное владение. В языке Java возможны все те же модели владения, но сборщик мусора снижает стоимость многих распространенных ошибок при совместном использовании ссылок, позволяя размышлять о владении без учёта мелких деталей.

Во многих случаях владение и инкапсуляция следуют рука об руку - объект инкапсулирует состояние, которым он владеет, и владеет состоянием, которое он инкапсулирует. Владелец переменной состояния принимает решение о протоколе блокировки, используемом для поддержания целостности состояния этой переменной. Владение подразумевает контроль, но как только вы опубликуете ссылку на изменяемый объект, у вас больше не будет эксклюзивного права контроля; в лучшем случае у вас может быть "совместное владение". Класс обычно не владеет объектами, переданными его методам или конструкторам, если только этот метод не предназначен для явного переноса права собственности на переданные объекты (например, фабричный метод обертки синхронизированной коллекции).

Классы коллекций часто демонстрируют форму "раздельного владения", в которой коллекции принадлежит состояние инфраструктуры коллекции, а клиентскому коду принадлежат объекты, хранящиеся в коллекции. Например, класс `ServletContext` из фреймворка сервлетов. Класс `ServletContext` предоставляет для сервлетов Мар-подобный контейнер объектов, с помощью которого они могут регистрировать и извлекать объекты приложения по имени, используя методы `setAttribute` и `getAttribute`. Объект `ServletContext`, реализованный контейнером сервлета, должен быть потокобезопасным, потому что к нему обязательно будут обращаться несколько потоков. Сервлетам нет необходимости использовать синхронизацию при вызове методов `setAttribute` и `getAttribute`, но им, возможно, придется использовать синхронизацию в моменты использования объектов, хранящихся в `ServletContext`. Эти объекты принадлежат приложению; они хранятся, для безопасного хранения, контейнером сервлетов от имени приложения. Как и прочие совместно используемые объекты, они должны использоваться безопасно; чтобы не допустить взаимного влияния от нескольких потоков обращающихся к одному и тому же объекту одновременно, они должны быть потокобезопасными, фактически неизменяемым или использовать явную блокировку⁵².

⁵² Интересно, что объект `HttpSession`, выполняющий аналогичную функцию в рамках сервлета, может иметь более строгие требования. Поскольку контейнер сервлета может обращаться к объектам в `httpsession`, чтобы их можно было сериализовать для репликации или пассивации, они должны быть потокобезопасными, так как контейнер будет обращаться к ним, а также к веб-приложению. (Мы говорим "может иметь", так как репликация и пассивация вне спецификации сервлета, но является общей особенностью контейнеров сервлета.)

4.2 Ограничение экземпляра

Если объект не является потокобезопасным, несколько подходов могут позволить безопасно использовать его в многопоточной программе. Вы можете гарантировать, что доступ к объекту будет осуществляться только из одного потока (ограничение потока) или что весь доступ к объекту будет должным образом защищен блокировкой.

Инкапсуляция упрощает создание потокобезопасных классов, способствуя ограничению экземпляра; часто это называется просто *ограничением* [CPJ 2.3.3]. Когда объект инкапсулируется в другой объект, все ветки кода, имеющие доступ к инкапсулированному объекту, известны и поэтому могут быть проанализированы значительно легче, чем в том случае, если бы этот объект был доступен из всех частей программы. Сочетание ограничения с соответствующей правилам блокировкой гарантирует, что не потокобезопасные объекты используются потокобезопасным способом.

Инкапсуляция данных в объекте ограничивает доступ к данным методами объекта, что упрощает обеспечение гарантии того, что доступ к данным будет всегда осуществляться с удержанием соответствующей блокировки.

Ограниченные объекты не должны покидать заданной области видимости. Объект может быть ограничен экземпляром класса (например, закрытым членом класса), лексической областью действия (например, локальной переменной) или потоком (например, объектом, который передается из метода в метод в потоке, но он не должен совместно использоваться другими потоками). Конечно, объекты не могут сбегать сами по себе - им нужна помощь разработчика, который может способствовать, опубликовав объект за пределами его области видимости.

Класс PersonSet в листинге 4.2 иллюстрирует, как ограничение и блокировка могут работать совместно, чтобы сделать класс потокобезопасным, даже если переменные состояния компонента таковыми не являются. Состояние класса PersonSet управляет классом HashSet, который не является потокобезопасным. Но, так как переменная mySet является приватной и не может сбежать, класс HashSet ограничен классом PersonSet. Единственные ветки кода, которые могут получить доступ к переменной mySet, это методы addPerson и containsPerson, и каждая из них захватывает блокировку на экземпляре PersonSet. Всё состояние класса защищено внутренней блокировкой, что делает класс PersonSet потокобезопасным.

```
@ThreadSafe
public class PersonSet {
    @GuardedBy("this")
    private final Set<Person> mySet = new HashSet<Person>();

    public synchronized void addPerson(Person p) {
        mySet.add(p);
    }

    public synchronized boolean containsPerson(Person p) {
        return mySet.contains(p);
    }
}
```

```
    }  
}
```

Листинг 4.2 Использование ограничения для обеспечения потокобезопасности

В этом примере не делается никаких предположений о потокобезопасности класса `Person`, но если класс изменяемый, при обращении к классу `Person`, полученному из `PersonSet`, будет необходимо использовать дополнительную синхронизацию. Самый надежный способ достичь этого – сделать класс `Person` потокобезопасным; менее надежным было бы защитить объекты `Person` с помощью блокировки и гарантировать, чтобы все клиенты будут следовать протоколу захвата соответствующей блокировки до осуществления доступа к экземпляру `Person`.

Ограничение экземпляра – один из простейших способов создания потокобезопасных классов. Оно позволяет проявить гибкость в выборе стратегии блокировки; класс `PersonSet` полагается на свою собственную внутреннюю блокировку, чтобы защитить свое состояние, но любая блокировка, используемая последовательно, будет работать же хорошо. Ограничение экземпляра также позволяет различным переменным состояния защищаться различными блокировками. (Пример класса, использующего несколько объектов блокировки для защиты своего состояния, см. в главе 11, в классе `ServerStatus`).

Существует множество примеров ограничений в библиотеках классов платформы, включая некоторые классы, которые существуют исключительно для того, чтобы превращать не потокобезопасные классы в потокобезопасные. Базовые классы коллекций, такие как `ArrayList` и `HashMap`, не являются потокобезопасными, но библиотека классов предоставляет фабричные методы-обертки (`Collections.synchronizedList` и другие), чтобы их можно было безопасно использовать в многопоточном окружении. Эти фабрики используют шаблон *Декоратор* (Гамма и другие, 1995)⁵³, чтобы обернуть коллекцию в синхронизированный объект-обёртку; обёртка реализует каждый метод соответствующего интерфейса как синхронизированный метод, который перенаправляет запрос базовому объекту коллекции. До тех пор, пока объект-обертка содержит единственную доступную ссылку на базовую коллекцию (то есть базовая коллекция ограничена обёрткой), объект-обертка также является потокобезопасным. В Javadoc для таких методов указано предупреждение, что весь доступ к базовой коллекции должен выполняться через объект-обёртку.

Конечно, по-прежнему возможно нарушить ограничение, опубликовав предположительно ограниченный объект; если объект должен быть ограничен определенной областью видимости, то позволение покинуть эту область видимости является ошибкой. Ограниченные объекты также могут сбегать после публикации другими объектами, такими как итераторы или внутренние экземпляры классов, которые могут косвенно публиковать ограниченные объекты.

Ограничение упрощает создание потокобезопасных классов, поскольку класс, ограничивающий своё состояние, может быть проанализирован на потокобезопасность без необходимости проверки всей программы.

⁵³ Речь о книге “Шаблоны проектирования” банды четырёх.

4.2.1 Шаблон Java “Монитор”

Следуя от принципа ограничения экземпляра к его логическому завершению, вы придёте к шаблону Java “монитор”⁵⁴. Объект, следующий шаблону Java “монитор”, инкапсулирует все своё изменяемое состояние и защищает его с помощью собственной внутренней блокировки.

Класс Counter в листинге 4.1 демонстрирует типичный пример использования этого шаблона. Он инкапсулирует одну переменную состояния, value, и весь доступ к этой переменной состояния осуществляется через методы класса Counter, которые являются синхронизированными.

Шаблон Java “монитор” используется многими библиотечными классами, такими как Vector и Hashtable. Иногда требуется более тонкая политика синхронизации; в главе 11 показано, как повысить масштабируемость за счет более утончённых стратегий блокировки. Основным преимуществом шаблона Java “монитор” является его простота.

Шаблон Java “монитор” - это просто соглашение; любой объект блокировки может использоваться для защиты состояния объекта, если он используется согласованно. В Листинге 4.3 иллюстрируется класс, который использует приватную блокировку для защиты своего состояния.

```
public class PrivateLock {  
    private final Object myLock = new Object();  
    @GuardedBy("myLock") Widget widget;  
  
    void someMethod() {  
        synchronized(myLock) {  
            // Access or modify the state of widget  
        }  
    }  
}
```

Листинг 4.3 Защита состояния с помощью приватной блокировки

Существуют преимущества использования приватного объекта блокировки вместо встроенной блокировки объекта (или любой другой общедоступной блокировки). При создании объекта блокировки приватным, блокировка инкапсулируется таким образом, что клиентский код не может ее получить, в то время как общедоступная блокировка позволяет клиентскому коду участвовать в политике синхронизации - правильно или неправильно. Клиенты, которые получают блокировку другого объекта некорректно, могут вызвать проблемы живучести, и проверка того, что общедоступная блокировка используется правильно, требует проверки всей программы, а не одного класса.

4.2.2 Пример: отслеживание парка автомобилей

Класс Counter в листинге 4.1 является кратким, но тривиальным примером использования шаблона Java “монитор”. Давайте создадим несколько менее

⁵⁴ Создание шаблона “монитор” было вдохновлено работой Хоара о мониторах (Hoare, 1974), хотя между этим шаблоном и истинным монитором существуют значительные различия. Даже инструкции байт-кода для входа и выхода из синхронизированного блока называются monitorenter и monitorexit, а встроенные (внутренние) блокировки Java иногда называют блокировкой монитора или монитором.

тривиальный пример: “трекер автомобилей” для управления парком транспортных средств, таких как такси, полицейские машины или грузовые автомобили. Сначала мы создадим его с помощью шаблона “монитор”, а затем рассмотрим, как ослабить некоторые требования к инкапсуляции, сохранив при этом потокобезопасность.

Каждое транспортное средство идентифицируется строкой и имеет местоположение, представленное координатами (x, y). Класс `VehicleTracker` инкапсулирует идентификационные данные и местоположения известных транспортных средств, что делает его хорошо подходящим для использования в качестве модели данных в шаблоне модель-представление-контроллер в GUI-приложении, где он может совместно использоваться потоком представления и несколькими потоками, обновляющими данные. Поток представления будет получать имена и местоположения транспортных средств, и отображать их на дисплее:

```
Map<String, Point> locations = vehicles.getLocations();
for (String key : locations.keySet())
    renderVehicle(key, locations.get(key));
```

Точно так же, обновляющие данные потоки будут изменять местоположения транспортных средств, внося данные полученные от устройств GPS, или введённые вручную диспетчером через графический интерфейс:

```
void vehicleMoved(VehicleMovedEvent evt) {
    Point loc = evt.getNewLocation();
    vehicles.setLocation(evt.getVehicleId(), loc.x, loc.y);
}
```

Поскольку поток представления и потоки обновлений будут обращаться к модели данных одновременно, она должна быть потокобезопасной. В листинге 4.4 показана реализация трекера транспортных средств с использованием шаблона Java “монитор”, который использует класс `MutablePoint` из листинга 4.5 для представления местоположения транспортных средств.

```
@ThreadSafe
public class MonitorVehicleTracker {
    @GuardedBy("this")
    private final Map<String, MutablePoint> locations;

    public MonitorVehicleTracker(
        Map<String, MutablePoint> locations) {
        this.locations = deepCopy(locations);
    }

    public synchronized Map<String, MutablePoint> getLocations() {
        return deepCopy(locations);
    }

    public synchronized MutablePoint getLocation(String id) {
        MutablePoint loc = locations.get(id);
        return loc == null ? null : new MutablePoint(loc);
```

```

}

public synchronized void setLocation(String id, int x, int y) {
    MutablePoint loc = locations.get(id);
    if (loc == null)
        throw new IllegalArgumentException("No such ID: " + id);
    loc.x = x;
    loc.y = y;
}

private static Map<String, MutablePoint> deepCopy(
    Map<String, MutablePoint> m) {
    Map<String, MutablePoint> result =
        new HashMap<String, MutablePoint>();
    for (String id : m.keySet())
        result.put(id, new MutablePoint(m.get(id)));
    return Collections.unmodifiableMap(result);
}
}

public class MutablePoint { /* ЛИСТИНГ 4.5 */ }

```

Листинг 4.4 Реализация трекера транспортных средств на основе монитора

```

@NotThreadSafe
public class MutablePoint {
    public int x, y;

    public MutablePoint() { x = 0; y = 0; }
    public MutablePoint(MutablePoint p) {
        this.x = p.x;
        this.y = p.y;
    }
}

```



Листинг 4.5 Класс изменяемой точки похожий на java.awt.Point

Даже не смотря на то, что класс `MutablePoint` не является потокобезопасным, класс трекера таковым является. Ни объект `Map`, ни какие-либо изменяемые точки (`points`), которые он содержит, никогда не публикуются. Когда нам необходимо вернуть местоположение транспортных средств вызывающим объектам, соответствующие значения копируются с использованием копирующего конструктора класса `MutablePoint`, либо с помощью метода `deepCopy`, который создает новый объект `Map`, значения которого являются копиями ключей и значений старого объекта `Map`⁵⁵.

Эта реализация частично поддерживает потокобезопасность, копируя изменяемые данные перед их возвратом клиенту. Как правило, это не является

⁵⁵ Обратите внимание, что метод `deepCopy` не может просто обернуть класс `Map` с помощью метода `unmodifiableMap`, поскольку последний защищает от модификации только коллекцию; это не мешает вызывающим объектам модифицировать изменяемые объекты, хранящиеся в ней. По той же причине заполнение класса `HashMap` в методе `deepCopy` с помощью копирующего конструктора также не сработало бы, потому что были скопированы только ссылки на объекты `Point`, а не сами объекты.

проблемой в плане производительности, но может стать ей, если набор транспортных средств станет очень большим⁵⁶. Другим следствием копирования данных при каждом вызове метода `getLocation` является то, что содержимое возвращаемой коллекции не изменяется, даже если местоположения изменяются в базовой коллекции. Хорошо это или плохо зависит от ваших требований. Это может быть полезно в том случае, если существуют внутренние требования к согласованности в наборе местоположений, и в этом случае возврат согласованного снимка имеет решающее значение, или недостатком, если вызывающему объекту требуется актуальная информация по каждому транспортному средству и, следовательно, возникает необходимость более частого обновления снимка.

4.3 Делегирование потокобезопасности

Все, кроме самых тривиальных объектов, являются составными объектами. Шаблон Java “монитор” полезен при создании классов с нуля или составлении классов из объектов, которые не являются потокобезопасными. Но что, если компоненты нашего класса уже потокобезопасны? Нужно ли вводить дополнительный уровень для обеспечения потокобезопасности? Ответ... “это зависит от обстоятельств”. В некоторых случаях композит, состоящий из потокобезопасных компонентов, является потокобезопасным (Листинги 4.7 и 4.9), а в других - просто является хорошей основой (листинг 4.10).

В классе `CountingFactorizer` мы добавили переменную типа `AtomicLong` к объекту без сохранения состояния, и полученный составной объект по-прежнему оставался потокобезопасным. Поскольку состояние `CountingFactorizer` определяется состоянием потокобезопасной переменной типа `AtomicLong`, и класс `CountingFactorizer` не накладывает никаких дополнительных ограничений на состояние счетчика, легко увидеть, что класс `CountingFactorizer` является потокобезопасным.

Мы могли бы сказать, что класс `CountingFactorizer` делегирует ответственность за обеспечение потокобезопасности переменной типа `AtomicLong`: класс `CountingFactorizer` потокобезопасен, потому что потокобезопасен класс `AtomicLong`.⁵⁷

4.3.1 Пример: транспортный трекер с использованием делегирования

В качестве более существенного примера делегирования, давайте создадим версию трекера транспортных средств, делегирующую обеспечение безопасности потокобезопасному классу. Мы храним местоположения в объекте `Map`, поэтому начнем с реализации потокобезопасной реализации класса `Map` – класса `ConcurrentHashMap`. Мы также храним местоположение с помощью неизменяемого класса `Point` вместо `MutablePoint`, как показано в листинге 4.6.

⁵⁶ Поскольку метод `deepCopy` вызывается из синхронизированного метода, внутренняя блокировка трекера удерживается на протяжении всего времени выполнения длительной операции копирования, что может привести к ухудшению отзывчивости пользовательского интерфейса при отслеживании множества транспортных средств.

⁵⁷ Если поле `count` не объявлено как `final`, провести анализ класса `CountingFactorizer` на предмет потокобезопасности будет сложнее. Если бы класс `CountingFactorizer` мог изменить поле `count`, с целью назначить ссылку на другой объект `AtomicLong`, мы должны были бы гарантировать, что это обновление было бы видимо всем потокам, которые могли бы обратиться к полю `count`, а также гарантировать, чтобы не возникало никаких условий гонки относительно ссылки на поле `count`. Это еще одна веская причина для того, чтобы использовать `final` поля там, где это практически целесообразно.

```
@Immutable
public class Point {
    public final int x, y;

    public Point(int x, int y) {
        this.x = x;
        this.y = y;
    }
}
```

Листинг 4.6 Неизменяемый класс Point используемый в классе DelegatingVehicleTracker.

Класс Point потокобезопасен, потому что он неизменяемый. Неизменяемые значения могут свободно распространяться и публиковаться, поэтому нам больше не нужно копировать местоположения при их возврате. Класс DelegatingVehicleTracker в листинге 4.7 не использует явную синхронизацию; весь доступ к состоянию управляется классом ConcurrentHashMap, а все ключи и значения объекта Map неизменямы.

```
@ThreadSafe
public class DelegatingVehicleTracker {
    private final ConcurrentHashMap<String, Point> locations;
    private final Map<String, Point> unmodifiableMap;

    public DelegatingVehicleTracker(Map<String, Point> points) {
        locations = new ConcurrentHashMap<String, Point>(points);
        unmodifiableMap = Collections.unmodifiableMap(locations);
    }

    public Map<String, Point> getLocations() {
        return unmodifiableMap;
    }

    public Point getLocation(String id) {
        return locations.get(id);
    }

    public void setLocation(String id, int x, int y) {
        if (locations.replace(id, new Point(x, y)) == null)
            throw new IllegalArgumentException(
                "invalid vehicle name: " + id);
    }
}
```

Листинг 4.7 Делегирование обеспечения потокобезопасности классу ConcurrentHashMap.

Если бы мы использовали исходный класс MutablePoint вместо Point, мы бы нарушили инкапсуляцию, позволив методу getLocations опубликовать ссылку на изменяемое состояние, которое не является потокобезопасным. Обратите внимание, что мы немного изменили поведение класса трекера транспортных

средств; в то время как версия с монитором возвращала снимок местоположений, делегирующая версия возвращает неизменяемое, но “живое” представление местоположений транспортного средства. Это означает, что если поток *A* вызывает метод `getLocations` и поток *B* позже изменяет местоположение некоторых точек, изменения отражаются в объекте `Map`, возвращаемом потоку *A*. Как мы отмечали ранее, в зависимости от ваших требований это может рассматриваться и как преимущество (более актуальные данные), и как недостаток (потенциально несогласованное представление о парке автомобилей).

Если требуется неизменяемое представление парка транспортных средств, метод `getLocations` может вместо него вернуть небольшую копию объекта `Map` переменной `locations`. Поскольку содержимое объекта `Map` является неизменяемым, необходимо скопировать только структуру объекта `Map`, а не его содержимое, как показано в листинге 4.8 (в котором возвращается простой объект `HashMap`, так как метод `getLocations` не давал обещание возвращать потокобезопасный объект `Map`).

```
public Map<String, Point> getLocations() {
    return Collections.unmodifiableMap(
        new HashMap<String, Point>(locations));
}
```

Листинг 4.8 Возврат статической копии набора местоположений вместо “живой”

4.3.2 Независимые переменные состояния

Примеры делегирования до сих пор касались делегирования единственной, потокобезопасной переменной состояния. Мы также можем делегировать потокобезопасность более чем одной базовой (нижележащей) переменной состояния, если базовые переменные состояния *независимы*; это означает, что составной класс не навязывает никаких инвариантов, включающих несколько переменных состояния.

Класс `VisualComponent` из листинга 4.9 представляет собой графический компонент, который позволяет клиентам регистрировать слушателей для событий мыши и нажатия клавиш. Он поддерживает список зарегистрированных слушателей каждого типа, поэтому, когда происходит событие, могут быть вызваны соответствующие слушатели. Но нет никакой связи между набором слушателей мыши и слушателей клавиш; оба слушателя являются независимыми, и поэтому класс `VisualComponent` может делегировать свои обязательства по обеспечению потокобезопасности обоим нижележащим потокобезопасным спискам.

```
public class VisualComponent {
    private final List<KeyListener> keyListeners
        = new CopyOnWriteArrayList<KeyListener>();
    private final List<MouseListener> mouseListeners
        = new CopyOnWriteArrayList<MouseListener>();

    public void addKeyListener(KeyListener listener) {
        keyListeners.add(listener);
    }
```

```

public void addMouseListener(MouseListener listener) {
    mouseListeners.add(listener);
}

public void removeKeyListener(KeyListener listener) {
    keyListeners.remove(listener);
}

public void removeMouseListener(MouseListener listener) {
    mouseListeners.remove(listener);
}
}

```

Листинг 4.9 Делегирование потокобезопасности множеству нижележащих переменных состояния

Класс `VisualComponent` использует метод `CopyOnWriteArrayList` для хранения каждого списка слушателей; это потокобезопасная реализация интерфейса `List`, особенно подходящая для управления списками слушателей (см. раздел [5.2.3](#)). Каждый экземпляр интерфейса `List` является потокобезопасным, и поскольку нет никаких ограничений, связывающих состояние одного списка с состоянием другого, класс `VisualComponent` может делегировать свои обязанности по обеспечению потокобезопасности нижележащим объектам `mouseListeners` и `keyListeners`.

4.3.3 Когда делегирование не работает

Большинство составных классов не так просты, как класс `VisualComponent`: у них есть инварианты, связанные с переменными состояния компонента. Класс `NumberRange` из листинга 4.10 использует два объекта `AtomicInteger`s для управления своим состоянием, но накладывает дополнительное ограничение – первое число должно быть меньше или равно второму.

```

public class NumberRange {
    // INVARIANT: lower <= upper
    private final AtomicInteger lower = new AtomicInteger(0);
    private final AtomicInteger upper = new AtomicInteger(0);

    public void setLower(int i) {
        // Warning -- unsafe check-then-act
        if (i > upper.get())
            throw new IllegalArgumentException(
                "can't set lower to " + i + " > upper");
        lower.set(i);
    }

    public void setUpper(int i) {
        // Warning -- unsafe check-then-act
        if (i < lower.get())
            throw new IllegalArgumentException(
                "can't set upper to " + i + " < lower");
        upper.set(i);
    }
}

```



```
}

public boolean isInRange(int i) {
    return (i >= lower.get() && i <= upper.get());
}
}
```

Листинг 4.10 Класс NumberRange недостаточно хорошо защищает свои инварианты. Не делайте так.

Класс NumberRange не является потокобезопасным; он не сохраняет инвариант, накладывающий ограничение на объекты lower и upper. Методы setLower и setUpUpper пытаются относиться к инварианту с уважением, но делают это плохо. Оба метода setLower и setUpUpper выполняют последовательность действий проверить-затем-выполнить, но они не используют блокировку достаточную для того, чтобы сделать операции атомарными. Если диапазон номеров удерживается в пределах (0, 10) и один поток вызывает метод setLower (5), а другой поток вызывает метод setUpUpper (4), то в некоторый неудачный момент временем оба метода будут проходить проверки в сеттерах, и будут применены обе модификации. В результате, диапазон теперь содержит значения (5, 4) - недопустимое состояние. Поэтому, в то время как базовые объекты AtomicIntegers являются потокобезопасными, составной класс таковым не является. Поскольку базовые переменные состояния lower и upper не являются независимыми, класс NumberRange не может просто взять и делегировать потокобезопасность своим переменным состояния.

Класс NumberRange можно сделать потокобезопасным, используя блокировку для поддержания сохранности его инвариантов, таких как защита переменных lower и upper с помощью общей блокировки. Также необходимо избегать публикации переменных lower и upper, чтобы предотвратить разрушение инвариантов клиентами.

Если класс имеет составные действия, как в примере с классом NumberRange, делегирование вновь не является должным подходом для обеспечения потокобезопасности. В таких случаях класс должен предоставлять свою собственную блокировку, чтобы гарантировать, что составные действия являются атомарными, если не всё составное действие может быть делегировано базовым переменным состояния.

Если класс состоит из нескольких независимых потокобезопасных переменных состояния, и не имеет операций, которые могут приводить к каким-либо недопустимым переходам состояния, тогда он может делегировать потокобезопасность базовым переменных состояния.

Проблема, которая не позволила классу NumberRange быть потокобезопасным, даже если бы его компоненты состояния были потокобезопасными, очень похожа на одно из правил о volatile переменных, описанных в разделе 3.1.4: переменная подходит для объявления как volatile только в том случае, если она не участвует в инвариантах, связанных с другими переменными состояния.

4.3.4 Публикация базовых переменных состояния

Когда вы делегируете потокобезопасность базовым переменным состояния объекта, при каких условиях вы сможете опубликовать эти переменные, чтобы другие классы могли также их изменять? Опять же, ответ зависит от того, какие инварианты ваш класс налагает на эти переменные. В то время как базовое поле `value` класса `Counter` может принимать любое целочисленное значение, класс `Counter` ограничивает его принятием только положительных значений, а операция инкремента ограничивает набор допустимых следующих состояний при заданном текущем состоянии. Если сделать поле `value` общедоступным, клиенты могут изменить его значение на недопустимое, поэтому публикация приведет к неверному представлению класса. С другой стороны, если переменная представляет текущую температуру или идентификатор последнего пользователя вошедшего в систему, то другой класс может изменить это значение в любое время, вероятно, это не будет нарушать какие-либо инварианты, поэтому публикация этой переменной может быть допустимым решением. (Это может быть не очень хорошей идеей, поскольку публикация изменяемых переменных затруднит разработку в будущем и возможность использования в подклассах, но это не обязательно приведет к тому, что класс будет не потокобезопасным.)

Если переменная состояния является потокобезопасной, она не участвует ни в каких инвариантах, которые ограничивают ее значение и не имеет запрещенных переходов состояний для любой из ее операций, тогда она может быть безопасно опубликована.

Например, было бы безопасно опубликовать поля `mouseListeners` или `keyListeners` из класса `VisualComponent`. Поскольку класс `VisualComponent` не накладывает никаких ограничений на допустимые состояния списков слушателей, эти поля могут быть сделаны публичными (`public`) или опубликованы иным образом без ущерба для потокобезопасности.

4.3.5 Пример: трекер транспортных средств публикующий своё состояние

Давайте создадим еще одну версию трекера транспортных средств, которая публикует его изменяемое базовое состояние. Нам вновь необходимо немного изменить интерфейс, чтобы приспособить это изменение, на этот раз, используя изменяемые, но потокобезопасные точки.

```
@ThreadSafe
public class SafePoint {
    @GuardedBy("this") private int x, y;

    private SafePoint(int[] a) { this(a[0], a[1]); }

    public SafePoint(SafePoint p) { this(p.get()); }

    public SafePoint(int x, int y) {
        this.x = x;
        this.y = y;
    }

    public int get() {
        return x;
    }
}
```

```

}

public synchronized int[] get() {
    return new int[] { x, y };
}

public synchronized void set(int x, int y) {
    this.x = x;
    this.y = y;
}
}

```

Листинг 4.11 Потокобезопасный изменяемый класс Point

Класс `SafePoint` из листинга 4.11 предоставляет геттер, возвращающий оба значения `x` и `y` за один раз, путём возврата двухэлементного массива⁵⁸. Если бы мы предоставили отдельные геттеры для `x` и `y`, то значения могли бы измениться в промежутке времени между моментом получения одной координаты и моментом получения другой, в результате чего, вызывающий код мог увидеть несогласованное значение: местоположение (`x, y`), где транспортное средство никогда не было. Используя класс `SafePoint`, мы можем создать трекер транспортных средств, который публикует базовое изменяемое состояние без ущерба для потокобезопасности, как показано в классе `PublishingVehicleTracker` в листинге 4.12.

```

@ThreadSafe
public class PublishingVehicleTracker {
    private final Map<String, SafePoint> locations;
    private final Map<String, SafePoint> unmodifiableMap;

    public PublishingVehicleTracker(
            Map<String, SafePoint> locations) {
        this.locations
            = new ConcurrentHashMap<String, SafePoint>(locations);
        this.unmodifiableMap
            = Collections.unmodifiableMap(this.locations);
    }

    public Map<String, SafePoint> getLocations() {
        return unmodifiableMap;
    }

    public SafePoint getLocation(String id) {
        return locations.get(id);
    }

    public void setLocation(String id, int x, int y) {
        if (!locations.containsKey(id))

```

⁵⁸ Приватный конструктор существует, чтобы избежать условия гонки, которое возникло бы, если бы копирующий конструктор был реализован как `this(p.x, p.y)`; это пример идиомы захвата приватного конструктора (Bloch and Gafter, 2005).

```
        throw new IllegalArgumentException( "invalid
                                         vehicle name: " + id);
    locations.get(id).set(x, y);
}
}
```

Листинг 4.12. Трекер транспортных средств безопасно публикующий базовое состояние

Класс `PublishingVehicleTracker` наследует потокобезопасность путём делегирования базовому классу `ConcurrentHashMap`, но на этот раз содержимым класса `Map` являются потокобезопасные изменяемые точки, а не неизменяемые. Метод `getLocations` возвращает неизменяемую копию базового объекта `Map`. Вызывающие не могут добавлять или удалять транспортные средства, но могут изменять местоположение одного из транспортных средств, изменения значения объектов `SafePoint` в возвращаемом объекте `Map`. Вновь, “живая” природа класса `Map` может быть как преимуществом, так и недостатком, в зависимости от ваших требований. Класс `PublishingVehicleTracker` является потокобезопасным, но это было бы не так, если бы он налагал какие-либо дополнительные ограничения на допустимые значения местоположений транспортных средств. Если возникает необходимость в наложении “вето” на изменение местоположения транспортных средств или принятии мер при изменении местоположения, подход, принятый в классе `PublishingVehicleTracker` не подойдёт.

4.4 Добавление функциональности в существующие потокобезопасные классы

Библиотека классов Java содержит множество полезных классов – “строительных блоков”. Повторное использование существующих классов часто предпочтительнее создания новых: повторное использование может снизить сложность разработки, риски разработки (поскольку существующие компоненты уже протестированы) и затраты на сопровождение. Иногда потокобезопасный класс, поддерживающий все операции, которые нам необходимы, уже существует, но часто лучшее, что мы можем найти, это класс, который поддерживает *почти* все необходимые операции, и нам остаётся только добавить к нему новую операцию, не нарушив его потокобезопасность.

В качестве примера предположим, что нам нужна потокобезопасная реализация интерфейса `List` с атомарной операцией *положить-если-отсутствует* (*put-if-absent*). Синхронизированные реализации интерфейса `List` почти полностью выполняют эту работу, так как они предоставляют методы `contains` и `add`, из которых мы можем построить операцию *положить-если-отсутствует*.

Концепция *положить-если-отсутствует* достаточно проста – проверьте, есть ли элемент в коллекции, прежде чем добавлять его, и не добавляйте, если он уже в ней. (Сейчас должен прекратиться предупреждающий звон колокольчиков операции *проверить-затем-выполнить*.) Требование, чтобы класс был потокобезопасным, неявно добавляет еще одно требование – чтобы такие операции, как *положить-если-отсутствует*, были *атомарными*. Любое разумное толкование предполагает, что если взять объект типа `List`, который не содержит объект `X`, и добавить объект `X` дважды с помощью операции *положить-если-отсутствует*, результирующая коллекция будет содержать только одну копию объекта `X`. Но, если операция *положить-если-отсутствует*, была не атомарной, в некоторый неудачный момент времени два потока могли увидеть, что объект `X`

отсутствует в списке и оба могли добавить объект X , в результате в списке будет находиться две копии объекта X .

Самый безопасный способ добавить новую атомарную операцию - изменить исходный класс для поддержки нужной операции, но это не всегда возможно, так как у вас может не быть доступа к исходному коду или вы не можете свободно изменять его. Если есть возможность изменить исходный класс, необходимо понять реализацию политики синхронизации, чтобы можно было изменить его в канве исходного дизайна. Добавление нового метода в класс означает, что весь код, который реализует политику синхронизации, для этого класса, все еще находится в одном исходном файле, что облегчает понимание и сопровождение.

Другой подход заключается в расширении класса, принимая допущение, что он был разработан для расширения. Класс `BetterVector` в листинге 4.13 расширяет класс `Vector`, чтобы добавить метод `putIfAbsent`. Расширение класса `Vector` достаточно простое, но не все классы предоставляют достаточный доступ к своему состоянию подклассам, чтобы позволять использование такого подхода.

Расширение более хрупко, чем добавление кода непосредственно в класс, поскольку реализация политики синхронизации теперь распространяется на несколько отдельно поддерживаемых файлов с исходным кодом. Если базовый класс изменит свою политику синхронизации, выбрав другой тип блокировки для защиты переменных состояния, подкласс незаметно и тихо сломается, так как он больше не сможет использовать правильную блокировку для управления параллельным доступом к состоянию базового класса. (Политика синхронизации класса `Vector` зафиксирована в его спецификации, поэтому класс `BetterVector` не пострадает от этой проблемы.)

```
@ThreadSafe
public class BetterVector<E> extends Vector<E> {
    public synchronized boolean putIfAbsent(E x) {
        boolean absent = !contains(x);
        if (absent)
            add(x);
        return absent;
    }
}
```

Листинг 4.13 Расширение класса `Vector` для добавления метода `putIfAbsent`

4.4.1 Блокировка на стороне клиента

Для класса `ArrayList`, обернутого методом `Collections.synchronizedList`, ни один из подходов - добавление метода к исходному классу или расширение класса - не работает, потому что клиентский код даже не знает что класс объекта `List`, возвращается синхронизирующим фабричным методом-оберткой. Третья стратегия заключается в расширении функциональности класса без расширения самого класса, путем размещения кода расширения во “вспомогательном” (*helper*) классе.

В листинге 4.14 показана неудачная попытка создать вспомогательный класс с атомарной операцией `положить-если-отсутствует` для работы с потокобезопасным объектом `List`.

```
@NotThreadSafe
public class ListHelper<E> {
```



```
public List<E> list =  
    Collections.synchronizedList(new ArrayList<E>());  
...  
public synchronized boolean putIfAbsent(E x) {  
    boolean absent = !list.contains(x);  
    if (absent)  
        list.add(x);  
    return absent;  
}  
}
```

Листинг 4.14 Непотокобезопасная попытка реализации операции *положить-если-отсутствует*

Почему бы это не сработало? Ведь метод `putIfAbsent` помечен как `synchronized`, верно? Проблема в том, что он синхронизируется на *неверной* блокировке. Какую бы блокировку объект типа `List` не использовал для защиты своего состояния, она не является собственной блокировкой класса `ListHelper`. Класс `ListHelper` предоставляет только *иллюзию синхронизации*; различные операции списка, несмотря на то, что все синхронизированы, используют разные блокировки, что означает, что метод `putIfAbsent` *не является* атомарным по отношению к другим операциям в объекте `List`. Таким образом, нет никакой гарантии, что другой поток не изменит список во время выполнения метода `putIfAbsent`.

Чтобы этот подход работал, мы должны использовать ту же блокировку, что и объект `List`, с помощью *блокировки на стороне клиента* (*client-side locking*) или *внешней блокировки* (*external locking*). Блокировка на стороне клиента влечёт за собой защиту клиентского кода, который использует некоторый объект `X` с использованием блокировки `X` для защиты своего собственного состояния. Чтобы использовать блокировку на стороне клиента необходимо знать, какую блокировку использует объект `X`.

Документация на класс `Vector` и классы синхронизирующих обёрток (*synchronized wrapper*) утверждает, хотя и косвенно, что они поддерживают блокировку на стороне клиента, используя встроенную блокировку для класса `Vector` или коллекцию классов-обёрток (не обёрнутую коллекцию). В листинге 4.15 показана операция `putIfAbsent` в потокобезопасном объекте `List`, корректно использующем клиентскую блокировку.

```
@ThreadSafe  
public class ListHelper<E> {  
    public List<E> list =  
        Collections.synchronizedList(new ArrayList<E>());  
    ...  
    public boolean putIfAbsent(E x) {  
        synchronized (list) {  
            boolean absent = !list.contains(x);  
            if (absent)  
                list.add(x);  
            return absent;  
        }  
    }  
}
```

Листинг 4.15 Реализация операции *положить-если-отсутствует* с блокировкой на стороне клиента

Если расширение класса для добавления другой атомарной операции является хрупким, так как он распределяет код блокировки класса по нескольким классам в иерархии объектов, блокировка на стороне клиента является еще более хрупкой, поскольку она влечет за собой размещение кода блокировки класса *C* в классы, которые полностью не связаны с классом *C*. Проявляйте осторожность при использовании блокировки на стороне клиента для классов, которые не фиксируют свою стратегию блокировки.

Блокировка на стороне клиента имеет много общего с расширением класса – они оба сочетают поведение производного класса с реализацией базового класса. Подобно тому, как расширение нарушает инкапсуляцию реализации [EJ Item 14], блокировка на стороне клиента нарушает инкапсуляцию политики синхронизации.

4.4.2 Композиция

Существует менее хрупкая альтернатива для добавления атомарной операции в существующий класс: *композиция*. Класс `ImprovedList` из листинга 4.16 реализует операции интерфейса `List`, путём делегирования их базовому экземпляру интерфейса `List` и добавляет атомарный метод *положить-если-отсутствует*. (Подобно фабричному методу `Collections.synchronizedList` и другим классам-обёрткам коллекций, класс `ImprovedList` предполагает, что однажды передав список в его конструктор, клиент не будет использовать базовый список напрямую, но получая к нему доступ только через класс `ImprovedList`.)

```
@ThreadSafe
public class ImprovedList<T> implements List<T> {
    private final List<T> list;

    public ImprovedList(List<T> list) { this.list = list; }

    public synchronized boolean putIfAbsent(T x) {
        boolean contains = list.contains(x);
        if (contains)
            list.add(x);
        return !contains;
    }

    public synchronized void clear() { list.clear(); }
    // ... similarly delegate other List methods
}
```

Листинг 4.16 Реализация операции *положить-если-отсутствует* с использованием композиции

Класс `ImprovedList` добавляет дополнительный уровень блокировки с использованием собственной встроенной блокировки. Не имеет значения, является ли базовый класс `List` потокобезопасным, поскольку он обеспечивает собственную согласованную блокировку, обеспечивающую безопасность потоков, даже если класс `List` не является потокобезопасным или его реализация блокировки изменится. Хотя дополнительный уровень синхронизации может привести к небольшому снижению производительности⁵⁹, реализация подхода как в

⁵⁹ Штраф будет небольшим, потому что синхронизация в базовом классе `List` гарантированно не будет затронута и, следовательно, будет быстрой; см. главу 11.

`ImprovedList` менее хрупкая, чем попытка имитировать стратегию блокировки другого объекта. Фактически, мы использовали шаблон Java монитор для инкапсуляции существующего объекта `List`, и это гарантирует потокобезопасность, пока наш класс содержит единственную выдающуюся ссылку на базовый объект `List`.

4.5 Документирование политики синхронизации

Документация является одним из самых мощных (и, к сожалению, наиболее сильно недоиспользуемых) инструментов для управления потокобезопасностью. Пользователи обращаются к документации, чтобы узнать, является ли класс потокобезопасным, а сопровождающие смотрят на документацию, чтобы понять стратегию реализации, чтобы они могли поддерживать его без непреднамеренной компрометации безопасности. К сожалению, обе группы клиентов обычно находят в документации меньше информации, чем им хотелось бы.

Документируйте гарантии потокобезопасности класса для его клиентов; документируйте политику синхронизации для сопровождающих его.

Каждое использование операторов `synchronized`, `volatile` или любого потокобезопасного класса отражается в *политике синхронизации*, определяющей стратегию обеспечения целостности данных перед лицом параллельного доступа. Эта политика является элементом дизайна вашей программы и должна быть задокументирована. Конечно, лучшее время для документирования проектных решений - этап проектирования. Недели или месяцы спустя, детали могут быть размыты - так что запишите их, прежде чем забыть.

Для создания политики синхронизации требуется принять ряд решений: какие переменные сделать `volatile`, какие переменные защитить с помощью блокировок, какие блокировки защищают какие переменные, какие переменные должны быть неизменяемыми или должны быть ограничены потоком, какие операции должны быть атомарными и т.д. Некоторые из них являются подробными сведениями о реализации и должны быть задокументированы для будущих сопровождающих, но некоторые из них влияют на публично наблюдаемое поведение блокировки вашего класса и должны быть задокументированы как часть его спецификации.

По крайней мере, задокументируйте гарантии потокобезопасности, сделанные классом. Это потокобезопасно? Делает ли он обратные вызовы, удерживая блокировку? Существуют ли какие-либо специфичные блокировки, влияющие на его поведение? Не заставляйте клиентов делать рискованные предположения. Если вы не хотите поддерживать блокировку на стороне клиента, это нормально, но так и скажите. Если вы хотите, чтобы клиенты могли создавать новые атомарные операции в вашем классе, как мы это делали в разделе 4.4, вам нужно задокументировать, какие блокировки они должны захватить, чтобы сделать это безопасно. Если вы используете блокировки для защиты состояния, задокументируйте это для будущих сопровождающих, потому что это так просто - аннотация `@GuardedBy` выполнит за вас этот трюк. Если вы используете более тонкие средства для обеспечения потокобезопасности, документируйте их, потому что они могут быть не очевидны для сопровождающих.

Текущее состояние дел в документации по потокобезопасности, даже в библиотеке классов платформы, не обнадеживает. Сколько раз вы смотрели в

Javadoc на описание класса и задавались вопросом, является ли он потокобезопасным?⁶⁰ Так или иначе, большинство классов не предлагают никаких подсказок. Многие официальные спецификации Java-технологий, такие как сервлеты и JDBC, удручающе относятся к документированию своих обещаний и требований к безопасности потоков.

Хотя благоразумие наводит на мысль, что мы не предполагаем поведения, не входящего в спецификацию, у нас есть работа, которая должна быть выполнена, и мы часто сталкиваемся с последствиями выбора плохих допущений. Должны ли мы предположить, что объект является потокобезопасным, потому как кажется, что он таковым должен быть? Должны ли мы предположить, что доступ к объекту можно сделать потокобезопасным, захватывая его блокировку? (Этот рискованный метод работает только в том случае, если мы контролируем весь код, который обращается к этому объекту; в противном случае, он нам дает лишь иллюзию потокобезопасности.) Ни один из вариантов не является удовлетворительным.

Усугубляя проблему, наша интуиция часто может ошибаться в том, какие классы “вероятно потокобезопасны”, а какие нет. Например, класс `java.text.SimpleDateFormat` является потокобезопасным, но в Javadoc забыли упомянуть об этом вплоть до JDK 1.4. То, что этот конкретный класс не является потокобезопасным, вызвало удивление многих разработчиков. Сколько программ ошибочно создают совместно используемый экземпляр объекта, не являющегося потокобезопасным, и используют его из нескольких потоков, не подозревая, что это может привести к ошибочным результатам при большой нагрузке?

Проблему с классом `SimpleDateFormat` можно обойти, если не считать класс потокобезопасным, если об этом не сказано прямо. С другой стороны, невозможно создать приложение на основе сервлета, не делая довольно сомнительных предположений о потокобезопасности объектов, предоставляемых контейнером, подобных `HttpSession`. Не заставляйте своих клиентов или коллег делать догадки, подобные этим.

4.5.1 Интерпретация расплывчатой документации

Многие спецификации технологий Java молчат или, по меньшей мере, неудовлетворительно рассказывают о гарантиях потокобезопасности и требованиях к таким интерфейсам, как `ServletContext`, `HttpSession` или `DataSource`⁶¹. Поскольку эти интерфейсы реализуются поставщиком контейнера или базы данных, часто невозможно просмотреть код, чтобы увидеть, что он делает. Кроме того, вы не хотите полагаться на детали реализации одного конкретного драйвера JDBC - вы хотите иметь совместимость со стандартом, так чтобы ваш код работал должным образом с любым драйвером JDBC. Но слова “поток” и “параллельный” вообще не появляются в спецификации JDBC и появляются удручающе редко в спецификации сервлета. Так чем вы займёесь?

Вам придется гадать. Один из способов улучшить качество вашей догадки - интерпретировать спецификацию с точки зрения того, кто ее *реализует* (например, поставщик контейнера или базы данных), в отличие от того, кто просто ее использует. Сервлеты всегда вызываются из потока, управляемого контейнером, и можно с уверенностью предположить, что если есть более одного такого потока, контейнер знает об этом. Контейнер сервлета предоставляет определенные

⁶⁰ Если вы никогда не задавались этим вопросом, мы восхищаемся вашим оптимизмом.

⁶¹ Мы находим особенно расстраивающим то, что эти упоминания сохраняются, несмотря на многочисленные пересмотры спецификаций.

объекты, которые обеспечивают обслуживание множества сервлетов, такие как `HttpSession` или `ServletContext`. Таким образом, контейнер сервлета должен рассчитывать, что эти объекты будут доступны одновременно, так как он создал несколько потоков и вызвал у них методы, подобные `Servlet.service`, что вполне ожидаемо при доступе к `ServletContext`.

Поскольку невозможно представить однопоточный контекст, в котором эти объекты были бы полезны, следует предположить, что они были созданы потокобезопасными, хотя спецификация явно не требует этого. Кроме того, если им требуется блокировка на стороне клиента, на какой блокировке должен синхронизироваться код клиента? Документация ничего не говорит об этом, и кажется абсурдным, что приходится гадать. Это “разумное предположение” далее подкрепляется примерами в спецификации и официальными руководствами, которые показывают, как обращаться к классам `ServletContext` или `HttpSession` и не использовать какую-либо клиентскую синхронизацию.

С другой стороны, объекты, помещенные в экземпляры `ServletContext` или `HttpSession` с помощью метода `setAttribute`, принадлежат веб-приложению, а не контейнеру сервлета. Спецификация сервлета не предполагает какого-либо механизма для координации параллельного доступа к совместно используемым атрибутам. Поэтому атрибуты, хранящиеся в контейнере от имени веб-приложения, должны быть потокобезопасными или эффективно неизменяемыми. Если бы контейнер хранил все атрибуты от имени веб-приложения, другим вариантом была бы гарантия, что они будут согласованно защищены блокировкой при доступе из кода приложения сервлета. Но поскольку контейнеру может потребоваться сериализовать объекты, содержащиеся в `HttpSession` для репликации или пассивации, а контейнер сервлета не может о реализуемом объектами протоколе блокировки, следует сделать их потокобезопасными.

Можно сделать аналогичный вывод об интерфейсе JDBC `DataSource`, который представляет собой пул повторно используемых соединений с базой данных. Интерфейс `DataSource` предоставляет приложению службу, и это не имеет большого смысла в контексте однопоточного приложения. Трудно представить себе вариант использования, который не включает вызов метода `getConnection` из нескольких потоков. И, как и в случае с сервлетами, примеры в спецификации JDBC не предполагают потребности в любой клиентской блокировке во множестве примеров кода, использующих объекты `DataSource`. Так что, хотя спецификация не утверждает, что источник данных является потокобезопасным и не требует, чтобы поставщик контейнера обеспечивал потокобезопасную реализацию, к тому же, учитывая аргумент “*было бы нелепо, если бы не*”, у нас нет выбора, кроме как предположить, что вызов `DataSource.getConnection` не требует дополнительной блокировки на стороне клиента.

С другой стороны, мы не будем приводить те же аргументы в отношении объектов JDBC `Connection`, порождаемых объектами `DataSource`, поскольку они не обязательно предназначены для совместного использования другими активностями (*activity*) до тех пор, пока не будут возвращены в пул соединений. Поэтому, если активность, которая получает объект JDBC `Connection`, охватывает несколько потоков, она должна нести ответственность за обеспечение надлежащего контроля доступа к соединению посредством синхронизации. (В большинстве приложений активности, использующие объекты JDBC `Connection`, реализованы таким образом, чтобы в любом случае ограничивать подключение определенным потоком.)

Глава 5 Строительные блоки

В предыдущей главе были рассмотрены несколько методов построения потокобезопасных классов, в том числе делегирование обеспечения потокобезопасности существующим потокобезопасным классам. В тех случаях, когда это целесообразно, делегирование является одной из наиболее эффективных стратегий создания потокобезопасных классов: просто позвольте существующим потокобезопасным классам управлять всем состоянием.

Библиотеки платформы включают в себя богатый набор параллельных строительных блоков, такие как потокобезопасные коллекции и разнообразные *синхронизаторы*, которые могут координировать потоки управления взаимодействующих потоков. В этой главе рассматриваются наиболее полезные параллельные строительные блоки, особенно те, которые представлены в Java 5.0 и Java 6, и некоторые шаблоны их использования способствующие структурированию параллельных приложений.

5.1 Синхронизированные коллекции

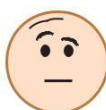
Классы синхронизированных коллекций включают в себя `Vector` и `Hashtable`, часть оригинального JDK, а также их “двоюродных братьев”, добавленных в JDK 1.2 и синхронизированные классы-обёртки, создаваемые фабричными методами `Collections.synchronizedXxx`. Эти классы обеспечивают потокобезопасность, инкапсулируя свое состояние и синхронизируя каждый открытый метод, чтобы одновременно только один поток мог получить доступ к состоянию коллекции.

5.1.1 Проблемы с синхронизированными коллекциями

Синхронизированные коллекции потокобезопасны, но иногда для защиты составных действий может потребоваться Дополнительная блокировка на стороне клиента. В общем случае, составные действия для коллекций включают в себя итерацию (многократную выборку элементов, пока коллекция не будет исчерпана), навигацию (поиск следующего элемента после текущего в соответствии с некоторым порядком) и условные операции, такие как *положить-если-отсутствует* (проверка, имеет ли объект Map сопоставление для ключа K , если нет, добавление сопоставления (K, V)). С синхронизированной коллекцией эти составные действия по-прежнему технически потокобезопасны, даже без блокировки на стороне клиента, но они могут вести себя не так, как можно было бы ожидать, когда другие потоки могут одновременно изменять коллекцию.

В листинге 5.1 показаны два метода, оперирующие классом `Vector`, `getLast` и `deleteLast`, оба из которых представляют собой последовательностями *проверить-затем-выполнить*. Каждый вызывает метод `size` для определения размера массива и использует полученное значение для извлечения или удаления последнего элемента.

```
public static Object getLast(Vector list) {  
    int lastIndex = list.size() - 1;  
    return list.get(lastIndex);  
}
```



```

public static void deleteLast(Vector list) {
    int lastIndex = list.size() - 1;
    list.remove(lastIndex);
}

```

Листинг 5.1 Составные действия над классом `Vector`, которые могут привести к запутанным результатам.

Эти методы кажутся безобидными, и в некотором смысле так и есть - они не могут повредить объект `Vector`, независимо от того, сколько потоков одновременно вызывает их. Но у вызвавшего эти методы может сложиться другое мнение. Если поток *A* вызывает метод `getLast` объекта `Vector` с десятью элементами, а поток *B* вызывает метод `deleteLast` у того же объекта `Vector`, и операции чередуются, как показано на рисунке 5.1, метод `getLast` бросит исключение `ArrayIndexOutOfBoundsException`. Между вызовом метода `size` и последующим вызовом метода `get` в методе `getLast`, размер объекта `Vector` уменьшился, и индекс, вычисленный на первом шаге, стал неправильным.

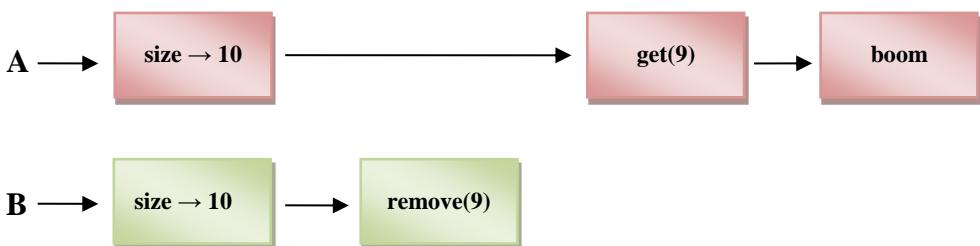


Рисунок 5.1 Чередование вызовов методов `getLast` и `deleteLast`, приводит к генерации исключения `ArrayIndexOutOfBoundsException`

Такое поведение полностью соответствует спецификации класса `Vector` - бросается исключение, если запрашивается несуществующий элемент. Это не тот результат, которого ожидал получить вызывающий объект от вызова метода `getLast`, даже на фоне параллельной модификации, если принять, что объект `Vector` не был пуст с самого начала.

Поскольку синхронизированные коллекции фиксируются в политике синхронизации, которая поддерживает блокировку на стороне клиента⁶², можно создавать новые операции, которые являются атомарными по отношению к другим операциям коллекций, при условии, что нам известно, какую блокировку использовать. Синхронизированные классы коллекций защищают каждый метод с помощью блокировки на самом синхронизированном объекте коллекции. Захватив блокировку коллекции, мы можем сделать методы `getLast` и `deleteLast` атомарными, гарантируя, что размер объекта `Vector` не изменится между вызовами методов `size` и `get`, как показано в листинге 5.2.

```

public static Object getLast(Vector list) {
    synchronized (list) {
        int lastIndex = list.size() - 1;
        return list.get(lastIndex);
    }
}

```

⁶² Это задокументировано только в Javadoc версии Java 5.0, как пример корректной идиомы итерации.

```
public static void deleteLast(Vector list) {  
    synchronized (list) {  
        int lastIndex = list.size() - 1;  
        list.remove(lastIndex);  
    }  
}
```

Листинг 5.2 Составные действия над классом `Vector` с использованием блокировки со стороны клиента

Риск того, что размер списка может измениться между вызовом `size` и соответствующим вызовом `get`, также присутствует при переборе элементов объекта `Vector`, как показано в листинге 5.3.

```
for (int i = 0; i < vector.size(); i++)  
    doSomething(vector.get(i));
```

Листинг 5.3 Итерация, которая может бросить исключение `ArrayIndexOutOfBoundsException`

Эта идиома итерации основывается на вере в то, что другие потоки не изменят объект `Vector` между вызовами методов `size` и `get`. В однопоточной среде это предположение вполне допустимо, но когда другие потоки могут параллельно изменять объект `Vector`, это может привести к проблемам. Так же, как с методом `getLast`, если другой поток удаляет элемент, в то время как вы перебираете элементы объекта `Vector`, и операции чередуются неудачно, эта идиома итерации бросает исключение `ArrayIndexOutOfBoundsException`.

Хотя итерация в листинге 5.3 может выдать исключение, это не означает, что класс `Vector` не является потокобезопасным. Состояние класса `Vector` остается действительным, и исключение фактически соответствует его спецификации. Тем не менее, что-то столь же обычное, как извлечение последнего элемента или итерация, бросает исключение, что явно нежелательно.

Проблема недостаточной итерации может быть вновь решена путем блокировки на стороне клиента, при некоторых дополнительных затратах на масштабируемость. Удерживая блокировку объекта `Vector` на протяжении всей итерации, как показано в листинге 5.4, мы не допускаем изменения объекта `Vector` другими потоками во время итерации. К сожалению, в течение этого времени мы также препятствуем другим потокам в получении доступа к нему, что ухудшает параллелизм.

```
synchronized (vector) {  
    for (int i = 0; i < vector.size(); i++)  
        doSomething(vector.get(i));  
}
```

Листинг 5.4 Итерация, с блокировкой на стороне клиента

5.1.2 Итераторы и `ConcurrentModificationException`

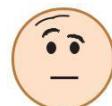
Мы используем класс `Vector` во многих наших примерах для большей наглядности, несмотря на то, что он считается “наследованным” (*legacy*) классом коллекций. Но более “современные” классы коллекций не устраниют проблему с составными действиями. Стандартный способ итерирования объекта `Collection` заключается в явном использовании объекта `Iterator` или неявном, с использованием синтаксиса цикла `for-each`, представленного в Java 5.0, но использование итераторов не устраниет необходимость в использовании блокировки коллекций во время итерации, если другие потоки могут параллельно

изменять их. Итераторы, возвращаемые синхронизированными коллекциями, не предназначены для работы в режиме параллельной модификации, они быстро падают (fail-fast) с ошибкой - это означает, что если они обнаруживают, что коллекция изменилась с начала итерации, они бросают исключение `ConcurrentModificationException`.

Эти быстро падающие с ошибкой итераторы не предназначены для надежной защиты - они спроектированы на основе принципа “добросовестных усилий” (*good-faith-effort*) для перехвата ошибок параллелизма и, таким образом, действуют только как индикаторы раннего предупреждения о проблемах параллелизма. Они реализуются путем связывания счетчика изменений с коллекцией: если счетчик изменений модифицируется во время итерации, методы `hasNext` или `next` бросают исключение `ConcurrentModificationException`. Однако эта проверка выполняется без синхронизации, поэтому существует риск увидеть устаревшее значение счетчика изменений и, следовательно, итератор может не узнать, что изменение было сделано. Это следствие специально принятого компромисса, для снижения влияния кода обнаружения параллельных изменений на производительность.⁶³

В листинге 5.5 иллюстрируется процесс итерирования коллекции с помощью синтаксиса цикла `for-each`. Внутренне `javac` генерирует код, который использует класс `Iterator`, повторно вызывающий методы `hasNext` и `next` для итерирования объекта `List`. Как и при итерировании класса `Vector`, способ предотвращения возникновения исключения `ConcurrentModificationException` заключается в том, чтобы удерживать на коллекции блокировку до полного завершения процесса итерации.

```
List<Widget> widgetList
    = Collections.synchronizedList(new ArrayList<Widget>());
...
// May throw ConcurrentModificationException
for (Widget w : widgetList)
    doSomething(w);
```



Листинг 5.5 Итерирование объекта `List` с помощью класса `Iterator`

Однако существует несколько причин, по которым блокировка коллекции во время итерации может быть нежелательной. Другие потоки, которым требуется доступ к коллекции, блокируются до завершения итерации; если коллекция большая или задача, выполняемая для каждого элемента, является длительной, потоки будут вынуждены долго ждать. Кроме того, если коллекция заблокирована, как показано в листинге 5.4, вызывается функция `doSomething` с блокировкой, которая является фактором риска возникновения взаимоблокировки (см. главу 10). Даже при отсутствии риска голода или взаимоблокировки, блокировка коллекций на значительные периоды времени ухудшает масштабируемость приложений. Чем дольше удерживается блокировка, тем больше вероятность того, что она будет оспорена, и если множество потоков заблокировано в ожидании освобождения блокировки, будет страдать пропускная способность и эффективность использования ЦП (см. главу 11).

Альтернативой блокировке коллекции во время итерации является клонирование коллекции и итерирования копии. Поскольку клон ограничен

⁶³ Исключение `ConcurrentModificationException` также может возникать и в однопоточном коде; это происходит, когда объекты удаляются из коллекции напрямую, а не с использованием метода `Iterator.remove`.

потоком, никакой другой поток не сможет изменить его во время итерации, таким образом, исключается возможность возбуждения исключения `ConcurrentModificationException`. (Коллекция по-прежнему должна быть заблокирована во время самой операции клонирования.) Клонирование коллекции оказывает очевидное влияние на производительность; является ли это выгодным компромиссом, зависит от многих факторов, включая размер коллекции, объем работы для каждого элемента, относительную частоту итераций по сравнению с другими операциями коллекции, а также требования к отзывчивости и пропускной способности.

5.1.3 Скрытые итераторы

Хотя блокировка может предотвращать возникновение исключения `ConcurrentModificationException`, необходимо помнить об использовании блокировки везде, где может выполняться итерирование совместно используемой коллекции. Пуще сказать, чем сделать, поскольку итераторы иногда скрыты, как в классе `HiddenIterator` из листинга 5.6. В классе `HiddenIterator` нет явной итерации, но код, выделенный жирным шрифтом, использует итерацию. Конкатенация строк преобразуется компилятором в вызов `StringBuilder.append(Object)`, который, в свою очередь, вызывает метод коллекции `toString`- и реализация `toString` в стандартных коллекциях итерирует коллекцию и вызывает `toString` на каждом элементе, чтобы сформировать хорошо отформатированное представление содержимого коллекции.

Метод `addTenThings` может бросить исключение `ConcurrentModificationException`, потому что коллекция итерируется при вызове метода `toString`, в процессе подготовки отладочного сообщения. Конечно, реальная проблема заключается в том, что класс `HiddenIterator` не является потокобезопасным; блокировка класса `HiddenIterator` должна захватываться перед использованием переменной `set` в вызове `println`, но код отладки и журналирования обычно пренебрегает этим.

Настоящий урок здесь в том, что чем больше расстояние между состоянием и синхронизацией, которая его защищает, тем больше вероятность того, что кто-то забудет согласованно использовать синхронизацию при доступе к этому состоянию. Если бы класс `HiddenIterator` обернул класс `HashSet` с помощью `synchronizedSet`, инкапсулируя синхронизацию, такая ошибка бы не возникла.

Подобно тому, как инкапсуляция состояния объекта упрощает сохранение его инвариантов, инкапсуляция его синхронизации упрощает реализацию политики синхронизации.

```
public class HiddenIterator {  
    @GuardedBy("this")  
    private final Set<Integer> set = new HashSet<Integer>();  
  
    public synchronized void add(Integer i) { set.add(i); }  
    public synchronized void remove(Integer i) { set.remove(i); }
```



```
public void addTenThings() {  
    Random r = new Random();  
    for (int i = 0; i < 10; i++)  
        add(r.nextInt());  
    System.out.println("DEBUG: added ten elements to " + set);  
}  
}
```

Листинг 5.6 Итерация скрытая внутри конкатенации строк. Не делайте так.

Итерация также косвенно вызывается методами коллекции `hashCode` и `equals`, которые могут вызываться, если коллекция используется как элемент или ключ другой коллекции. Аналогично, методы `containsAll`, `removeAll` и `keepAll`, а также конструкторы, которые принимают коллекции в качестве аргументов, также выполняют итерацию коллекции. Все эти косвенные применения итераций могут приводить к возбуждению исключения `ConcurrentModificationException`.

5.2 Параллельные коллекции

В Java 5.0 синхронизированные коллекции были улучшены, путём ввода нескольких классов параллельных коллекций. Синхронизированные коллекции обеспечивают потокобезопасность путём сериализации всех обращений к состоянию коллекции. Цена такого подхода – низкая степень параллелизма; когда несколько потоков конкурируют за блокировку всей коллекции, страдает пропускная способность.

Параллельные коллекции, с другой стороны, предназначены для одновременного доступа из нескольких потоков. В Java 5.0 добавлен класс `ConcurrentHashMap`, в качестве замены для синхронизированных основанных на хэше реализаций интерфейса `Map`, и класс `CopyOnWriteArrayList`, в качестве замены для синхронизированных реализаций интерфейса `List`, в случаях, когда обход элементов является доминирующей операцией. Новый интерфейс `ConcurrentMap` добавляет поддержку общих составных действий, таких как положить-если-отсутствует, замена и удаление по условию.

Замена синхронизированных коллекций параллельными коллекциями может предложить значительное улучшение масштабируемости с небольшим риском.

В Java 5.0 также было добавлено два новых типа коллекции, `Queue` и `BlockingQueue`. Класс `Queue` предназначен для временного хранения набора элементов в ожидании обработки. Предлагается несколько реализаций, в том числе класс `ConcurrentLinkedQueue`, традиционная очередь FIFO, класс `PriorityQueue`, (не параллельная) приоритетная упорядоченная очередь. Операции очередей не блокируются; если очередь пуста, операция извлечения возвращает значение `null`. Хотя вы можете имитировать поведение класса `Queue` с помощью интерфейса `List` - фактически, класс `LinkedList` также является реализацией очереди - классы `Queue` были добавлены, из-за того, что устранение требований произвольного доступа интерфейса `List` допускает более эффективные параллельные реализации.

Класс `BlockingQueue` расширяет класс `Queue`, добавляя блокировки операциям вставки и извлечения данных. Если очередь пуста, извлечение данных блокируется

до тех пор, пока элемент не будет доступен, и если очередь заполнена (для ограниченных очередей), вставка блокируется до того момента, пока не появится свободное место. Блокирующие очереди чрезвычайно полезны в шаблоне производитель-потребитель (*producer-consumer*) и более подробно рассматриваются в разделе 5.3.

Также как класс `ConcurrentHashMap` является параллельной заменой для синхронизированного основанного на хэше класса `Map`, Java 6 добавляет классы `ConcurrentSkipListMap` и `ConcurrentSkipListSet`, которые выступают в качестве параллельной замены для синхронизированных классов `SortedMap` или `SortedSet` (например, классов `TreeMap` или `TreeSet`, обёрнутых с помощью метода `synchronizedMap`).

5.2.1 Класс `ConcurrentHashMap`

Синхронизированные классы коллекций удерживают блокировку во время выполнения каждой операции. Некоторые операции, такие как `HashMap.get` или `List.contains`, могут требовать больше работы, чем предполагалось изначально: перемещение хэш-корзины (*hash bucket*) или поиск конкретного объекта в списке влечёт за собой вызов метода `equals` (который сам по себе может вызывать значительное количество вычислений) для нескольких объектов-кандидатов. Если метод `hashCode`, в коллекции, основанной на хэше, плохо распределяет хэш-значения, элементы могут быть неравномерно распределены между корзинами; в вырожденном случае, плохая хеш-функция превратит хэш-таблицу в связанный список. Обход длинного списка и вызов метода `equals` для некоторых или всех элементов может занять много времени, и в течение этого времени ни один другой поток не может получить доступ к коллекции.

Класс `ConcurrentHashMap` – это основанная на хэше реализация интерфейса `Map`, подобная классу `HashMap`, но использующая совершенно другую стратегию блокировки, которая обеспечивает лучшую параллельность и масштабируемость. Вместо синхронизации каждого метода на общей блокировке, с ограничением доступа одним потоком за раз, он использует более детализированный механизм блокировки, называемый *чередуемой блокировкой* (см. раздел [11.4.3](#)), позволяющий обеспечить большую степень совместного доступа. Множество потоков-читателей может произвольно и параллельно получать доступ к переменной `map`⁶⁴, читатели могут получать доступ к переменной `map` одновременно с писателями, и ограниченное число писателей может параллельно изменять переменную `map`. Результатом является гораздо более высокая пропускная способность при параллельном доступе с небольшим снижением производительности при однопоточном доступе.

Класс `ConcurrentHashMap`, наряду с другими параллельными коллекциями, далее улучшает синхронизированные классы коллекций, предоставляя итераторы, которые не бросают исключение `ConcurrentModificationException`, таким образом, устранив необходимость в блокировке коллекции во время итерирования. Итераторы, возвращаемые классом `ConcurrentHashMap`, слабо согласованы (*weakly consistent*), в отличие от итераторов, быстро падающих с ошибкой, возвращаемых синхронизированными коллекциями. Слабо согласованный итератор может допускать параллельную модификацию, обход элементов так, как они

⁶⁴ В контексте – переменная `map` является экземпляром класса `ConcurrentHashMap`.

существовали на момент построения итератора, и может (но не гарантируется) отражать изменения в коллекции после построения итератора.

Как и со всеми прочими улучшениями, пришлось пойти на несколько компромиссов. Семантика методов, оперирующих всем классом Map, таких как `size` и `isEmpty`, была немного ослаблена, чтобы отразить параллельный характер коллекции. Поскольку результат вызова метода `size` может быть устаревшим к моменту его вычисления, фактически выполняется только оценка, поэтому метод `size` может возвращать вместо точного подсчета аппроксимированное значение. Хотя на первый взгляд это может вызывать беспокойство, на самом деле такие методы, как `size` и `isEmpty`, гораздо менее полезны в параллельных средах, поскольку эти величины являются изменяющимися значениями. Таким образом, требования к этим операциям были ослаблены, чтобы обеспечить оптимизацию производительности для наиболее важных операций, в первую очередь `get`, `put`, `containsKey` и `remove`.

Одна особенность, предлагаемая синхронизированными реализациями интерфейса Map, но не предлагаемая классом `ConcurrentHashMap`, - это возможность блокировки переменной `map` для монопольного доступа. С помощью класса `Hashtable` и метода `synchronizedMap` захват блокировки объекта Map предотвращает доступ к нему любого другого потока. Это может быть необходимо в нестандартных случаях, таких как атомарное добавление нескольких отображений или итерирование объекта Map несколько раз и необходимость видеть одни и те же элементы в одном и том же порядке. В целом, это разумный компромисс: ожидается, что параллельные коллекции будут постоянно изменять свое содержимое.

Поскольку он имеет так много преимуществ и так мало недостатков по сравнению с классом `Hashtable` или методом `synchronizedMap`, замена синхронизированных реализаций Map на класс `ConcurrentHashMap` в большинстве случаев приводит к улучшению масштабируемости. Класс `ConcurrentHashMap` не подходит для полной замены только в том случае, если приложению необходимо заблокировать переменную `map` для монопольного доступа⁶⁵.

5.2.2 Дополнительные атомарные операции интерфейса Map

Поскольку класс `ConcurrentHashMap` не может быть заблокирован для монопольного доступа, мы не можем использовать блокировку на стороне клиента для создания новых атомарных операций, таких как положить-если-отсутствует, как мы сделали для класса `Vector` в разделе 4.4.1. Вместо этого, ряд общих составных операций, таких как положить-если-отсутствует, удалить-если-равно и заменить-если-равно реализованы как атомарные операции и определены в интерфейсе `ConcurrentMap`, приведённом в листинге 5.7.

```
public interface ConcurrentMap<K,V> extends Map<K,V> {
    // Insert into map only if no value is mapped from K
    V putIfAbsent(K key, V value);

    // Remove only if K is mapped to V
```

⁶⁵ Или в случае, если вы полагаетесь на побочные эффекты от синхронизации синхронизированной реализации интерфейса Map.

```
boolean remove(K key, V value);

// Replace value only if K is mapped to oldValue
boolean replace(K key, V oldValue, V newValue);

// Replace value only if K is mapped to some value
V replace(K key, V newValue);
}
```

Листинг 5.7 Интерфейс ConcurrentMap

Если вы найдёте для себя возможным добавить такую функциональность в существующую синхронизированную реализацию интерфейса Map, это, вероятно, будет признаком того, что вместо нее следует рассмотреть возможность использования интерфейса ConcurrentMap.

5.2.3 Класс CopyOnWriteArrayList

Класс CopyOnWriteArrayList является параллельной заменой для синхронизированной реализации интерфейса List, предлагает лучшую степень параллелизма в распространенных ситуациях и устраняет необходимость использования блокировки или копирования коллекции в процессе итерирования. (Аналогичным образом, класс CopyOnWriteArraySet является параллельной заменой синхронизированной реализации интерфейса Set.)

Коллекции типа “скопировать-при-записи” (*copy-on-write*) обеспечивают потокобезопасность благодаря тому, что при правильной публикации фактически неизменяемого объекта дальнейшая синхронизация при доступе к нему не требуется. Они реализуют изменяемость, создавая и публикую новую копию коллекции при каждом ее изменении. Итераторы для коллекций “скопировать-при-записи” сохраняют ссылку на резервный массив, который был текущим на момент начала итерации, и так как он никогда не изменяется, им необходима только краткая синхронизация, чтобы обеспечить видимость содержимого массива. В результате несколько потоков могут выполнять итерирование коллекции без влияния друг на друга или влияния потоков, желающих изменить коллекцию. Итераторы, возвращаемые коллекциями типа “скопировать-при-записи”, не возбуждают исключение ConcurrentModificationException и возвращают элементы точно такими же, какими они были на момент создания итератора, независимо от последующих изменений.

Очевидно, что копирование резервного массива при каждом изменении коллекции сопряжено с определенными затратами, особенно если коллекция велика; коллекции типа “скопировать-при-записи” целесообразно использовать только тогда, когда итерация используется гораздо чаще, чем модификация. Этот критерий точно описывает принцип работы многих систем уведомления о событиях: для доставки уведомления требуется итерация по списку зарегистрированных слушателей и вызов каждого из них, и, в большинстве случаев, регистрация или отмена регистрации слушателя событий происходит гораздо реже, чем получение уведомления о событии. (См. [CPJ 2.4.4] для получения дополнительной информации о шаблоне “скопировать-при-записи”.)

5.3 Блокирующие очереди и шаблон производитель-потребитель

Блокирующие очереди предоставляют блокирующие методы `put` и `take`, а также их эквиваленты с ограничением времени `offer` и `poll`. Если очередь заполнена, выполнение метода `put` блокируется до тех пор, пока в очереди не появится свободное место; если очередь пуста, метод `take` блокируется, до того момента, пока не появится доступный элемент. Очереди могут быть ограниченными или неограниченными; неограниченные очереди никогда не заполняются, поэтому метод `put` в неограниченной очереди никогда не блокируется.

Блокирующие очереди поддерживают шаблон проектирования “производитель-потребитель” (*producer-consumer*). Модель “производитель-потребитель” отделяет идентификацию работы, которая должна быть выполнена, от выполнения этой работы, помещая элементы работы в список “сделать это” (*to do*) для последующей обработки, вместо того, чтобы обрабатывать их сразу, по мере идентификации. Шаблон “производитель-потребитель” упрощает разработку, так как удаляет код зависимостей между классами производителя и потребителя, и упрощает управление нагрузкой путем разделения активностей, которые могут производить или потреблять данные с различной или переменной скоростью.

В схеме ‘производитель-потребитель’, построенной вокруг блокирующей очереди, производители помещают данные в очередь по мере их поступления, а потребители извлекают данные из очереди, когда они готовы предпринять соответствующие действия. Производителям не нужно ничего знать об идентичности или количестве потребителей, или даже о том, являются ли они единственным производителем - всё, что им нужно сделать, это поместить элементы данных в очередь. Точно так же потребители не должны знать, кто является производителем или откуда взялась работа. Класс `BlockingQueue` упрощает реализацию дизайна “производитель-потребитель” с любым количеством производителей и потребителей. Одним из наиболее распространенных дизайнов реализации шаблона “производитель-потребитель” является пул потоков в сочетании с рабочей очередью; этот шаблон воплощен в фреймворке выполнения задач `Executor`, что является темой глав 6 и 8.

Разделение труда для двух человек, моющих посуду, является знакомым примером дизайна ‘производителя-потребителя’: один человек моет посуду и помещает её в стойку для посуды, а другой человек извлекает посуду из стойки и сушит её. В этом случае стойка с тарелками действует как блокирующая очередь; если в стойке нет посуды, потребитель ждет, пока она появится, чтобы высушить её, и если стойка заполняется, производитель должен прекратить мойку, до того момента, пока в шкафу не освободится место. Эта аналогия распространяется на несколько производителей (хотя это может привести к конкуренции за доступ к мойке) и нескольких потребителей; каждый рабочий (*worker*) взаимодействует только с посудой. Никто не должен знать, сколько производителей или потребителей существует, или кто создал полученный элемент работы.

Ярлыки “производитель” и “потребитель” относительны; активность, которая выступает в качестве потребителя в одном контексте, может выступать в качестве производителя в другом. Сушка посуды “потребляет” чистую влажную посуду и “производит” чистую сухую посуду. Третий человек, желающий помочь, может убирать сухую посуду, и в этом случае сушильщик будет выступать и как потребитель, как и производитель, и теперь есть две общие рабочие очереди (каждая из которых может блокировать работу сушильщика от продолжения.)

Блокирующие очереди упрощают кодирование потребителей, так как принимают блоки до тех пор, пока данные доступны. Если производители генерируют работу не достаточно быстро, чтобы занять потребителей, потребители просто ждут, пока не станет доступно больше работы. Иногда это вполне приемлемо (как в серверном приложении, когда ни один клиент не обращается к службе), а иногда это указывает на то, что отношение потоков-производителей к потокам-потребителям должно быть скорректировано для достижения лучшего использования ресурсов (как в веб-сканере или другом приложении, в котором работа фактически бесконечна).

Если производители последовательно генерируют работу быстрее, чем потребители могут ее обработать, это в конечном итоге приведёт к тому, что у приложения закончится память, потому что рабочие элементы будут добавляться в очередь без ограничений. Опять же, блокирующая природа метода `put` значительно упрощает кодирование производителей; если мы используем *ограниченную очередь*, то, когда очередь заполняет блок производителей, потребителям предоставляется время нагнать их, потому что заблокированный производитель больше не может генерировать элементы работы.

Блокирующие очереди также предоставляют метод `offer`, который возвращает статус сбоя, если элемент не может быть помещён в очередь. Это позволяет создавать более гибкие политики для работы с перегрузкой, такие как сброс нагрузки, сериализация избыточных рабочих элементов и запись их на диск, уменьшение числа потоков производителя или регулирование работы производителей каким-либо другим способом.

Ограниченные очереди являются мощным инструментом управления ресурсами для создания надежных приложений: они делают вашу программу более устойчивой к перегрузке, регулируя работу активностей, рискующих произвести больше работы, чем может быть обработано.

В то время как шаблон “производитель-потребитель” позволяет отделить код производителя и потребителя друг от друга, их поведение по-прежнему косвенно связано через общую рабочую очередь. Представляется заманчивым предположить, что потребители всегда будут следить за тем, чтобы вам не нужно было размещать какие-либо ограничения на размер рабочих очередей, но это рецепт для последующего перепроектирования вашей системы. *Стройте управление ресурсами в своем дизайне, заблаговременно используя блокирующие очереди - намного проще это сделать сначала, чем позже модифицировать его.* Блокирующие очереди упрощают дело в ряде ситуаций, но если блокирующие очереди не вписываются в ваш дизайн, вы можете создавать другие блокирующие структуры данных, используя класс `Semaphore` (см. раздел [5.5.3](#)).

Библиотека классов содержит несколько реализаций интерфейса `BlockingQueue`. Классы `LinkedBlockingQueue` и `ArrayBlockingQueue` являются очередями FIFO, аналогичными классам `LinkedList` и `ArrayList`, но с лучшей параллельной производительностью, чем синхронизированные реализации интерфейса `List`. Класс `PriorityBlockingQueue` является очередью с приоритетом, которая полезна в случае, когда вы хотите обработать элементы в порядке, отличном от FIFO. Как и другие сортированные коллекции, `PriorityBlockingQueue` может сравнивать элементы в соответствии с их

естественным порядком (если они реализуют интерфейс Comparable) или с помощью интерфейса Comparator.

Последняя реализация интерфейса BlockingQueue, класс SynchronousQueue, на самом деле не очередь вообще, в том смысле, что он не поддерживает пространство для хранения элементов в очереди. Вместо этого он поддерживает список *потоков* в очереди, ожидающих постановки в очередь или удаления элемента из очереди. По аналогии с мытьём посуды, это подобно тому, что шкафа для посуды нет, вместо этого, вымытая посуда вручается непосредственно следующему свободному сушильщику. Хотя это может показаться странным способом реализации очереди, он уменьшает задержку, связанную с перемещением данных от производителя к потребителю, поскольку элемент работы может быть передан непосредственно. (В традиционной очереди операции постановки в очередь и изъятия из очереди должны выполняться последовательно, прежде чем можно будет передать единицу работы.) Прямая передача также возвращает производителю больше информации о состоянии задачи; когда передача принята, он знает, что потребитель взял на себя ответственность за нее, а не просто позволяет ей находиться где-то в очереди - это очень похоже на разницу между тем, чтобы передать документ коллеге напрямую или просто положить его в её почтовый ящик и надеяться, что она получит его в ближайшее время. Поскольку класс SynchronousQueue не имеет емкости хранения, методы put и take будут блокироваться, за исключением ситуации, когда другой поток уже ожидает участия в передаче. Синхронные очереди, как правило, подходят только тогда, когда есть достаточное количество потребителей, которые почти всегда будут готовы принять передачу.

5.3.1 Пример: поиск в компьютере

Одним из типов программ, поддающихся декомпозиции на производителей и потребителей, является агент, который сканирует локальные диски на наличие документов и индексирует их для последующего поиска, подобно Google Desktop или службе Windows Indexing. В классе DiskCrawler в листинге 5.8, показана задача-производитель, которая ищет иерархию файлов удовлетворяющих критерию индексации и помещает их имена в рабочую очередь; класс Indexer в листинге 5.8 демонстрирует пример задачи-потребителя, которая принимает имена файлов из очереди и индексирует их.

```
public class FileCrawler implements Runnable {
    private final BlockingQueue<File> fileQueue;
    private final FileFilter fileFilter; private
    final File root;
    ...
    public void run() {
        try {
            crawl(root);
        } catch (InterruptedException e) {
            Thread.currentThread().interrupt();
        }
    }

    private void crawl(File root) throws InterruptedException {
        File[] entries = root.listFiles(fileFilter);
```

```

        if (entries != null) {
            for (File entry : entries)
                if (entry.isDirectory())
                    crawl(entry);
                else if (!alreadyIndexed(entry))
                    fileQueue.put(entry);
            }
        }
    }

public class Indexer implements Runnable {
    private final BlockingQueue<File> queue;

    public Indexer(BlockingQueue<File> queue) {
        this.queue = queue;
    }

    public void run() {
        try {
            while (true)
                indexFile(queue.take());
        } catch (InterruptedException e) {
            Thread.currentThread().interrupt();
        }
    }
}

```

Листинг 5.8 Задачи производителя и потребителя в приложении поиска на компьютере

Шаблон “производитель-потребитель” предлагает дружественный к потокам способ декомпозиции проблемы поиска в компьютере в более простые компоненты. Факторинг⁶⁶ сканера файлов и индексирования в отдельные активности порождает лучше читаемый и повторно используемый код, чем в том случае, если бы действие было монолитным и выполняло обе операции; каждая из активностей отвечает за выполнение только одной задачи, и блокирующая очередь берёт на себя всё управление потоками, таким образом, код каждой активности становится проще и яснее.

Шаблон “производитель-потребитель” также включает несколько преимуществ в повышении производительности. Производители и потребители могут выполняться одновременно; если один ограничен операциями ввода/вывода, а другой ограничен центральным процессором, одновременное выполнение улучшает общую пропускную способность, по сравнению с последовательным выполнением. Если активности производителя и потребителя параллелизуются в разной степени, их тесная связь снижает параллелизуемость по сравнению с менее параллелизуемыми активностями.

В листинге 5.9 запускается несколько сканеров и индексаторов, каждый в своем потоке. Как и написано, потоки потребителя никогда не завершаются, что предотвращает завершение программы; мы рассмотрим несколько методов

⁶⁶ В **математике** факторизация или **факторинг** — это декомпозиция объекта (например, числа, полинома или матрицы) в произведение других объектов или факторов, которые, будучи перемноженными, дают исходный объект.

решения этой проблемы в главе 7. Хотя в этом примере используются явно управляемые потоки, многие проекты вида “производитель-потребитель” могут быть выражены с помощью фреймворка выполнения задач Executor, который сам использует шаблон “производитель-потребитель”.

```
public static void startIndexing(File[] roots) {  
    BlockingQueue<File> queue = new LinkedBlockingQueue<File>(BOUND);  
    FileFilter filter = new FileFilter() {  
        public boolean accept(File file) { return true; }  
    };  
  
    for (File root : roots)  
        new Thread(new FileCrawler(queue, filter, root)).start();  
  
    for (int i = 0; i < N_CONSUMERS; i++)  
        new Thread(new Indexer(queue)).start();  
}
```

Листинг 5.9 Запуск поиска в компьютере

5.3.2 Последовательное ограничение потока

Все реализации блокирующих очередей в пакете java.util.concurrent содержат достаточную внутреннюю синхронизацию для безопасной публикации объектов от потока-производителя к потоку-потребителю.

Для изменяемых объектов шаблон “производитель-потребитель” и блокирующие очереди облегчают *последовательное ограничение потоков*, в целях передачи прав владения на объекты от производителей к потребителям. Объект, ограниченный потоком, принадлежит исключительно одному потоку, но это право собственности может быть “передано” путем безопасной публикации, когда только один поток получит доступ к нему и гарантирует, что публикующий поток не получит доступ к объекту после передачи. Безопасная публикация гарантирует, что состояние объекта будет видно новому владельцу, и поскольку исходный владелец больше не будет к нему прикасаться, теперь объект будет ограничен новым потоком. Новый владелец может изменять его свободно, так как имеет к нему эксклюзивный доступ.

Пулы объектов используют последовательное ограничение потока, “одолживая” объект запрашивающему потоку. Пока пул содержит достаточно внутренней синхронизации для безопасной публикации объекта пула, и пока клиенты сами не опубликуют объект пула или не будут использовать его после возвращения в пул, право владения может безопасно передаваться из потока в поток.

Можно также использовать другие механизмы публикации для передачи права собственности на изменяемый объект, но необходимо убедиться, что только один поток получает передаваемый объект. Блокирующие очереди упрощают этот процесс; затратив немного усилий, тоже самое можно сделать с помощью атомарного метода remove класса ConcurrentMap или метода compareAndSet класса AtomicReference.

5.3.3 Двусторонние очереди и кража работы

В Java 6 также было добавлено еще два типа коллекций: интерфейсы Deque (произносится как «дек») и BlockingDeque, которые расширяют интерфейсы

`Queue` и `BlockingQueue`. Интерфейс `Deque` - это двусторонняя очередь, которая позволяет эффективно вставлять и удалять элементы, как с головы, так и с хвоста. Реализации включают классы `ArrayDeque` и `LinkedBlockingDeque`.

Подобно тому, как блокирующие очереди “одалживают” себя с помощью шаблона “производитель-потребитель”, двусторонние очереди “одалживают” себя подобным образом, с помощью шаблона называемого “*кражей работы*” (*work stealing*). В дизайне “производитель-потребитель” используется одна общая рабочая очередь для всех потребителей; в дизайне “кража работы” у каждого потребителя есть своя собственная двусторонняя очередь. Если потребитель исчерпывает работу в своей собственной двусторонней очереди, он может “украсть” работу с *хвоста* чужой двусторонней очереди. “Кража работы” может быть более масштабируемой, чем традиционный дизайн “производитель-потребитель”, поскольку работники не конкурируют за общую рабочую очередь; большую часть времени они получают доступ только к своей собственной двусторонней очереди, уменьшая конкуренцию. Когда один работник должен получить доступ к очереди другого работника, он делает это с хвоста, а не с головы, что еще больше уменьшает конкуренцию.

“Кража работы” хорошо подходит для проблем, в которых потребители также являются производителями - когда выполнение единицы работы, возможно, приведет к идентификации большего количества работы. Например, обработка страницы в веб-сканере обычно приводит к идентификации новых страниц для обхода. Аналогичным образом, многие алгоритмы исследования графов, например, такой как пометка кучи (*heap*) во время сборки мусора, могут быть эффективно распараллелены с помощью “кражи работы”. Когда работник определяет новую единицу работы, он помещает ее в конец своей собственной двусторонней очереди (или, в качестве альтернативы, в дизайне “совместного использования работы” (*work sharing*), перекидывает на другого работника); когда его двусторонняя очередь пустеет, он ищет работу в конце чужой двусторонней очереди, гарантируя таким образом, что каждый работник будет занят.

5.4 Методы блокирования и прерывания

Потоки могут блокироваться или приостанавливаться по нескольким причинам: ожидание завершения операций ввода/вывода, ожидание получения блокировки, ожидание пробуждения из метода `Thread.sleep` или ожидание результатов вычислений других потоков. Когда поток блокируется, он обычно приостанавливается и переводится в одно из состояний блокировки (`BLOCKED`, `WAITING` или `TIMED_WAITING`). Различие между блокирующей операцией и обычной, только требующей длительного времени для завершения выполнения, состоит в том, что заблокированный поток должен ожидать события, которое находится вне его контроля, прежде чем сможет продолжить своё выполнение – например, завершение операции ввода/вывода, блокировка становится доступной или завершаются внешние (по отношению к потоку) вычисления. При возникновении такого внешнего события поток переходит в состояние `RUNNABLE` и снова становится доступен для планирования.

Методы `put` и `take` интерфейса `BlockingQueue` бросают проверяемое исключение `InterruptedException`, как и несколько других библиотечных методов, таких как `Thread.sleep`. Когда метод может бросить исключение `InterruptedException`, это говорит вам о том, что это - блокирующий метод, и

что в дальнейшем, если его выполнение будет *прервано*, это приведёт к раннему прекращению действия блокировки.

Класс `Thread` предоставляет метод `interrupt` для прерывания потока и для запроса, был ли поток прерван. Каждый поток имеет свойство типа `boolean`, содержащее информацию о том, был ли поток прерван; прерывание потока устанавливает его состояние.

Прерывание - *кооперативный* механизм. Один поток не может заставить другой прекратить то, что он делает, и сделать что-то иное; когда поток *A* прерывает поток *B*, поток *A* просто отправляет запрос, чтобы поток *B* остановил то, что он делает, когда достигнет удобной точки для остановки. Хотя в API или спецификации языка нет ничего, что требовало бы какой-либо конкретной семантики уровня приложения для выполнения прерывания, наиболее разумным использованием прерывания является отмена некоторого действия. Методы блокировки, реагирующие на прерывания, упрощают своевременную отмену операций, выполняющихся длительное время.

Когда ваш код вызывает метод, который бросает исключение `InterruptedException`, ваш метод также будет являться блокирующими, и должен иметь план реакции на прерывание. Для библиотечного кода, в основном, существует два варианта:

Передача `InterruptedException` выше по стеку. Часто, наиболее разумной политикой будет передача дальше – просто распространите исключение `InterruptedException` до вызывающего кода. Это может подразумевать отсутствие перехвата исключения `InterruptedException`, или перехват и повторное возбуждение, после выполнения некоторой краткой, специфичной для активности, чистки.

Восстановление прерывания. Иногда вы не можете бросить исключение `InterruptedException`, например, в том случае, когда ваш код является частью интерфейса `Runnable`. В такой ситуации необходимо перехватить исключение `InterruptedException` и восстановить состояние прерывания, путём вызова метода `interrupt` в текущем потоке, чтобы код выше по стеку вызовов мог увидеть, что было выполнено прерывание, как показано в листинге 5.10.

Вы можете получить гораздо более сложные варианты с прерыванием, но эти два подхода должны работать в подавляющем большинстве ситуаций. Есть одна вещь, которую вы *не* должны делать с исключением `InterruptedException` – перехватывать его и ничего не делать в ответ. Это лишает код, находящийся выше в стеке вызовов, возможности действовать при прерывании, поскольку теряется доказательство того, что поток был прерван. *Единственная ситуация, в которой приемлемо проглотить прерывание, – это когда вы расширяете класс `Thread` и, следовательно, контролируете весь код выше по стеку вызовов.* Отмена и прерывание более подробно рассматриваются в главе 7.

```
public class TaskRunnable implements Runnable {  
    BlockingQueue<Task> queue;  
    ...  
    public void run() {  
        try {
```

```

        processTask(queue.take());
    } catch (InterruptedException e) {
        // restore interrupted status
        Thread.currentThread().interrupt();
    }
}
}

```

Листинг 5.10 Восстановление состояния `interrupted`, чтобы не “проглотить” прерывание

5.5 Синхронизаторы

Блокирующие очереди уникальны среди классов коллекций: они не только действуют как контейнеры для объектов, но и могут координировать управление потоками производителя и потребителя, поскольку методы `take` и `put` блокируются до тех пор, пока очередь не перейдет в желаемое состояние (не пустая или не полная).

Синхронизатор является любым объектом, который координирует поток управления (*control flow*) потоками на основе своего состояния. Блокирующие очереди могут действовать как синхронизаторы; другие типы синхронизаторов включают семафоры, барьеры и защелки. В библиотеке платформы есть несколько классов синхронизатора; если они не соответствуют вашим потребностям, вы можете создать свой собственный синхронизатор, используя механизмы, описанные в главе 14.

Все синхронизаторы разделяют определенные структурные свойства: они инкапсулируют состояние, которое определяет, должны ли потоки, поступающие в синхронизатор, проходить или принудительно ожидать, предоставляя методы для манипулирования этим состоянием, и предоставляя методы для эффективного ожидания синхронизатора, с целью ввода желаемого состояния.

5.5.1 Защелки

Защелка (*latches*) - это синхронизатор, который может задерживать ход выполнения потоков до достижения ими *конечного* (*terminal*) состояния [CPI 3.4.2]. Защелка действует как затвор: до тех пор, пока защелка не достигнет конечного состояния, затвор будет закрыт и никакие потоки не смогут пройти, в конечном состоянии затвор открываются, позволяя пройти всем потокам. Как только защелка достигнет конечного состояния, она более не сможет изменять своё состояние, поэтому она навсегда останется открытой. Защелки можно использовать для обеспечения того, чтобы определенные действия не выполнялись до завершения других разовых действий, таких как:

- Обеспечение того, чтобы вычисление не выполнялось до тех пор, пока необходимые ему ресурсы не будут инициализированы. Простую бинарную (двухпозиционную) защелку можно использовать, чтобы указать, что “Ресурс *R* был инициализирован”, и любая активность, которая нуждается в ресурсе *R*, будет сначала ожидать на этой защелке.
- Обеспечение того, чтобы служба не запускалась до тех пор, пока не запустятся другие службы, от которых она зависит. У каждой службы есть ассоциированная бинарная защелка; запуск службы *S* сперва, будет включать ожидание на защелках других служб, от которых зависит служба *S*, а затем

освобождение защелки службы S после завершения запуска других служб, чтобы затем все службы, зависящие от S , могли продолжить выполнение.

- Ожидание, пока все стороны, участвующие в некоторой деятельности, например игроки в многопользовательской игре, будут готовы продолжить. В этом случае защелка достигает конечного состояния после того, как все игроки готовы.

Класс `CountDownLatch` представляет собой гибкую реализацию защелки, которая может использоваться в любой из этих ситуаций; он позволяет одному или нескольким потокам ожидать возникновения набора событий. Состояние защелки состоит из счетчика, инициализированного положительным числом, представляющего количество ожидающих событий. Метод `countDown` уменьшает счетчик, указывая на то, что событие произошло, и методы `await` ожидают достижения счетчиком значения ноль, что является следствием возникновения всех ожидаемых событий. Если на входе счетчик отличен от нуля, выполнение метода `await` блокируется до тех пор, пока счётчик не достигнет нуля, ожидающий поток будет прерван или истечёт время ожидания.

Класс `TestHarness` из листинга 5.11, иллюстрирует два распространённых случая применения защелок. Класс `TestHarness` создает несколько потоков, которые одновременно выполняют переданное задание. Он использует две защелки, «начальный затвор» и «конечный затвор». Начальный затвор инициализируется значением «один», конечный затвор инициализируется счетчиком, равным количеству рабочих потоков. Первое, что делает каждый рабочий поток - выполняет ожидание на стартовом затворе; это гарантирует, что ни один из потоков не начнет выполнять работу, пока все они не будут готовы к запуску. Последнее, что делает каждый из них, - это уменьшение счётчика на конечном затворе; это позволяет главному потоку дождаться того момента, когда последний рабочий поток завершится, и он сможет вычислить прошедшее с момента запуска потоков время.

```
public class TestHarness {
    public long timeTasks(int nThreads, final Runnable task)
        throws InterruptedException {
        final CountDownLatch startGate = new CountDownLatch(1);
        final CountDownLatch endGate = new CountDownLatch(nThreads);

        for (int i = 0; i < nThreads; i++) {
            Thread t = new Thread() {
                public void run() {
                    try {
                        startGate.await();
                        try {
                            task.run();
                        } finally {
                            endGate.countDown();
                        }
                    } catch (InterruptedException ignored) { }
                }
            };
            t.start();
        }
    }
}
```

```
        }

        long start = System.nanoTime();
        startGate.countDown();
        endGate.await();
        long end = System.nanoTime();
        return end-start;
    }
}
```

Листинг 5.11 Использование класса CountDownLatch для запуска и остановки потоков в тестах на время

Почему мы побеспокоились о защелках в классе TestHarness вместо того, чтобы просто запускать потоки сразу после их создания? По видимому, мы хотели измерить, сколько времени потребуется для параллельного выполнения задачи n раз. Если бы мы просто создали и запустили потоки, потоки, начатые ранее, имели бы "фору" перед потоками, запущенными позднее, и степень конкуренции изменялась бы с течением времени по мере увеличения или уменьшения числа активных потоков. Использование стартового затвора позволяет главному потоку сразу отпустить все рабочие потоки, а конечный затвор позволяет главному потоку ждать завершения только последнего потока, а не ждать последовательного завершения каждого потока.

5.5.2 Класс FutureTask

Класс `FutureTask` также действует как защелка. (Класс `FutureTask` реализует интерфейс `Future`, описывающий абстрактный опорный результат (*result-bearing*) вычисления [CPJ 4.3.3].) Вычисление, представленное классом `FutureTask`, реализуется с использованием интерфейса `Callable`, возвращающим результат эквивалентом интерфейса `Runnable` и может находиться в одном из трех состояний: *ожидание выполнения*, *выполнение* или *завершено*. Под завершением понимаются все способы выполнения вычисления, включая обычное завершение, отмену и исключение. Как только экземпляр класса `FutureTask` переходит в завершенное состояние, он остается в этом состоянии навсегда.

Поведение метода `Future.get` зависит от состояния задачи. Если она завершена, метод `get` немедленно возвращает результат, в противном случае метод блокируется до тех пор, пока задача не перейдет в завершенное состояние, а затем не вернет результат или не бросит исключение. Класс `FutureTask` передает результат потока, выполняющего вычисление, в поток (и), получающий(е) результат; спецификация класса `FutureTask` гарантирует, что эта передача основана на безопасной публикации результата.

Класс `FutureTask` используется фреймворком `Executor` для представления асинхронных задач и может также использоваться для представления любых потенциально длительных вычислений, которые могут быть запущены до того, как их результаты понадобятся. Класс `Preloader` из листинга 5.12 использует класс `FutureTask` для выполнения долгостоящих вычислений, результаты которых понадобятся позже; начав вычисление на раннем этапе, вы сократите время, которое вам придется подождать, когда вам позже действительно понадобятся результаты.

```
public class Preloader {
```

```

private final FutureTask<ProductInfo> future =
    new FutureTask<ProductInfo>(new Callable<ProductInfo>() {
        public ProductInfo call() throws DataLoadException {
            return loadProductInfo();
        }
    });
private final Thread thread = new Thread(future);

public void start() { thread.start(); }

public ProductInfo get()
    throws DataLoadException, InterruptedException {
    try {
        return future.get();
    } catch (ExecutionException e) {
        Throwable cause = e.getCause();
        if (cause instanceof DataLoadException)
            throw (DataLoadException) cause;
        else
            throw launderThrowable(cause);
    }
}
}

```

Листинг 5.12 Использование `FutureTask` для предварительной загрузки данных, которые понадобятся позднее

Класс `Preloader` создает экземпляр `FutureTask`, который описывает задачу загрузки информации о продукте из базы данных и поток, в котором будет выполняться вычисление. Он предоставляет метод `start` для запуска потока, так как нецелесообразно запускать поток из конструктора или статического инициализатора. Когда программе позже понадобится экземпляр `ProductInfo`, она может вызвать метод `get`, который вернёт загруженные данные, если они готовы, или будет ожидать завершения загрузки, если они ещё не готовы.

Задачи, описанные интерфейсом `Callable`, могут бросать проверяемые и непроверяемые исключения, и любой код может бросать исключение `Error`. Независимо от того, какое исключение может бросить код задачи, оно обрамчивается исключением `ExecutionException` и перебрасывается из метода `Future.get`. Это усложняет код вызывающий метод `get`, и не только потому, что он должен иметь дело с возможностью возбуждения исключения `ExecutionException` (и непроверяемого исключения `CancellationException`), но также и потому, что причина исключения `ExecutionException` возвращается как исключение `Throwable`, что неудобно в использовании.

Когда метод `get` класса `Preloader` бросает исключение `ExecutionException`, причина исключения попадает в одну из трех категорий: проверяемое исключение брошенное интерфейсом `Callable`, `RuntimeException` или `Error`. Мы должны обрабатывать каждый из этих случаев отдельно, но в листинге 5.13 мы будем использовать метод-утилиту `launderThrowable`, чтобы инкапсулировать “грязную” часть логики обработки исключений. Перед вызовом метода `launderThrowable`, класс `Preloader` тестирует исключения на соответствие известным проверяемым исключениям и, в случае успеха, пробрасывает их дальше. Таким образом, остаются только непроверяемые исключения, которые класс

`Preloader` обрабатывает, вызывая метод `launderThrowable` и бросая результат. Если исключение `Throwable`, переданное методу `launderThrowable`, является экземпляром `Error`, метод `launderThrowable` пробрасывает его напрямую; если это не исключение `RuntimeException`, он бросает исключение `IllegalStateException`, чтобы указать на логическую ошибку. Таким образом, остаётся только исключение `RuntimeException`, которое метод `launderThrowable` возвращает вызывающему методу и которое вызывающий метод, как правило, пробрасывает дальше.

```
/* If the Throwable is an Error, throw it; if it is a
 * RuntimeException return it, otherwise throw IllegalStateException
 */
public static RuntimeException launderThrowable(Throwable t) {
    if (t instanceof RuntimeException)
        return (RuntimeException) t;
    else if (t instanceof Error)
        throw (Error) t;
    else
        throw new IllegalStateException("Not unchecked", t);
}
```

Листинг 5.13 Приведение непроверяемого исключения `Throwable` к типу `RuntimeException`.

5.5.3 Семафоры

Подсчет семафоров используется для управления количеством активностей, которые могут обращаться к определенному ресурсу или выполнять определенное действие в некоторый момент времени[CPJ 3.4.1]. Подсчет семафоров может использоваться для реализации пулов ресурсов или для наложения ограничений на коллекцию.

Класс `Semaphore` управляет набором виртуальных разрешений; начальное количество разрешений передается в конструкторе класса `Semaphore`. Активности могут приобретать разрешения (до тех пор, пока есть в наличии) и освобождать разрешения, когда выполняют свою работу. Если разрешение не доступно, метод `acquire` блокируется до тех пор, пока не получит его (или до тех пор, пока не будет прерван или пока не истечёт время выполнения операции). Метод `release` возвращает разрешение семафору⁶⁷.

Вырожденным случаем счетного семафора является бинарный семафор, то есть, семафор с начальным значением счётчика равным единице. Бинарный семафор может использоваться как мьютекс с семантикой нереентабельной блокировки; тот, кто имеет единственное разрешение, держит мьютекс.

Семафоры полезны для реализации пулов ресурсов, таких как пулы соединений с базами данных. Довольно легко построить пул фиксированного размера, который упадёт с ошибкой, если вы запросите ресурс из пустого пула, но что вы действительно хотите, так это заблокировать операцию получения, если пул пуст и вновь разблокировать её, как только появятся данные. Если вы инициализируете класс `Semaphore` размером пула, с помощью метода `acquire` приобретите

⁶⁷ Реализация не имеет реальных объектов разрешений, и класс `Semaphore` не ассоциирует выданные разрешения с потоками, поэтому разрешение, приобретённое в одном потоке, может быть освобождено из другого потока. Вы можете думать о методе `acquire` как о потреблении разрешения и методе `release`, как о его создании; класс `Semaphore` не ограничен количеством разрешений, с которым он был создан.

разрешение, прежде чем попытается извлечь ресурс из пула, и отпустите разрешение с помощью метода `release`, после помещения ресурса в пул, метод `acquire` блокируется до тех пор, пока пул остаётся пустым. Этот подход используется в классе ограниченного буфера, в главе 12. (Более простой способ создания блокирующего пула объектов - использовать класс `BlockingQueue` для хранения полученных ресурсов.)

Аналогичным образом, вы можете использовать класс `Semaphore`, чтобы превратить любую коллекцию в блокирующую ограниченную коллекцию, как показано в классе `BoundedHashSet`, в листинге 5.14. Семафор инициализируется желаемым максимальным значением размера коллекции. Операция `add` приобретает разрешение перед добавлением элемента в базовую коллекцию. Если базовая операция `add` фактически ничего не добавляет, она немедленно освобождает разрешение. Аналогичным образом, успешное выполнение операции `remove` освобождает разрешение, позволяя добавлять дополнительные элементы. Базовая реализация интерфейса `Set` ничего не знает об ограничении; всё это ложится на класс `BoundedHashSet`.

```
public class BoundedHashSet<T> {
    private final Set<T> set;
    private final Semaphore sem;

    public BoundedHashSet(int bound) {
        this.set = Collections.synchronizedSet(new HashSet<T>());
        sem = new Semaphore(bound);
    }

    public boolean add(T o) throws InterruptedException {
        sem.acquire();
        boolean wasAdded = false;
        try {
            wasAdded = set.add(o);
            return wasAdded;
        }
        finally {
            if (!wasAdded)
                sem.release();
        }
    }

    public boolean remove(Object o) {
        boolean wasRemoved = set.remove(o);
        if (wasRemoved)
            sem.release();
        return wasRemoved;
    }
}
```

Листинг 5.14 Использование класса `Semaphore` для ограничения коллекции

5.5.4 Барьеры

Мы видели, как защёлки могут облегчить запуск исполнения группы связанных активностей или ожидание завершения исполнения группы связанных активностей. Защёлки - одноразовые объекты; как только защелка переходит в терминальное состояние, ее нельзя сбросить.

Барьеры похожи на защелки тем, что они блокируют группу потоков до тех пор, пока не произойдет какое-либо событие [CPJ 4.4.3]. Ключевое отличие барьера от защёлки состоит в том, что все потоки должны совместно подойти к барьере *в некоторый момент времени*, только после этого они смогут продолжить выполнение. Защелки используются для ожидания событий; барьеры используются для ожидания других потоков. Барьер реализует протокол, который некоторые семьи используют для встречи в течение дня в торговом центре: “все встречаются в Макдональдсе в 6:00; как только вы доберетесь туда, оставайтесь там, пока все не появятся, а затем мы решим, что будем делать дальше”.

Класс `CyclicBarrier` позволяет фиксированному числу участников неоднократно встречаться у *барьера* и полезен в параллельных итерационных алгоритмах, которые разбивают задачу на фиксированное число независимых подзадач. При достижении барьера потоки вызывают метод `await`, и метод `await` блокирует выполнение потоков до тех пор, пока все из них не достигнут барьера. Если все потоки приходят к барьеру, то барьер переходит в состояние “успешно пройден”, в этом случае все потоки освобождаются и барьер сбрасывается, поэтому его можно использовать вновь. Если истекает время вызова метода `await` или поток заблокированный вызовом метода `await` прерывается, то барьер считается сломанным и все исходящие вызовы метода `await` будут завершены с возбуждением исключения `BrokenBarrierException`. Если барьер успешно пройден, метод `await` возвращает уникальный индекс прихода для каждого потока, который можно использовать для “выбора” лидера, выполняющего особые действия в следующей итерации. Класс `CyclicBarrier` также позволяет передать конструктору *действие барьера* (*barrier action*); это экземпляр интерфейса `Runnable`, который выполняется (в одной из подзадач потоков), когда барьер успешно пройден, но до того, как заблокированные потоки будут освобождены.

Барьеры часто используются в моделировании, где работа по вычислению одного шага может выполняться параллельно, но вся работа, связанная с заданным шагом, должна завершиться до перехода на следующий шаг. Например, при n -мерном моделировании частиц, каждый шаг вычисляет изменения в местоположении каждой частицы на основе местоположений и других атрибутов прочих частиц. Ожидание у барьера гарантирует, что все обновления для шага k завершаться прежде, чем станет доступной возможность перейти к шагу $k + 1$.

Класс `CellularAutomata` в листинге 5.15 демонстрируют использование барьера для компьютерного моделирования клеточного автомата, например такого, как игра Конвея в жизнь (Gardner, 1970). При распараллеливании моделирования, как правило, нецелесообразно назначать отдельный поток для каждого элемента (в случае игры “Жизнь”, это клетка); это потребует слишком много потоков, и накладные расходы на их координацию сильно замедлят вычисления. Вместо этого имеет смысл *разбить* проблему на несколько частей, позволить каждому потоку заниматься только своей частью, а затем объединить результаты. Класс `CellularAutomata` разделяет доску на N_{CPU} частей, где N_{CPU} - это количество

доступных процессоров, и назначает каждую часть отдельному потоку⁶⁸. На каждом шаге рабочие потоки вычисляют новые значения для всех ячеек в своей части доски. Когда все рабочие потоки достигают барьера, действие барьера фиксирует новые значения в модели данных. После выполнения действия барьера, рабочие потоки освобождаются для вычисления следующего шага вычисления, который включает в себя сверку с результатом выполнения метода `isDone`, с целью определения, требуются ли дальнейшие итерации.

```
public class CellularAutomata {  
    private final Board mainBoard;  
    private final CyclicBarrier barrier;  
    private final Worker[] workers;  
  
    public CellularAutomata(Board board) {  
        this.mainBoard = board;  
        int count = Runtime.getRuntime().availableProcessors();  
        this.barrier = new CyclicBarrier(count,  
            new Runnable() {  
                public void run() {  
                    mainBoard.commitNewValues();  
                }});  
        this.workers = new Worker[count];  
        for (int i = 0; i < count; i++)  
            workers[i] = new Worker(mainBoard.getSubBoard(count, i));  
    }  
  
    private class Worker implements Runnable {  
        private final Board board;  
  
        public Worker(Board board) { this.board = board; }  
        public void run() {  
            while (!board.hasConverged()) {  
                for (int x = 0; x < board.getMaxX(); x++)  
                    for (int y = 0; y < board.getMaxY(); y++)  
                        board.setNewValue(x, y, computeValue(x, y));  
                try {  
                    barrier.await();  
                } catch (InterruptedException ex) {  
                    return;  
                } catch (BrokenBarrierException ex) {  
                    return;  
                }  
            }  
        }  
    }  
}
```

⁶⁸ Для вычислительных задач, подобных этой, не осуществляющих операции ввода/вывода и не обращающихся к общим данным, количество потоков равное N_{cpu} или $N_{cpu} + 1$ даёт оптимальную пропускную способность; большее количество потоков не помогает и, фактически, может ухудшить производительность, поскольку потоки начинают конкурировать за ресурсы ЦП и памяти.

```
public void start() {  
    for (int i = 0; i < workers.length; i++)  
        new Thread(workers[i]).start();  
    mainBoard.waitForConvergence();  
}  
}
```

Листинг 5.15 Координация вычислений в клеточном автомате с использованием класса CyclicBarrier.

Другой формой барьера является класс Exchanger - двусторонний барьер, в котором стороны обмениваются данными у барьера [CPJ 3.4.3]. Обменники полезны, когда стороны выполняют асимметричные действия, например, когда один поток заполняет буфер данными, а другой поток использует данные из буфера; эти потоки могут использовать класс Exchanger для встречи и обмена полного буфера на пустой. Когда два потока обмениваются объектами через класс Exchanger, обмен представляет собой безопасную публикацию обоих объектов другой стороне.

Время обмена зависит от требований к скорости отклика приложения. Простейший подход заключается в том, что задача заполнения обменивается, когда буфер заполнен, и задача очистки обменивается, когда буфер пуст; это минимизирует количество обменов, но может задержать обработку некоторых данных, если скорость прибытия новых данных непредсказуема. Другой подход заключается в том, что заполнитель обменивается, когда буфер заполнен, но также и в том случае, когда буфер частично заполнен и прошло определенное время.

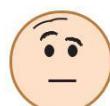
5.6 Создание эффективного и масштабируемого кэша результатов

Почти каждое серверное приложение использует некоторую форму кэширования. Повторное использование результатов предыдущих вычислений позволяет уменьшить задержки и увеличить пропускную способность, за счет дополнительного использования памяти.

Подобно многим другим часто изобретаемым велосипедам, кэширование обычно выглядит проще, чем есть на самом деле. Наивная реализация кэша, вероятно, превратит узкое место производительности в узкое место масштабируемости, даже если это улучшит однопоточную производительность. В этом разделе мы разрабатываем эффективный и масштабируемый кэш результатов для ресурсоемкой функции. Начнем с очевидного подхода – простого объекта HashMap – а затем рассмотрим некоторые недостатки его параллелизма и способы их устранения.

Интерфейс Computable<A, V> в листинге 5.16 описывает функцию с входными данными типа A и результатом типа V. Класс ExpensiveFunction, реализующий интерфейс Computable, тратит значительное количество времени на вычисление результатов; мы хотели бы создать обёртку для интерфейса Computable, которая бы запомнила результаты предыдущих вычислений и инкапсулировала процесс кэширования. (Эта техника называется *мемоизацией* (*memorization*))

```
public interface Computable<A, V> {  
    V compute(A arg) throws InterruptedException;  
}
```



```

public class ExpensiveFunction
    implements Computable<String, BigInteger>
{ public BigInteger compute(String arg) {
    // after deep thought...
    return new BigInteger(arg);
}
}

public class Memoizer1<A, V> implements Computable<A, V>
{ @GuardedBy("this")
private final Map<A, V> cache = new HashMap<A,
V>(); private final Computable<A, V> c;

public Memoizer1(Computable<A, V> c)
{ this.c = c;
}

public synchronized V compute(A arg) throws InterruptedException {
    V result = cache.get(arg);
    if (result == null) {
        result = c.compute(arg);
        cache.put(arg, result);
    }
    return result;
}
}

```

Листинг 5.16 Начальная попытка кэширования с использованием класса `HashMap` и синхронизации.

В классе `Memoizer1` в листинге 5.16 показана первая попытка: использование класса `HashMap` для хранения результатов предыдущих вычислений. Метод `compute` сначала проверяет, закэширован ли желаемый результат, и возвращает предварительно вычисленное значение, если это так. В противном случае результат вычисляется и перед возвратом кэшируется в экземпляре `HashMap`.

Класс `HashMap` не является потокобезопасным, поэтому для обеспечения того, чтобы два потока не обращались к экземпляру `HashMap` одновременно, класс `Memoizer1` использует консервативный подход, заключающийся в синхронизации всего метода `compute`. Такой подход гарантирует потокобезопасность, но имеет очевидную проблему с масштабируемостью: в общем случае, только один поток за раз может выполнять вычисления. Если другой поток занят вычислением результата, другие потоки, вызывающие метод `compute`, могут быть заблокированы на длительное время. Если несколько потоков находятся в очереди, ожидая вычисления значений, которые еще не вычислены, вычисление может занять больше времени, чем без использования мемоизации. На рисунке 5.2 иллюстрируется, что может произойти, когда несколько потоков пытаются использовать функцию, мемоизированную с помощью такого подхода. Это явно не тот сорт повышения производительности, которого мы надеялись достичь благодаря кэшированию.

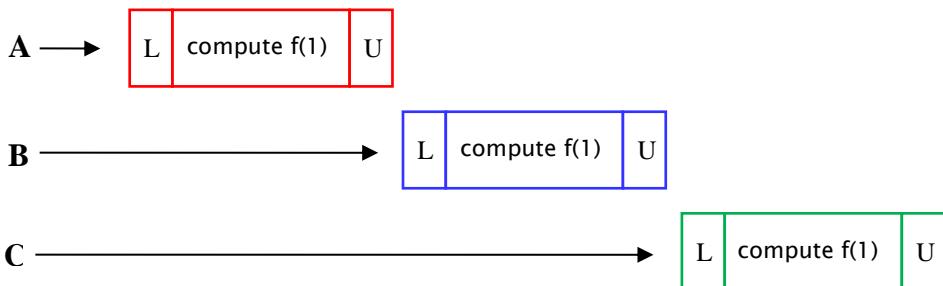


Рисунок 5.2 Плохой параллелизм в классе Memoizer1

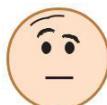
Класс Memoizer2 из листинга 5.17, улучшает ужасное параллельное поведение Memoizer1, путём замены класса HashMap на класс ConcurrentHashMap. Поскольку класс ConcurrentHashMap потокобезопасен, нет необходимости обеспечивать синхронизацию при доступе к кэширующему объекту Map, тем самым устраняется сериализация, вызванная синхронизацией метода compute в классе Memoizer1.

Класс Memoizer2, безусловно, обеспечивает лучшее параллельное поведение, чем класс Memoizer1: несколько потоков могут использовать его одновременно. Но он все еще, выступая в качестве кэша, обладает некоторыми недостатками - существует окно уязвимости, в котором два потока, одновременно вызывающие метод compute, могут вычислить одно и то же значение. В случае использования принципа мемоизации это просто неэффективно - предназначение кэша состоит в том, чтобы предотвратить многократное вычисление одинаковых данных. Для более универсального механизма кэширования дела обстоят намного хуже; для кэша объектов, который должен обеспечивать однократную инициализацию, эта уязвимость также представляет угрозу безопасности.

```
public class Memoizer2<A, V> implements Computable<A, V> {
    private final Map<A, V> cache = new ConcurrentHashMap<A, V>();
    private final Computable<A, V> c;

    public Memoizer2(Computable<A, V> c) { this.c = c; }

    public V compute(A arg) throws InterruptedException {
        V result = cache.get(arg);
        if (result == null) {
            result = c.compute(arg);
            cache.put(arg, result);
        }
        return result;
    }
}
```



Листинг 5.17 Замена класса HashMap классом ConcurrentHashMap

Проблема с классом Memoizer2 заключается в том, что если один поток начинает выполнение дорогостоящего вычисления, другие потоки не знают, что вычисление выполняется, и поэтому могут начать то же вычисление, как показано

на рисунке 5.3. Мы хотели бы как-то отметить, что “поток X в данный момент вычисляет значение $f(27)$ ”, так чтобы если другой поток прибудет для поиска значения $f(27)$, он знал, что самый эффективный способ найти значение, это встать в очередь за потоком X, дождаться завершения выполнения X, а потом спросить “Эй, что же вы получили при вычислении значения $f(27)$?”.

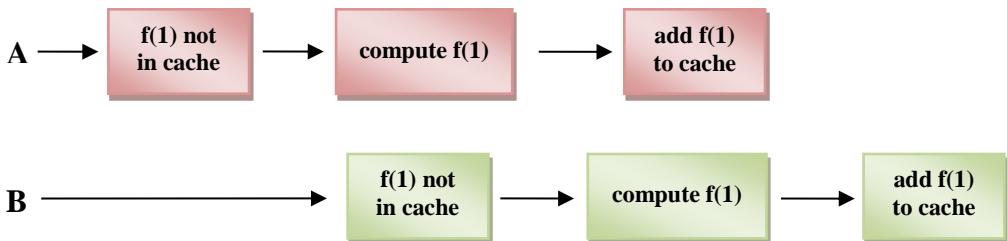


Рисунок 5.3 Два потока вычисляют одно и тоже значение, используя класс Memoizer2

Мы уже ранее встречали класс, который делает почти именно то, что мы хотим получить: `FutureTask`. Класс `FutureTask` представляет собой вычислительный процесс, который может находиться в завершенном или незавершённом состоянии. Вызов `FutureTask.get` немедленно возвращает результат вычисления, если он доступен; в противном случае вызов блокируется до тех пор, пока результат не будет вычислен, а затем возвращает его.

Класс `Memoizer3` из листинга 5.18 переопределяет экземпляр интерфейса `Map`, используемый для хранения кэша, как `ConcurrentHashMap<A, Future<V>>` вместо `ConcurrentHashMap<A, V>`. Класс `Memoizer3` сначала проверяет, было ли запущено соответствующее вычисление (в отличие от принятого по умолчанию как “завершённое”, в классе `Memoizer2`). Если это не так, он создает экземпляр класса `FutureTask`, регистрирует его в экземпляре `Map` и запускает вычисление; в противном случае он ожидает результата существующего вычисления. Результат может быть доступен немедленно или находиться в процессе вычисления - но он прозрачен для объекта, вызывающего метод `Future.get`.

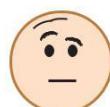
Реализация класса `Memoizer3` почти идеальна: она демонстрирует очень хороший параллелизм (в основном, полученный за счёт превосходной реализации параллелизма в классе `ConcurrentHashMap`), результат возвращается эффективно, если он уже известен, и если вычисление выполняется другим потоком, вновь прибывающие потоки терпеливо ожидают результата. У него есть только один недостаток - все еще существует небольшое окно уязвимости, в котором два потока могли бы начать вычисление одного и того же значения.

```

public class Memoizer3<A, V> implements Computable<A, V> {
    private final Map<A, Future<V>> cache
        = new ConcurrentHashMap<A, Future<V>>();
    private final Computable<A, V> c;

    public Memoizer3(Computable<A, V> c) { this.c = c; }

    public V compute(final A arg) throws InterruptedException {
        Future<V> f = cache.get(arg);
        if (f == null) {
            Callable<V> eval = new Callable<V>() {
  
```



```

        public V call() throws InterruptedException {
            return c.compute(arg);
        }
    };
    FutureTask<V> ft = new FutureTask<V>(eval);
    f = ft;
    cache.put(arg, ft);
    ft.run(); // call to c.compute happens here
}
try {
    return f.get();
} catch (ExecutionException e) {
    throw launderThrowable(e.getCause());
}
}
}

```

Листинг 5.18 Мемоизированная обёртка использующая класс FutureTask

Это окно намного меньше, чем в классе Memoizer2, но поскольку блок `if` в методе `compute` по прежнему является неатомарной последовательностью операций `проверить-затем-выполнить`, два потока могут вызвать метод `compute` с одинаковым значением приблизительно в одно и то же время, оба могут увидеть, что кэш не содержит желаемого значения, и оба могут начать вычисление. Пример такого неудачного момента показан на рисунке 5.4.

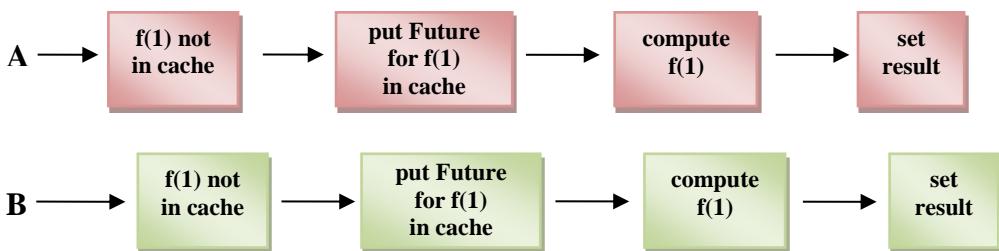


Рисунок 5.4 Неудачный момент времени, который может складываться в классе Memoizer3 и быть причиной двойного вычисления некоторого значения

Класс Memoizer3 уязвим для этой проблемы, потому что составное действие (*положить-если-отсутствует*) выполняется над базовым⁶⁹ экземпляром Map и не может быть сделано атомарным с помощью блокировки. Класс Memoizer из листинга 5.19 использует атомарный метод `putIfAbsent` интерфейса `ConcurrentMap`, закрывая окно уязвимости в классе Memoizer3.

Кэширование экземпляра Future вместо значения создает предпосылки к загрязнению кэша (*cache pollution*): если вычисление отменяется или завершается неудачей, дальнейшие попытки вычислить результат также будут указывать на отмену или сбой. Чтобы избежать этого, класс Memoizer удаляет экземпляр Future из кэша, если он обнаруживает, что вычисление было отменено; также может быть желательным удалять экземпляр Future при обнаружении возникновения исключения `RuntimeException`, в том случае, если есть шанс, что последующая попытка вычисления будет успешной.

⁶⁹ Речь идёт об экземпляре класса Map, содержащемся в классе ConcurrentHashMap.

Класс `Memoizer` также не рассматривает вопросы устаревания значений кэша (*cache expiration*), но это может быть реализовано с помощью подкласса `FutureTask`, который связывает время устаревания с каждым результатом и периодически просматривает кэш на наличие устаревших записей. (Подобным образом, не рассматривается вопрос вытеснения значений кэша, при котором старые записи удаляются, чтобы освободить место для новых, чтобы кэш не потреблял слишком много памяти.)

```
public class Memoizer<A, V> implements Computable<A, V> {
    private final ConcurrentMap<A, Future<V>> cache
        = new ConcurrentHashMap<A, Future<V>>();
    private final Computable<A, V> c;

    public Memoizer(Computable<A, V> c) { this.c = c; }

    public V compute(final A arg) throws InterruptedException {
        while (true) {
            Future<V> f = cache.get(arg);
            if (f == null) {
                Callable<V> eval = new Callable<V>() {
                    public V call() throws InterruptedException {
                        return c.compute(arg);
                    }
                };
                FutureTask<V> ft = new FutureTask<V>(eval);
                f = cache.putIfAbsent(arg, ft);
                if (f == null) { f = ft; ft.run(); }
            }
            try {
                return f.get();
            } catch (CancellationException e) {
                cache.remove(arg, f);
            } catch (ExecutionException e) {
                throw launderThrowable(e.getCause());
            }
        }
    }
}
```

Листинг 5.19 Финальная реализация класса `Memoizer`.

С завершением работы над нашей параллельной реализацией кэша, мы теперь можем добавить реальное кэширование к сервлету факторизации из главы 2, как и было обещано. Класс `Factorizer` из листинга 5.20 использует класс `Memoizer` для эффективного и масштабируемого кэширования ранее вычисленных значений.

```
@ThreadSafe
public class Factorizer implements Servlet {
    private final Computable<BigInteger, BigInteger[]> c =
        new Computable<BigInteger, BigInteger[]>() {
```

```
public BigInteger[] compute(BigInteger arg) {
    return factor(arg);
}
};

private final Computable<BigInteger, BigInteger[]> cache
= new Memoizer<BigInteger, BigInteger[]>(c);

public void service(ServletRequest req,
                    ServletResponse resp) {
    try {
        BigInteger i = extractFromRequest(req);
        encodeIntoResponse(resp, cache.compute(i));
    } catch (InterruptedException e) {
        encodeError(resp, "factorization interrupted");
    }
}
}
```

Листинг 5.20 Сервлет факторизации, кэширующий результаты с использованием класса Memoizer

Итоги части I

Мы рассмотрели очень много материала! В следующей “шпаргалке по параллелизму” кратко излагаются основные понятия и правила, представленные в Части 1.

- *Это изменчивое состояние, тутика.*⁷⁰
Все проблемы параллелизма сводятся к координации доступа к изменяющему состоянию. Чем менее доступно изменяющееся состояние, тем легче обеспечить потокобезопасность.
- *Сделайте поля final, если они не должны быть изменяющимися.*
Неизменяющиеся объекты значительно упрощают параллельное программирование. Они проще и безопаснее, и могут свободно использоваться без использования блокировки или защитного копирования.
- *Инкапсуляция позволяет управлять сложностью.*
Можно написать потокобезопасную программу со всеми данными, хранящимися в глобальных переменных, но зачем вам это? Инкапсуляция данных внутри объектов упрощает сохранение их инвариантов; инкапсуляция синхронизации внутри объектов упрощает соблюдение политики синхронизации.
- *Защищайте каждую изменяющуюся переменную с помощью блокировки.*
- *Защищайте все переменные инварианта одной и той же блокировкой.*
- *Удерживайте блокировки во время выполнения составных действий.*
- *Программа, которая обращается к изменяющейся переменной из нескольких потоков без использования синхронизации, является поврежденной программой.*
- *Не полагайтесь на умные рассуждения о том, почему вам не нужно использовать синхронизацию.*
- *Учитывайте потокобезопасность в процессе проектирования или явно задокументируйте, что класс не является потокобезопасным.*
- *Задокументируйте вашу политику синхронизации.*

⁷⁰ During the 1992 U.S. presidential election, electoral strategist James Carville hung a sign in Bill Clinton's campaign headquarters reading "The economy, stupid", to keep the campaign on message.

Часть II Структурирование параллельных приложений

Глава 6 Выполнение задач

Большинство параллельных приложений организованы вокруг выполнения *задач*: абстрактных, дискретных единиц работы. Разделение работы выполняемой приложением на задачи, упрощает организацию программы, облегчает восстановление после возникновения ошибок путём предоставления естественных границ транзакций, а также способствует параллелизму, предоставляя естественную структуру для распараллеливания работы.

6.1 Выполнение задач в потоках

Первым шагом в организации программы вокруг выполнения задачи является определение разумных *границ задачи* (*task boundaries*). В идеале задачи представляют собой независимые активности: работа, которая не зависит от состояния, результата или побочных эффектов, возникающих при выполнении других задач. Независимость облегчает реализацию параллелизма, так как независимые задачи могут выполняться параллельно при наличии достаточных вычислительных ресурсов. Для большей гибкости при планировании и балансировке нагрузки, каждая задача также должна представлять собой небольшую часть вычислительной мощности приложения.

Серверные приложения должны демонстрировать хорошую пропускную способность и высокую скорость отклика при нормальной нагрузке. Поставщики приложений хотят, чтобы приложения поддерживали как можно больше пользователей, чтобы снизить затраты на обеспечение работы для каждого пользователя; в свою очередь, пользователи хотят быстро получать ответы. Кроме того, приложения должны демонстрировать плавную деградацию по мере перегрузки, а не просто падать под большой нагрузкой. Грамотный выбор границ задачи в сочетании с разумной политикой выполнения задачи (см. раздел [6.2.2](#)) может помочь в достижении этих целей.

Большинство серверных приложений предлагают естественный выбор границ задач: индивидуальные запросы клиентов. Веб-серверы, почтовые серверы, файловые серверы, контейнеры EJB и серверы баз данных принимают запросы от удаленных клиентов через сетевые подключения. Использование отдельных запросов в качестве границ задачи обычно определяет как независимость, так и соответствующий размер задачи. Например, на результат отправки сообщения почтовому серверу не влияют другие сообщения, обрабатываемые в тот же момент времени, и обработка одного сообщения обычно требует очень малого процента от общей вычислительной емкости сервера.

6.1.1 Последовательное выполнение задач

Существует ряд возможных политик для планирования задач в приложении, некоторые из которых лучше других используют потенциал параллелизма. Проще всего выполнять задачи последовательно в одном потоке. Однопоточный Веб-сервер из листинга 6.1 обрабатывает свои задачи - HTTP-запросы поступающие на 80 порт - последовательно. Детали обработки запроса не так важны; нас интересуют характеристики параллелизма при использовании различных политик планирования.

```
class SingleThreadWebServer {  
    public static void main(String[] args) throws IOException {  
        ServerSocket socket = new ServerSocket(80);  
        while (true) {  
            Socket connection = socket.accept();  
            handleRequest(connection);  
        }  
    }  
}
```



Листинг 6.1. Веб-сервер с последовательной обработкой входящих запросов

Класс `SingleThreadedWebServer` прост и теоретически корректен, но плохо работает в продуктиве, поскольку может обрабатывать за раз только один запрос. Основной поток чередует операции принятия соединений и обработки связанных запросов. В то время как сервер обрабатывает запрос, новые соединения вынуждены ожидать, пока он не завершит обработку текущего запроса и вновь не вызовет метод `accept`. Такой подход мог бы иметь право на существование, если бы обработка запроса была настолько быстрой и эффективной, что метод `handleRequest` возвращал бы управление немедленно, но это не соответствует ни единому описанию веб-сервера в реальном мире⁷¹.

Обработка веб-запроса включает в себя сочетание вычислений и операций ввода/вывода. Сервер должен выполнить операции ввода/вывода с сокетом (`socket`) для чтения запроса и записи ответа, сокет может быть заблокирован из-за перегрузки сети или проблем с подключением. Он также может выполнять файловый ввод/вывод или выполнять запросы к базе данных, которые также могут блокироваться. В однопоточном сервере блокировка не только задерживает выполнение текущего запроса, но и препятствует обработке отложенных запросов. Если один запрос блокируется в течение необычно длительного промежутка времени, пользователи могут подумать, что сервер недоступен, поскольку он не отвечает. В то же время ресурсы используются крайне неэффективно, так как ЦП простояивает, в то время как единственный поток ждет завершения операций ввода/вывода для завершения обработки запроса.

В серверных приложениях последовательная обработка редко обеспечивает хорошую пропускную способность или высокую скорость отклика. Есть исключения - например, когда задач мало и они “долгоиграющие”, или когда сервер обслуживает одного клиента, который делает только один запрос за раз - но большинство серверных приложений не используют такой подход в своей работе⁷².

6.1.2 Явное создание потоков для задач

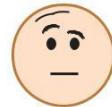
Более эффективный подход заключается в создании нового потока для обслуживания каждого входящего запроса, как показано в классе `ThreadPerTaskWebServer`, в листинге 6.2.

```
class ThreadPerTaskWebServer {
```

⁷¹ Речь идёт о продуктиве. Домашние странички в расчёт не берутся.

⁷² В некоторых ситуациях последовательная обработка может предложить преимущества простоты или безопасности; большинство фреймворков GUI обрабатывают задачи последовательно, используя единственный поток. Мы вернемся к последовательной модели в главе 9.

```
public static void main(String[] args) throws IOException {
    ServerSocket socket = new ServerSocket(80);
    while (true) {
        final Socket connection = socket.accept();
        Runnable task = new Runnable() {
            public void run() {
                handleRequest(connection);
            }
        };
        new Thread(task).start();
    }
}
```



Листинг 6.2 Веб сервер, запускающий новый поток для обработки каждого входящего запроса

Класс ThreadPerTaskWebServer схож по своей структуре с однопоточной версией - основной поток по-прежнему чередует прием входящих соединений и отправку запросов. Разница заключается в том, что для каждого входящего соединения главный цикл создает новый поток для обработки запроса, а не обрабатывает его в основном потоке. Это позволяет сделать три основных вывода:

- Обработка запросов выводится за рамки основного потока, позволяя главному циклу быстрее возвращаться к ожиданию последующих входящих соединений. Это позволяет принимать новые подключения до завершения обработки предыдущих запросов, что приводит к улучшению отзывчивости.
- Задачи могут обрабатываться параллельно, что позволяет обслуживать несколько запросов одновременно. Это может привести к повышению пропускной способности, если имеется несколько процессоров или в случае, если задачи по любым причинам зависят от блокировок, например завершение операции ввода/вывода, захват блокировки или ожидание доступности ресурсов.
- Код обработки задач должен быть потокобезопасным, так как он может вызываться одновременно для нескольких задач.

В свете умеренной нагрузки, подход “один поток на одну задачу” имеет преимущества по сравнению с последовательным выполнением. До тех пор, пока скорость поступления запросов не превысит возможности сервера по обработке запросов, этот подход обеспечивает лучшую отзывчивость и пропускную способность.

6.1.3 Недостатки неограниченного создания потоков

Однако, использование в продуктиве подхода “один поток на одну задачу” имеет некоторые практические недостатки, особенно при создании огромного числа потоков:

Накладные расходы на поддержание жизненного цикла. Создание потока и его завершение не даются даром. Фактические накладные расходы варьируются в зависимости от платформы, но создание потока занимает время, добавляю задержку ко времени обработки запроса, и требует некоторой обработки

активностями JVM и ОС. Если запросы выполняются часто и являются легковесными, как в большинстве серверных приложений, создание нового потока для каждого запроса может привести к потреблению значительных вычислительных ресурсов сервера.

Потребление ресурсов. Активные потоки потребляют системные ресурсы, особенно память. Если количество запускаемых потоков превышает количество доступных процессоров, потоки бездействуют. Множество свободных потоков может связать много оперативной памяти, приводя к повышению нагрузки на механизм сборщика мусора, а также, располагая множеством потоков, конкурирующих за процессоры, можно нарваться на другие эксплуатационные расходы. Если количество потоков достаточно для загрузки всех процессоров, создание большего количества потоков не только не поможет, а напротив, может даже ухудшить ситуацию.

Масштабируемость. Существует ограничение на количество создаваемых потоков. Ограничение зависит от платформы и зависит от таких факторов, как параметры вызова JVM, запрошенный размер стека в конструкторе потоков и ограничения на потоки, накладываемые базовой операционной системой⁷³. При достижении этого предела, наиболее вероятным результатом является возбуждение исключения `OutOfMemoryError`. Попытка восстановления после такой ошибки очень рискованна; гораздо проще структурировать программу так, чтобы избежать выхода за доступные пределы.

До определенного момента увеличение количества потоков может повышать пропускную способность, но при превышении некоторого уровня, создание большего числа потоков приведёт к замедлению работы приложения, а создание слишком большого числа потоков может привести к краху всего приложения. Чтобы избежать опасности, необходимо установить некоторые ограничения на количество потоков, создаваемых приложением, и тщательно протестировать приложение, чтобы гарантировать, что даже при достижении лимита ресурсы не закончатся.

Проблема подхода “один поток на одну задачу” заключается в том, что ничто не ограничивает количество созданных потоков, кроме скорости, с которой удаленные пользователи могут отправлять HTTP-запросы. Подобно прочим угрозам параллелизму, создание неограниченного количества потоков может работать очень хорошо во время прототипирования и в процессе разработки, но проблемы, как правило, возникают только при развертывании приложения в продуктиве и под большой нагрузкой. Таким образом, как злонамеренные пользователи, так и вполне рядовые, могут обрушить веб-сервер, если трафик когда-либо достигнет определенного порогового значения. Для серверного приложения, которое должно обеспечивать высокую доступность и плавное снижение нагрузки, это очень серьезная проблема.

⁷³ В 32-разрядных машинах, основным ограничивающим фактором является размер адресного пространства для стеков потока. Каждый поток поддерживает два стека выполнения: один для кода Java и один для машинного кода. Типичные значения по умолчанию для JVM порождают общий размер стека около 0.5 мегабайта (Вы можете изменить это значение с помощью флага JVM `-Xss` или с использованием конструктора потока). Если вы разделите размер стека потока на 2^{32} , то получите ограничение в несколько тысяч или десятков тысяч потоков. Другие факторы, такие как ограничения ОС, могут налагать более строгие ограничения.

6.2 Фреймворк Executor

Задачи - это логические единицы работы, а потоки - это механизм, с помощью которого задачи могут выполняться асинхронно. Мы рассмотрели две политики для выполнения задач с помощью потоков - выполнение задач последовательно в одном потоке и выполнение каждой задачи в своем собственном потоке. У обоих подходов есть серьезные ограничения: последовательный подход страдает от плохой отзывчивости и пропускной способности, а подход “один поток на одну задачу” страдает от плохого управления ресурсами.

В главе 5 мы увидели, как использовать *ограниченные очереди*, для предотвращения переполнения памяти приложения. Пулы потоков предлагают те же преимущества, только для управления потоками, а пакет `java.util.concurrent` обеспечивает гибкую реализацию пула потоков в рамках фреймворка Executor. Основной абстракцией для выполнения задач в библиотеках классов Java является *не поток*, а интерфейс `Executor`, показанный в листинге 6.3.

```
public interface Executor {  
    void execute(Runnable command);  
}
```

Листинг 6.3 Интерфейс Executor

Может быть, `Executor` и является простым интерфейсом, но он формирует основу гибкого и мощного фреймворка для асинхронного выполнения задач, поддерживающего широкий спектр политик выполнения задач. Он предоставляет стандартное средство, отделяющее *отправку задачи* от *выполнения задачи*, описывая задачи с помощью интерфейса `Runnable`. Реализации интерфейса `Executor` также обеспечивают поддержку жизненного цикла и хуки для добавления сбора статистики, управления приложениями и мониторинга.

Реализация интерфейса `Executor` основана на шаблоне производитель-потребитель, где активности, которые предоставляют задачи, являются производителями (производят единицы работы, которые должны быть выполнены), а выполняющие задачи потоки, являются потребителями (потребители этих единиц работы). *Использование интерфейса Executor обычно является самым простым способом реализации архитектуры производитель-потребитель в вашем приложении.*

6.2.1 Пример: веб-сервер использующий фреймворк Executor

Создать веб-сервер с помощью фреймворка `Executor` очень просто. Класс `TaskExecutionWebServer` из листинга 6.4 заменяет жестко закодированное создание потока использованием экземпляра интерфейса `Executor`. В этом случае мы используем одну из стандартных реализаций интерфейса `Executor` - пул фиксированного размера на 100 потоков.

```
class TaskExecutionWebServer {  
    private static final int NTHREADS = 100;  
    private static final Executor exec  
        = Executors.newFixedThreadPool(NTHREADS);
```

```
public static void main(String[] args) throws IOException {
    ServerSocket socket = new ServerSocket(80);
    while (true) {
        final Socket connection = socket.accept();
        Runnable task = new Runnable() {
            public void run() {
                handleRequest(connection);
            }
        };
        exec.execute(task);
    }
}
```

Листинг 6.4 Веб-сервер использующий пул потоков

В классе `TaskExecutionWebServer` отправка задачи обрабатывающей запрос отделена от её выполнения с помощью экземпляра `Executor`, и ее поведение можно изменить, просто подставив другую реализацию интерфейса `Executor`. Изменение реализаций или конфигурации экземпляра `Executor` гораздо менее агрессивный способ, чем изменение способа отправки задач; конфигурирование экземпляра `Executor`, как правило, является одноразовым действием и может быть легко введено в конфигурацию во время развертывания, в то время как код отправки задачи, как правило, разбросан по всей программе и, следовательно, изменить его значительно сложнее.

Мы можем легко изменить поведение класса `TaskExecutionWebServer` так, чтобы он начал вести себя также, как класс `ThreadPerTaskWebServer`, просто заменив реализацию интерфейса `Executor` вариантом, создающим новый поток для каждого запроса. Написать такую реализацию интерфейса `Executor` тривиальная задача, как показано в классе `ThreadPerTaskExecutor` в листинге 6.5.

```
public class ThreadPerTaskExecutor implements Executor {
    public void execute(Runnable r) {
        new Thread(r).start();
    }
}
```

Листинг 6.5 Реализация интерфейса `Executor`, запускающая для каждой задачи новый поток

Точно так же легко написать реализацию интерфейса `Executor`, которая заставляла бы класс `TaskExecutionWebServer` вести себя как же, как однопоточная версия, синхронно выполняющая каждую задачу перед возвратом из метода `execute`, как показано в классе `WithinThreadExecutor`, в листинге 6.6.

```
public class WithinThreadExecutor implements Executor {
    public void execute(Runnable r) {
        r.run();
    }
}
```

Листинг 6.6 Реализация интерфейса `Executor`, синхронно выполняющая задачи в вызывающем потоке

6.2.2 Политики выполнения

Значение отделения отправки задачи от её выполнения заключается в том, что оно позволяет легко определять и впоследствии без особых трудностей изменять политику выполнения для данного класса задач. Политика выполнения определяет "что, где, когда и как" выполняет задачи, в том числе:

- В каком потоке будет выполняться задача?
- В каком порядке должны выполняться задачи (FIFO, LIFO, в порядке приоритета)?
- Сколько задач может выполняться параллельно?
- Сколько задач может быть поставлено в очередь ожидания выполнения?
- Если задача должна быть отклонена из-за перегрузки системы, какая задача должна быть выбрана на роль жертвы и каким образом должно быть уведомлено приложение?
- Какие действия следует предпринять до или после выполнения задачи?

Политики выполнения представляют собой инструмент управления ресурсами, и оптимальная политика зависит от доступных вычислительных ресурсов и требований к качеству обслуживания. Ограничиваая число параллельно выполняющихся задач, можно гарантировать, что приложение не потерпит краха в связи с исчерпанием доступных ресурсов или не столкнётся с проблемами с производительностью из-за конкуренции за доступ к дефицитным ресурсам⁷⁴. Разделение спецификации политики выполнения и механизма отправки задач позволяет выбрать политику выполнения во время развертывания, соответствующую доступному оборудованию.

Всякий раз, когда вы видите код типа:

```
new Thread(runnable).start()
```

и вы думаете, что в какой-то момент вам может понадобиться более гибкая политика выполнения, серьезно подумайте о ее замене с использованием интерфейса Executor.

6.2.3 Пулы потоков

Пул потоков, как следует из названия, управляет однородным пулом рабочих потоков. Пул потоков тесно связан с *рабочей очередью*, содержащей задачи, ожидающие выполнения. Рабочие потоки живут достаточно просто: запрашивают следующую задачу из рабочей очереди, выполняют ее и возвращаются к ожиданию других задач.

Выполнение задач в пule потоков предлагает ряд преимуществ по сравнению с подходом "один поток на одну задачу". Повторное использование существующего потока вместо создания нового, амортизирует создание потока и снижает затраты,

⁷⁴Такое поведение соответствует одной из ролей монитора транзакций в корпоративном приложении: он может регулировать скорость, с которой транзакции могут выполняться, чтобы не исчерпывать или не перегружать ограниченные ресурсы.

в случае обработки множества запросов. В качестве дополнительного бонуса, поскольку рабочий поток часто уже существует на момент получения запроса, затраты времени, связанные с созданием потока, не задерживают выполнение задачи, тем самым повышая отзывчивость приложения. Правильно настроив размер пула потоков, можно получить достаточно потоков, чтобы процессоры были загружены, но при этом не так много, чтобы приложение исчерпывало доступную память или терпело крах из-за конкуренции за ресурсы между потоками.

Библиотека классов предоставляет гибкую реализацию пула потоков, наряду с некоторыми полезными предопределенными конфигурациями. Пул потоков можно создать, вызвав один из статических фабричных методов класса `Executors`:

newFixedThreadPool. Пул потоков фиксированного размера, создает потоки по мере отправки задач, вплоть до максимального размера пула, а затем старается сохранять размер пула неизменным (добавляя новые потоки, если поток умирает из-за возникновения неожиданного исключения `Exception`).

newCachedThreadPool. Кэшированный пул потоков обладает большой гибкостью в утилизации простаивающих потоков, когда текущий размер пула превышает потребность в обработчиках, и в добавления новых потоков, при увеличении потребности, но не накладывает ограничений на размер пула.

newSingleThreadExecutor. Однопоточная реализация интерфейса `Executor` создает единственный рабочий поток для обработки задач, заменяя его, если он неожиданно умирает. Гарантируется, что задачи будут обрабатываться последовательно в соответствии с порядком, определяемым очередью задач (FIFO, LIFO, в порядке приоритета).⁷⁵

newScheduledThreadPool. Пул потоков фиксированного размера, поддерживающий отложенное и периодичное выполнение задач, подобно классу `Timer` (см. раздел [6.2.5](#)).

Фабричные методы `newFixedThreadPool` и `newCachedThreadPool` возвращают экземпляры класса общего назначения - `ThreadPoolExecutor`, которые могут непосредственно использоваться для создания более специализированных экземпляров интерфейса `Executor`. Подробнее параметры конфигурации пула потоков обсудим в главе 8.

Веб-сервер в классе `TaskExecutionWebServer` использует экземпляр интерфейса `Executor` с ограниченным пулом рабочих потоков. Вызов метода `execute` отправляет задачу в рабочую очередь, а рабочие потоки непрерывно забирают задачи из рабочей очереди и выполняют их.

Переход от политики "один поток на одну задачу" к политике на основе пула потоков, оказывает большое влияние на стабильность приложения: веб-сервер больше не будет отказывать под большой нагрузкой⁷⁶. Как следствие, это приводит

⁷⁵ Однопоточные реализации интерфейса `Executor` также обеспечивают достаточную внутреннюю синхронизацию, чтобы гарантировать, что любая запись в память, сделанная задачей, будет видна последующим задачам; это означает, что объекты могут быть безопасно ограничены "потоком задачи", даже если этот поток может временно отменяться другим.

⁷⁶ Сервер не может отказать (упасть) из-за создания слишком большого количества потоков, но, если скорость поступления задач превышает скорость их обработки достаточно длительное время, все еще возможно (просто сложнее) исчерпать память из-за растущей очереди экземпляров `Runnable`,

к плавному снижению производительности, поскольку не создаются тысячи потоков, конкурирующих за ограниченные ресурсы ЦП и памяти. А использование фреймворка Executor открывает двери для всех видов дополнительных возможностей по настройке, управлению, мониторингу, ведению журнала, отчетах об ошибках и других, которые было бы гораздо сложнее добавить без фреймворка выполнения задач.

6.2.4 Жизненный цикл экземпляра Executor

Ранее мы уже рассматривали, как создать экземпляр интерфейса Executor, но не как завершить его. Реализация Executor, вероятнее всего, создаст потоки для обработки задач. Но JVM не сможет завершить свою работу до тех пор, пока все потоки (не демоны) не завершаться, поэтому неудача с закрытием экземпляра Executor может помешать завершению работы JVM.

Поскольку экземпляр Executor обрабатывает задачи асинхронно, то есть в любой момент времени, состояние ранее переданных на обработку задач не сразу становится очевидным. Некоторые из них могут быть завершены, некоторые могут быть запущены в текущий момент, а другие могут быть поставлены в очередь ожидания выполнения. Существует целый спектр вариантов завершения работы приложения - от мягкого выключения (завершается то, что было ранее начато, но новая работа не принимается) до внезапного выключения (выключается питание в комнате с компьютером), а также различные варианты между ними. Поскольку экземпляры Executor предоставляют приложениям службы, они также должны иметь возможность завершать свою работу, как плавно, так и внезапно, и передавать приложениям сведения о состоянии задач, на которые повлияло завершение работы.

Рассматривая вопросы жизненного цикла службы выполнения отметим, что интерфейс ExecutorService расширяет интерфейс Executor, добавляя ряд методов для управления жизненным циклом (а также некоторые удобные методы для отправки задач). Методы интерфейса ExecutorService, используемые для управления жизненным циклом, показаны в листинге 6.7.

```
public interface ExecutorService extends Executor {  
    void shutdown();  
    List<Runnable> shutdownNow();  
    boolean isShutdown();  
    boolean isTerminated();  
    boolean awaitTermination(long timeout, TimeUnit unit)  
        throws InterruptedException;  
    // ... additional convenience methods for task submission  
}
```

Листинг 6.7 Методы управления жизненным циклом интерфейса ExecutorService

Жизненный цикл, подразумеваемый интерфейсом ExecutorService, включает в себя три состояния: *запущено (running)*, *завершение работы (shutting down)* и *завершено (terminated)*. Экземпляры интерфейса ExecutorService изначально создаются в состоянии *запущено*. Метод shutdown инициирует мягкое завершение

ожидающих выполнения. Эта проблема может быть решена в рамках фреймворка Executor с помощью ограниченной рабочей очереди - см. раздел [8.3.2](#).

работы: новые задачи не принимаются к обработке, но ранее переданные завершаются, включая те, которые, что еще не принимались на выполнение. Метод `shutdownNow` производит внезапное завершение работы: он пытается отменить выполняющиеся задачи и не запускает какие-либо задачи, которые находятся в очереди, но не запущены.

Задачи, передаваемые экземпляру `ExecutorService` после того, как его работа была завершена, обрабатываются *обработчиком отклонённых задач* (*rejected execution handler*, см. раздел [8.3.3](#)), который может тихо избавляться от задач или может бросать непроверяемое исключение `RejectedExecutionException`. Как только все задачи будут выполнены, экземпляр `ExecutorService` перейдёт в состояние *завершено*. Можно дождаться, пока служба `ExecutorService` достигнет завершенного состояния с помощью метода `awaitTermination`, или опрашивать, завершилась ли она с помощью метода `isTerminated`. Как правило, следом за вызовом метода `shutdown` следует вызов метода `awaitTermination`, тем самым создавая эффект синхронного завершения работы экземпляра `ExecutorService` (Завершение работы экземпляра `Executor` и отмена задач подробно рассматриваются в главе 7).

Класс `LifecycleWebServer` из листинга 6.8, расширяет функциональность нашего веб-сервера, добавляя поддержку управления жизненным циклом. Сервер может быть погашен двумя способами: программным путем, с помощью вызова метода `stop`, и с помощью запроса со стороны клиента, путём отправки веб-серверу специально отформатированного HTTP-запроса.

```
class LifecycleWebServer {
    private final ExecutorService exec = ...;

    public void start() throws IOException {
        ServerSocket socket = new ServerSocket(80);
        while (!exec.isShutdown()) {
            try {
                final Socket conn = socket.accept();
                exec.execute(new Runnable() {
                    public void run() { handleRequest(conn); }
                });
            } catch (RejectedExecutionException e) {
                if (!exec.isShutdown())
                    log("task submission rejected", e);
            }
        }
    }

    public void stop() { exec.shutdown(); }

    void handleRequest(Socket connection) {
        Request req = readRequest(connection);
        if (isShutdownRequest(req))
            stop();
        else
            dispatchRequest(req);
    }
}
```

Листинг 6.8 Веб-сервер с поддержкой завершения работы

6.2.5 Отложенные и периодические задачи

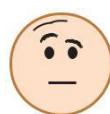
Класс `Timer` предоставляет возможность управления выполнением отсроченных (“выполнить эту задачу через 100мс”) и периодических (“выполнять задачу каждые 10 мс”) задач. Однако `Timer` имеет некоторые недостатки, и класс `ScheduledThreadPoolExecutor` следует рассматривать как его замену⁷⁷. Вы можете создать экземпляр класса `ScheduledThreadPoolExecutor` с помощью его конструктора или с помощью фабричного метода `newScheduledThreadPool`.

Класс `Timer` создает только один поток, выполняющий задания таймера. Если задание таймера выполняется слишком долго, может пострадать точность синхронизации других экземпляров `TimerTask`. Если экземпляр `TimerTask` запланирован для повторного выполнения каждые 10 мс, а другой экземпляр `TimerTask` на повтор каждые 40 мс, повторяющаяся задача либо (в зависимости от того, была ли она запланирована с фиксированной скоростью или с фиксированной задержкой) вызывается четыре раза подряд после завершения длительной задачи, либо “пропускает” четыре вызова полностью. Запланированные пулы потоков устраниют это ограничение, позволяя предоставлять несколько потоков для выполнения отложенных и периодических задач.

Другая проблема с классом `Timer` заключается в том, что он ведет себя плохо, если экземпляр `TimerTask` бросает непроверяемое исключение. Поток таймера не перехватывает исключение, поэтому непроверяемое исключение, брошенное экземпляром `TimerTask`, завершает поток таймера. Класс `Timer` также не воскрешает поток в этой ситуации; вместо этого он ошибочно предполагает, что весь экземпляр `Timer` был отменен. В этом случае задачи таймера, которые уже были запланированы, но еще не выполнены, никогда не выполняются, и новые задачи не смогут быть запланированы (эта проблема, называемая “утечкой потоков”, описана в разделе 7.3 вместе с методами ее предотвращения).

Класс `OutOfTime` из листинга 6.9 иллюстрирует, каким образом класс `Timer` может запутаться, и как путаница любит компанию, как класс `Timer` разделяет своё замешательство со следующим незадачливым вызывающим объектом, который пытается отправить экземпляр `TimerTask`. Можно ожидать, что программа будет работать в течение шести секунд и завершит работу, но на самом деле происходит так, что она завершается через одну секунду с возбуждением исключения `IllegalStateException`, текст которого звучит, как “Таймер уже отменен”. Класс `ScheduledThreadPoolExecutor` должным образом имеет дело с некорректными задачами; существует мало причин для использования класса `Timer` в Java 5.0 и выше.

```
public class OutOfTime {
    public static void main(String[] args) throws Exception {
        Timer timer = new Timer();
        timer.schedule(new ThrowTask(), 1);
        SECONDS.sleep(1);
        timer.schedule(new ThrowTask(), 1);
```



⁷⁷ Класс `Timer` осуществляет планирование на основе абсолютного, а не относительного времени, так что на задачи может оказываться влияние при изменении в системных часах; класс `ScheduledThreadPoolExecutor` использует в своей работе только относительное время.

```
        SECONDS.sleep(5);
    }

    static class ThrowTask extends TimerTask {
        public void run() { throw new RuntimeException(); }
    }
}
```

Листинг 6.9 Класс, иллюстрирующий запутанное поведение класса Timer

Если вам необходимо создать свою собственную службу планировщик, вы все равно можете воспользоваться библиотекой, используя класс `DelayQueue` - реализацию `BlockingQueue`, обеспечивающую функциональность планировщика с помощью класса `ScheduledThreadPoolExecutor`. Класс `DelayQueue` управляет коллекцией объектов `Delayed`. Класс `Delayed` содержит время задержки: класс `DelayQueue` позволяет получить элемент, только если его задержка истекла. Объекты возвращаются из экземпляра `DelayQueue` в порядке, определённом временем их задержки.

6.3 Поиск уязвимостей параллелизма

Фреймворк `Executor` упрощает определение политики выполнения, но для того, чтобы его использовать, вы должны быть в состоянии описать свою задачу как экземпляр интерфейса `Runnable`. В большинстве серверных приложений существует очевидная граница задачи: один запрос от клиента. Но иногда хорошие границы задач не так очевидны, как во многих десктопных приложениях. В серверных приложениях также может использоваться параллелизм в рамках одного запроса клиента, как это иногда бывает в серверах баз данных (для дальнейшего обсуждения конкурирующих тенденций проектирования при выборе границ задач см. [CPJ 4.4.1.1]).

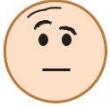
В этом разделе мы разрабатываем несколько версий компонента, допускающих различную степень параллелизма. Наш пример компонента - это часть отрисовки страницы в приложении браузера, которая берет страницу HTML и отображает ее в буфер изображений. Для простоты мы предполагаем, что HTML состоит только из размеченного текста с элементами изображения с заранее заданными размерами и URL.

6.3.1 Пример: последовательный генератор страниц

Простейший подход заключается в последовательной обработке HTML-документа. При обнаружении текстовой разметки отрисуйте ее в буфер изображений; при обнаружении ссылок на изображения извлеките изображения по сети и отрисуйте их в буфер изображений. Этот подход легко реализовать и он требует касания каждого элемента ввода только один раз (буферизации документа даже не требуется), но, вероятнее всего он будет раздражать пользователя, которому возможно придется долго ждать, прежде чем весь текст будет отображен.

Менее раздражающий, но все же последовательный подход включает в себя сначала отрисовку текстовых элементов, оставляя прямоугольные заполнители для изображений, а после завершения начального прохода по документу, возврат, загрузку изображений и их отрисовку в соответствующие заполнители. Пример такого подхода представлен в классе `SingleThreadRenderer`, в листинге 6.10.

```
public class SingleThreadRenderer {  
    void renderPage(CharSequence source) {  
        renderText(source);  
        List<ImageData> imageData = new ArrayList<ImageData>();  
        for (ImageInfo imageInfo : scanForImageInfo(source))  
            imageData.add(imageInfo.downloadImage());  
        for (ImageData data : imageData)  
            renderImage(data);  
    }  
}
```



Листинг 6.10 Последовательная отрисовка элементов страницы

В основном, процесс загрузки изображения включает в себя ожидание завершения операций ввода/вывода, и в течение этого времени процессор выполняет достаточно мало работы. Таким образом, последовательный подход может привести к недоиспользованию процессора, а также заставляет пользователя ожидать дольше, чем необходимо, для того, чтобы увидеть готовую страницу. Мы можем добиться более эффективного использования и отзывчивости, разбив проблему на независимые задачи, которые могут выполняться одновременно.

6.3.2 Задачи возвращающие результат: Callable и Future

Фреймворк Executor использует интерфейс `Runnable` в качестве базового представления задачи. Интерфейс `Runnable` представляет собой довольно ограниченную абстракцию; метод `run` не может возвращать значения или бросать проверяемые исключения, однако он может оказывать сторонние воздействия, такие как запись в файл журнала или размещение результата в общей структуре данных.

Многие задачи, по своей сути, являются отложенными вычислениями - выполнение запроса к базе данных, получение ресурса по сети или вычисление сложной функции. Для этих типов задач интерфейс `Callable` является лучшей абстракцией: он ожидает, что основная точка входа, метод `call`, вернет значение и ожидает, что он может бросить исключение⁷⁸. Класс `Executors` включает в себя несколько методов-утилит, предназначенных для упаковки других типов задач, в том числе экземпляров `Runnable` и `java.security.PrivilegedAction`, в экземпляры интерфейса `Callable`.

Интерфейсы `Runnable` и `Callable` описывают абстрактные вычислительные задачи. Задачи обычно конечны: у них есть четкая отправная точка, и они в конечном итоге завершаются. Жизненный цикл задачи, выполняемой экземпляром `Executor`, состоит из четырех фаз: *создано (created)*, *отправлено (submitted)*, *запущено (started)* и *завершено (completed)*. Поскольку выполнение задач может занимать много времени, мы также хотим иметь возможность отменять задачу. В фреймворке `Executor` задачи, которые были отправлены, но еще не запущены, всегда могут быть отменены, а выполняющиеся задачи иногда могут быть отменены, если они реагируют на прерывание. Отмена ранее завершенной задачи не оказывает никакого эффекта (отмена более подробно рассматривается в главе 7).

⁷⁸ Для описания задачи, не возвращающей значения, с помощью интерфейса `Callable`, используйте `Callable<Void>`

Интерфейс `Future` отражает жизненный цикл задачи и предоставляет методы для проверки того, была ли задача завершена или отменена, получения результата выполнения и отмены задачи. Интерфейсы `Callable` и `Future` представлены в листинге 6.11. В спецификации интерфейса `Future` неявно подразумевается, что жизненный цикл задачи может продвигаться только вперед, а не назад - как и жизненный цикл интерфейса `ExecutorService`. Как только задача выполнена, она остается в этом состоянии навсегда.

```
public interface Callable<V> {  
    V call() throws Exception;  
}  
  
public interface Future<V> {  
    boolean cancel(boolean mayInterruptIfRunning);  
    boolean isCancelled();  
    boolean isDone();  
    V get() throws InterruptedException, ExecutionException,  
              CancellationException;  
    V get(long timeout, TimeUnit unit)  
        throws InterruptedException, ExecutionException,  
              CancellationException, TimeoutException;  
}
```

Листинг 6.11 Интерфейсы `Callable` и `Future`

Поведение метода `get` зависит от состояния задачи (еще не запущено, выполняется, завершено). Если задача уже выполнена метод, немедленно возвращает результат или бросает исключение `Exception`, иначе блокируется до завершения задачи. Если задача завершается путем возбуждения исключения, метод `get` оборачивает его в экземпляр `ExecutionException` и пробрасывает дальше; если задача была отменена, метод `get` бросает исключение `CancellationException`. Если метод `get` бросает исключение `ExecutionException`, базовое исключение может быть получено с помощью метода `getCause`.

Существует несколько способов создания экземпляра `Future`, описывающего задачу. Все методы `submit` интерфейса `ExecutorService` возвращают экземпляры `Future`, так что вы можете отправить экземпляры `Runnable` или `Callable` исполнителю и получить назад экземпляр `Future`, который может быть использован для получения результата или отмены задачи. Можно также явно создать экземпляр `FutureTask` для переданного экземпляра `Runnable` или `Callable` (поскольку класс `FutureTask` реализует интерфейс `Runnable`, он может быть передан экземпляру `Executor` для выполнения или непосредственно выполнен, путем вызова метода `run`).

Начиная с Java 6, реализации интерфейса `ExecutorService` имеют возможность переопределять метод `newTaskFor` класса `AbstractExecutorService`, и таким образом получают возможность контролировать инстанцирование (*instantiation*) экземпляров `Future`, соответствующих представленным экземплярам `Callable` или `Runnable`. Реализация по умолчанию просто создает новый экземпляр `FutureTask`, как показано в листинге 6.12.

```
protected <T> RunnableFuture<T> newTaskFor(Callable<T> task) {
    return new FutureTask<T>(task);
}
```

Листинг 6.12 Реализация “по умолчанию” метода newTaskFor класса ThreadPoolExecutor

Отправка экземпляров `Runnable` или `Callable` экземпляру `Executor` основывается на безопасной публикации (см. раздел 3.5) экземпляров `Runnable` или `Callable` от отправляющего потока к потоку, который, в конечном счете, выполнит задачу. Аналогично, установка значения результата экземпляру `Future` представляет собой безопасную публикацию результата из потока, в котором он был вычислен, в любой поток, который получает результат через вызов метода `get`.

6.3.3 Пример: отрисовка страниц с помощью Future

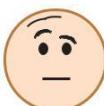
В качестве первого шага в повышении степени параллелизма отрисовщика страниц, давайте разделим его на две задачи, одна из которых будет отрисовывать текст, а другая будет загружать все изображения (поскольку одна задача в значительной степени ограничена ЦП, а другая задача в большей степени ограничена операциями ввода/вывода, такой подход может привести к улучшению производительности даже в системах с одним ЦП).

Интерфейсы `Callable` и `Future` могут помочь нам отразить взаимодействие между этими кооперирующимися задачами. В классе `FutureRenderer` в листинге 6.13, мы создаем экземпляр интерфейса `Callable` для загрузки всех изображений и отправляем его экземпляру `ExecutorService`. Он возвращает экземпляр `Future`, описывающий выполнение задачи; когда основная задача доходит до точки, в которой ей нужны изображения, она ожидает результата вызова `Future.get`. Если нам повезет, результаты к моменту запроса уже будут готовы; в противном случае, мы, по крайней мере, получили фору при загрузке изображений.

```
public class FutureRenderer {
    private final ExecutorService executor = ...;

    void renderPage(CharSequence source) {
        final List<ImageInfo> imageInfos = scanForImageInfo(source);
        Callable<List<ImageData>> task =
            new Callable<List<ImageData>>() {
                public List<ImageData> call() {
                    List<ImageData> result
                        = new ArrayList<ImageData>();
                    for (ImageInfo imageInfo : imageInfos)
                        result.add(imageInfo.downloadImage());
                    return result;
                }
            };
        Future<List<ImageData>> future = executor.submit(task);
        renderText(source);

        try {
```



```

        List<ImageData> imageData = future.get();
        for (ImageData data : imageData)
            renderImage(data);
    } catch (InterruptedException e) {
        // Re-assert the thread's interrupted status
        Thread.currentThread().interrupt();
        // We don't need the result, so cancel the task too
        future.cancel(true);
    } catch (ExecutionException e) {
        throw launderThrowable(e.getCause());
    }
}
}

```

Листинг 6.13 Ожидание загрузки изображений с использованием экземпляра Future

По своей природе метод `get` зависит от состояния, это значит, что вызывающий объект не должен ничего знать о состоянии задачи, а безопасная публикация, свойственная отправке задачи и извлечению результата, делает этот подход потокобезопасным. Код обработки исключений, окружающий метод `Future.get`, может столкнуться с двумя возможными проблемами: задача в процессе выполнения породила исключение `Exception`, или поток вызванный методом `get` был прерван, прежде чем стали доступны результаты (см. разделы [5.5.2](#) и [5.4](#)).

Класс `FutureRenderer` позволяет тексту отрисовываться одновременно с загрузкой данных изображения. Когда все изображения загружены, они отображаются на странице. Это улучшение заключается в том, что пользователь быстрее видит результат и в том, что используется некоторый параллелизм, но мы можем всё сделать значительно лучше. Нет никакой необходимости в том, чтобы заставлять пользователей ожидать загрузки *всех* изображений; они, вероятнее всего, предпочтут видеть отдельные изображения, отрисовывающиеся по мере доступности.

6.3.4 Ограничения распараллеливания разнородных задач

В предыдущем примере, мы попытались параллельно выполнить два разных типа задач - загрузку изображений и отрисовку страницы. Но получить значительные улучшения производительности, пытаясь распараллелить последовательные гетерогенные задачи, может быть сложно.

Два человека могут разделить работу по уборке посуды достаточно эффективно: один человек моет её, а другой сушит. Тем не менее, назначение различных типов задач каждому работнику масштабируется не очень хорошо; если появятся еще несколько человек, не очевидно, как они смогут помочь, не вставая на пути или без проведения значительной реструктуризации разделения труда. Если не найти меньших параллелизуемых частей среди подобных задач, такой подход приведёт к снижению отдачи.

Еще одна проблема связана с разделением разнородных задач между несколькими работниками заключается в том, что задачи могут иметь несопоставимые размеры. Если вы разделите задачи *A* и *B* между двумя работниками, но *A* выполняется в десять раз больше времени, чем *B*, вы ускорите общий процесс только на 9%. Наконец, разделение задачи между несколькими

работниками всегда сопряжено с определенными накладными расходами на координацию; для того чтобы разделение было оправданным, эти накладные расходы должны быть более чем компенсированы повышением производительности за счет параллелизма.

Класс `FutureRenderer` использует две задачи: одну для отрисовки текста, а другую для загрузки изображений. Если отрисовка текста происходит гораздо быстрее, чем загрузка изображений, как это часто бывает на практике, производительность будет не сильно отличаться от последовательной версии, но код станет намного сложнее. Лучшее, что мы можем сделать с двумя потоками, это получить ускорение в два раза. Таким образом, попытка увеличить параллелизм за счет распараллеливания разнородных действий может быть очень трудоемкой, и существуют пределы того, сколько дополнительного параллелизма можно получить за счёт распараллеливания. (См. разделы [11.4.2](#) и [11.4.3](#), в них приводится другой пример того же явления.)

Реальная отдача от разделения рабочей нагрузки программы на задачи возникает тогда, когда существует большое количество независимых однородных задач, которые могут обрабатываться параллельно.

6.3.5 CompletionService: Executor пересекается с BlockingQueue

Если у вас есть пакет вычислений для отправки экземпляру `Executor` и вы хотите получать результаты по мере их доступности, вы можете сохранить экземпляры интерфейса `Future`, связанные с каждой задачей, и периодически опрашивать их на предмет завершения выполнения, с помощью вызова метода `get` с нулевым таймаутом. Так делать можно, но это утомительно. К счастью, есть лучший способ: *служба завершения* (*completion service*).

Интерфейс `CompletionService` объединяет в себе функционал интерфейсов `Executor` и `BlockingQueue`. Вы можете отправлять ему задачи для выполнения в виде экземпляров `Callable` и использовать методы подобные тем, что имеются у очередей - методы `take` и `poll` - для получения завершенных результатов, упакованных в экземпляры интерфейса `Future`, по мере их доступности. Класс `ExecutorCompletionService` реализует интерфейс `CompletionService`, делегируя вычисления интерфейсу `Executor`.

Реализация класса `ExecutorCompletionService` довольно проста. Конструктор создает экземпляр `BlockingQueue` для хранения завершенных результатов. Класс `FutureTask` имеет метод `done`, который вызывается после завершения вычислений. Когда задача отправляется, она оборачивается классом `QueueingFuture`, являющимся подклассом `FutureTask`, переопределяющим метод `done` с целью помещать результат выполнения в экземпляр `BlockingQueue`, как показано в листинге 6.14. Методы `take` и `poll` делегируются экземпляру `BlockingQueue`, блокируясь, если результаты еще недоступны.

```
private class QueueingFuture<V> extends FutureTask<V> {
    QueueingFuture(Callable<V> c) { super(c); }
    QueueingFuture(Runnable t, V r) { super(t, r); }
```

```
    protected void done() {
        completionQueue.add(this);
    }
}
```

Листинг 6.14 Класс QueueingFuture используемый классом ExecutorCompletionService.

6.3.6 Пример: рендер страниц с использованием CompletionService

Мы можем использовать интерфейс `CompletionService` для повышения производительности рендера страниц двумя способами: сократив общее времени выполнения и улучшив отзывчивость. Мы можем создавать отдельные задачи для загрузки *каждого* изображения и выполнять их в пуле потоков, превращая последовательную загрузку в параллельную: такой подход сократит общее время загрузки всех изображений. Извлекая результаты из экземпляра `CompletionService` и отрисовывая каждое изображение, как только оно станет доступно, мы можем предоставить пользователю более динамичный и отзывчивый пользовательский интерфейс. Пример реализации показан в классе `Renderer`, в листинге 6.15.

```
public class Renderer {
    private final ExecutorService executor;

    Renderer(ExecutorService executor) { this.executor = executor; }

    void renderPage(CharSequence source) {
        List<ImageInfo> info = scanForImageInfo(source);
        CompletionService<ImageData> completionService =
            new ExecutorCompletionService<ImageData>(executor);
        for (final ImageInfo imageInfo : info)
            completionService.submit(new Callable<ImageData>() {
                public ImageData call() {
                    return imageInfo.downloadImage();
                }
            });
        renderText(source);

        try {
            for (int t = 0, n = info.size(); t < n; t++) {
                Future<ImageData> f = completionService.take();
                ImageData imageData = f.get();
                renderImage(imageData);
            }
        } catch (InterruptedException e) {
            Thread.currentThread().interrupt();
        } catch (ExecutionException e) {
            throw launderThrowable(e.getCause());
        }
    }
}
```

```
}
```

Листинг 6.15 Использование интерфейса CompletionService для рендеринга элементов страницы по мере того, как они становятся доступны.

Несколько экземпляров ExecutorCompletionServices могут совместно использовать один и тот же экземпляр Executor, так что будет вполне разумно создавать экземпляры ExecutorCompletionServices, закрытыми для конкретных вычислений, пока экземпляр Executor используется совместно. Когда интерфейс CompletionService используется подобным образом, он выступает в качестве обработчика для пакета вычислений, во многом повторяя поведение, интерфейса Future, выступающего в качестве обработчика для отдельного вычисления. Запомнив, сколько задач было отправлено экземпляру CompletionService, и подсчитывав, сколько завершенных результатов было получено, можно узнать, когда были получены все результаты для данного пакета, даже если экземпляр Executor используется совместно.

6.3.7 Ограничение времени выполнения задач

Иногда, если активность не завершает свою работу в течение определенного периода времени, результат становится больше не нужен, и работа активности может быть прервана. Например, веб-приложение может получать объявления с внешнего сервера объявлений, но если объявление недоступно в течение двух секунд, оно отображает объявление по умолчанию, чтобы отсутствие объявлений не нарушало требований к отзывчивости сайта. Аналогично, сайт портала может параллельно получать данные из нескольких источников, но может быть готов ожидать, пока данные станут доступны, только определенное время, прежде чем отобразить страницу без них.

Основной проблемой, при выполнении задач в рамках выделенного бюджета времени, является возможность убедиться, что для получения ответа или выяснения его отсутствия, вам не придется ожидать дольше, чем позволяет бюджет времени. Версия Future.get с параметром “время ожидания” поддерживает это требование: она возвращает результат, как только он готов, в случае если результат не был подготовлен в течение заданного периода ожидания - бросает исключение TimeoutException.

Вторичной проблемой при использовании ограниченных временем задач является их остановка по истечении выделенного периода времени, чтобы они не тратили вычислительные ресурсы впустую, продолжая вычислять результат, который не будет использован. Этого можно добиться сделав так, чтобы задача строго управляла собственным бюджетом времени и прерывала своё выполнение, когда у нее заканчивается время, или, отменяя задачу по истечении времени ожидания. Опять же, интерфейс Future может помочь и в этом случае; если ограниченная по времени (*timed*) версия get завершается возбуждением исключения TimeoutException, вы можете отменить задачу с использованием экземпляра Future. Если задача написана как “отменяемая” (см. главу 7), ее можно завершить досрочно, чтобы не использовать излишние ресурсы. Такой подход используется в Листингах 6.13 и 6.16.

В листинге 6.16 показано типичное применение “временного” метода Future.get. Метод renderPageWithAd генерирует составную веб-страницу, содержащую запрошенное содержимое и рекламу, получаемую с сервера объявлений. Он отправляет задачу, извлекающую рекламу, исполнителю, обрабатывает остальную часть содержимого страницы, а затем ожидает получения

рекламы, пока не закончится выделенный ему бюджет времени⁷⁹. Если время ожидания вызова метода `get` истекло, задача извлечения рекламы отменяется⁸⁰ и вместо нее используется объявление по умолчанию.

```
Page renderPageWithAd() throws InterruptedException {
    long endNanos = System.nanoTime() + TIME_BUDGET;
    Future<Ad> f = exec.submit(new FetchAdTask());
    // Render the page while waiting for the ad
    Page page = renderPageBody();
    Ad ad;
    try {
        // Only wait for the remaining time budget
        long timeLeft = endNanos - System.nanoTime();
        ad = f.get(timeLeft, NANOSECONDS);
    } catch (ExecutionException e) {
        ad = DEFAULT_AD;
    } catch (TimeoutException e) {
        ad = DEFAULT_AD;
        f.cancel(true);
    }
    page.setAd(ad);
    return page;
}
```

Листинг 6.16 Получение рекламы в рамках выделенного на процедуру бюджета времени

6.3.8 Пример: портал бронирования путешествий

Подход к составлению бюджета времени, описанный в предыдущем разделе, можно легко обобщить на произвольное число задач. Рассмотрим портал бронирования путешествий: пользователь вводит даты поездки и требования, а портал извлекает и отображает предложения от ряда авиакомпаний, отелей или компаний по прокату автомобилей. В зависимости от компании, получение предложения может включать использование вызова веб-службы, консультации с базой данных, выполнение транзакции EDI или использование другого механизма. Вместо того чтобы время отклика страницы привязывать к скорости ответа самого медленного запроса, может быть предпочтительнее представлять только ту информацию, которая была получена в рамках выделенного бюджета времени. Провайдеров, которые не ответили вовремя на запрос, страница могла бы полностью опустить или отобразить заставку, вроде такой “ответ от Air Java не получен”.

Получение предложения от одной компании не зависит от получения предложения от другой, поэтому получение одного предложения является разумной границей задачи, которая позволяет получать предложения параллельно. Было бы достаточно легко создать n задач, отправить их в пул потоков, сохранить

⁷⁹ Значение тайм-аута, переданное методу `get`, вычисляется путем вычитания текущего времени из крайнего срока (*deadline*); на практике может быть вычислено отрицательное число, но все “временные” методы в `java.util.concurrent` рассматривают отрицательные значения тайм-аутов как ноль, поэтому никакой дополнительный код для обработки этого случая не требуется.

⁸⁰ Параметр `true` в вызове метода `Future.cancel` означает, что поток задачи может быть прерван, если задача в данный момент выполняется; см. главу 7.

экземпляры `Future` и использовать “временную” версию метода `get` для последовательной выборки каждого результата с помощью экземпляра `Future`, но есть еще более простой способ – вызов метода `invokeAll`.

В листинге 6.17, ограниченная по времени версия метода `invokeAll` используется для отправки нескольких задач экземпляру `ExecutorService` и получения результатов. Метод `invokeAll` принимает коллекцию задач и возвращает коллекцию экземпляров `Future`. Обе коллекции имеют идентичные структуры; метод `invokeAll` добавляет экземпляры `Future` к возвращаемой коллекции в порядке, установленном итератором коллекции задач, позволяя, таким образом, вызывающему объекту связывать экземпляры `Future` с экземплярами `Callable`, который он представляет. Ограниченнная по времени версия метода `invokeAll` возвратит управление тогда, когда будут завершены все задачи, будет прерван вызывающий поток или когда истечёт время ожидания. Все задачи, которые не были завершены по истечении тайм-аута, отменяются. При возврате из метода `invokeAll`, каждая задача будет либо выполнена в обычном режиме, либо отменена; клиентский код может вызвать метод `get` или метод `isCancelled`, чтобы выяснить, в каком состоянии находится задача. Listing 6.17 uses the timed version of `invokeAll` to submit multiple tasks to an `ExecutorService` and retrieve the results. The `invokeAll` method takes a collection of tasks and returns a collection of `Futures`. The two collections have identical structures; `invokeAll` adds the `Futures` to the returned collection in the order imposed by the task collection’s iterator, thus allowing the caller to associate a `Future` with the `Callable` it represents. The timed version of `invokeAll` will return when all the tasks have completed, the calling thread is interrupted, or the timeout expires. Any tasks that are not complete when the timeout expires are cancelled. On return from `invokeAll`, each task will have either completed normally or been cancelled; the client code can call `get` or `isCancelled` to find out which.

```
private class QuoteTask implements Callable<TravelQuote> {
    private final TravelCompany company;
    private final TravelInfo travelInfo;
    ...
    public TravelQuote call() throws Exception {
        return company.solicitQuote(travelInfo);
    }
}

public List<TravelQuote> getRankedTravelQuotes(
    TravelInfo travelInfo, Set<TravelCompany> companies,
    Comparator<TravelQuote> ranking, long time, TimeUnit unit)
    throws InterruptedException {
    List<QuoteTask> tasks = new ArrayList<QuoteTask>();
    for (TravelCompany company : companies)
        tasks.add(new QuoteTask(company, travelInfo));

    List<Future<TravelQuote>> futures =
        exec.invokeAll(tasks, time, unit);

    List<TravelQuote> quotes =
```

```
    new ArrayList<TravelQuote>(tasks.size());
Iterator<QuoteTask> taskIter = tasks.iterator();
for (Future<TravelQuote> f : futures) {
    QuoteTask task = taskIter.next();
    try {
        quotes.add(f.get());
    } catch (ExecutionException e) {
        quotes.add(task.getFailureQuote(e.getCause()));
    } catch (CancellationException e) {
        quotes.add(task.getTimeoutQuote(e));
    }
}
Collections.sort(quotes, ranking);
return quotes;
}
```

Листинг 6.17 Запрос расценок на путешествия в рамках операции с ограниченным временем

6.4 Итоги

Структурирование приложений вокруг выполнения задач может упростить разработку и облегчить использование параллелизма. Фреймворк Executor позволяет отделить отправку задач от политики выполнения и поддерживает множество различных политик выполнения; всякий раз, когда вы создаете потоки для выполнения задач, рассмотрите возможность использования фреймворка Executor. Чтобы извлечь максимальную пользу из разбиения приложения на задачи, необходимо определить разумные границы задачи. В некоторых приложениях очевидные границы задач работают хорошо, в то время как в других может потребоваться некоторый анализ для выявления деталей эксплуатируемого параллелизма.

Глава 7 Отмена и завершение

Запускать задачи и потоки довольно легко. Большую часть времени мы позволяем им самим принимать решение о том, когда остановиться, позволяя работать до завершения. Однако, иногда мы хотим остановить выполнение задач или потоков раньше, чем они остановились бы сами по себе, возможно, в связи с тем, что пользователь отменил операцию или приложение должно быстро завершиться.

Безопасно, быстро и надежно остановить задачу или поток не всегда легко. Язык Java не предоставляет никакого механизма для безопасной принудительной остановки потока, вместо этого Java предоставляет *прерывание (interruption)* - механизм взаимодействия, позволяющий одному потоку попросить другой поток прекратить выполнять то, что он делает⁸¹.

Кооперативный подход необходим, потому что достаточно редко возникает необходимость в том, чтобы задача, поток или служба *немедленно* останавливались, так как это может привести к тому, что совместно используемые структуры данных могут остаться в несогласованном состоянии. Вместо этого задачи и службы можно закодировать так, чтобы по запросу они “очищали” все выполняемые в данный момент работы, а *затем* завершали их. Такой подход обеспечивает большую гибкость, поскольку сам код задачи, как правило, лучше способен оценить требуемые действия по очистке используемых ресурсов, чем код, запрашивающий отмену.

Проблемы, связанные с завершением жизненного цикла, могут усложнить проектирование и реализацию задач, служб и приложений, и этот важный элемент разработки программ слишком часто игнорируется. Работа со сбоями, завершением работы и отменой - это одна из характеристик, отличающих корректное приложение от просто работающего приложения. В этой главе рассматриваются механизмы отмены и прерывания, а также кодирование задач и служб для перехвата и обработки запросов на отмену.

7.1 Отмена задачи

Выполнение активности может быть отменено, если внешний код может подвести её к завершению до нормального завершения. Существует ряд причин, по которым может потребоваться отменить выполнение активности:

Отмена, инициированная пользователем. Пользователь нажал на кнопку "отмена" в приложении с GUI или запросил отмену через интерфейс управления, такой как JMX (Java Management Extensions).

Активности, ограниченные по времени. Приложение занимается поиском решения в пространстве проблем на протяжении конечного промежутка времени и выбирает лучшее решение, найденное за это время. По истечении таймера все задачи поиска отменяются.

События приложения. Приложение ищет решение в пространстве проблем путем разложения его на отдельные задачи так, что различные задачи выполняют

⁸¹ Устаревшие методы `Thread.stop` и `suspend` являлись попыткой реализовать такой механизм, но довольно скоро были осознаны серьезные недостатки, связанные с ними, поэтому их использования следует избегать. За разъяснениями обратитесь к статье <http://java.sun.com/j2se/1.5.0/docs/guide/misc/threadPrimitiveDeprecation.html>, в которой рассматриваются проблемы связанные с использованием этих методов.

поиск в различных регионах пространства проблем. Когда одна задача находит решение, все остальные задачи, которые еще находятся в состоянии поиска, отменяются.

Ошибки. Веб-сканер выполняет поиск соответствующих страниц, сохраняя страницы или сводные данные на диске. Когда задача сканера сталкивается с ошибкой (например, диск переполнен), другие задачи обхода отменяются, возможно, с записью своего текущего состояния, чтобы их можно было перезапустить позже.

Завершение работы. Когда приложение или служба завершает свою работу, что-то должно быть сделано с работой, выполняющейся в текущий момент или ожидающей в очереди на выполнение. При плавном завершении работы, задачам, выполняющимся в текущий момент, может быть разрешено завершить свою работу; при немедленном завершении работы, текущие задачи могут быть отменены.

В Java нет безопасного способа превентивной остановки потока, и поэтому нет безопасного способа превентивной остановки задачи. Существуют лишь механизмы кооперации, в рамках которых задача и код запрашивающий отмену, следуют согласованному протоколу.

Одним из таких механизмов кооперации является установка флага "запрос отмены", периодически проверяемого задачей; если задача обнаруживает, что флаг установлен, задача завершается раньше. Этот подход иллюстрируется классом `PrimeGenerator` из листинга 7.1, который перечисляет простые числа до тех пор, пока его работа не будет отменена. Метод `cancel` устанавливает флаг `cancelled`, и главный цикл опрашивает состояние этого флага перед поиском следующего простого числа. (Для надежной работы переменная `cancelled` должна быть объявлена с ключевым словом `volatile`.)

```
@ThreadSafe
public class PrimeGenerator implements Runnable {
    @GuardedBy("this")
    private final List<BigInteger> primes
        = new ArrayList<BigInteger>();
    private volatile boolean cancelled;

    public void run() {
        BigInteger p = BigInteger.ONE;
        while (!cancelled) {
            p = p.nextProbablePrime();
            synchronized (this) {
                primes.add(p);
            }
        }
    }

    public void cancel() { cancelled = true; }

    public synchronized List<BigInteger> get() {
        return new ArrayList<BigInteger>(primes);
    }
}
```

```
    }  
}
```

Листинг 7.1 Использование поля типа volatile для хранения состояния отмены

В листинге 7.2 показан пример использования этого класса, позволяющий генератору простых чисел работать в течение одной секунды перед отменой выполнения. Генератор не обязательно остановит работу ровно через одну секунду, так как может возникнуть некоторая задержка между моментом времени запроса на отмену и моментом времени следующей проверки состояния флага, управляющего отменой выполнения цикла. Метод `cancel` вызывается из блока `finally`, чтобы гарантированно прервать работу генератора простых чисел, даже в том случае, если вызов метода `sleep` будет прерван. Если метод `cancel` не будет вызван, основной поисковый поток будет работать вечно, потребляя циклы ЦП и предотвращая завершение работы JVM.

```
List<BigInteger> aSecondOfPrimes() throws InterruptedException {  
    PrimeGenerator generator = new PrimeGenerator();  
    new Thread(generator).start();  
    try {  
        SECONDS.sleep(1);  
    } finally {  
        generator.cancel();  
    }  
    return generator.get();  
}
```

Листинг 7.2 Секундная генерация простых чисел

Задача, которая хочет быть отменяемой, должна иметь политику отмены, которая определяет "как", "когда" и "что" отменять - как другой код может запросить отмену задачи, когда задача проверяет, была ли запрошена отмена, и какие действия задача выполняет в ответ на запрос отмены.

Рассмотрим реальный пример остановки платежа по чеку. У банков есть правила, описывающие как подать запрос на остановку платежа, какая оперативность гарантируется при обработке таких запросов и каким процедурам они следуют при фактической остановке платежа (например, уведомление другого банка, участвующего в транзакции, и оценка комиссии по счету плательщика). Вместе взятые, эти процедуры и гарантии составляют политику отмены для чека на оплату.

Класс `PrimeGenerator` использует простую политику отмены: клиентский код запрашивает отмену задачи с помощью вызова метода `cancel`, экземпляр `PrimeGenerator` проверяет флаг отмены каждый раз, как находит простое число и завершает свою работу, когда обнаруживает, что была запрошена отмена выполнения.

7.1.1 Прерывание

Механизм отмены в классе `PrimeGenerator` в конечном итоге приведет к завершению задачи поиска простого числа, но это может занять некоторое время. Однако, если задача, которая использует этот подход, вызывает блокирующий метод, такой как `BlockingQueue.put`, у нас может возникнуть более серьезная

проблема - задача может никогда не проверить флаг отмены и, следовательно, никогда не завершиться.

Класс `BrokenPrimeProducer` из листинга 7.3 иллюстрирует эту проблему. Поток производителя создаёт простые числа и помещает их в блокирующую очередь. Если производитель опередит потребителя, очередь заполнится и вызов метода `put` будет заблокирован. Что произойдет, если потребитель попытается отменить задачу производителя, пока заблокирован вызов метода `put`? Он может вызвать метод `cancel`, который в свою очередь установит флагу `cancelled` значение `true`, но производитель никогда не проверит состояние флага, потому что он никогда не сможет выйти из заблокированного метода `put` (потому что потребитель остановил процесс получения простых чисел из очереди).

```
class BrokenPrimeProducer extends Thread {  
    private final BlockingQueue<BigInteger> queue;  
    private volatile boolean cancelled = false;  
  
    BrokenPrimeProducer(BlockingQueue<BigInteger> queue) {  
        this.queue = queue;  
    }  
  
    public void run() {  
        try {  
            BigInteger p = BigInteger.ONE;  
            while (!cancelled)  
                queue.put(p = p.nextProbablePrime());  
        } catch (InterruptedException consumed) { }  
    }  
  
    public void cancel() { cancelled = true; }  
}  
  
void consumePrimes() throws InterruptedException {  
    BlockingQueue<BigInteger> primes = ...;  
    BrokenPrimeProducer producer = new BrokenPrimeProducer(primes);  
    producer.start();  
    try {  
        while (needMorePrimes())  
            consume(primes.take());  
    } finally {  
        producer.cancel();  
    }  
}
```



Листинг 7.3 Ненадежный механизм отмены, который может оставить задачу производителя застрявшей в заблокированной операции. *Не делайте так.*

Как мы уже упоминали в главе 5, некоторые библиотечные методы блокировки поддерживают прерывание. Прерывание потока является кооперативным механизмом, позволяя одному потоку сигнализировать другому, что он должен, по своему усмотрению и если для него это допустимо, остановить то, что он делает, и сделать что-то другое.

В спецификации API или языка нет ничего, что связывало бы прерывание с какой-либо конкретной семантикой механизма отмены, но на практике использование прерывания для чего-либо, кроме механизма отмены, является хрупким и сложным для сопровождения, особенно в более крупных приложениях.

Каждый поток имеет *статус прерывания* (*interrupted status*) типа `boolean`; прерывание потока устанавливает его статус прерывания в `true`. Класс `Thread` содержит методы для прерывания потока и запроса статуса прерывания, как показано в листинге 7.4. Метод `interrupt` прерывает целевой поток, а метод `isInterrupted` позволяет запросить статус прерывания целевого потока. Плохо именованный статический метод `interrupted` сбрасывает статус прерывания текущего потока и возвращает его предыдущее значение; это единственный способ сбросить статус прерывания.

```
public class Thread {  
    public void interrupt() { ... }  
    public boolean isInterrupted() { ... }  
    public static boolean interrupted() { ... }  
    ...  
}
```

Листинг 7.4 Методы прерывания класса `Thread`

Блокирующие библиотечные методы, такие как `Thread.sleep` и `Object.wait`, пытаются определить, когда поток был прерван, и вернуть управление раньше. Они отвечают на прерывание, сбрасывая статус прерывания и бросая исключение `InterruptedException`, что указывает на то, что блокирующая операция была завершена раньше из-за прерывания. Среда JVM не даёт гарантий, насколько быстро блокирующий метод обнаружит прерывание, но на практике это происходит достаточно быстро.

Если поток прерывается, когда он не заблокирован, его статус прерывания устанавливается, и дальнейшее зависит от отменяемой активности, которая должна для обнаружения прерывания опрашивать статус прерывания. Таким образом, прерывание является "липким" - если не будет возбуждено исключение `InterruptedException`, свидетельство о прерывания сохранится, пока кто-нибудь специально не сбросит статус прерывания.

Вызов метода `interrupt` не обязательно приводит к остановке операции, выполняемой целевым потоком; он просто доставляет сообщение о том, что было запрошено прерывание.

Хороший способ думать о прерывании – это представлять себе, что оно фактически не прерывает выполняющийся поток; оно просто запрашивает, чтобы поток прервал сам себя при следующей удобной возможности (эти возможности называются *точками отмены* (*cancellation points*)). Некоторые методы, такие как `wait`, `sleep` и `join` воспринимают такие запросы серьезно, бросая исключение, когда они получают запрос на прерывание или сталкиваются с уже установленным состоянием прерывания на входе. Хорошо ведут себя методы, которые могут полностью игнорировать такие запросы, оставляя запросы на прерывание на месте,

чтобы вызывающий код мог что-то с ними сделать. Плохо ведут себя методы, “проглатывающие” запросы на прерывание, и тем самым лишающие код, расположенный выше по стеку вызовов, возможности реагировать на них.

Статический метод `interrupted` следует использовать с осторожностью, так как он сбрасывает текущий статус прерывания потока. Если вы вызываете метод `interrupted` и он возвращает `true`, в случае если вы не планируете “проглатывать” прерывание, вы должны что-то с этим сделать - либо бросить исключение `InterruptedException` или восстановить статус прерывания потока, вызвав метод `interrupt` еще раз, как показано в листинге 5.10.

Класс `BrokenPrimeProducer` иллюстрирует тот момент, что пользовательские механизмы отмены не всегда хорошо взаимодействуют с библиотечными методами блокировки. Если вы кодируете задачи так, чтобы они реагировали на прерывание, вы можете использовать прерывание в качестве механизма отмены и воспользоваться поддержкой прерывания, предоставляемой множеством библиотечных классов.

Прерывание, как правило, является наиболее разумным способом реализации механизма отмены.

Класс `BrokenPrimeProducer` можно легко исправить (и упростить), используя прерывание вместо логического флага запроса на отмену, как показано в листинге 7.5. В каждой итерации цикла есть две точки, в которых прерывание может быть обнаружено: в блокирующем вызове метода `put` и явным опросом статуса прерывания в заголовке цикла.

```
class PrimeProducer extends Thread {
    private final BlockingQueue<BigInteger> queue;

    PrimeProducer(BlockingQueue<BigInteger> queue) {
        this.queue = queue;
    }

    public void run() {
        try {
            BigInteger p = BigInteger.ONE;
            while (!Thread.currentThread().isInterrupted())
                queue.put(p = p.nextProbablePrime());
        } catch (InterruptedException consumed) {
            /* Allow thread to exit */
        }
    }

    public void cancel() { interrupt(); }
}
```

Листинг 7.5 Использование прерывания для отмены

Явная проверка здесь не является строго необходимой из-за блокирующего вызова `put`, но она делает класс `PrimeProducer` более отзывчивым к прерыванию, потому что она проверяет прерывание *перед* началом выполнения длительной задачи

поиска простого числа, а не после. Если вызовы прерываемых блокирующих методов недостаточно часты для обеспечения требуемой скорости отклика, может помочь явная проверка статуса прерывания.

7.1.2 Политики прерывания

Так же, как задачи должны иметь политику отмены, также и потоки должны иметь *политику прерывания* (*interruption policy*). Политика прерывания определяет, как поток интерпретирует запрос на прерывание - что он делает (если есть что) при обнаружении, какие единицы работы считаются атомарными по отношению к прерыванию и как быстро он реагирует на прерывание.

Наиболее разумной политикой прерывания является некоторая форма отмены на уровне потока или на уровне службы: скорейшее завершение работы, очистка ресурсов при необходимости и, возможно, уведомление некоторой сущности-владельца о завершении потока. Можно установить и другие политики прерывания, например приостановку или возобновление работы службы, но потоки или пулы потоков с нестандартными политиками прерывания, возможно, потребуется ограничить задачами, которые были написаны с учетом таких политик.

Важно различать, как задачи и потоки должны реагировать на прерывания. Один запрос на прерывание может быть адресован нескольким получателям - прерывание рабочего потока в пуле потоков может означать как "отмена текущей задачи", так и "завершение рабочего потока".

Задачи не выполняются в собственных потоках; они заимствуют потоки, принадлежащие другим службам, например пулу потоков. Код, не являющийся владельцем потока (для пула потоков таким является любой код за пределами реализации пула потоков), должен заботиться о сохранении статуса прерывания, чтобы код-владелец мог в конечном итоге обработать прерывание, даже если "гостевой" код также обработает прерывание. (Если вы по чьей-то просьбе ожидаете у кого-то дома, вы не выбрасываете почту, которая приходит, пока хозяев нет дома, - вы сохраняете ее и позволяете им разобраться с ней, когда они возвращаются домой, даже если вы читаете их журналы.)

Вот почему большинство блокирующих библиотечных методов в ответ на прерывание просто бросают исключение `InterruptedException`. Они никогда не будут выполняться в потоке, которым они владеют, поэтому они реализуют наиболее разумную политику отмены для библиотечного кода или кода задач: как можно быстрее отойти в сторону и передать прерывание обратно вызывающему объекту, чтобы код, расположенный выше по стеку вызовов, мог предпринять дальнейшие действия.

Нет необходимости в том, чтобы задача обязательно что-нибудь бросала, когда обнаруживает запрос на прерывание - она может отложить это до более благоприятного момента, запомнив, что она была прервана, дождаться завершения операции, которую она выполняла, а затем бросив исключение `InterruptedException` или иным образом указав на прерывание. Этот метод может защитить структуры данных от повреждения, когда выполнение активности прерывается в середине операции обновления.

Задача не должна делать предположений о политике прерывания выполняющегося потока, только если она явно не предназначена для выполнения в службе, имеющей определенную политику прерывания. Независимо от того,

интерпретирует ли задача прерывание как отмену или выполняет какое-либо другое действие при прерывании, она должна позаботиться о сохранении состояния прерывания выполняющегося потока. Если она не сможет распространить исключение `InterruptedException` до вызывающего объекта, она должна будет восстановить статус прерывания после перехвата исключения `InterruptedException`:

```
Thread.currentThread().interrupt();
```

Так же, как код задачи не должен делать предположений о том, что прерывание значит для выполняющего его потока, код отмены не должен делать предположений о политике прерывания произвольных потоков. Поток должен прерываться только своим владельцем; владелец может инкапсулировать знание о политике прерывания потока в соответствующем механизме отмены, таком как метод завершения работы.

Поскольку каждый поток имеет свою собственную политику прерывания, не следует прерывать поток, если вы не знаете, что прерывание означает для этого потока.

Критики высмеивали средства прерывания в Java, потому что они не предоставляют возможность упреждающего прерывания, но, в то же время, заставляет разработчиков обрабатывать исключение `InterruptedException`. Однако возможность отложить запрос на прерывание, позволяет разработчикам создавать гибкие политики прерывания, балансирующие отзывчивость и надежность в зависимости от требований приложения.

7.1.3 Ответ на прерывание

Как было отмечено ранее в разделе 5.4, когда вы вызываете прерывание блокирующих методов, таких как `Thread.sleep` или `BlockingQueue.put`, существуют две практические стратегии для обработки исключения `InterruptedException`:

- Распространить исключение (возможно, после некоторой очистки в рамках конкретной задачи), фактически сделав ваш метод прерываемым блокирующим методом; или
- Восстановить статус прерывания, чтобы код, находящийся выше по стеку вызовов, мог с ним работать.

Распространение исключения `InterruptedException` может быть так же просто, как добавление исключения `InterruptedException` в выражение `throws`, как показано в методе `getNextTask`, в листинге 7.6.

```
BlockingQueue<Task> queue;  
...  
public Task getNextTask() throws InterruptedException {  
    return queue.take();  
}
```

Листинг 7.6 Распространение исключения `InterruptedException` до вызывающего кода

Если Вы не хотите или не можете распространять исключение `InterruptedException` (возможно потому, что ваша задача определена с помощью интерфейса `Runnable`), вам нужно найти другой способ сохранить запрос на прерывание. Стандартным способом сделать это, является восстановление статуса прерывания путем повторного вызова метода `interrupt`. Чего вы *не* должны делать, так это проглатывать исключение `InterruptedException`, путём его перехвата и невыполнения никаких действий в блоке `catch`, если только ваш код не реализует политику прерывания для потока. Класс `PrimeProducer` проглатывает прерывание, но делает это со знанием того, что поток собирается завершить свою работу, и что в связи с этим нет никакого кода выше по стеку вызовов, который должен знать о прерывании. Большая часть кода не знает, в каком потоке будет выполняться, поэтому статус прерывания следует сохранять.

Только код, реализующий политику прерывания потока, может проглатывать запрос на прерывание. Задачи общего назначения и библиотечный код никогда не должны проглатывать запросы на прерывание.

Активности, которые не поддерживают отмену, но по-прежнему вызывают прерываемые блокирующие методы, должны вызывать их в цикле, повторяя попытку при обнаружении прерывания. В этом случае они должны локально сохранить статус прерывания и восстановить его непосредственно перед возвратом, как показано в листинге 7.7, а не сразу после перехвата исключения `InterruptedException`. Слишком ранняя установка статуса прерывания может привести к бесконечному циклу, поскольку большинство прерываемых блокируемых методов проверяет состояние прерывания на входе и немедленно кидает исключение `InterruptedException`, если статус прерывания установлен. (Прерываемые методы обычно проводят опрос на предмет установки флага прерывания перед блокировкой или выполнением какой-либо значительной работы, чтобы быть максимально отзывчивыми к прерыванию.)

```
public Task getNextTask(BlockingQueue<Task> queue) {
    boolean interrupted = false;
    try {

        while (true) {
            try {
                return queue.take();
            } catch (InterruptedException e) {
                interrupted = true;
                // fall through and retry
            }
        }
    } finally {
        if (interrupted)
            Thread.currentThread().interrupt();
    }
}
```

Листинг 7.7 Неотменяемые задачи, восстанавливающие прерывание до своего завершения

Если код не вызывает прерываемые блокирующие методы, он по-прежнему может реагировать на прерывание путем опроса статуса прерывания текущего потока по всему коду задачи. Выбор частоты опроса является компромиссом между эффективностью и отзывчивостью. При наличии высоких требований к быстродействию нельзя вызывать потенциально долговременные методы, которые сами не реагируют на прерывания, что потенциально ограничивает возможности по вызову библиотечного кода.

Отмена может включать в себя состояние, отличное от состояния прерывания; прерывание может использоваться для привлечения внимания потока, а информация, хранящаяся в другом месте, может быть использована для предоставления дальнейших инструкций прерванному потоку (при получении доступа к этой информации обязательно используйте синхронизацию). Например, когда рабочий поток, принадлежащий экземпляру `ThreadPoolExecutor`, обнаруживает прерывание, он проверяет, завершает ли пул свою работу. Если это так, он выполняет некоторую очистку пула перед завершением; в противном случае он может создать новый поток, для восстановления численности потоков в пуле до заданных значений.

7.1.4 Пример: запуск по времени

Для решения многих проблем может потребоваться вечность (например, перечисление всех простых чисел); для других ответ может быть найден достаточно быстро, но также может занять вечность. Возможность сказать “потратьте до десяти минут на поиск ответа” или “перечислите все ответы, полученные за десять минут”, может быть крайне полезна в таких ситуациях.

Метод `aSecondOfPrimes` в листинге 7.2 запускает на выполнение экземпляр `PrimeGenerator` и прерывает его выполнение через секунду. Несмотря на то, что остановка экземпляра `PrimeGenerator` может занять несколько дольше времени, чем одна секунда, в конечном счёте, прерывание будет замечено и выполнение остановится, позволив потоку завершиться. Но другой аспект выполнения задачи заключается в том, что необходимо выяснить, бросает ли задача исключение. Если экземпляр `PrimeGenerator` бросит непроверяемое исключение до истечения времени ожидания, это, вероятнее всего, останется незамеченным, так как генератор простых чисел работает в отдельном потоке, в котором обработка исключений в явном виде не прописана.

В листинге 7.8 показана попытка выполнения произвольного экземпляра `Runnable` в течение заданного промежутка времени. Метод запускает задачу в вызывающем потоке и планирует запуск отменяющей задачи, чтобы прервать ее после истечения заданного интервала времени. Это решает проблему непроверяемых исключений, которые могут быть брошены задачей, поскольку они могут быть перехвачены объектами, вызвавшими метод `timedRun`.

```
private static final ScheduledExecutorService cancelExec = ...;

public static void timedRun(Runnable r,
                           long timeout, TimeUnit unit) {
    final Thread taskThread = Thread.currentThread();
    cancelExec.schedule(new Runnable() {
        public void run() { taskThread.interrupt(); }
    }, timeout, unit);
```



```
r.run();  
}
```

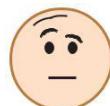
Листинг 7.8 Планирование прерывания в заимствованном потоке. Не делайте так

Такой подход умоляюще прост, но он нарушает правила: вы должны знать о политике прерывания потока, прежде чем прервать его. Так как метод `timedRun` может быть вызван из произвольного потока, он не может знать о политике прерывания вызывающего потока. Если задача завершается до истечения времени ожидания, то отменяющая задача, предназначенная для прерывания работы потока в котором был вызван метод `timedRun`, может завершиться *после* возвращения управления от метода `timedRun` вызывающему объекту. Мы не знаем, какой код будет выполняться, когда это произойдет, но результат хорошим точно не назовешь. (Устранить риск возникновения такой ситуации возможно, хотя и на удивление сложно, если для отмены отменяющей задачи использовать экземпляр интерфейса `ScheduledFuture`, возвращаемый методом `schedule`.)

Кроме того, если задача не реагирует на прерывания, метод `timedRun` не вернет управление до тех пор, пока задача не завершится, что может сильно превысить указанное время ожидания (или даже не завершится вовсе). Служба, ограниченная по времени, но не возвращающая управление по истечении указанного промежутка времени, скорее всего, будет вызывать раздражение у вызвавших её пользователей.

В листинге 7.9 рассматриваются проблемы обработки исключений, возникающих в методе `aSecondOfPrimes`, а также проблемы с предыдущей попыткой. Поток, созданный для выполнения задачи, может иметь собственную политику выполнения, и даже если задача не отвечает на прерывание, метод ограниченный по времени все равно может вернуть управление вызывающему объекту. После запуска потока задачи, метод `timedRun` вызывает метод `join` с параметром ограничения по времени у вновь созданного потока. После того, как метод `join` возвращает управление, проверяется, было ли из задачи брошено исключение и если было, то пробрасывает его в поток, вызвавший метод `timedRun`. Сохраненный экземпляр `Throwable` совместно используется двумя потоками, и поэтому для его безопасной публикации из потока задачи в поток метода `timedRun`, объявляется с ключевым словом `volatile`.

```
public static void timedRun(final Runnable r,  
                           long timeout, TimeUnit unit)  
                           throws InterruptedException {  
  
    class RethrowableTask implements Runnable {  
        private volatile Throwable t;  
        public void run() {  
            try { r.run(); }  
            catch (Throwable t) { this.t = t; }  
        }  
        void rethrow() {  
            if (t != null)  
                throw launderThrowable(t);  
        }  
    }  
  
    RethrowableTask task = new RethrowableTask();  
    final Thread taskThread = new Thread(task);  
    taskThread.start();
```



```
cancelExec.schedule(new Runnable() {
    public void run() { taskThread.interrupt(); }
}, timeout, unit);
taskThread.join(unit.toMillis(timeout));
task.rethrow();
}
```

Листинг 7.9 Прерывание задачи в выделенном потоке

Эта версия устраняет проблемы, описанные в предыдущих примерах, но так как она основана на использовании метода `join` с параметром ограничения по времени, она имеет тот же недостаток, что и при использовании простого метода `join`: мы не знаем, было ли управление возвращено в результате нормального завершения потока или в связи с истечением времени ожидания метода `join`⁸².

7.1.5 Выполнение отмены с помощью интерфейса Future

Мы уже использовали абстракцию для управления жизненным циклом задачи, обработки исключений и облегчения операции отмены - интерфейс `Future`. Следуя общему принципу, заключающемуся в том, что лучше использовать существующие библиотечные классы, чем разворачивать свои собственные, давайте построим метод `timedRun` с использованием интерфейса `Future` и фреймворка выполнения задач.

Вызов `ExecutorService.submit` возвращает экземпляр `Future` описывающий задачу. Интерфейс `Future` имеет метод `cancel`, принимающий аргумент `mayInterruptIfRunning` типа `boolean`, и возвращает значение, указывающее, была ли попытка отмены операции успешной. (Результат выполнения метода расскажет вам только о том, удалось ли выполнить прерывание, а не о том, была ли обнаружена задача, и затронуло ли её прерывание.) Если параметр `mayInterruptIfRunning` имеет значение `true` и задача в данный момент выполняется в некотором потоке, то этот поток прерывается. Установка этого параметра в `false` означает “не выполнять эту задачу, если она еще не запущена” и должна применяться к задачам, которые изначально не были спроектированы для обработки прерывания.

Поскольку вы не должны прерывать поток, если не знаете о его политике прерывания, каким образом приемлемо вызвать метод `cancel` с аргументом `true`? Потоки, выполняющие задачи, созданные стандартными реализациями интерфейса `Executor`, реализуют политику прерывания, которая позволяет задачам быть отмененными с помощью прерывания, таким образом, вполне безопасно установить значение параметру `mayInterruptIfRunning` при отмене задач через их экземпляры `Future`, когда они выполняются стандартной реализацией интерфейса `Executor`. При попытке отменить задачу, вы не должны напрямую прерывать поток, находящийся в пуле, так как вы не знаете, какая задача будет выполняться в момент доставки запроса на прерывание - делайте это только через связанный с задачей экземпляр `Future`. Это еще одна причина для кодирования задач, рассматривающих прерывание как запрос на отмену: в таком случае они могут быть отменены через их экземпляры `Future`.

⁸² Это является недостатком API потоков, потому что независимо от того, завершается ли вызов метода `join` успешно или нет, согласно модели памяти Java для видимости памяти наступают последствия, но метод `join` всё равно не возвращает статус успешности завершения выполнения.

В листинге 7.10 показана версия метода `timedRun`, которая отправляет задачу экземпляру `ExecutorService` и получает результат с помощью ограниченной по времени версии `Future.get`. Если вызов метода `get` завершается возбуждением исключения `TimeoutException`, задача отменяется с помощью собственного экземпляра `Future`. (В целях упрощения кодирования, приведённая версия безусловно вызывает метод `Future.cancel` в блоке `finally`, пользуясь тем преимуществом, что отмена завершенной задачи не оказывает никакого влияния.) Если базовое вычисление⁸³ бросает исключение до момента отмены, оно прорывается дальше из метода `timedRun`, что позволяет вызывающему объекту иметь дело с исключением наиболее удобным способом. В листинге 7.10 также иллюстрируется другая хорошая практика: отмена задач, результат выполнения которых больше не нужен.

```
public static void timedRun(Runnable r,
                            long timeout, TimeUnit unit)
                            throws InterruptedException {
    Future<?> task = taskExec.submit(r);
    try {
        task.get(timeout, unit);
    } catch (TimeoutException e) {
        // task will be cancelled below
    } catch (ExecutionException e) {
        // exception thrown in task; rethrow
        throw launderThrowable(e.getCause());
    } finally {
        // Harmless if task already completed
        task.cancel(true); // interrupt if running
    }
}
```

Листинг 7.10 Отмена задачи с использованием `Future`

Когда вызов `Future.get` бросает исключение `InterruptedException` или `TimeoutException`, и вы знаете, что результат программе больше не нужен, отменяйте задачу с помощью вызова метода `Future.cancel`.

7.1.6 Работа с непрерываемыми блокировками

Многие блокирующие библиотечные методы реагируют на прерывание, возвращая управление как можно раньше и бросая исключение `InterruptedException`, что упрощает построение задач, реагирующих на отмену. Однако не все блокирующие методы или механизмы блокировки реагируют на прерывание; если поток блокируется при выполнении операций синхронного ввода/вывода сокета или ожидает захвата встроенной блокировки, прерывание не окажет никакого эффекта, кроме установки статуса прерывания потока. Иногда мы можем убедить потоки, заблокированные в непрерываемых действиях, чтобы они остановились, с помощью средств, аналогичных прерыванию, но это требует большего понимания того, по какой причине поток был заблокирован.

⁸³ Код, выполняемый в экземпляре `Runnable`.

Синхронные операции ввода/вывода сокета в java.io. Распространенной формой блокировки ввода/вывода в серверных приложениях является чтение или запись в сокет. К сожалению, методы `read` и `write` в классах `InputStream` и `OutputStream` не реагируют на прерывание, но закрытие лежащего в основе сокета приводит к тому, что любые потоки, заблокированные в методах `read` или `write`, бросают исключение `SocketException`.

Синхронные операции ввода/вывода в java.nio. Прерывание потока ожидающего канал `InterruptibleChannel` приводит к возбуждению исключения `ClosedByInterruptException` и закрытию канала (а также вынуждает все прочие потоки, заблокированные на этом канале бросить исключение `ClosedByInterruptException`). Закрытие экземпляра `InterruptibleChannel` `AsynchronousCloseException` заставляет потоки, заблокированные на канальных операциях, бросать исключение `AsynchronousCloseException`. Большая часть стандартных экземпляров `Channel` реализует интерфейс `InterruptibleChannel`.

Асинхронные операции ввода/вывода с классом Selector. Если поток заблокирован в вызове метода `Selector.select` (в языке `java.nio.channels`), вызов методов `close` или `wakeup` приводит к преждевременному возврату.

Захват блокировки. Если поток блокируется в ожидании встроенной блокировки, вы ничего не можете сделать для того, чтобы вскоре остановить его, за исключением того, что вы могли бы привлечь его внимание каким-то другим способом, когда он, в конечном итоге, захватит блокировку и достигнет достаточного прогресса в выполнении задачи. Однако явно определённые классы `Lock` предлагают метод `lockInterruptibly`, который позволяет ожидать захвата блокировки и, при этом, оставаться отзывчивым к прерываниям - см. главу [13](#).

В классе `ReaderThread` из листинга 7.11, демонстрируется метод инкапсуляции нестандартной отмены. Класс `ReaderThread` управляет одиночным соединением сокета, осуществляя синхронное чтение из сокета и передавая все полученные данные методу `processBuffer`. Для упрощения разрыва пользовательского соединения или завершения работы сервера, класс `ReaderThread` переопределяет метод `interrupt`, используя его как для доставки стандартного прерывания, так и для закрытия нижележащего сокета; таким образом, прерывание экземпляра `ReaderThread` приводит к остановке того, что он делает, независимо от того, заблокирован ли он в вызове метода `read` или в прерываемом блокирующем методе

```
public class ReaderThread extends Thread {  
    private final Socket socket;  
    private final InputStream in;  
  
    public ReaderThread(Socket socket) throws IOException {  
        this.socket = socket;  
        this.in = socket.getInputStream();  
    }  
}
```

```

public void interrupt() {
    try {
        socket.close();
    }
    catch (IOException ignored) { }
    finally {
        super.interrupt();
    }
}

public void run() {
    try {
        byte[] buf = new byte[BUFSZ];
        while (true) {
            int count = in.read(buf);
            if (count < 0)
                break;
            else if (count > 0)
                processBuffer(buf, count);
        }
    } catch (IOException e) { /* Allow thread to exit */ }
}
}

```

Листинг 7.11 Инкапсулирование нестандартной отмены в потоке путем переопределения метода `interrupt`

7.1.7 Инкапсуляция нестандартной отмены с помощью метода `newTaskFor`

Подход, используемый классом `ReaderThread` для инкапсуляции нестандартной отмены, может быть усовершенствован с помощью метода-ловушки `newTaskFor`, добавленного классу `ThreadPoolExecutor` в Java 6. При отправке экземпляра `Callable` службе `ExecutorService`, метод `submit` возвращает экземпляр `Future`, который можно использовать для отмены задачи. Метод-ловушка `newTaskFor` является фабричным методом, который создает экземпляр `Future`, представляющий задачу. Он возвращает экземпляр `RunnableFuture` - интерфейс, расширяющий интерфейсы `Future` и `Runnable` (и реализуемый классом `FutureTask`).

Настройка экземпляра `Future`, связанного с задачей, позволяет переопределить метод `Future.cancel`. Пользовательский код отмены может выполнять ведение журнала или собирать статистику по отменам, а также может использоваться для отмены действий, которые не реагируют на прерывание. Класс `ReaderThread` инкапсулирует отмену работы с сокетом – используя потоки переопределяет метод `interrupt`; то же самое может быть сделано для задач, путём переопределения метода `Future.cancel`.

В листинге 7.12, определяется интерфейс `CancellableTask`, расширяющий интерфейс `Callable` и добавляются методы `cancel` и фабричный метод `newTask` для построения экземпляра `RunnableFuture`. Класс `CancellingExecutor` расширяет класс `ThreadPoolExecutor` и переопределяет метод `newTaskFor`, чтобы

позволить интерфейсу CancellableTask создать свой собственный экземпляр Future.

```
public interface CancellableTask<T> extends Callable<T> {
    void cancel();
    RunnableFuture<T> newTask();
}

@ThreadSafe
public class CancellingExecutor extends ThreadPoolExecutor {
    ...
    protected<T> RunnableFuture<T> newTaskFor(Callable<T> callable) {
        if (callable instanceof CancellableTask)
            return ((CancellableTask<T>) callable).newTask();
        else
            return super.newTaskFor(callable);
    }
}

public abstract class SocketUsingTask<T>
    implements CancellableTask<T> {
    @GuardedBy("this") private Socket socket;

    protected synchronized void setSocket(Socket s) { socket = s; }

    public synchronized void cancel() {
        try {
            if (socket != null)
                socket.close();
        } catch (IOException ignored) { }
    }

    public RunnableFuture<T> newTask() {
        return new FutureTask<T>(this) {
            public boolean cancel(boolean mayInterruptIfRunning) {
                try {
                    SocketUsingTask.this.cancel();
                } finally {
                    return super.cancel(mayInterruptIfRunning);
                }
            }
        };
    }
}
```

Листинг 7.12 Инкапсуляция нестандартной отмены в задачу с помощью метода newTaskFor

Класс SocketUsingTask реализует интерфейс CancellableTask и определяет метод Future.cancel, используемый для закрытия сокета, а также вызова метода super.cancel. Если класс SocketUsingTask отменяется через собственный экземпляр Future, сокет закрывается и выполнение потока прерывается. Это

повышает отзывчивость задачи к запросам на отмену: она может не только безопасно вызывать прерываемые блокирующие методы, оставаясь отзывчивой к запросам на отмену, но также может вызывать блокирующие методы ввода/вывода сокета.

7.2 Остановка служб основанных на потоках

Приложения обычно создают службы, владеющие потоками, например пулы потоков, и время жизни этих служб обычно превышает время существования создавших их методов. Если приложение должно завершиться корректно, потоки, принадлежащие этим службам, также должны быть завершены. Поскольку нет упреждающего способа остановить поток, их необходимо убедить самостоятельно завершить свою работу.

Здравомыслящие практики инкапсуляции диктуют, что не следует манипулировать потоком - прерывать его, изменять его приоритет и т. д. - если только вы не являетесь его владельцем. Формально API потока не имеет концепции “владения потоком”: поток представлен объектом `Thread`, который может свободно совместно использоваться, как и любой другой объект. Однако, имеет смысл думать о потоке как об имеющем владельца, и обычно таковым оказывается класс, создавший поток. Таким образом, пул потоков владеет своими рабочими потоками, и если эти потоки должны быть прерваны, пул потоков должен позаботиться об этом.

Как и в случае с любым другим инкапсулированным объектом, владение потоком не является транзитивным: приложение может владеть службой, а служба может владеть рабочими потоками, но приложение не владеет рабочими потоками и поэтому не должно пытаться остановить их напрямую. Вместо этого служба должна предоставлять методы жизненного цикла для завершения работы, которые также завершают работу собственных потоков; затем приложение может завершить работу службы, а служба может завершить работу потоков. Интерфейс `ExecutorService` предоставляет методы `shutdown` и методы `shutdownNow`; другие службы, владеющие потоками, должны обеспечивать аналогичный механизм завершения работы.

Предоставляйте методы жизненного цикла всякий раз, когда срок жизни службы, владеющей потоком, превышает срок жизни метода, её создавшего.

7.2.1 Пример: Сервис логирования

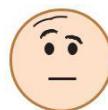
Большинство серверных приложений использует логирование, которое может быть реализовано так же просто, как вставка в код операторов `println`. Поточные классы, такие как `PrintWriter`, потокобезопасны, так что этот простой подход не требует выполнения явной синхронизации⁸⁴. Однако, как мы увидим в разделе 11.6, встроенный в код механизм ведения журнала может приводить к заметным затратам производительности в приложениях большого размера. Другой

⁸⁴ Если вы логируете несколько строк как части одного сообщения журнала, может потребоваться дополнительная блокировка на стороне клиента для предотвращения нежелательного чередования вывода из нескольких потоков. Если два потока логируют многострочные трассировки стека (**stack traces**) в один и тот же поток с вызовом метода `println` на каждую строку, результаты будут непредсказуемо чередоваться и вполне могут сойти за одну большую, но бессмысленную трассировку стека.

альтернативой при вызове метода `log`, является помещение сообщения в очередь для обработки другим потоком.

Класс `LogWriter` из листинга 7.13 демонстрирует пример простой службы логирования, в которой операция логирования перемещена в отдельный логирующий поток. Вместо того чтобы иметь поток, который производит сообщения и тут же их записывает непосредственно в выходной поток, класс `LogWriter` передает сообщения логирующему потоку с помощью класса `BlockingQueue`, и логирующий поток записывает их. Этот дизайн подразумевает несколько производителей и одного потребителя: любая активность, вызывающая метод `log` выступает в качестве производителя, а фоновый логирующий поток выступает в качестве потребителя. Если логирующий поток отстает, класс `BlockingQueue`, в конечном счёте, заблокирует производителей, пока логирующий поток их не нагонит.

```
public class LogWriter {  
    private final BlockingQueue<String> queue;  
    private final LoggerThread logger;  
  
    public LogWriter(Writer writer) {  
        this.queue = new LinkedBlockingQueue<String>(CAPACITY);  
        this.logger = new LoggerThread(writer);  
    }  
  
    public void start() { logger.start(); }  
  
    public void log(String msg) throws InterruptedException {  
        queue.put(msg);  
    }  
  
    private class LoggerThread extends Thread {  
        private final PrintWriter writer;  
        ...  
        public void run() {  
            try {  
                while (true)  
                    writer.println(queue.take());  
            } catch(InterruptedException ignored) {  
            } finally {  
                writer.close();  
            }  
        }  
    }  
}
```



Листинг 7.13 Служба логирования производитель-потребитель без поддержки завершения работы

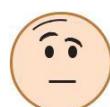
Чтобы служба подобная классу `LogWriter` была полезна в продуктиве, нам нужен способ завершения работы логирующего потока, такой, чтобы он не препятствовал нормальному завершению работы JVM. Остановить логирующий поток достаточно просто, поскольку он многократно вызывает чувствительный к прерыванию метод `take`; если логирующий поток будет изменён и будет завершать

свою работу при перехвате исключения `InterruptedException`, то прерывание логирующего потока остановит работу службы.

Однако простой выход из логирующего потока является не очень удовлетворительным механизмом завершения работы. Такое резкое завершение работы приведёт к сбросу логируемых сообщений, ожидающих записи в лог, но, что ещё более важно, потоки, заблокированные в методе `log`, никогда не будут разблокированы, так как очередь заполнена. Отмена выполнения активности производителя-потребителя требует отмены и производителей, и потребителей. Прерывание логирующего потока имеет дело с потребителем, но поскольку производители в данном случае не являются отдельными потоками, их отменить значительно сложнее.

Другой подход к завершению работы класса `LogWriter` заключается в том, чтобы установить значение флага "запрос на завершение работы", чтобы предотвратить отправку дальнейших сообщений, как показано в листинге 7.14. Получив уведомление о запросе на завершение работы, потребитель мог бы тогда "осушить" очередь, записывая любые ожидающие сообщения в лог и разблокировав всех производителей, заблокированных в методе `log`. Однако такой подход содержит в себе предпосылки для возникновения условия гонок, что делает его ненадежным. Реализация метода `log` представляет собой последовательность операций проверить-затем-выполнить: производители могли бы наблюдать, что работа службы ещё не была завершена, и продолжать клать сообщения в очередь, даже после завершения работы потребителя, таким образом, вновь возникает риск того, что производители могли бы заблокироваться в методе `log` и никогда не разблокироваться. Есть приёмы, которые уменьшают вероятность возникновения такой ситуации (например, потребитель ждёт несколько секунд, прежде чем объявить очередь опустошённой), но они не меняют механизмов, лежащих в основе проблемы, просто уменьшают вероятность того, что по этой причине произойдёт сбой.

```
public void log(String msg) throws InterruptedException {
    if (!shutdownRequested)
        queue.put(msg);
    else
        throw new IllegalStateException("logger is shut down");
}
```



Листинг 7.14 Ненадежный способ добавить поддержку завершения работы в службу логирования

Способ обеспечения надежного завершения работы класса `LogWriter` заключается в том, чтобы исправить условия, при которых возникает состояние гонки, что означает, что отправка новых сообщений в журнал должна происходить атомарно. Но мы не хотим удерживать блокировку при попытке поставить сообщение в очередь, так как метод `put` может быть заблокирован. Вместо этого мы можем атомарно проверить, была ли работа службы завершена и по условию, если служба работает, увеличить значение счетчика, чтобы "зарезервировать" право на отправку сообщения, как показано в классе `LogService` в листинге 7.15.

```
public class LogService {
    private final BlockingQueue<String> queue;
    private final LoggerThread loggerThread;
```

```

private final PrintWriter writer;
@GuardedBy("this") private boolean isShutdown;
@GuardedBy("this") private int reservations;

public void start() { loggerThread.start(); }

public void stop() {
    synchronized (this) { isShutdown = true; }
    loggerThread.interrupt();
}

public void log(String msg) throws InterruptedException {
    synchronized (this) {
        if (isShutdown)
            throw new IllegalStateException(...);
        ++reservations;
    }
    queue.put(msg);
}

private class LoggerThread extends Thread {
    public void run() {
        try {
            while (true) {
                try {
                    synchronized (LogService.this) {
                        if (isShutdown && reservations == 0)
                            break;
                    }
                    String msg = queue.take();
                    synchronized (LogService.this) {
                        --reservations;
                    }
                    writer.println(msg);
                } catch (InterruptedException e) { /* retry */ }
            }
        } finally {
            writer.close();
        }
    }
}
}

```

Листинг 7.15 Добавление надёжного механизма отмены классу LogWriter

7.2.2 Завершение работы ExecutorService

Ранее, в разделе 6.2.4, мы уже видели, что класс ExecutorService предлагает два способа завершения работы: плавное завершение работы с помощью метода `shutdown` и резкое завершение работы с помощью метода `shutdownNow`. При резком завершении работы, метод `shutdownNow`, после попытки отменить все

активно выполняемые задачи, возвращает список задач, которые еще не были запущены.

Два различных варианта завершения предлагают компромисс между безопасностью и отзывчивостью: резкое завершение быстрее, но более рискованно, потому что задачи могут быть прерваны в процессе выполнения, в свою очередь нормальное завершение медленнее, но безопаснее, потому что класс `ExecutorService` не завершит свою работу до тех пор, пока все поставленные в очередь задачи не будут обработаны. Другим службам-владельцам потоков следует рассмотреть возможность предоставления аналогичного набора режимов завершения работы.

Простые программы могут легко запускать и завершать работу глобального экземпляра службы `ExecutorService` из метода `main`. Более сложные программы, вероятнее всего, инкапсулируют экземпляр `ExecutorService` в службу более высокого уровня, которая предоставляет собственные методы управления жизненным циклом, например такую, как вариант реализации класса `LogService` в листинге 7.16, который делегирует управление потоками экземпляру `ExecutorService`, вместо самостоятельного управления ими. Инкапсуляция экземпляра `ExecutorService` расширяет цепочку владения от приложения к потоку службы, добавив ещё одну ссылку; каждый член цепочки управляет жизненным циклом служб или потоков, которыми он владеет.

```
public class LogService {
    private final ExecutorService exec = newSingleThreadExecutor();
    ...
    public void start() { }

    public void stop() throws InterruptedException {
        try {
            exec.shutdown();
            exec.awaitTermination(TIMEOUT, UNIT);
        } finally {
            writer.close();
        }
    }
    public void log(String msg) {
        try {
            exec.execute(new WriteTask(msg));
        } catch (RejectedExecutionException ignored) { }
    }
}
```

Листинг 7.16 Служба логирования использующая интерфейс `ExecutorService`

7.2.3 Ядовитые пилюли

Еще один способ убедить службу производитель-потребитель завершить свою работу, это использовать **ядовитую пиллюлю** (*poison pill*): узнаваемый объект, помещенный в очередь, что трактуется как “остановитесь, когда вы это получите”. В случае очереди FIFO ядовитые пиллюли гарантируют, что потребители закончат работу со своей очередью перед закрытием, так как любая работа, представленная

до отправки ядовитой пилюли, будет извлечена перед пилюлей; производители не должны отправлять какую-либо работу после помещения ядовитой таблетки в очередь. Класс `IndexingService`, представленный в листингах 7.17, 7.18 и 7.19, демонстрирует версию примера поиска на рабочем столе из листинга 5.8, реализованную с помощью одного производителя и одного потребителя, в которой для выключения службы используется ядовитая пилюля.

```
public class IndexingService {  
    private static final File POISON = new File("");  
    private final IndexerThread consumer = new IndexerThread();  
    private final CrawlerThread producer = new CrawlerThread();  
    private final BlockingQueue<File> queue;  
    private final FileFilter fileFilter;  
    private final File root;  
  
    class CrawlerThread extends Thread { /* Listing 7.18 */ }  
    class IndexerThread extends Thread { /* Listing 7.19 */ }  
  
    public void start() {  
        producer.start();  
        consumer.start();  
    }  
  
    public void stop() { producer.interrupt(); }  
  
    public void awaitTermination() throws InterruptedException {  
        consumer.join();  
    }  
}
```

Листинг 7.17 Завершение работы с помощью ядовитой пилюли

```
public class CrawlerThread extends Thread {  
    public void run() {  
        try {  
            crawl(root);  
        } catch (InterruptedException e) { /* fall through */}  
        finally {  
            while (true) {  
                try {  
                    queue.put(POISON);  
                    break;  
                } catch (InterruptedException e1) { /* retry */}  
            }  
        }  
    }  
  
    private void crawl(File root) throws InterruptedException {  
        ...  
    }  
}
```

```
}
```

Листинг 7.18 Поток производителя в классе IndexingService

```
public class IndexerThread extends Thread {  
    public void run() {  
        try {  
            while (true) {  
                File file = queue.take();  
                if (file == POISON)  
                    break;  
                else  
                    indexFile(file);  
            }  
        } catch (InterruptedException consumed) {}  
    }  
}
```

Листинг 7.19 Поток потребителя в классе IndexingService

Ядовитые пилюли работают только тогда, когда известно количество производителей и потребителей. Подход, принятый в классе `IndexingService` может быть распространен на несколько производителей, если каждый производитель будет помещать пилюлю в очередь и потребитель остановится только тогда, когда он получит таблетки от всех $N_{\text{производителей}}$. Также этот подход может быть распространен и на нескольких потребителей; каждый производитель должен помещать в очередь пилюли на $N_{\text{потребителей}}$, хотя такая структура может стать очень громоздкой при большом количестве производителей и потребителей. Ядовитые пилюли надежно работают только с неограниченными очередями.

7.2.4 Пример: одноразовая служба выполнения

Если метод должен обработать пакет задач и не возвращать управление, пока все задачи не будут выполнены, можно упростить управление жизненным циклом службы с помощью приватного экземпляра `Executor`, время жизни которого ограничено временем жизни метода (методы `invokeAll` и `invokeAny` часто могут быть полезны в таких ситуациях).

Метод `checkMail` в листинге 7.20 параллельно проверяет наличие свежей почты на нескольких хостах. Он создает приватный экземпляр исполнителя (`executor`) и отправляет ему по одной задаче для каждого хоста: затем он инициирует завершение работы исполнителя и ожидает завершения, которое происходит, когда все задачи проверки почты полностью выполняются⁸⁵.

```
boolean checkMail(Set<String> hosts, long timeout, TimeUnit unit)  
    throws InterruptedException {  
    ExecutorService exec = Executors.newCachedThreadPool();  
    final AtomicBoolean hasNewMail = new AtomicBoolean(false);  
    try {
```

⁸⁵ Причина, по которой вместо типа `volatile boolean` был использован тип `AtomicBoolean` заключается в том, что для доступа к флагу `hasNewMail` из внутреннего экземпляра `Runnable` он должен быть объявлен как `final`, что исключает его изменение.

```

        for (final String host : hosts)
            exec.execute(new Runnable() {
                public void run() {
                    if (checkMail(host))
                        hasNewMail.set(true);
                }
            });
    } finally {
        exec.shutdown();
        exec.awaitTermination(timeout, unit);
    }
    return hasNewMail.get();
}

```

Листинг 7.20 Использование приватного экземпляра Executor, чей жизненный цикл ограничен вызовом метода

7.2.5 Ограничения метода shutdownNow

При внезапном завершении работы службы ExecutorService с помощью метода `shutdownNow`, она пытается отменить текущие задачи и возвращает список задач, которые были отправлены службе, но никогда не запускались, чтобы их можно было отметить в логе или сохранить для последующей обработки⁸⁶. Однако, нет общего способа узнать, какие задачи были запущены, но не были завершены. Это означает, что нет способа узнать состояние выполняемых задач во время завершения работы, только если сами задачи не выполняют какую-либо контрольную точку. Чтобы узнать, какие задачи не были выполнены, нужно знать не только то, какие задачи не запускались, но и какие задачи выполнялись в момент завершения работы исполнителя⁸⁷.

Класс `TrackingExecutor` в листинге 7.21 демонстрирует подход к определению того, какие задачи выполнялись в процессе завершения работы. Инкапсулируя экземпляр `ExecutorService` и инструментируя метод `execute` (а так же метод `submit`, не показанный здесь) для сохранения информации о том, какие задачи были отменены после завершения работы, класс `TrackingExecutor` может определять, какие задачи запускались, но нормально не завершились. После завершения работы исполнителя метод `getCancelledTasks` возвращает список отмененных задач. Для того чтобы этот метод работал, задачи, когда возвращают управление, должны сохранять статус прерывания потока, как в любом случае и ведут себя хорошо написанные задачи.

```

public class TrackingExecutor extends AbstractExecutorService {
    private final ExecutorService exec;
    private final Set<Runnable> tasksCancelledAtShutdown =
        Collections.synchronizedSet(new HashSet<Runnable>());
    ...
    public List<Runnable> getCancelledTasks() {
        if (!exec.isTerminated())

```

⁸⁶ Объекты `Runnable`, возвращаемые методом `shutdownNow`, могут не совпадать с объектами, отправленными службе `ExecutorService`: они могут быть *обернутыми* (*wrapped*) экземплярами отправленных задач.

⁸⁷ К сожалению, нет опции завершения работы, в которой задачи, еще не запущенные на выполнение, возвращались бы вызывающему объекту, но незавершенные задачи могли бы быть завершены; такая опция устранила бы это неопределенное промежуточное состояние.

```

        throw new IllegalStateException(...);
    return new ArrayList<Runnable>(tasksCancelledAtShutdown);
}

public void execute(final Runnable runnable) {
    exec.execute(new Runnable() {
        public void run() {
            try {
                runnable.run();
            } finally {
                if (isShutdown()
                    && Thread.currentThread().isInterrupted())
                    tasksCancelledAtShutdown.add(runnable);
            }
        }
    });
}

// delegate other ExecutorService methods to exec
}

```

Листинг 7.21 Экземпляр ExecutorService отслеживающий отмененные задачи после завершения работы

Класс `WebCrawler`, из листинга 7.22, демонстрирует применение класса `TrackingExecutor`. Работа веб-сканера часто не ограничена, поэтому, если сканер должен быть выключен, мы можем захотеть сохранить его состояние, таким образом, он может быть перезапущен позже. Класс `CrawlTask` предоставляет метод `getPage`, позволяющий узнать, с какой страницей он работает. Когда сканер завершает свою работу, то задачи на сканирование, которые не были запущены и те, что были отменены, а также их URL записываются, чтобы задачи сканирования страниц для этих URL-адресов можно было добавить в очередь при перезапуске сканера.

```

public abstract class WebCrawler {
    private volatile TrackingExecutor exec;
    @GuardedBy("this")
    private final Set<URL> urlsToCrawl = new HashSet<URL>();
    ...
    public synchronized void start() {
        exec = new TrackingExecutor(
            Executors.newCachedThreadPool());
        for (URL url : urlsToCrawl) submitCrawlTask(url);
        urlsToCrawl.clear();
    }

    public synchronized void stop() throws InterruptedException {
        try {
            saveUncrawled(exec.shutdownNow());
            if (exec.awaitTermination(TIMEOUT, UNIT))
                saveUncrawled(exec.getCancelledTasks());
        } finally {

```

```

        exec = null;
    }
}

protected abstract List<URL> processPage(URL url);

private void saveUncrawled(List<Runnable> uncrawled) {
    for (Runnable task : uncrawled)
        urlsToCrawl.add(((CrawlTask) task).getPage());
}

private void submitCrawlTask(URL u) {
    exec.execute(new CrawlTask(u));
}

private class CrawlTask implements Runnable {
    private final URL url;
    ...
    public void run() {
        for (URL link : processPage(url)) {
            if (Thread.currentThread().isInterrupted())
                return;
            submitCrawlTask(link);
        }
    }
    public URL getPage() { return url; }
}
}

```

Листинг 7.22 Использование класса `TrackingExecutorService` для сохранения незавершённых задач с целью последующего выполнения

Класс `TrackingExecutor` содержит неизбежное условие гонки, которое может привести к ложным срабатываниям: задачи могут идентифицироваться как отмененные, но фактически будут завершены. Это происходит потому, что пул потоков может быть закрыт в момент времени между выполнением последней инструкции задачи и записью задачи как завершенной. Это не проблема, если задачи *идемпотентны*⁸⁸ (если их повторное выполнение оказывает тот же эффект, что и выполнение в первый раз), как это обычно и бывает в веб-сканере. В противном случае, приложение, извлекающее отмененные задачи, должно учитывать риск возникновения такой ситуации и быть готовым к работе с ложными срабатываниями.

7.3 Обработка аварийного завершения потока

Очевидно, что когда однопоточное консольное приложение завершает свою работу из-за не перехваченного исключения - программа прекращает работу и выводит в консоль трассировку стека, которая сильно отличается от типичного вывода программы. Сбой потока в параллельном приложении не всегда так очевиден. Трассировка стека может быть напечатана в консоли, но за консолью может никто

⁸⁸ **Идемпотентность** - свойство объекта или операции при повторном применении операции к объекту давать тот же результат, что и при первом (из вики).

не наблюдать. Кроме того, при сбое потока приложение может продолжать работать, поэтому его сбой может остаться незамеченным. К счастью, существуют средства обнаружения и предотвращения “утечки” (*leaking*) потоков из приложения.

Основной причиной преждевременной смерти потока является возбуждение исключения `RuntimeException`. Поскольку эти исключения указывают на программную ошибку или другую неисправимую проблему, они обычно не перехватываются. Вместо этого они распространяются вверх по стеку, и в такой ситуации поведение по умолчанию заключается в печати трассировки стека в консоль и завершении потока.

Последствия аварийного завершения потока лежат в диапазоне от “безвредно” до “катастрофа”, в зависимости от роли потока в приложении. Потеря потока из пула потоков может иметь последствия для производительности, но приложение, которое хорошо работает с пулом из 50 потоков, вероятно, будет также хорошо работать и с пулом из 49 потоков. С другой стороны, потеря потока диспетчеризующего события в приложении с графическим интерфейсом, была бы весьма заметна - приложение бы прекратило обработку событий, и в результате графический интерфейс как бы “замерз”. В классе `OutOfTime` были продемонстрированы серьезные последствия, к которым привела утечка потока: служба, представленная классом `Timer`, постоянно выходит из строя.

Почти любой код может бросить исключение `RuntimeException`. Всякий раз, когда вы вызываете другой метод, вы верите, что метод нормально вернёт управление или бросит одно из проверяемых исключений, объявленных в его сигнатуре. Чем меньше вы знакомы с вызываемым кодом, тем более скептически вы должны относиться к его поведению.

Потоки, обрабатывающие задачи, такие как рабочие потоки в пуле потоков или поток диспетчеризации событий `Swing`, проводят всю свою жизнь, вызывая неизвестный код через абстрактные барьеры, такие как экземпляры `Runnable`, и эти потоки должны быть очень скептически настроены по отношению к тому, что код, который они вызывают, будет хорошо себя вести. Было бы очень плохо, если бы служба, подобная потоку событий `Swing`, упала бы (*failed*) только потому, что какой-то плохо написанный обработчик событий бросил бы исключение `NullPointerException`. Соответственно, эти средства⁸⁹ должны вызывать задачи в блоке `try-catch`, который перехватывает непроверяемые исключения, или в блоке `try-finally`, чтобы гарантировать, что, если поток завершает работу аварийно, фреймворк будет проинформирован об этом и сможет предпринять корректирующие действия. Это один из немногих случаев, когда может потребоваться перехват исключения `RuntimeException` - при вызове неизвестного, ненадежного кода с использование некой абстракции, такой как `Runnable`⁹⁰.

В листинге 7.23 иллюстрируется способ структурирования рабочего потока в пуле потоков. Если задача бросает непроверяемое исключение, она позволяет потоку умереть, но не ранее момента уведомления фреймворка о том, что поток умер. Затем фреймворк может принять решение о замене рабочего потока новым потоком или не заменять его, так как пул потоков уже завершил работу или уже имеется достаточное количество рабочих потоков, полностью удовлетворяющее текущий спрос. Класс `ThreadPoolExecutor` и фреймворк `Swing` используют этот метод, чтобы гарантировать, что плохо работающая задача не создаст препятствий

⁸⁹ Имеются в виду потоки, обрабатывающие задачи.

⁹⁰ До сих пор идут споры по поводу безопасности такого подхода; когда поток бросает непроверяемое исключение, всё приложение может быть скомпрометировано. Но альтернатива – завершение работы приложения - обычно непрактична.

для выполнения последующих задач. Если вы пишете класс рабочего потока, который выполняет отправленные задачи, или вызываете ненадежный внешний код (например, динамически загружаемые плагины), используйте один из представленных подходов для предотвращения завершения потока, в котором происходит вызов, плохо написанной задачи или плагина.

```
public void run() {
    Throwable thrown = null;
    try {
        while (!isInterrupted())
            runTask(getTaskFromWorkQueue());
    } catch (Throwable e) {
        thrown = e;
    } finally {
        threadExited(this, thrown);
    }
}
```

Листинг 7.23 Типичная структура рабочего потока в пуле потоков

7.3.1 Обработчики неперехваченных исключений

В предыдущем разделе предлагалось использовать упреждающий подход к проблеме непроверяемых исключений. Поточное API также предоставляет средство `UncaughtExceptionHandler`, которое позволяет обнаруживать ситуацию, когда поток умирает из-за неперехваченного исключения.

Когда поток завершает свою работу из-за неперехваченного исключения, JVM сообщает о возникновении этого события экземпляру `UncaughtExceptionHandler`, предоставленному приложением (см. листинг 7.24); если обработчик не существует, поведение по умолчанию будет заключаться в печати трассировки стека в поток `System.err`⁹¹.

```
public interface UncaughtExceptionHandler {
    void uncaughtException(Thread t, Throwable e);
}
```

Листинг 7.24 Интерфейс `UncaughtExceptionHandler`

Что обработчик должен сделать с неперехваченным исключением, зависит от требований к качеству обслуживания (*quality-of-service*). Наиболее распространенным ответом является запись сообщения об ошибке и трассировки стека в лог приложения, как показано в листинге 7.25. Обработчики также могут

⁹¹ До Java 5.0 единственным способом управления обработчиком `UncaughtExceptionHandler` было создание подклассов `ThreadGroup`. В Java 5.0 и более поздних версиях можно установить обработчик `UncaughtExceptionHandler` каждому потоку с помощью метода `Thread.setUncaughtExceptionHandler`, а также можно установить обработчик `UncaughtExceptionHandler` по умолчанию, с помощью метода `Thread.setDefaultUncaughtExceptionHandler`. Однако вызван будет только один из этих обработчиков - сначала JVM ищет обработчик для каждого потока, а затем обработчик класса `ThreadGroup`. Реализация обработчика по умолчанию в `ThreadGroup` делегирует ответственность за обработку родительской группы потоков и так далее по цепочке до тех пор, пока один из обработчиков `ThreadGroup` не обработает неперехваченное исключение или оно не всплывет до группы потоков верхнего уровня. Обработчик группы потоков верхнего уровня делегирует обработку системному обработчику по умолчанию (если он существует; по умолчанию - нет), в противном случае выводит трассировку стека в консоль.

выполнять и более непосредственные действия, такие как попытка перезапуска потока, завершение работы приложения, отправка сообщения оператору на пейджер⁹² или другие корректирующие или диагностические действия.

```
public class UEHLogger implements Thread.UncaughtExceptionHandler {  
    public void uncaughtException(Thread t, Throwable e) {  
        Logger logger = Logger.getAnonymousLogger();  
        logger.log(Level.SEVERE,  
                   "Thread terminated with exception: " + t.getName(),  
                   e);  
    }  
}
```

Листинг 7.25 Обработчик UncaughtExceptionHandler логирующий информацию об исключении

В долговременных приложениях для всех потоков всегда используйте обработчики неперехваченных исключений, которые должны, по меньшей мере, записывать информацию о возникшем исключении в лог.

Для того, чтобы пулу потоков установить экземпляр обработчика `UncaughtExceptionHandler`, предоставьте экземпляр `ThreadFactory` конструктору `ThreadPoolExecutor`. (Как и во всех прочих манипуляциях с потоками, только владелец потока должен изменять обработчик `UncaughtExceptionHandler`.) Стандартные пулы потоков позволяют брошенному задачей неперехваченному исключению завершать работу потока в пуле, но используют блок *try-finally* для уведомления, когда это происходит, чтобы поток можно было заменить. Без обработчика неперехваченных исключений или другого механизма уведомления о сбоях, может оказаться так, что задачи завершаются со сбоем в фоновом режиме (тихо), что может привести к путанице. Если вы хотите получать уведомления, когда задача “падает” из-за возбуждённого исключения, так что бы вы могли выполнить специфичные для задачи действия по восстановлению, либо оберните задачу с помощью экземпляров `Runnable` или `Callable`, которые перехватывают исключения, либо переопределите хук (*hook*, ловушку)⁹³ `afterExecute` экземпляра `ThreadPoolExecutor`.

Несколько сбивает с толку, что исключения, бросаемые задачами, превращаются в неперехваченные обработчиком исключений только для задач, отправленных с помощью метода `execute`; для задач, отправленных с помощью `submit`, любое брошенное исключение, проверяемое оно или нет, считается частью состояния возврата задачи. Если задача, отправленная с помощью метода `submit`, завершается с исключением, исключение, обёрнутое в экземпляр `ExecutionException`, пробрасывается дальше с помощью метода `Future.get`.

7.4 Завершение работы JVM

Среда JVM может завершать свою работу в *плановой* (*orderly*) или *внезапной* (*abrupt*) манере. Плановое завершение работы инициируется, когда последний “нормальный” (не демон) поток завершается, кто-то вызывает метод `System.exit`

⁹² Кто сейчас вспомнит, что это такое, а когда-то была очень полезная и распространённая шутка.

⁹³ Хук – метод-перехватчик, выражаясь другими словами – обработчик, срабатывающий при наступлении некоторого события.

или другими, специфичными для платформы, средствами (например, отправка сигнала SIGINT или нажатие комбинации клавиш Ctrl-C). Хотя это стандартный и предпочтительный способ завершения работы виртуальной машины Java, ее работу также можно завершить внезапно, вызвав метод `Runtime.halt` или убив процесс JVM средствами операционной системы (например, отправив процессу сигнал SIGKILL).

7.4.1 Завершающие хуки

При плановом завершении работы, среда JVM в первую очередь запускает все зарегистрированные *завершающие хуки* (*shutdown hooks*). Завершающие хуки представляют собой незапущенные потоки, зарегистрированные с помощью метода `Runtime.addShutdownHook`. Среда JVM не даёт никаких гарантий относительно порядка, согласно которому будут запускаться завершающие хуки. Если потоки какого-либо приложения (демоны или не демоны) все еще выполняются в момент начала завершения работы, они продолжат выполнение параллельно с процессом завершения работы. Когда все завершающие хуки выполняются, JVM может принять решение о выполнении финализаторов, если метод `runFinalizersOnExit` вернёт значение `true`, и затем остановится. Среда JVM не предпринимает никаких попыток остановить или прервать потоки приложения, которые все еще продолжают выполняться во время завершения работы; они будут резко завершены, когда JVM, в конечном счете, остановится. Если завершающие хуки или финализаторы ещё не закончили выполнение, то плановый процесс завершения работы “зависает” и работа JVM должна быть завершена внезапно. При внезапном завершении работы, от среды JVM не требуется делать ничего, кроме остановки JVM; завершающие хуки выполняться не будут.

Завершающие хуки должны быть потокобезопасными: они должны использовать синхронизацию в процессе доступа к совместно используемым данным и должны заботиться о том, чтобы избегать взаимоблокировок, как, впрочем, и любой другой параллельный код. Кроме того, они не должны строить предположений о состоянии приложения (например, о том, завершили ли уже свою работу другие службы или выполнились ли все обычные потоки) или о том, почему завершает работу JVM, и поэтому должны быть закодированы с чрезвычайной степенью защиты. Наконец, они должны завершать своё выполнение как можно быстрее, так как их существование задерживает завершение работы JVM в то время, когда пользователь может ожидать, что JVM завершится быстро.

Завершающие хуки можно использовать для очистки служб или приложений, например для удаления временных файлов или очистки ресурсов, которые автоматически не очищаются ОС. В листинге 7.26 демонстрируется, как класс `LogService` из листинга 7.16 может зарегистрировать в методе `start` завершающий хук, чтобы обеспечить закрытие файла журнала после завершения работы.

```
public void start() {
    Runtime.getRuntime().addShutdownHook(new Thread() {
        public void run() {
            try { LogService.this.stop(); }
            catch (InterruptedException ignored) {}
        }
    })
}
```

```
});  
}
```

Листинг 7.26 Регистрация завершающего хука для остановки службы логирования.

Поскольку все завершающие хуки выполняются одновременно, закрытие файла лога может привести к проблемам с другими завершающими хуками, которые также хотят использовать службу логирования. Чтобы избежать этой проблемы, завершающие хуки не должны полагаться на службы, которые могут быть закрыты приложением или другими завершающими хуками. Один из способов достигнуть этого - использовать один завершающий хук для всех служб, а не по одному хуку для каждой службы, и обязать его последовательно вызывать действия по завершению работы. Это гарантирует, что действия по завершению работы выполняются последовательно в одном потоке, что позволит избежать возникновения условий гонки или взаимоблокировок между действиями, завершающими работу. Этот метод можно использовать независимо от того, используете ли вы завершающие хуки или нет; последовательное, а не параллельное выполнение действий по завершению работы, позволяет устраниить множество потенциальных источников проблем. В приложениях, которые поддерживают явные сведения о зависимостях между службами, этот метод также может гарантировать, что действия по завершению работы выполняются в правильном порядке.

7.4.2 Потоки демоны

Иногда вы хотите создать поток, который выполняет некоторую вспомогательную функцию, но вы не хотите, чтобы существование этого потока препятствовало завершению работы JVM. Для этого предназначены потоки демоны (*daemon threads*).

Потоки делятся на два типа: обычные потоки и потоки демоны. При запуске виртуальной машины Java, все создаваемые ею потоки (например, сборщик мусора и другие служебные потоки) являются потоками демонами, за исключением основного потока. При создании нового потока он наследует статус демона, от создавшего потока, поэтому по умолчанию все потоки, созданные основным потоком, также являются обычными потоками.

Обычные потоки и потоки демоны различаются только тем, что происходит при завершении работы. Когда поток завершает работу, JVM выполняет инвентаризацию выполняющихся потоков, и если единственные оставшиеся потоки являются потоками демонами, среда инициирует упорядоченное плановое завершение работы. Когда JVM останавливается, работа всех оставшихся потоков демонов резко прерывается, - блоки `finally` не выполняются, стеки не разматываются - JVM просто завершает свою работу.

Потоки демоны должны использоваться крайне умеренно - выполнение нескольких выполняющихся активностей может быть прервано в любой момент времени без очистки. В частности, опасно использовать потоки демоны для задач, которые могут выполнять какие-либо операции ввода/вывода. Потоки демоны лучше всего приберегать для задач "очистки", таких как фоновый поток, который периодически удаляет устаревшие записи из кэша в памяти.

Потоки демоны не являются хорошим вариантом для замены обычных потоков в контексте правильного управления жизненным циклом служб приложения.

7.4.3 Финализаторы

Сборщик мусора выполняет полезную работу по освобождению ресурсов памяти, когда они больше не используются, но некоторые ресурсы, такие как дескрипторы файлов или сокетов, должны быть явно возвращены операционной системе, когда больше не нужны. Чтобы помочь этому процессу, сборщик мусора обрабатывает объекты, имеющие специальный нетривиальный метод `finalize`: после их освобождения сборщиком мусора, вызывается метод `finalize`, позволяющий освободить удерживаемые ресурсы.

Поскольку финализаторы могут работать в потоке управляемом JVM, любое состояние, доступ к которому получается финализатор, будет доступно более чем в одном потоке и, следовательно, доступ к нему должен осуществляться только с использованием синхронизации. Финализаторы не предоставляют гарантий относительно момента времени, в который они будут выполнены, или даже самого факта выполнения и приводят к значительным затратам производительности на объекты обладающие нетривиальными финализаторами. Кроме того, их довольно сложно написать корректно⁹⁴. В большинстве случаев сочетание блоков `finally` и явных вызовов методов `close` лучше справляется с управлением ресурсами, чем финализаторы; единственным исключением является необходимость управления объектами, содержащими ресурсы, полученные native методами. По этим и другим причинам старайтесь избегать написания или использования классов с финализаторами (кроме классов библиотеки платформы) [EJ Item 6].

Избегайте использования финализаторов.

7.5 Итоги

Вопросы, касающиеся завершения жизненного цикла (*end-of-lifecycle*) задач, потоков, служб и приложений могут усложнить их разработку и реализацию. Язык Java не предоставляет упреждающего механизма для отмены выполняющихся активностей или завершения потоков. Вместо этого Java предоставляет механизм кооперативного прерывания, который можно использовать для облегчения процесса отмены, но создание протоколов для отмены и их согласованное использование зависит только от вас. Использование класса `FutureTask` и фреймворка `Executor` упрощает построение отменяемых задач и сервисов.

⁹⁴ См. (Boehm, 2005) о некоторых проблемах, связанных с написанием финализаторов.

Глава 8 Применение пулов потоков

В главе 6 был представлен фреймворк выполнения задач, упрощающий управление жизненным циклом задач и потоков и предоставляющий простые и гибкие средства для отделения отправки задач от политики выполнения. В главе 7 рассматривались некоторые запутанные детали жизненного цикла службы, возникающие при использовании фреймворка выполнения задач в реальных приложениях. В этой главе рассматриваются дополнительные параметры для конфигурирования и настройки пулов потоков, описываются опасности, которые следует учитывать при использовании фреймворка выполнения задач, и предлагаются некоторые более продвинутые примеры использования фреймворка Executor.

8.1 Неявные связи между задачами и политиками выполнения

Ранее мы утверждали, что фреймворк выполнения задач отделяет отправку задачи от ее выполнения. Подобно многим прочим попыткам развязать комплексные процессы, это было небольшим преувеличением. Хотя фреймворк Executor и предоставляет значительную гибкость в определении и изменении политик выполнения, не все задачи совместимы со всеми политиками выполнения. Типы задач, которым требуются определенные политики выполнения, включают в себя:

Зависимые задачи. Наилучшим образом проявляют себя независимые задачи: те, которые не зависят от момента времени, результатов выполнения или побочных эффектов, возникающих при выполнении других задач. Когда в пуле потоков выполняются независимые задачи, можно свободно изменять размер пула и его конфигурацию, не влияя ни на что, кроме производительности. С другой стороны, когда вы отправляете задачи, зависящие от других задач, в пул потоков, вы неявно накладываете ограничения на политику выполнения, которыми необходимо тщательно управлять, чтобы избежать проблем с живучестью (см. раздел [8.1.1](#)).

Задачи, использующие ограничение потока. Однопоточные исполнители (*executors*) берут на себя более строгие обязательства, касающиеся обеспечения параллелизма, чем произвольные пулы потоков. Они гарантируют, что задачи не выполняются параллельно, что позволяет вам несколько ослабить обеспечение потокобезопасности в коде задачи. Объекты могут быть ограничены потоком задачи, что позволяет задачам, спроектированным для выполнения в этом потоке, получать доступ к этим объектам без синхронизации, даже если эти ресурсы не являются потокобезопасными. Это приводит к формированию неявной связи между задачей и её политикой выполнения - задачи определяют требование, чтобы их исполнитель был однопоточным⁹⁵. В этом случае, при изменении однопоточной реализации интерфейса Executor на пул потоков, свойство потокобезопасности может быть потеряно.

⁹⁵ Требование не такое уж сильное; было бы достаточно гарантии того, что задачи не выполняются параллельно и обеспечивается достаточный уровень синхронизации, так, чтобы воздействие, оказанное на память одной задачей, было гарантированно видно следующей задаче – именно такие гарантии и предлагает фабричный метод Executors.newSingleThreadExecutor.

Задачи, чувствительные ко времени ответа. Приложения GUI чувствительны к времени отклика: как правило, пользователей раздражает длительная задержка между нажатием кнопки и соответствующим визуальным откликом. Отправка долговременной задачи однопоточному исполнителю или отправка нескольких долговременных задач в пул потоков небольшого размера⁹⁶, может ухудшить отзывчивость службы, управляемой этой реализацией Executor.

Задачи, использующие класс ThreadLocal. Класс ThreadLocal позволяет каждому потоку иметь свою собственную "версию" переменной. Однако исполнители могут свободно использовать потоки по своему усмотрению. Стандартные реализации интерфейса Executor могут уничтожать неиспользуемые потоки, когда спрос низкий и добавлять новые, когда спрос высокий, а также заменять рабочие потоки свежими, если задачей было брошено непроверяемое исключение. Класс ThreadLocal имеет смысл использовать с пулом потоков, только если локальное по отношению к потоку значение имеет жизненный цикл, ограниченный этой задачей; класс ThreadLocal не должен использоваться в пуле потоков для передачи значений между задачами.

Пулы потоков лучше всего работают, когда задачи *гомогенны и независимы*. Смешивание длительных и кратковременных задач может привести к "засорению" пула, если он не сильно велик; отправка задач, зависящих от других задач, может привести к взаимоблокировке, если пул не является неограниченным. К счастью, запросы в типичных сетевых серверных приложениях - веб-серверах, почтовых серверах, файловых серверах - обычно соответствуют этим рекомендациям.

Некоторые задачи имеют характеристики, требующие или исключающие определенную политику выполнения. Задачи, зависящие от других задач, требуют, чтобы пул потоков был достаточно большим, чтобы задачи никогда не ставились в очередь или не отклонялись; задачи, использующие ограничение потока, требуют последовательного выполнения.

Документируйте эти требования, чтобы будущие сопровождающие не подрывали безопасность или живучесть, устанавливая несовместимую политику выполнения.

8.1.1 Взаимоблокировка потоков, вызванная голоданием

Если задачи, зависящие от других задач, выполняются в пуле потоков, они могут попадать в состояние взаимоблокировки. В случае использования однопоточного исполнителя, задача, отправляющая другую задачу некоторому исполнителю и ожидающая результата её выполнения, всегда будет попадать в состояние взаимоблокировки. Вторая задача будет находиться в рабочей очереди до завершения первой задачи, но первая не будет завершена, так как она ожидает результата второй задачи. То же самое может произойти и в больших пулах потоков, если все потоки выполняют задачи, которые заблокированы в ожидании других задач, все еще находящихся в рабочей очереди. Это называется

⁹⁶ Пул поток с небольшим количеством рабочих потоков, выполняющих задачи.

взаимоблокировкой потоков, вызванной голоданием (*thread starvation deadlock*), и может возникать всякий раз, когда пул задач инициирует неограниченное ожидание блокировки какого-нибудь ресурса или условия, которое может быть выполнено лишь на основе действия другой задачи из пула потоков, такого как ожидание возвращаемого значения или побочного эффекта, вызванного другой задачей, если вы не можете гарантировать, что пул достаточно большой.

Класс `ThreadDeadlock` из листинга 8.1, иллюстрирует ситуацию голодания потока. Метод `RenderPageTask` отправляет две дополнительные задачи исполнителю для извлечения верхнего и нижнего колонтитулов страницы, отображает тело страницы, ожидает результатов задач получения верхнего и нижнего колонтитулов, а затем объединяет верхний и нижний колонтитулы в готовую страницу. С однопоточным исполнителем в классе `ThreadDeadlock` всегда будет возникать взаимоблокировка. Аналогично, задачи осуществляющие координацию друг с другом с помощью барьера, также могут приходить к состоянию взаимоблокировки связанной с голоданием потока, если пул потоков недостаточно велик.

```
public class ThreadDeadlock {  
    ExecutorService exec = Executors.newSingleThreadExecutor();  
  
    public class RenderPageTask implements Callable<String> {  
        public String call() throws Exception {  
            Future<String> header, footer;  
            header = exec.submit(new LoadFileTask("header.html"));  
            footer = exec.submit(new LoadFileTask("footer.html"));  
            String page = renderBody();  
            // Will deadlock -- task waiting for result of subtask  
            return header.get() + page + footer.get();  
        }  
    }  
}
```



Листинг 8.1 Задача, попадающая в состояние взаимоблокировки в однопоточном экземпляре `Executor`. Не делайте так.

Всякий раз, когда вы отправляете исполнителю задачи, которые не являются независимыми, знайте о возможности возникновения взаимоблокировки, вызванной голоданием потока, и документируйте любые ограничения, налагаемые на размер пула или конфигурацию в коде или файле конфигурации, где конфигурируется экземпляр `Executor`.

Помимо явных ограничений, накладываемых на размер пула потоков, также могут быть и неявные ограничения, возникающие из-за ограничений на другие ресурсы. Если приложение использует пул соединений JDBC с десятью соединениями, и каждой задаче требуется подключение к базе данных, то пул потоков будет иметь только десять потоков, поскольку задачи, количество которых превышает десять, будут блокироваться в ожидании получения соединения.

8.1.2 Долговременные задачи

У пулов потоков могут возникать проблемы с отзывчивостью, если задачи могут блокироваться на длительные периоды времени, даже если взаимоблокировка невозможна. Пул потоков может “засоряться” длительными задачами, что приводит к увеличению времени обслуживания даже для коротких задач. Если размер пула слишком мал по сравнению с ожидаемым устойчивым количеством долговременных задач, в конечном итоге все потоки в пуле будут выполнять долговременные задачи, и отзывчивость ухудшится.

Одним из способов, позволяющих смягчить неблагоприятные последствия от длительного выполнения задач заключается в том, чтобы задачи использовали ограниченное по времени ожидание ресурсов, взамен неограниченного по времени. Большинство блокирующих методов в библиотеках платформы, таких как `Thread.join`, `BlockingQueue.put`, `CountDownLatch.await`, и `Selector.select`, поставляются в обеих, ограниченной и неограниченной по времени, версиях, если время ожидания истекло, можно пометить задачу как завершившуюся с ошибкой и прервать ее или заново поместить в очередь для последующего выполнения. Такой подход гарантирует, что каждая задача в конечном итоге будет продвигаться к успешному или неудачному завершению, освобождая потоки для задач, которые смогут завершиться быстрее. Если часто возникает такая ситуация, что пул потоков заполнен заблокированными задачами, это может быть свидетельством того, что размер пула слишком мал.

8.2 Размеры пулов потоков

Идеальный размер пула потоков зависит от типов отправляемых задач и характеристик системы развертывания. Размеры пула потоков редко когда должны быть жестко заданы; вместо этого размеры пула должны предоставляться механизмом конфигурации или вычисляться динамически, с помощью метода `Runtime.availableProcessors`.

Определение размеров пулов потоков не является точной наукой, но, к счастью, вам только и нужно, что избегать крайностей “слишком большой” и “слишком маленький”. Если пул потоков слишком велик, потоки конкурируют за ограниченные ресурсы ЦП и памяти, что приводит к увеличению использования памяти и возможному исчерпанию ресурсов. Если размер пула слишком мал, пропускная способность падает, поскольку процессорные мощности не используются в полном объёме, несмотря на доступную работу.

Чтобы правильно определить размер пула потоков, необходимо понимать вычислительную среду, бюджет ресурсов и характер задач. Сколько процессоров в системе развертывания? Сколько памяти? Выполняют ли задачи в основном вычисления, операции ввода/вывода или какую-то их комбинацию? Они требуют дефицитного ресурса, такого как соединения JDBC? Если у вас в наличии различные категории задач с очень разным поведением, рекомендуется использовать несколько пулов потоков, чтобы каждый из них можно было настроить в соответствии с определённым типом рабочей нагрузки.

Для ресурсоемких задач N_{cpu} -процессорная система обычно достигает оптимального использования ресурсов с размером пула $N_{cpu} + 1$ потоков. (Даже ресурсоемкие потоки иногда натыкаются на ошибку выделения страницы памяти или встают на паузу по какой-либо другой причине, поэтому “дополнительный”

запускаемый поток предотвращает простаивание циклов CPU в такой ситуации.) Для задач, которые также включают операции ввода/вывода или другие блокирующие операции, требуется больший пул, поскольку не все потоки будут доступны для планирования. Чтобы правильно определить размер пула, необходимо оценить отношение времени ожидания ко времени вычисления задач; эта оценка не обязательно должна быть точной и может быть получена с помощью профилирования или с помощью других измерительных средств. Кроме того, размер пула потоков можно настроить, запустив приложение с использованием нескольких различных размеров пула потоков при тестовой нагрузке и наблюдая за уровнем загрузки ЦП.

Примем такие определения:

$$\begin{aligned} N_{cpu} &= \text{количество } CPU \\ U_{cpu} &= \text{целевая загрузка } CPU, 0 \leq U_{cpu} \leq 1 \\ \frac{W}{C} &= \text{отношение времени ожидания к времени вычисления} \end{aligned}$$

Оптимальный размер пула для поддержания желаемого уровня использования процессора:

$$N_{threads} = N_{cpu} * U_{cpu} * \left(1 + \frac{W}{C} \right)$$

Количество процессоров можно определить с помощью класса `Runtime`:

```
int N_CPUS = Runtime.getRuntime().availableProcessors();
```

Конечно, циклы CPU не единственный ресурс, которым можно управлять с помощью пулов потоков. Другие ресурсы, такие как размер памяти, дескрипторы файлов, дескрипторы сокетов, подключения к базе данных, также могут способствовать в определении размеров ограничений. Вычисление ограничений размера пула для этих типов ресурсов проще: просто сложите количество единиц этого ресурса, требуемого каждой задаче, и разделите его на общее доступное количество единиц ресурса. В результате будет получена верхняя граница размера пула потоков.

Когда задачам требуется ресурс, находящийся в пуле ресурсов, например, такой как соединение с базой данных, размер пула потоков и размер пула ресурсов оказывают друг на друга взаимное влияние. Если для каждой задачи требуется соединение, эффективный размер пула потоков ограничивается размером пула соединений. Аналогично, если единственными потребителями соединений являются задачи пула потоков, эффективный размер пула соединений ограничен размером пула потоков.

8.3 Конфигурирование класса `ThreadPoolExecutor`

Класс `ThreadPoolExecutor` предоставляет базовую реализацию для исполнителей, возвращаемых фабричными методами `newCachedThreadPool`, `newFixedThreadPool`, и `newScheduledThreadPool` класса `Executors`. Класс `ThreadPoolExecutor` представляет собой гибкую и надежную реализацию пула, которая позволяет выполнять различные настройки.

Если политика выполнения по умолчанию не соответствует вашим потребностям, вы можете создать экземпляр `ThreadPoolExecutor` с помощью его конструктора и настроить его по своему усмотрению; вы можете обратиться к

исходному коду класса `Executors`, чтобы увидеть политики выполнения для конфигураций по умолчанию и использовать их в качестве отправной точки. Класс `ThreadPoolExecutor` имеет несколько конструкторов, наиболее общий из которых показан в листинге 8.2.

```
public ThreadPoolExecutor(int corePoolSize,
                          int maximumPoolSize,
                          long keepAliveTime,
                          TimeUnit unit,
                          BlockingQueue<Runnable> workQueue,
                          ThreadFactory threadFactory,
                          RejectedExecutionHandler handler) { ... }
```

Листинг 8.2 Общий конструктор класса `ThreadPoolExecutor`.

8.3.1 Создание и удаление потока

Корневой размер пула, максимальный размер пула и время жизни определяют создание и удаление потоков. Корневой размер является целевым размером; реализация пытается поддерживать пул в пределах этого размера, даже если нет никаких задач для выполнения⁹⁷, и не будет создавать потоков больше, чем заданное количество, пока рабочая очередь не заполнится⁹⁸. Максимальный размер пула представляет собой верхнюю границу количества одновременно активных потоков в пуле. Поток, который простоявает дольше установленного времени ожидания активности (*keep-alive times*), становится кандидатом на поглощение и может быть завершен, если текущий размер пула превысит корневой размер.

Настраивая размер корневого пула и время ожидания активности, можно поощрять пул к освобождению ресурсов, используемых простояющими потоками, что делает их доступными для более полезной работы. (Подобно прочему, это компромисс: поглощение бездействующих потоков приводит к дополнительной задержке в связи с необходимостью создания потоков, когда позже потоки должны будут быть созданы при увеличении спроса.)

Фабричный метод `newFixedThreadPool` устанавливает и корневой, и максимальный размер пула до требуемых значений, при этом определяя бесконечное время ожидания; фабричный метод `newCachedThreadPool` устанавливает максимальный размер пула в значение `Integer.MAXVALUE` и корневой размер пула равным нулю, с таймаутом в одну минуту, создавая эффект бесконечно расширяемого пула потоков, который будет сжиматься снова при

⁹⁷ Когда класс `ThreadPoolExecutor` изначально создается, корневые потоки запускаются не сразу, а по мере отправки задач, если только вы не вызовите метод `prestartAllCoreThreads`.

⁹⁸ Иногда у разработчиков возникает соблазн установить корневой размер пула равным нулю, чтобы рабочие потоки в конечном итоге поглощались, и, как следствие, не препятствовали завершению работы JVM, но это может вызвать странное поведение в пулах потоков, которые не используют класс `SynchronousQueue` для своей рабочей очереди (например, как это делает `newCachedThreadPool`). Если пулу уже установлен корневой размер, класс `ThreadPoolExecutor` создаст новый поток только тогда, когда рабочая очередь заполнится. Таким образом, задачи, отправленные в пул потоков с рабочей очередью, имеющей любую емкость и нулевой размер ядра, не будут выполняться до тех пор, пока очередь не заполнится, что обычно является нежелательным. Начиная с Java 6, метод `allowCoreThreadTimeOut` позволяет отправить всем потокам в пуле указание на завершение выполнения по тайм-ауту; включите эту функцию с корневым размером пула равным нулю, если вам требуется ограниченный пул потоков с ограниченной рабочей очередью, но при этом все потоки будут поглощаться, когда для них не будет работы.

уменьшении спроса. Другие комбинации можно скомпоновать с помощью явного использования конструктора `ThreadPoolExecutor`.

8.3.2 Управление задачами в очереди

Ограниченные пулы потоков устанавливают предельное количество одновременно выполняемых задач. (Однопоточные исполнители являются примечательным частным случаем: они гарантируют, что никакие задачи не будут выполняться одновременно, предлагая возможность достижения потокобезопасности через ограничение потока.)

Ранее, в разделе 6.1.2, мы видели, что создание неограниченного количества потоков может привести к нестабильности, и решили эту проблему, используя пул потоков фиксированного размера вместо создания нового потока для каждого поступающего запроса. Тем не менее, это лишь частичное решение проблемы; приложение все еще может исчерпать ресурсы при большой нагрузке, только сделать это будет сложнее. Если скорость поступления новых запросов превышает скорость, с которой они могут быть обработаны, запросы все равно будут помещаться в очередь. В случае использования пула потоков, запросы ожидают в очереди экземпляров `Runnable`, управляемой экземпляром `Executor`, вместо того, чтобы ожидать в очереди в качестве потоков, конкурирующих за ЦП. Представление ожидающей задачи с помощью экземпляра `Runnable` и узла списка, безусловно, обходится намного дешевле, чем представление с помощью потока, но риск исчерпания ресурсов по-прежнему остается, так как клиенты могут отправлять запросы на сервер быстрее, чем он сможет их обрабатывать.

Запросы часто поступают пакетами, даже если в среднем скорость поступления запросов достаточно стабильна. Очереди могут помочь в сглаживании временных переходов между всплесками задач, но если задачи продолжат поступать слишком быстро, вам в конечном итоге придется регулировать скорость прибытия, чтобы избежать нехватки памяти⁹⁹. Ещё до того, как закончится доступная память, время ответа будет постепенно увеличиваться, по мере роста количества задач в очереди.

Класс `ThreadPoolExecutor` позволяет вам использовать экземпляр `BlockingQueue` для содержания задач, ожидающих выполнения. Существует три основных подхода к организации очереди задач: неограниченная очередь, ограниченная очередь и синхронная передача. Выбор типа очереди оказывает влияние на другие параметрами конфигурации, например, на размер пула.

По умолчанию, фабричные методы `newFixedThreadPool` и `newSingleThreadExecutor` используют неограниченную реализацию очереди `LinkedBlockingQueue`. Задачи будут помещаться в очередь, если все рабочие потоки будут заняты, в случае, если задачи продолжат поступать быстрее, чем они могут быть выполнены, очередь будет расти без ограничений,

Более стабильной стратегией управления ресурсами является использование ограниченной очереди, например класса `ArrayBlockingQueue` или ограниченной версии `LinkedBlockingQueue` или `PriorityBlockingQueue`. Ограниченные очереди помогают предотвратить исчерпание ресурсов, но поднимают вопрос о том, что делать с новыми задачами при заполнении очереди. (Существует ряд возможных *политик насыщения* (*saturation policies*) для решения этой проблемы;

⁹⁹ Можно провести аналогию с управлением потоком в сетях связи: вы можете захотеть буферизировать определенный объем данных, но в конечном итоге вы будете вынуждены найти способ заставить другую сторону прекратить отправку данных или начать отбрасывать лишние данные и надеяться, что отправитель позже повторно отправит их вновь, когда вы будете свободнее.

см. раздел 8.3.3.) С ограниченной рабочей очередью размер очереди и размер пула должны настраиваться вместе. Большая очередь в сочетании с пулом небольшого размера может снизить использование памяти, ЦП и переключение контекста за счет потенциального ограничения пропускной способности.

Для очень больших или неограниченных пулов потоков также можно полностью обойтись без помещения задач в очередь и передавать задачи непосредственно от производителей рабочим потокам, с помощью класса `SynchronousQueue`. Класс `SynchronousQueue`, на самом деле, является не очередью, а механизмом управления передачей между потоками. Следуя порядку размещения элементов в экземпляре `SynchronousQueue`, должен существовать поток, ожидающий принятия передачи. Если нет ожидающих потоков, но текущий размер пула меньше, чем максимально заданный, экземпляр `ThreadPoolExecutor` создает новый поток; иначе задача отбрасывается согласно политике насыщения. Использование прямой передачи является более эффективным механизмом взаимодействия, поскольку задача может быть напрямую передана потоку, который ее выполнит, а не помещена в очередь, из которой её извлечёт рабочий поток. Класс `SynchronousQueue` является оправданным выбором только в том случае, если пул потоков является неограниченным или если отклонение избыточных задач приемлемо. Фабричный метод `newCachedThreadPool` использует класс `SynchronousQueue`.

Использование очередей FIFO, например, таких как `LinkedListQueue` или `ArrayBlockingQueue`, приводит к запуску задач в том порядке, в котором они были приняты. Для большего контроля над порядком выполнения задач можно использовать класс `PriorityBlockingQueue`, который упорядочивает задачи в соответствии с заданным приоритетом. Приоритет может быть определен естественным порядком (если задачи реализуют интерфейс `Comparable`) или с помощью интерфейса `Comparator`.

Фабричный метод `newCachedThreadPool` является хорошим выбором по умолчанию для реализации `Executor`, обеспечивая лучшую производительность очереди, чем фиксированный пул потоков¹⁰⁰. Пул потоков фиксированного размера является хорошим выбором, когда необходимо ограничить число параллельных задач в целях управления ресурсами, как в случае серверного приложения, которое принимает запросы от клиентов по сети, и в ином случае было бы уязвимо для перегрузки.

Решение с ограничением пула потоков или рабочей очереди подходит только в том случае, если задачи независимы. С задачами, которые зависят от других задач, ограниченные пулы потоков или очереди могут привести к взаимоблокировке потоков вызванной голоданием; вместо этого используйте конфигурацию неограниченного пула, например фабричный метод `newCachedThreadPool`¹⁰¹.

¹⁰⁰ Это различие в производительности происходит от использования класса `SynchronousQueue` вместо класса `LinkedListQueue`. В Java 6 реализация класса `SynchronousQueue` была заменена новым неблокирующими алгоритмом, который улучшил пропускную способность экземпляра `Executor` в тестах на факторизацию более чем в три раза, по сравнению с реализацией класса `SynchronousQueue` из Java 5.0 (Scherer et al., 2006).

¹⁰¹ Альтернативной конфигурацией для задач, отправляющих другие задачи и ожидающих результатов их выполнения, является использование ограниченного пула потоков, класса `SynchronousQueue` в качестве рабочей очереди и политики насыщения вызывающего объекта.

8.3.3 Политики насыщения

Когда ограниченная рабочая очередь заполняется, вступает в игру *политика насыщения* (*saturation policy*). Политика насыщения для класса `ThreadPoolExecutor` может быть изменена путем вызова метода `setRejectedExecutionHandler`. (Политика насыщения также используется, когда задача передается экземпляру `Executor`, который был выключен.) Предоставляется несколько реализаций интерфейса `RejectedExecutionHandler`, реализующих ту или другую политику насыщения: `AbortPolicy`, `CallerRunsPolicy`, `DiscardPolicy`, и `DiscardOldestPolicy`.

Политика по умолчанию, прерывание (*abort*), заставляет метод `execute` бросить непроверяемое исключение `RejectedExecutionException`;зывающий объект может перехватить это исключение и реализовать собственную обработку переполнения пула по своему усмотрению. Политика *сброса* (*discard*) автоматически отбрасывает только что отправленную задачу, если она не может быть поставлена в очередь на выполнение; политика сброса старейшего элемента (*discard-oldest*) сбрасывает задачу, которая в противном случае была бы выполнена следующей, и пытается повторно отправить новую задачу. (Если рабочая очередь представляет собой очередь с приоритетом, это приводит к сбросу элемента с наивысшим приоритетом, таким образом, комбинация политики насыщения путём сброса старейшего элемента и очереди с приоритетом не является хорошей идеей.)

Политика *выполнение вызывающего объекта* (*caller-runs*) реализует форму регулирования, при которой задачи не сбрасываются, исключения не бросаются, а наоборот прилагаются все усилия по замедлению поступления потока новых задач, путём возврата некоторой части работы обратно вызывающему объекту. Только что отправленная задача выполняется не в потоке пула, а в потоке, вызвавшем метод `execute`. Если мы внесём изменения в класс `WebServer`, представленный в качестве примера, чтобы использовать ограниченную очередь и политику “выполнение вызывающего объекта”, после того, как все пулы потоков будут заполнены, и рабочая очередь также будет заполнена, следующая задача будет выполняться в главном потоке, в процессе вызова на выполнение. Поскольку обработка, вероятно, займет некоторое время, основной поток не сможет отправлять задачи, по крайней мере, некоторое время, давая рабочим потокам какое-то время на то, чтобы сократить разрыв. Основной поток также не будет вызывать метод `accept` в течение этого времени, поэтому входящие запросы будут попадать в очередь на уровне TCP, а не на уровне приложения. Если перегрузка будет сохраняться, в конечном счёте, уровень TCP примет решение о том, что он поставил в очередь достаточное количество запросов на подключение и начнёт отбрасывать запросы на подключение. По мере того как сервер будет перегружаться, перегрузка будет постепенно выталкиваться наружу – из пула потоков в рабочую очередь приложения, далее на уровень TCP и, в конечном счете, клиенту, обеспечивая таким образом более плавное снижение производительности под нагрузкой.

Выбор политики насыщения или внесение других изменений в политику выполнения может быть выполнен в процессе создания экземпляра `Executor`. В листинге 8.3 показано создание пула потоков фиксированного размера с помощью политики “выполнение вызывающего объекта”.

`ThreadPoolExecutor executor`

```
= new ThreadPoolExecutor(N_THREADS, N_THREADS,
    0L, TimeUnit.MILLISECONDS,
    new LinkedBlockingQueue<Runnable>(CAPACITY));
executor.setRejectedExecutionHandler(
    new ThreadPoolExecutor.CallerRunsPolicy());
```

Листинг 8.3 Создание пула потоков фиксированного размера с ограниченной очередью и политикой насыщения “Выполнение вызывающего объекта”

Не существует предопределённой политики насыщения, выполняющей блокировку метода `execute`, когда очередь заполняется. Однако, похожий эффект может быть достигнут с помощью семафора, для ограничения частоты внедрения задач, как показано в классе `BoundedExecutor` из листинга 8.4. При таком подходе используйте неограниченную очередь (нет причин ограничивать и размер очереди, и частоту внедрения) и установите границу на семафоре равную размеру пула *плюс* количеству поставленных в очередь задач, которые вы хотите обработать, так как семафор ограничивает количество задач, как выполняемых в данный момент, так и ожидающих выполнения.

```
@ThreadSafe
public class BoundedExecutor { private
    final Executor exec; private final
    Semaphore semaphore;

    public BoundedExecutor(Executor exec, int bound) {
        this.exec = exec;
        this.semaphore = new Semaphore(bound);
    }

    public void submitTask(final Runnable command)
        throws InterruptedException {
        semaphore.acquire();
        try {
            exec.execute(new Runnable() {
                public void run() {
                    try {
                        command.run();
                    } finally {
                        semaphore.release();
                    }
                }
            });
        } catch (RejectedExecutionException e) {
            semaphore.release();
        }
    }
}
```

Листинг 8.4 Использование класса `Semaphore` для управления частотой отправки задач

8.3.4 Фабрики задач

Всякий раз, когда пул потоков должен создать поток, он делает это с помощью *фабрики потоков* (*thread factory*, см. листинг 8.5). Фабрика потоков по умолчанию создает новый поток, не демон, без специальной конфигурации. Определяемые фабрики потоков позволяет вам настроить конфигурацию пула потоков. Интерфейс *ThreadFactory* имеет один метод, *newThread*, который вызывается всякий раз, когда пулу потоков необходимо создать новый поток.

```
public interface ThreadFactory {  
    Thread newThread(Runnable r);  
}
```

Листинг 8.5 Интерфейс *ThreadFactory*

Существует ряд причин для использования настраиваемой версии (*custom*) фабрики потоков. Вы можете захотеть определить для пула потоков исключение *UncaughtExceptionHandler* или создать экземпляр настраиваемой версии класса *Thread*, например такого, который выполняет ведение журнала отладки. Вы можете захотеть изменить приоритет (обычно не очень хорошая идея; см. раздел [10.3.1](#)) или установить статус демона (снова, не всё в этой идеи хорошо; см. раздел [7.4.2](#)) потокам в пуле. Или, может быть, вы просто хотите дать потокам пула более осмысленные имена, чтобы упростить интерпретацию дампов потоков и логов с ошибками.

В классе *MyThreadFactory* из листинга 8.6, иллюстрируется пример настраиваемой фабрики потоков. Фабрика создаёт новый экземпляр *MyAppThread*, передавая конструктору специфичное для экземпляра пула потоков имя, так, чтобы потоки от каждого пула могли быть легко различимы в дампах потоков и журналах ошибок. Класс *MyAppThread* также можно использовать и в других частях приложения, так что все потоки могут воспользоваться предоставляемыми возможностями для облегчения отладки.

```
public class MyThreadFactory implements ThreadFactory {  
    private final String poolName;  
  
    public MyThreadFactory(String poolName) {  
        this.poolName = poolName;  
    }  
  
    public Thread newThread(Runnable runnable) {  
        return new MyAppThread(runnable, poolName);  
    }  
}
```

Листинг 8.6 Настраиваемая реализация фабрики потоков

В классе *MyAppThread* показанном в листинге 8.7, выполняется интересная настройка, которая позволяет вам указать имя потока, устанавливает настраиваемое исключение *UncaughtExceptionHandler*, записывающее сообщение в экземпляр класса *Logger*, ведет статистику о том, сколько потоков было создано и уничтожено, и, дополнительно, записывает отладочное сообщение в лог, когда поток создается или завершается.

```
public class MyAppThread extends Thread {
    public static final String DEFAULT_NAME = "MyAppThread";
    private static volatile boolean debugLifecycle = false;
    private static final AtomicInteger created = new AtomicInteger();
    private static final AtomicInteger alive = new AtomicInteger();
    private static final Logger log = Logger.getAnonymousLogger();

    public MyAppThread(Runnable r) { this(r, DEFAULT_NAME); }

    public MyAppThread(Runnable runnable, String name) {
        super(runnable, name + "-" + created.incrementAndGet());
        setUncaughtExceptionHandler(
            new Thread.UncaughtExceptionHandler() {
                public void uncaughtException(Thread t,
                                              Throwable e) {
                    log.log(Level.SEVERE,
                            "UNCAUGHT in thread " + t.getName(), e);
                }
            });
    }

    public void run() {
        // Copy debug flag to ensure consistent value throughout.
        boolean debug = debugLifecycle;
        if (debug) log.log(Level.FINE, "Created "+getName());
        try {
            alive.incrementAndGet();
            super.run();
        } finally {
            alive.decrementAndGet();
            if (debug) log.log(Level.FINE, "Exiting "+getName());
        }
    }

    public static int getThreadsCreated() { return created.get(); }
    public static int getThreadsAlive() { return alive.get(); }
    public static boolean getDebug() { return debugLifecycle; }
    public static void setDebug(boolean b) { debugLifecycle = b; }
}
```

Листинг 8.7 Настраиваемый класс, построенный на основе потока

Если приложение использует преимущества политик безопасности в целях предоставления разрешений определенным кодовым базам, для создания фабрики потоков вы можете использовать фабричный метод `privilegedThreadFactory` класса `Executors`. Метод создает пул потоков, имеющий те же разрешения, экземпляры `AccessControlContext` и `contextClassLoader` с теми же правами, что и поток, создаваемый фабричным методом `privilegedThreadFactory`. В противном случае потоки, созданные пулом потоков, наследуют разрешения от любого клиента вызвавшего методы `execute` или `submit` во время создания

нового потока, что может привести к путанице в исключениях, связанных с безопасностью.

8.3.5 Настройка экземпляра ThreadPoolExecutor после построения

Большинство параметров, передаваемых конструкторам ThreadPoolExecutor, также можно изменить после построения с помощью сеттеров (например, корневой размер пула потоков, максимальный размер пула потоков, время ожидания, фабрику потоков и обработчик выполнения отклонённых запросов). Если экземпляр Executor создается одним из фабричных методов класса Executors (исключая фабричный метод newSingleThreadExecutor), вы можете привести результат к классу ThreadPoolExecutor, для доступа к сеттерам, как показано в листинге 8.8.

```
ExecutorService exec = Executors.newCachedThreadPool();
if (exec instanceof ThreadPoolExecutor)
    ((ThreadPoolExecutor) exec).setCorePoolSize(10);
else
    throw new AssertionError("Oops, bad assumption");
```

Листинг 8.8 Изменение экземпляра Executor, созданного с помощью стандартной фабрики

Класс Executors включает в себя фабричный метод unconfigurableExecutorService, который принимает существующий экземпляр ExecutorService и обертывает его, раскрывая только методы интерфейса ExecutorService, без возможности настройки в дальнейшем. В отличие от реализаций с пулами потоков, фабричный метод newSingleThreadExecutor возвращает экземпляр ExecutorService, обернутый подобным образом, а не “сырой” (raw) экземпляр ThreadPoolExecutor. В то время как однопоточный исполнитель фактически реализован как пул потоков с одним потоком, он также берёт на себя обязательство не выполнять задачи параллельно. Если какой-то ошибочный код увеличит размер пула в однопоточном исполнителе, это нарушит предполагаемую семантику выполнения.

Вы можете использовать этот подход с собственными исполнителями, для предотвращения внесения изменений в политику выполнения. Если вы будете предоставлять экземпляр ExecutorService коду, которому не доверяете, без внесения изменений в него, вы можете обернуть экземпляр с помощью метода unconfigurableExecutorService.

8.4 Расширение класса ThreadPoolExecutor

Класс ThreadPoolExecutor был разработан для расширения, предоставляя несколько “хуков” для переопределения в подклассах - beforeExecute, afterExecute, и terminated – что может быть использовано для расширения поведения класса ThreadPoolExecutor.

Хуки beforeExecute и afterExecute вызываются в потоке, который выполняет задачу, и могут использоваться для добавления логирования, отчёта времени выполнения, мониторинга или сбора статистики. Хук afterExecute вызывается независимо от того, завершается ли задача обычным образом, путём возврата управления из метода run или путём возбуждения исключения Exception. (Если задача завершается путём возбуждения исключения Error, хук

`afterExecute` не будет вызван.) Если хук `beforeExecute` бросит исключение `RuntimeException`, задача не будет выполняться, и хук `afterExecute` также не будет вызван.

Хук `terminated` вызывается, когда выполнен процесс завершения работы пула потоков, после завершения всех задач и завершения выполнения всех рабочих потоков. Он может быть использован для высвобождения ресурсов, выделенных исполнителем в течение его жизненного цикла, отправки уведомлений или логирования, или завершения сбора статистики.

8.4.1 Пример: добавление статистики в пул потоков

В классе `TimingThreadPool` из листинга 8.9, показан настраиваемый пул потоков, который использует хуки `beforeExecute`, `afterExecute`, и `terminated` для добавления возможностей логирования и сбора статистики. Чтобы измерить время выполнения задачи, хук `beforeExecute` должен записать время начала и сохранить его в таком месте, в котором хук `afterExecute` сможет найти его. Поскольку выполняемые хуки вызываются в том же потоке, который выполняет задачу, значение, помещенное в экземпляр `ThreadLocal` хуком `beforeExecute`, может быть получено хуком `afterExecute`. Класс `TimingThreadPool` использует пару переменных типа `AtomicLong`, чтобы отслеживать общее количество обработанных задач и общее время количества времени, потраченное на обработку, и использует хук `terminated`, чтобы вывести сообщение лога, отражающее среднее время выполнения задачи.

```
public class TimingThreadPool extends ThreadPoolExecutor {
    private final ThreadLocal<Long> startTime
        = new ThreadLocal<Long>();
    private final Logger log = Logger.getLogger("TimingThreadPool");
    private final AtomicLong numTasks = new AtomicLong();
    private final AtomicLong totalTime = new AtomicLong();

    protected void beforeExecute(Thread t, Runnable r) {
        super.beforeExecute(t, r);
        log.fine(String.format("Thread %s: start %s", t, r));
        startTime.set(System.nanoTime());
    }

    protected void afterExecute(Runnable r, Throwable t) {
        try {
            long endTime = System.nanoTime();
            long taskTime = endTime - startTime.get();
            numTasks.incrementAndGet();
            totalTime.addAndGet(taskTime);
            log.fine(String.format("Thread %s: end %s, time=%dns",
                t, r, taskTime));
        } finally {
            super.afterExecute(r, t);
        }
    }
}
```

```
protected void terminated() {
    try {
        log.info(String.format("Terminated: avg time=%dns",
                               totalTime.get() / numTasks.get()));
    } finally {
        super.terminated();
    }
}
```

Листинг 8.9 Пул потоков, расширенный механизмами логирования и учёта времени выполнения задач

8.5 Распараллеливание рекурсивных алгоритмов

Пример рендеринга страницы из раздела 6.3, подвергся серии доработок, в поисках подходящего для использования механизма параллелизма. Первая попытка была полностью последовательной; вторая использовала два потока, но все же выполняла все операции загрузки изображений последовательно; окончательная версия, для достижения большей степени параллелизма, рассматривала каждую операцию загрузки изображения как отдельную задачу. Циклы, тела которых содержат нетривиальные вычисления или выполняют потенциально подверженные блокировке операции ввода/вывода, часто являются хорошими кандидатами для распараллеливания, если итерации независимы.

Если у нас есть цикл, итерации которого независимы, и для продолжения нам не нужно ожидать завершения выполняемых в них операций, мы можем использовать экземпляр `Executor` для преобразования последовательного цикла в параллельный, как показано в методах `processSequentially` и `processInParallel` из листинга 8.10.

```
void processSequentially(List<Element> elements) {
    for (Element e : elements)
        process(e);
}

void processInParallel(Executor exec, List<Element> elements) {
    for (final Element e : elements)
        exec.execute(new Runnable() {
            public void run() { process(e); }
        });
}
```

Листинг 8.10. Преобразование последовательного выполнения в параллельное

Вызов метода `processInParallel` завершается быстрее, чем вызов метода `processSequentially`, потому что он возвращает управление, как только все задачи будут поставлены в очередь экземпляра `Executor`, вместо того чтобы ждать, пока все они будут завершены. Если вы хотите отправить набор задач и дождаться их завершения, можно использовать метод `ExecutorService.invokeAll`; для получения результатов по мере их поступления, можно использовать экземпляр `CompletionService`, как в классе `Renderer` из раздела 6.3.6.

Последовательные итерации цикла подходят для распараллеливания, когда каждая итерация независима от других и работа, выполняемая в каждой итерации тела цикла, достаточно значительна, чтобы компенсировать затраты на управление новой задачей.

Распараллеливание циклов также может быть применено в случае некоторых рекурсивных дизайнов; в рекурсивном алгоритме часто встречаются последовательные циклы, к распараллеливанию которых можно применить такой же подход, как в листинге 8.10. Простейший случай - это когда итерации не ожидают возвращения результатов выполнения рекурсивных итераций, которые они вызывают. Например, метод `sequentialRecursive` в листинге 8.11 выполняет обход дерева “в глубину”, выполняя вычисления на каждом узле и помещая результат в коллекцию. Преобразованная версия, приведённая в методе `parallelRecursive`, также выполняет обход в глубину, но вместо вычисления результата при посещении каждого узла, она отправляет задачу для вычисления результата соответствующего узлу.

```
public<T> void sequentialRecursive(List<Node<T>> nodes,
                                    Collection<T> results) {
    for (Node<T> n : nodes) {
        results.add(n.compute());
        sequentialRecursive(n.getChildren(), results);
    }
}

public<T> void parallelRecursive(final Executor exec,
                                List<Node<T>> nodes,
                                final Collection<T> results) {
    for (final Node<T> n : nodes) {
        exec.execute(new Runnable() {
            public void run() {
                results.add(n.compute());
            }
        });
        parallelRecursive(exec, n.getChildren(), results);
    }
}
```

Листинг 8.11 Преобразование последовательной хвостовой рекурсии в распараллеленную рекурсию

К моменту возвращения управления методом `parallelRecursive`, каждый узел дерева уже был посещен (обход по-прежнему осуществляется последовательно: параллельно выполняются только вызовы метода `compute`) и вычисления, выполняемые для каждого узла, уже были поставлены в очередь экземпляра `Executor`. Объекты,зывающие метод `parallelRecursive`, могут ожидать получения всех результатов, создав специфичный, предназначенный для обхода, экземпляр `Executor` и используя методы `shutdown` и `awaitTermination`, как показано в листинге 8.12.

```
public<T> Collection<T> getParallelResults(List<Node<T>> nodes)
    throws InterruptedException {
```

```

ExecutorService exec = Executors.newCachedThreadPool();
Queue<T> resultQueue = new ConcurrentLinkedQueue<T>();
parallelRecursive(exec, nodes, resultQueue);
exec.shutdown();
exec.awaitTermination(Long.MAX_VALUE, TimeUnit.SECONDS);
return resultQueue;
}

```

Листинг 8.12 Ожидание результатов параллельных вычислений

8.5.1 Пример: фреймворк для решения головоломок

Такой подход привлекателен для поиска решений в головоломках, которые включают в себя поиск последовательности преобразований из некоторого начального состояния в некоторое целевое состояние, например, такие головоломки как «головоломки на перемещение элементов»¹⁰², «Hi-Q», «Instant Insanity» и прочие головоломки-пасьянсы.

Мы определяем “головоломку” как комбинацию начальной позиции, целевой позиции и набора правил, определяющих допустимые ходы. Набор правил состоит из двух частей: вычислённого списка возможных, с заданной позиции, ходов и вычисления результата применения хода к позиции. В интерфейсе `Puzzle` из листинга 8.13, показана абстракция нашей головоломки; параметры типа `P` и `M` представляют собой классы позиции и хода. Основываясь на приведённом интерфейсе, мы можем написать простой последовательный решатель, который просматривает пространство шагов головоломки до тех пор, пока не будет найдено решение или пространство шагов головоломки не будет исчерпано.

```

public interface Puzzle<P, M> {
    P initialPosition();
    boolean isGoal(P position);
    Set<M> legalMoves(P position);
    P move(P position, M move);
}

```

Листинг 8.13 Абстракция для головоломки, подобной “головоломка на перемещение элементов”

Класс `Node` из листинга 8.14, представляет собой позицию, которая была достигнута с помощью некоторой серии шагов, путём удержания ссылки на перемещение, создавшее позицию, и предыдущий узел. Переход по ссылкам назад, начиная с текущего экземпляра `Node`, позволяет восстановить последовательность ходов, которые привели к текущей позиции.

```

@Immutable
static class Node<P, M> {
    final P pos;
    final M move;
    final Node<P, M> prev;

    Node(P pos, M move, Node<P, M> prev) {...}

    List<M> asMoveList() {

```

¹⁰² See <http://www.puzzleworld.org/SlidingBlockPuzzles>.

```

        List<M> solution = new LinkedList<M>();
        for (Node<P, M> n = this; n.move != null; n = n.prev)
            solution.add(0, n.move);
        return solution;
    }
}

```

Листинг 8.14 Связующий узел в фреймворке решения загадок

В классе SequentialPuzzleSolver из листинга 8.15, приведён последовательный решатель, используемый в фреймворке решения головоломок и выполняющий поиск решения “в глубину” в пространстве шагов головоломки. Он завершает свою работу, когда находит решение (которое не обязательно является кратчайшим решением).

```

public class SequentialPuzzleSolver<P, M> {
    private final Puzzle<P, M> puzzle;
    private final Set<P> seen = new HashSet<P>();

    public SequentialPuzzleSolver(Puzzle<P, M> puzzle) {
        this.puzzle = puzzle;
    }

    public List<M> solve() {
        P pos = puzzle.initialPosition();
        return search(new Node<P, M>(pos, null, null));
    }

    private List<M> search(Node<P, M> node) {
        if (!seen.contains(node.pos)) {
            seen.add(node.pos);
            if (puzzle.isGoal(node.pos))
                return node.asMoveList();
            for (M move : puzzle.legalMoves(node.pos)) {
                P pos = puzzle.move(node.pos, move);
                Node<P, M> child = new Node<P, M>(pos, move, node);
                List<M> result = search(child); if (result != null)

                    return result;
            }
        }
        return null;
    }

    static class Node<P, M> { /* Listing 8.14 */ }
}

```

Листинг 8.15 Последовательная версия решателя головоломок

Переписав решатель с использованием параллелизма, мы сможем параллельно вычислять последующие шаги, а также, параллельно, проводить оценку на соответствие целевому условию, так как процесс оценки одного хода в основном

не зависит от оценки других ходов. (Мы говорим "в основном", потому что задачи имеют некоторое совместно используемое состояние, такое как набор посещённых позиций.)

Параллельная версия решателя головоломок, представленная в листинге 8.16, использует внутренний класс `SolverTask` расширяющий класс `Node` и реализующий интерфейс `Runnable`. Большая часть работы выполняется в методе `run`: оценка набора возможных следующих позиций, отсечение ранее посещённых позиций, проверка факта достижения успеха (этой задачей или какой-то другой) и отправка ранее не посещённых позиций экземпляру `Executor`.

```
public class ConcurrentPuzzleSolver<P, M> {
    private final Puzzle<P, M> puzzle;
    private final ExecutorService exec;
    private final ConcurrentMap<P, Boolean> seen;
    final ValueLatch<Node<P, M>> solution
        = new ValueLatch<Node<P, M>>();
    ...
    public List<M> solve() throws InterruptedException {
        try {
            P p = puzzle.initialPosition();
            exec.execute(newTask(p, null, null));
            // block until solution found
            Node<P, M> solnNode = solution.getValue();
            return (solnNode == null) ? null : solnNode.asMoveList();
        } finally {
            exec.shutdown();
        }
    }

    protected Runnable newTask(P p, M m, Node<P,M> n) {
        return new SolverTask(p, m, n);
    }

    class SolverTask extends Node<P, M> implements Runnable {
        ...
        public void run() {
            if (solution.isSet()
                || seen.putIfAbsent(pos, true) != null)
                return; // already solved or seen this position
            if (puzzle.isGoal(pos))
                solution.setValue(this);
            else
                for (M m : puzzle.legalMoves(pos))
                    exec.execute(
                        newTask(puzzle.move(pos, m), m, this));
        }
    }
}
```

Листинг 8.16 Параллельная версия решателя головоломок

Чтобы избежать бесконечных циклов, последовательная версия поддерживала экземпляр `Set`, с набором ранее посещённый позиций; класс `ConcurrentPuzzleSolver` использует той же цели экземпляр `ConcurrentHashMap`. Такой подход обеспечивает потокобезопасность и позволяет избежать возникновения состояния гонки, присущего условному обновлению совместно используемой коллекции, за счёт использования метода `putIfAbsent` для атомарного добавления позиции только в том случае, если она не была ранее посещена. Для хранения состояния поиска, параллельная версия решателя головоломок использует, вместо стека вызовов, внутреннюю рабочую очередь потоков.

Параллельный подход также приводит к подмене ограничения одной формы на другую, которая может быть более подходящей для данной проблемной области. Последовательная версия выполняет поиск “в глубину”, поэтому поиск ограничен доступным размером стека. Параллельная версия выполняет поиск “в ширину” и поэтому свободна от ограничения размера стека (но все еще может столкнуться с проблемой нехватки памяти, если набор позиций для посещения или уже посещённых позиций превысит доступную память).

Для того, чтобы остановить поиск, когда решение найдено, нам нужен способ определения, нашел ли какой-либо поток решение. Если мы хотим принять первое найденное решение, нам также необходимо обеспечить обновление решения только в том случае, если ни одна другая задача ранее не нашла его. Эти требования описывают некоторую разновидность защёлки (см. раздел [5.5.1](#)) и в частности, *защелку с результатом* (*result-bearing latch*). Мы могли бы легко создать блокирующую защёлку с результатом, используя методы, описанные в главе 14, но часто проще и менее подвержено ошибкам использовать существующие библиотечные классы, а не низкоуровневые механизмы языка. Класс `ValueLatch` из листинга 8.17 использует класс `CountDownLatch` для обеспечения необходимого фиксирующего поведения и использует блокировку, чтобы гарантировать факт того, что решение будет установлено только один раз.

```
@ThreadSafe
public class ValueLatch<T> {
    @GuardedBy("this") private T value = null;
    private final CountDownLatch done = new CountDownLatch(1);

    public boolean isSet() {
        return (done.getCount() == 0);
    }

    public synchronized void setValue(T newValue) {
        if (!isSet()) {
            value = newValue;
            done.countDown();
        }
    }

    public T getValue() throws InterruptedException {
        done.await();
        synchronized (this) {
            return value;
        }
    }
}
```

```
    }
}
}
```

Листинг 8.17 Защёлка с результатом, используемая классом ConcurrentPuzzleSolver

Каждая задача сначала проверяет защёлку с решением и прекращает выполнение, если решение уже найдено. Главный поток должен ожидать, пока решение не будет найдено; метод `getValue` экземпляра `ValueLatch` блокируется, пока какой либо поток не установит значение. Класс `ValueLatch` обеспечивает способ хранения значения таким образом, что только первый вызов фактически устанавливает значение, вызывающие объекты могут проверить, было ли значение установлено, и вызывающие объекты могут быть заблокированы в ожидании его установки. Решение обновляется при первом вызове метода `setValue`, а счётчик экземпляра `CountDownLatch` уменьшается, освобождая, таким образом, основной поток решателя от ожидания в методе `getValue`.

Первый нашедший решение поток также завершает работу экземпляра `Executor`, для предотвращения отправки новых задач. Чтобы избежать необходимости иметь дело с исключением `RejectedExecutionException`, обработчик отклонённых задач должен быть настроен на то, чтобы отбрасывать отправляемые задачи. Затем все незавершенные задачи, в конечном итоге, завершают своё выполнение, и любые последующие попытки выполнить новые задачи тихо прерываются, позволяя исполнителю завершить свою работу. (Если задачи выполняются слишком долго, мы могли бы прервать их, вместо того, чтобы позволить им завершить своё выполнение.)

Класс `ConcurrentPuzzleSolver` плохо справляется со случаем, когда у головоломки нет решения: если все возможные ходы и позиции были оценены и решение не найдено, метод `solve` зависнет в вечном ожидании результата вызова метода `getSolution`. Последовательная версия завершает свою работу, когда исчерпает всё пространство поиска, но процесс завершения работы параллельных программ иногда может быть сложнее. Одним из возможных решений является сохранение количества активных задач решателя и присвоение решению значения `null` при уменьшении количества активных задач до нуля, как показано в листинге 8.18.

```
public class PuzzleSolver<P,M> extends ConcurrentPuzzleSolver<P,M> {
    ...
    private final AtomicInteger taskCount = new AtomicInteger(0);

    protected Runnable newTask(P p, M m, Node<P,M> n) {
        return new CountingSolverTask(p, m, n);
    }

    class CountingSolverTask extends SolverTask {
        CountingSolverTask(P pos, M move, Node<P, M> prev) {
            super(pos, move, prev);
            taskCount.incrementAndGet();
        }
        public void run() {
            try {
```

```
        super.run();
    } finally {
        if (taskCount.decrementAndGet() == 0)
            solution.setValue(null);
    }
}
}
```

Листинг 8.18 Решатель, с механизмом распознавания отсутствия решения

Поиск решения также может занять больше времени, чем мы готовы ожидать; есть несколько дополнительных условий завершения, которые мы могли бы наложить на решатель. Одним из них является ограничение по времени; его легко добавить, реализовав временную версию метода `getValue` в классе `ValueLatch` (который будет использовать временную версию метода `await`) и, завершая работу экземпляра `Executor` с объявлением сбоя в том случае, если истечёт время ожидания, установленное методу `getValue`. Другой подход - это своего рода метрика, специфичная для головоломки, такая как поиск только до определенного количества позиций. Или мы можем предусмотреть механизм отмены и позволить клиенту принять собственное решение о том, когда следует прекратить поиск.

8.6 Итоги

Фреймворк `Executor` представляет собой мощный и гибкий фреймворк для обеспечения параллельного выполнения задач. Он предлагает ряд настроек параметров, таких как политика создания и завершения потоков, обработка задач в очереди и настройка действий по отношению к избыточным задачам, а также предоставляет несколько хуков, расширяющих его поведение. Однако, как и в большинстве других мощных фреймворков, существуют комбинации параметров, которые совместно работают плохо; некоторые типы задач требуют определенных политик выполнения, а некоторые комбинации настроек параметров могут привести к странным результатам.

Глава 9 Приложения GUI

Если вы когда-либо пытались написать даже простое приложение GUI с помощью Swing, вы знаете, что GUI-приложения имеют собственные своеобразные проблемы с потоками. Для обеспечения безопасности, определённые задачи должны выполняться в потоке событий Swing. Но вы не можете выполнять длительные задачи в потоке событий из-за боязни того, что пользовательский интерфейс перестанет отвечать на запросы. И структуры данных Swing не являются потокобезопасными, поэтому необходимо быть осторожными, ограничивая их рамками потока событий.

Почти все наборы инструментов GUI, включая Swing и SWT, реализованы в виде *однопоточных подсистем*, в которых все действия GUI ограничены одним потоком. Если вы не планируете писать полностью однопоточную программу, в ней будут активности, которые будут частично выполняться в потоке приложения, а частично - в потоке событий. Подобно многим другим ошибкам многопоточности, неправильное разделение между этими частями не обязательно приведёт к немедленному возникновению аварии в вашей программе; вместо этого, программа может вести себя странно при сложных условиях. Несмотря на то, что фреймворки GUI сами по себе являются однопоточными подсистемами, ваше приложение может не являться таковым, и вам все равно придётся тщательно рассматривать проблемы с потоками при написании кода GUI.

9.1 Почему фреймворки GUI однопоточны?

В прошлом, приложения GUI были однопоточными, и события GUI обрабатывались из "основного цикла обработки событий". Современные фреймворки GUI используют немного отличающуюся модель: они создают выделенный поток диспетчеризации событий (EDT, *event dispatch thread*) для обработки событий GUI.

Однопоточные фреймворки GUI не являются уникальным явлением, присущим только Java; Qt, NextStep, MacOS Cocoa, X Windows, и многие другие также однопоточны. Сложившаяся ситуация является таковой не из-за отсутствия попыток; было предпринято множество попыток написать многопоточные фреймворки GUI, но из-за постоянных проблем с условиями гонок и взаимоблокировками, все они, в конечном итоге, пришли к однопоточной модели с очередью событий, в которой выделенный поток извлекает события из очереди и отправляет их обработчикам событий, определённым приложением. (Фреймворк AWT изначально пытался в большей степени поддерживать многопоточный доступ, и решение сделать Swing однопоточным, было, в основном, основано на опыте, полученном при работе с AWT.)

Многопоточные фреймворки GUI, как правило, особенно чувствительны к взаимоблокировкам, частично из-за неудачного взаимодействия между механизмом обработки входящих событий и объектно-ориентированной моделью компонентов GUI. Действия, инициируемые пользователем, как правило "всплывают" от ОС к приложению — щелчок мыши обнаруживается ОС, превращается с помощью инструментария в событие "щелчок мыши" и, в конечном итоге, доставляется слушателю приложения как событие более высокого уровня, такое как "нажатие кнопки". С другой стороны, действия инициируемые приложением "проваливаются" от приложения в ОС - изменение цвета фона

компонента происходит в приложении и отправляется определенному классу компонента, и, в конечном итоге, в ОС для визуализации. Сочетание тенденции получения активностями доступа к одним и тем же объектам GUI в противоположном порядке, с требованием сделать каждый объект потокобезопасным, порождает рецепт возникновения несогласованного порядка блокировок, что приводит к возникновению взаимоблокировок (см. главу 10). И это именно тот опыт, который, всякий раз в процессе разработки, переоткрывают для себя почти все инструментарии GUI.

Другим фактором, приводящим к возникновению взаимоблокировок в многопоточных фреймворках GUI, является преобладание шаблона модель-представление-контроллер (model-view-controller, MVC). Факторинг взаимодействий пользователя в кооперацию объектов, представляющих модель, представление и контроллер, значительно упрощает реализацию приложений GUI, но вновь повышает риск несогласованного упорядочивания блокировок. Контроллер отправляет вызов в модель, которая уведомляет представление о том, что что-то изменилось. Но контроллер также может отправить вызов в представление, которое, в свою очередь, может отправить вызов в модель для запроса состояния модели. В результате будет получен несогласованный порядок блокировки, с сопутствующим риском возникновения взаимоблокировки.

В своём блоге¹⁰³, Sun VP Грэм Гамильтон отлично подводит итоги по вышеупомянутым проблемам, описывая, почему многопоточный инструментарий GUI является повторяющейся "провальной мечтой" компьютерной науки.

Я считаю, что вы можете успешно программировать с использованием многопоточного инструментария GUI, если инструментарий спроектирован очень тщательно; если инструментарий во всех деталях раскрывает используемую методологию блокировки; если вы очень умны, очень осторожны, и у вас есть глобальное понимание всей структуры инструментария. Если вы хотя бы немного ошибётесь в одной из этих вещей, все будет в основном работать, но вы получите случайные зависания (из-за взаимоблокировок) или глюки (из-за условий гонок). Такой многопоточный подход лучше всего подходит для людей, которые принимали непосредственное участие в разработке инструментария.

К сожалению, я не думаю, что перечисленный набор характеристик подходит для широкого коммерческого использования. Вы создаёте приложение, которое работает не совсем надежно по причинам, которые вовсе не очевидны, силами нормальных умных программистов. Таким образом, авторы приложения очень недовольны и разочарованы и употребляют плохие слова по адресу бедного невинного инструментария.

Однопоточные фреймворки GUI обеспечивают потокобезопасность посредством ограничения потока; все объекты GUI, включая визуальные компоненты и модели данных, доступны исключительно из потока событий. Конечно, это просто переносит часть нагрузки по обеспечению потокобезопасности обратно на плечи разработчика приложения, который должен убедиться, что объекты ограничены правильным образом.

¹⁰³ <http://weblogs.java.net/blog/kgh/archive/2004/10>

9.1.1 Последовательная обработка событий

Приложения GUI ориентированы на обработку небольших событий, таких как щелчок мыши, нажатие клавиши или истечение времени таймера. Событие представляет собой некоторый вид задачи; механизм обработки событий, предоставляемый AWT и Swing, структурно похож на механизм `Executor`.

Поскольку существует только один поток для обработки задач GUI, они обрабатываются последовательно - одна задача завершается до начала выполнения следующей, и никакие две задачи не перекрываются. Это знание упрощает написание кода задачи – вам не нужно беспокоиться о вмешательстве других задач.

Недостатком последовательной обработки задач является то, что если выполнение одной задачи занимает много времени, другие задачи должны ждать ее завершения. Если эти другие задачи ответственны за реакцию на ввод данных пользователем или обеспечение визуальной обратной связи, приложение будет казаться замороженным. Если в потоке событий выполняется длительная задача, пользователь даже не сможет нажать кнопку "Отмена", поскольку слушатель кнопки отмены не будет вызван до завершения длительной задачи. Поэтому задачи, выполняемые в потоке событий, должны как можно быстрее возвращать управление потоку событий. Чтобы запустить длительную задачу, такую как проверка орфографии большого документа, поиск в файловой системе или получение ресурса по сети, необходимо запустить эту задачу в другом потоке, чтобы элемент управления мог быстро вернуться в поток событий. Для обновления состояния индикатора выполнения во время обработки длительной задачи или обеспечения визуальной обратной связи после ее завершения, необходимо вновь выполнить код в потоке событий. Всё это ведёт к быстрому усложнению программы.

9.1.2 Ограничение потока в Swing

Все компоненты Swing (такие как `JButton` и `JTable`) и объекты модели данных (такие как `TableModel` и `TreeModel`) ограничены потоком событий, таким образом, любой код, который обращается к этим объектам, должен работать в потоке событий. Объекты GUI хранятся согласованно без использования синхронизации, но с помощью ограничения потока. Преимуществом такого подхода является то, что задачи, выполняемые в потоке событий, не должны беспокоиться о синхронизации, когда получают доступ к объектам презентации; недостатком такого подхода является то, что вы вообще не можете получить доступ к объектам презентации извне потока событий.

Правило однопоточности Swing: компоненты и модели Swing должны создаваться, изменяться и запрашиваться только из потока диспетчеризации событий.

Как и во всех правилах, имеется несколько исключений. Небольшое количество методов Swing может быть безопасно вызвано из любого потока; они явно определены в Javadoc как потокобезопасные. Другие исключения из правила однопоточности включают:

- Метод `SwingUtilities.isEventDispatchThread`, определяющий, является ли текущий поток потоком событий;

- Метод `SwingUtilities.invokeLater`, планирующий выполнение экземпляра `Runnable` в потоке событий (вызываемый из любого потока);
- Метод `SwingUtilities.invokeAndWait`, планирующий задачу выполнения экземпляра `Runnable` в потоке событий и блокирующий текущий поток до её завершения (вызывается *только* не из потоков GUI);
- Методы, помещающие запросы на перерисовку или повторную проверку в очередь событий (вызываются из любого потока); и
- Методы, добавляющие и удаляющие слушателей (можно вызывать из любого потока, но слушатели всегда будут вызываться в потоке событий).

Методы `invokeLater` и `invokeAndWait` функционируют подобно механизму `Executor`. На практике, достаточно просто реализовать связанные с потоками методы класса `SwingUtilities`, используя однопоточную реализацию `Executor`, как показано в листинге 9.1.

```
public class SwingUtilities {  
    private static final ExecutorService exec =  
        Executors.newSingleThreadExecutor(new SwingThreadFactory());  
    private static volatile Thread swingThread;  
  
    private static class SwingThreadFactory implements ThreadFactory {  
        public Thread newThread(Runnable r) {  
            swingThread = new Thread(r);  
            return swingThread;  
        }  
    }  
  
    public static boolean isEventDispatchThread() {  
        return Thread.currentThread() ==  
            swingThread;  
    }  
  
    public static void invokeLater(Runnable task)  
    { exec.execute(task);  
    }  
  
    public static void invokeAndWait(Runnable task)  
        throws InterruptedException, InvocationTargetException {  
        Future f = exec.submit(task);  
        try {  
            f.get();  
        } catch (ExecutionException e) {  
            throw new InvocationTargetException(e);  
        }  
    }  
}
```

Листинг 9.1 Реализация класса `SwingUtilities` с использованием экземпляра `Executor`

Описанное выше не отражает того, как класс `SwingUtilities` реализован фактически, поскольку фреймворк Swing предшествует фреймворку Executor, но, вероятно отражает, как он мог бы быть построен, если бы Swing реализовывался сегодня.

Поток событий Swing можно рассматривать как однопоточную реализацию Executor, обрабатывающую задачи из очереди событий. Как и в случае пулов потоков, рабочий поток иногда умирает и заменяется новым, но это должно происходить прозрачно для задач. Последовательное однопоточное выполнение является разумной политикой выполнения, когда задачи кратковременные, предсказуемость планирования не имеет значения или крайне важно, чтобы задачи не выполнялись параллельно.

Класс `GuiExecutor` из листинга 9.2 представляет собой реализацию Executor, делегирующую задачи на выполнение классу `SwingUtilities`. Подобный подход мог бы быть реализован и в терминах других GUI-фреймворков; например, фреймворк SWT предоставляет метод `Display.asyncExec`, подобный методу `invokeLater` фреймворка Swing.

```
public class GuiExecutor extends AbstractExecutorService {  
    // Singletions have a private constructor and a public factory  
    private static final GuiExecutor instance = new GuiExecutor();  
  
    private GuiExecutor() { }  
  
    public static GuiExecutor instance() { return instance; }  
  
    public void execute(Runnable r) {  
        if (SwingUtilities.isEventDispatchThread())  
            r.run();  
        else  
            SwingUtilities.invokeLater(r);  
    }  
  
    // Plus trivial implementations of lifecycle methods  
}
```

Листинг 9.2 Реализация Executor построенная на основе класса `SwingUtilities`.

9.2 Кратковременные задачи GUI

В приложении GUI события происходят в потоке событий и всплывают до уровня слушателей, предоставляемых приложением, которые, вероятно, выполнят некоторые вычисления, влияющие на объекты уровня представления. Для простых кратковременных (*short-running*) задач, действие может целиком оставаться в потоке событий; для долговременных (*longer-running*) задач, часть обработки должна быть выгружена в отдельный поток.

В простом случае, ограничение объектов представления потоком событий вполне естественно. В листинге 9.3 создаётся кнопка, цвет которой в момент нажатия изменяется случайным образом. Когда пользователь нажимает на кнопку, инструментарий доставляет экземпляр события `ActionEvent` всем

зарегистрированным в потоке событий слушателям действий (*action listeners*). В ответ слушатель действий выбирает новый цвет и изменяет цвет фона кнопки.

```
final Random random = new Random();
final JButton button = new JButton("Change Color");
...
button.addActionListener(new ActionListener() {
    public void actionPerformed(ActionEvent e) {
        button.setBackground(new Color(random.nextInt()));
    }
});
```

Листинг 9.3 Простой слушатель событий

Таким образом, событие берёт своё начало в инструментарии GUI и доставляется приложению, а приложение изменяет GUI в ответ на действие пользователя. Как показано на рис. 9.1, элемент управления никогда не покидает поток событий.

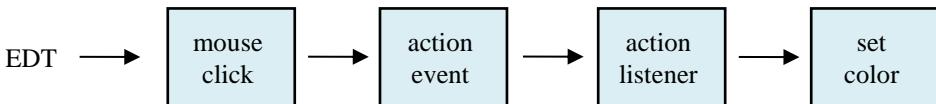


Рисунок 9.1 Контроль управления в случае простого клика по кнопке

Этот тривиальный пример характеризует большинство взаимодействий между GUI-приложениями и инструментарием GUI. До тех пор, пока задачи кратковременны и имеют доступ только к объектам GUI (или другим ограниченным потоком или потокобезопасным объектам приложений), вы можете почти полностью игнорировать проблемы с потоками и делать все из потока событий, и всё будет происходить правильно.

Несколько более сложная версия того же сценария, иллюстрируемая на рис. 9.2, предполагает использование формальной модели данных, такой как `TableModel` или `TreeModel`. Swing разделяет большинство визуальных компонентов на два объекта: модель (*model*) и представление (*view*). Отображаемые данные находятся в модели, а правила, регулирующие их отображение, - в представлении.

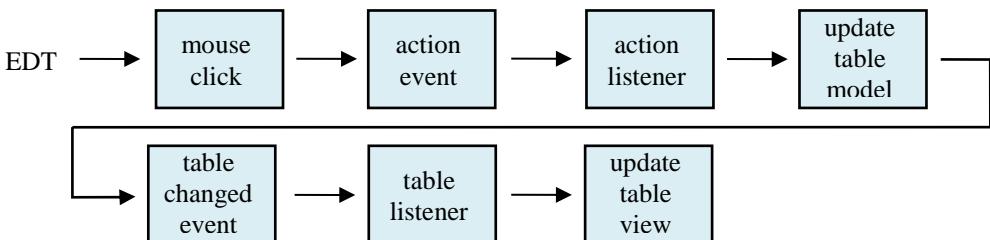


Рисунок 9.2 Контроль управления в случае разделения объектов модели и представления

Объекты модели могут инициировать события, указывающие на изменение данных в модели, и представления подписываются на эти события. Когда представление получает событие, указывающее, что данные модели, возможно, изменились, оно запрашивает у модели новые данные и обновляет отображение. Таким образом, в слушателе кнопки, изменяющей содержимое таблицы, слушатель действий обновит модель и вызовет один из методов `fireXXX`, который, в свою очередь, вызовет слушателя табличной модели представления, который обновит представление. И вновь, элемент управления никогда не покидает поток событий.

(Методы `fireXXX` модели данных Swing всегда напрямую вызывают слушателей модели, вместо того, чтобы передать новое событие в очередь событий, таким образом, методы `fireXXX` должны вызывать только в потоке событий.)

9.3 Долговременные задачи GUI

Если бы все задачи были кратковременными (и у приложения не было бы существенной части кода без GUI), то всё приложение могло бы работать в потоке событий, и вам в принципе не пришлось бы обращать внимание на потоки. Однако сложные приложения с графическим интерфейсом могут выполнять задачи, выполнение которых может занять больше времени, чем пользователь готов ожидать, например, проверка орфографии, фоновая компиляция или извлечение удаленных ресурсов. Эти задачи должны выполняться в другом потоке, чтобы графический интерфейс оставался отзывчивым во время их выполнения.

Фреймворк Swing упрощает выполнение задачи в потоке событий, но (до Java версии 6) не предоставляет механизма, оказывающего помощь потокам GUI в выполнении кода в других потоках. В принципе, в этом вопросе мы можем обойтись и без помощи со стороны Swing: мы можем создать свою собственную реализацию Executor для обработки долговременных задач. Кэшированный пул потоков является хорошим выбором для долговременных задач; достаточно редко приложения GUI инициируют большое количество долговременных задач, таким образом, существующий риск роста пула потоков без ограничений невелик.

Мы начнём с простых задач, которые не поддерживают отмену или индикацию прогресса и не обновляют GUI по завершению выполнения, а затем будем постепенно, по одной, добавлять эти функции. В листинге 9.4 показан слушатель действий, ограниченный визуальным компонентом, который отправляет долговременную задачу экземпляру Executor. Несмотря на два уровня внутренних классов, инициировать выполнение задачи с использованием приведённого выше подхода, при наличии задачи GUI, довольно просто: слушатель действий пользовательского интерфейса вызывается в потоке событий и отправляет экземпляр Runnable на выполнение в пул потоков.

```
ExecutorService backgroundExec = Executors.newCachedThreadPool();  
...  
button.addActionListener(new ActionListener() {  
    public void actionPerformed(ActionEvent e) {  
        backgroundExec.execute(new Runnable() {  
            public void run() { doBigComputation(); }  
        });  
    }});
```

Листинг 9.4 Привязка долговременной задачи к визуальному компоненту

Этот пример извлекает долговременную задачу из потока событий способом "воспламенить и забыть", что, вероятно, не очень полезно. Обычно, при выполнении долговременной задачи, возникает визуальная обратная связь. Но вы не можете получить доступ к презентационным объектам из фонового потока, поэтому, после завершения выполнения, задача должна отправить другую задачу для запуска в потоке событий, с целью обновления состояния пользовательского интерфейса.

В листинге 9.5 иллюстрируется очевидный способ реализации такого подхода, начавший постепенно усложняться; теперь мы имеем уже три уровня внутренних классов. Слушатель действий сначала затемняет кнопку и устанавливает метку, указывающую, что выполняется вычисление, а затем отправляет задачу фоновому исполнителю. Когда текущая задача завершается, она помещает в очередь на выполнение в потоке событий другую задачу, которая повторно включает кнопку и восстанавливает текст метки.

```
button.addActionListener(new ActionListener() {
    public void actionPerformed(ActionEvent e) {
        button.setEnabled(false);
        label.setText("busy");
        backgroundExec.execute(new Runnable() {
            public void run() {
                try {
                    doBigComputation();
                } finally {
                    GuiExecutor.instance().execute(new Runnable() {
                        public void run() {
                            button.setEnabled(true);
                            label.setText("idle");
                        }
                    });
                }
            }
        });
    }
});
```

Листинг 9.5 Долговременные пользовательские задачи с обратной связью

Задача, запускаемая при нажатии кнопки, состоит из трех последовательных подзадач, выполнение которых чередуется между потоком событий и фоновым потоком. Первая подзадача обновляет пользовательский интерфейс, показывая, что началось выполнение длительной операции, и запускает вторую подзадачу в фоновом потоке. По завершении вторая подзадача помещает третью подзадачу в очередь потока событий для повторного выполнения, целью обновления пользовательского интерфейса, для отражения того, что операция завершена. Этот вид “скаккообразной перестройки потока” типичен для обработки длительных задач в приложениях GUI.

9.3.1 Отмена

Любая задача, запуск которой на выполнение в другом потоке занимает достаточно много времени, вероятно, также потребует достаточно много времени на её отмену пользователем. Вы можете реализовать отмену напрямую, используя прерывание потока, но гораздо проще использовать интерфейс `Future`, который был специально разработан для управления отменяемыми задачами.

При вызове метода `cancel` экземпляра `Future`, с параметром `mayInterruptIfRunning` установленным в `true`, реализация `Future` прерывает поток, выполняющий задачу, если он выполняется в данный момент. Если ваша задача написана отзывчивой на прерывание, в случае инициации процесса отмены,

она может вернуть управление раньше. В листинге 9.6 иллюстрируется задача, которая опрашивает статус прерывания потока и при прерывании возвращает управление раньше.

```
Future<?> runningTask = null; // thread-confined
...
startButton.addActionListener(new ActionListener() {
    public void actionPerformed(ActionEvent e) {
        if (runningTask == null) {
            runningTask = backgroundExec.submit(new Runnable() {
                public void run() {
                    while (moreWork()) {
                        if (Thread.currentThread().isInterrupted()) {
                            cleanUpPartialWork();
                            break;
                        }
                        doSomeWork();
                    }
                }
            });
        }
    }
});

cancelButton.addActionListener(new ActionListener() {
    public void actionPerformed(ActionEvent event) {
        if (runningTask != null)
            runningTask.cancel(true);
    }
});
```

Листинг 9.6 Отмена долговременной задачи

Поскольку переменная `runningTask` ограничена потоком событий, при установке или проверке значения синхронизация не требуется, а слушатель кнопки запуска гарантирует, что одновременно выполняется только одна фоновая задача. Однако лучше получать уведомления о завершении задачи, чтобы, например, можно было отключить кнопку отмены. Мы рассмотрим этот подход в следующем разделе.

9.3.2 Индикатор прогресса и завершения

Использование экземпляра `Future` для представления долговременной задачи значительно упрощает реализацию механизма отмены. Класс `FutureTask` имеет хук `done`, что так же облегчает уведомление о завершении. После завершения фонового выполнения экземпляра `Callable`, вызывается метод `done`. Выполняя метод `done` в потоке событий, по завершении задачи, мы можем построить класс `BackgroundTask`, предоставляющий хук `onCompletion`, который вызывается в потоке событий, как показано в листинге 9.7.

```
abstract class BackgroundTask<V> implements Runnable, Future<V> {
    private final FutureTask<V> computation = new Computation();
```

```

private class Computation extends FutureTask<V> {
    public Computation() {
        super(new Callable<V>() {
            public V call() throws Exception { return
                BackgroundTask.this.compute();
            }
        });
    }
    protected final void done() {
        GuiExecutor.instance().execute(new Runnable() {
            public void run() {
                V value = null;
                Throwable thrown = null;
                boolean cancelled = false;
                try {
                    value = get();
                } catch (ExecutionException e) {
                    thrown = e.getCause();
                } catch (CancellationException e) {
                    cancelled = true;
                } catch (InterruptedException consumed) {
                } finally {
                    onCompletion(value, thrown, cancelled);
                }
            };
        });
    }
    protected void setProgress(final int current, final int max) {
        GuiExecutor.instance().execute(new Runnable() {
            public void run() { onProgress(current, max); }
        });
    }
    // Called in the background thread
    protected abstract V compute() throws Exception;
    // Called in the event thread
    protected void onCompletion(V result, Throwable exception,
                               boolean cancelled) { }
    protected void onProgress(int current, int max) { }
    // Other Future methods forwarded to computation
}

```

Листинг 9.7 Класс фоновой задачи, поддерживающий механизм отмены, уведомление о завершении задачи и уведомление о ходе выполнения

Класс `BackgroundTask` также поддерживает индикацию хода выполнения. Метод `compute` может вызывать метод `setProgress`, указывая прогресс в числовом выражении. Хук `onProgress` вызывающийся из потока событий, может обновить пользовательский интерфейс для визуального отображения состояния прогресса.

Для реализации абстрактного класса `BackgroundTask`, достаточно реализовать метод `compute`, вызывающийся в фоновом потоке. У вас также есть возможность переопределения хуков `onCompletion` и `onProgress`, которые вызываются в потоке событий.

Класс `BackgroundTask` базирующийся на классе `FutureTask`, также упрощает отмену. Вместо того, чтобы опрашивать статус прерывания потока, метод `compute` может вызвать метод `Future.isCancelled`. В Листинге 9.8 демонстрируется пример из листинга 9.6, переписанный с использованием класса `BackgroundTask`.

```
startButton.addActionListener(new ActionListener() {
    public void actionPerformed(ActionEvent e) {
        class CancelListener implements ActionListener {
            BackgroundTask<?> task;
            public void actionPerformed(ActionEvent event) {
                if (task != null)
                    task.cancel(true);
            }
        }

        final CancelListener listener = new CancelListener();
        listener.task = new BackgroundTask<Void>() {
            public Void compute() {
                while (moreWork() && !isCancelled())
                    doSomeWork();
                return null;
            }
            public void onCompletion(boolean cancelled, String s,
                                   Throwable exception) {
                cancelButton.removeActionListener(listener);
                label.setText("done");
            }
        };
        cancelButton.addActionListener(listener);
        backgroundExec.execute(listener.task);
    }
});
```

Листинг 9.8 Инициация долговременной, отменяемой задачи с помощью класса `BackgroundTask`

9.3.3 Класс `SwingWorker`

Мы построили простой фреймворк, с использованием класса `FutureTask` и интерфейса `Executor`, для выполнения долговременных задач в фоновых потоках, без ущерба для отзывчивости GUI. Эти методы могут быть применены к любому однопоточному фреймворку GUI, а не только Swing. В Swing, многие ранее разработанные функции, предоставляются классом `SwingWorker`, включая отмену, уведомление о завершении и индикацию хода выполнения. Различные версии `SwingWorker` были опубликованы в *The Swing Connection* и *The Java Tutorial*, а их обновленная версия была включена в Java 6.

9.4 Совместно используемые модели данных

Объекты представления Swing, включая объекты модели данных, такие как `TableModel` или `TreeModel`, ограничены потоком событий. В простых программах с GUI все изменяемое состояние хранится в объектах презентации, и единственный поток, кроме потока событий, является основным потоком. В этих программах легко применить правило однопоточности: не обращаться к компонентам модели данных или представления из основного потока. Более сложные программы могут использовать другие потоки для перемещения данных в постоянное хранилище или из него, например, в файловую систему или базу в данных, чтобы не ухудшать скорость отклика.

В простейшем случае, данные в модели данных вводятся пользователем или загружаются статически из файла или другого источника данных при запуске приложения, и в этом случае к данным никогда не обращаются потоки, отличные от потока событий. Но иногда объект модели представления (*presentation model object*) выступает только представлением (*view*) другого источника данных, например базы данных, файловой системы или удаленной службы. В этом случае более чем один поток может получить доступ к данным в момент входа или в момент выхода из приложения.

Например, можно отобразить содержимое удаленной файловой системы с помощью древовидного элемента управления. Вы не хотели бы перечислять всю файловую систему, прежде чем вы смогли бы отобразить древовидный элемент управления - это займет слишком много времени и памяти. Вместо этого, дерево можно заполнять по мере раскрытия узлов. Перечисление даже одного каталога на удаленном томе может занять довольно много времени, поэтому может потребоваться выполнить перечисление в фоновой задаче. После завершения фоновой задачи, необходимо каким-то образом поместить данные в модель дерева. Это можно сделать с помощью потокобезопасной модели дерева, "проталкивая" данные из фоновой задачи в поток событий, путем публикации задачи с помощью метода `invokeLater` или путем опроса состояния потоком событий, чтобы увидеть, доступны ли данные.

9.4.1 Потокобезопасные модели данных

До тех пор, пока блокировка не оказывает чрезмерного влияния на отзывчивость, проблему нескольких потоков, работающих с данными, можно решить с помощью потокобезопасной модели данных. Если модель данных поддерживает разбиение на небольшие параллельные операции, поток событий и фоновые потоки должны иметь возможность совместно использовать её без проблем с откликом. Например, класс `DelegatingVehicleTracker` использует базовый класс `ConcurrentHashMap`, операции извлечения данных которого обеспечивают высокую степень параллелизма. Недостатком является то, что он не предлагает согласованный моментальный снимок данных (*snapshot of the data*), который может требоваться или не требоваться. Потокобезопасные модели данных также должны генерировать события при обновлении модели, чтобы представления могли обновляться при изменении данных.

Иногда возможно получить потокобезопасность, согласованность и хорошую отзывчивость с помощью версионной модели данных (*versioned data model* such), такой как класс `CopyOnWriteArrayList` [CPJ 2.2.3.3]. Когда вы захватываете итератор для коллекции скопировать-и-записать, итератор проходит по коллекции

в том виде, в котором она существовала на момент создания итератора. Однако коллекции скопировать-и-записать обеспечивают хорошую производительность только в том случае, если количество проходов значительно превышает количество модификаций, что, вероятно, не будет иметь место, например, в приложении отслеживания транспортных средств. Более специализированные версионные структуры данных могут избежать этого ограничения, но создание версионных структур данных, которые обеспечивают эффективный параллельный доступ и не сохраняют старые версии данных дольше, чем это необходимо, непросто, и поэтому такой вариант следует рассматривать только тогда, когда другие подходы непрактичны.

9.4.2 Разделение моделей данных

С точки зрения GUI, классы табличной модели Swing, такие как `TableModel` и `TreeModel`, являются официальным репозиторием для отображаемых данных. Однако объекты модели часто сами являются "представлениями" других объектов, управляемых приложением. О программе, которая имеет как домен представления (*presentation-domain*), так и домена приложения (*applicationdomain*), говорят, что она спроектирована в виде разделённой модели (*split-model*, Fowler, 2005).

В дизайне с разделённой моделью, модель представления ограничена потоком событий и другой моделью - совместно используемой моделью (*shared model*), которая является потокобезопасной и доступ к ней может осуществляться как из потока событий, так и из потоков приложения. Модель представления регистрирует слушателя в совместно используемой модели, чтобы получать уведомления об обновлениях. Затем модель представления можно обновить из совместно используемой модели путем встраивания моментального снимка соответствующего состояния в сообщение об обновлении или путем получения моделью представления данных непосредственно из совместно используемой модели, при получении события обновления.

Подход с моментальными снимками прост, но имеет ограничения. Он хорошо работает, когда модель данных мала, обновления не слишком часты, и структура обеих моделей подобна. Если модель данных имеет большой размер или обновления происходят очень часто, или если одна или обе стороны разделения содержат информацию, которая не видна другой стороне, более эффективным подходом может быть отправка инкрементных обновлений вместо полных снимков. Этот подход приводит к эффекту сериализации обновлений в совместно используемой модели и воссоздании их в потоке событий для модели представления. Другим преимуществом инкрементных обновлений является то, что более детальная информация о том, что изменилось, может улучшить качество восприятия информации с дисплея - если движется только один автомобиль, нам не нужно перекрашивать весь дисплей, а только те регионы, что были затронуты обновлениями.

Рассмотрите возможность проектирования с разделённой моделью, когда модель данных должна совместно использоваться несколькими потоками, и реализация потокобезопасной модели данных нецелесообразна по причинам блокировки, согласованности или сложности.

9.5 Другие формы однопоточных подсистем

Ограничение потока не ограничивается GUI: его можно использовать всякий раз, когда объект реализуется как однопоточная подсистема. Иногда ограничение потока принудительно вводится разработчиком по причинам, которые не имеют ничего общего с желанием избежать синхронизации или взаимоблокировок. Например, некоторые нативные (*native*) библиотеки требуют, чтобы весь доступ к библиотеке, даже загрузка библиотеки с помощью вызова `System.loadLibrary`, осуществлялся из одного и того же потока.

Заимствуя подход, принятый в GUI фреймворках, для доступа к нативной библиотеке вы можете легко создать выделенный поток или однопоточный исполнитель и предоставить прокси-объект, который будет перехватывать вызовы, поступающие к ограниченному потоком объекту, и отправлять их в качестве задач выделенному потоку. Для упрощения реализации такого подхода, экземпляр `Future` и метод `newSingleThreadExecutor` работают в связке; прокси-метод может отправить задачу и немедленно вызвать метод `Future.get` для ожидания результата. (Если ограниченный потоком класс реализует интерфейс, можно автоматизировать процесс отправки экземпляра `Callable` фоновому исполнителю потока и ожидать получение результата с помощью динамического прокси.)

9.6 Итоги

Фреймворки GUI почти всегда реализуются как однопоточные подсистемы, в которых весь связанный с представлением код выполняется в качестве задач в потоке событий. Поскольку существует только один поток событий, долговременные задачи могут существенно снижать скорость отклика и поэтому должны выполняться в фоновых потоках. Вспомогательные классы, такой как `SwingWorker` или приведённый в главе класс `BackgroundTask`, обеспечивающие поддержку отмены, индикации прогресса и индикации завершения, могут упростить разработку долговременных задач, имеющих как GUI, так и компоненты без GUI.

**Часть 3 III Живучесть,
производительность и тестирование**

Глава 10 Предотвращение возникновения угроз живучести

Между требованиями безопасности и живучести часто возникает противоречие. Мы используем блокировку для обеспечения потокобезопасности, но беспорядочное использование блокировок может привести к взаимоблокировкам, вызванным порядком блокировок (*lock-ordering deadlocks*). Точно так же мы используем пулы потоков и семафоры для ограничения потребления ресурсов, но непонимание механизмов ограничения активностей может привести к взаимоблокировкам ресурсов (*resource deadlocks*). Приложения Java не восстанавливаются после возникновения взаимоблокировок, поэтому стоит убедиться, что ваш проект исключает условия, которые могут привести к их появлению. В этой главе рассматриваются некоторые из причин возникновения сбоев живучести и подходы для их предотвращения.

10.1 Взаимоблокировки

Взаимоблокировка иллюстрируется классической, хотя и немного антисанитарной, проблемой “обедающих философов”. Пять философов выходят за китайской едой и садятся за круглый стол. Есть пять палочек для еды (не пять пар), по одной между каждой парой обедающих. Философы чередуют процессы размышления и приёма пищи. Каждый из них должен захватить две палочки для еды достаточно надолго, но затем может положить палочки обратно и вернуться к размышлению. Существует несколько алгоритмов управления доступностью палочек, результатом применения которых является то, что каждый философ ест более или менее своевременно (голодный философ пытается схватить обе соседние палочки, но если одна уже используется, кладёт обратно ту, которую взял, и ждет минуту или около того, прежде чем вновь попытаться взять палочки), это может привести к тому, что некоторые или все философы умрут от голода (каждый философ немедленно хватается за палочку слева от него и ожидает, пока станет доступна палочка справа, прежде чем вернуть палочку слева). В последующей ситуации, при которой каждый из них имеет ресурс, необходимый другому, и ожидает возможности захвата ресурса, удерживаемого другим, и не освободит тот, который у него есть, пока не приобретёт тот, которого у него нет, иллюстрирует принцип взаимоблокировки.

Когда поток удерживает блокировку навсегда, другие потоки, пытающиеся получить ту же блокировку, также блокируются навсегда. Когда поток *A* удерживает блокировку *L* и пытается получить блокировку *M*, но в то же время поток *B* удерживает блокировку *M* и пытается получить блокировку *L*, оба потока будут ждать вечно. Приведённая ситуация является самым простым случаем взаимоблокировки (или *смертельного объятия*, *deadly embrace*), когда несколько потоков находятся в вечном ожидании из-за циклической зависимости блокировок. (Представим, что потоки - это узлы ориентированного графа, ребра которого представляют собой отношение “поток *A* ожидает ресурс, удерживаемый потоком *B*”. Если получившийся график цикличен, значит, произошла взаимоблокировка.)

Системы баз данных проектируются для обнаружения состояния взаимоблокировки и восстановления при выходе из него. Транзакция может затребовать наложения множества блокировок, и блокировки будут удерживаться до момента фиксации транзакции. Таким образом, вполне возможна ситуация, и на

самом деле не является редкостью, при которой две транзакции попадают в состояние взаимоблокировки. Без внешнего вмешательства они будут вечно ожидать завершения друг друга (удерживая блокировки на ресурсах, которые, вероятно, требуются и другим транзакциям). Но сервер баз данных не позволит этому случиться. Когда сервер обнаруживает, что набор транзакций находится в состоянии взаимоблокировки (это осуществляется путем поиска циклов в графе ожидания (*is-waiting-for*)), он выбирает жертву и прерывает выбранную транзакцию. Это приводит к освобождению блокировок, удерживаемых жертвой, позволяя другим транзакциям продолжить своё выполнение. Позже, после завершения конкурирующих транзакций, приложение может повторить выполнение прерванной транзакции.

В отличие от серверов баз данных, JVM не оказывает помощи в разрешении ситуаций возникновения взаимоблокировок. Когда набор Java потоков попадает в состояние взаимоблокировки, это конец игры - эти потоки мгновенно выходят из строя. В зависимости от того, что делают эти потоки, приложение может полностью зависнуть или может зависнуть определенная подсистема, или может пострадать производительность. Единственный способ восстановить работоспособность приложения, это прервать его выполнение и перезапустить - и надеяться, что то же самое впредь не повторится.

Как и многие другие угрозы параллелизма, взаимоблокировки редко проявляются сразу. Тот факт, что класс имеет потенциальные взаимоблокировки, не означает, что он всегда *будет* попадать в состояние взаимоблокировки, а только то, что это может произойти. Когда взаимоблокировки все-таки проявляют себя, то зачастую это случается в самый неподходящий момент - под большой производственной нагрузкой.

10.1.1 Взаимоблокировки, вызванные порядком наложения блокировок

Класс `LeftRightDeadlock`, представленный в листинге 10.1, подвержен риску возникновения взаимоблокировки. Каждый из методов `leftRight` и `rightLeft` пытается захватить блокировки на объектах `left` и `right`. Если один поток вызывает метод `leftRight`, а другой поток вызывает метод `rightLeft`, и их действия чередуются, как показано на рисунке 10.1, они попадут в состояние взаимоблокировки.

```
// Warning: deadlock-prone!
public class LeftRightDeadlock {
    private final Object left = new Object();
    private final Object right = new Object();

    public void leftRight() {
        synchronized (left) {
            synchronized (right) {
                doSomething();
            }
        }
    }

    public void rightLeft() {
        synchronized (right) {

```



```

        synchronized (left) {
            doSomethingElse();
        }
    }
}

```

Листинг 10.1 Простая взаимоблокировка, вызванная порядком наложения блокировок. Не делайте так.

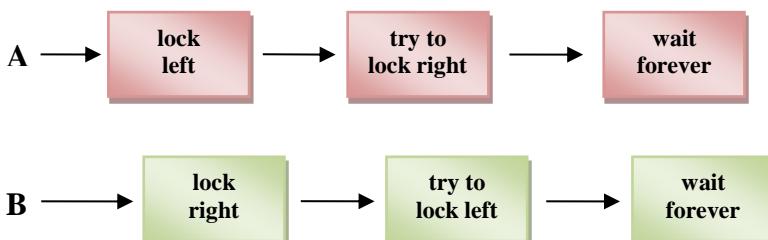


Рисунок 10.1 Неудачный момент времени для класса LeftRightDeadlock

Взаимоблокировка в классе `LeftRightDeadlock` произошла потому, что два потока пытались захватить блокировки в *различном порядке* (*different order*). Если бы они запрашивали блокировки в одном и том же порядке, не было бы никакой циклической зависимости между блокировками и, следовательно, никакой взаимоблокировки бы не произошло. Если вы сможете гарантировать, что каждый поток, которому одновременно нужны блокировки *L* и *M*, всегда захватывает *L* и *M* в одном и том же порядке, взаимоблокировки возникать не будут.

Программа будет свободна от взаимоблокировок вызванных порядком наложения блокировок, если все потоки будут получать необходимые им блокировки в глобально зафиксированном порядке.

Проверка согласованности порядка блокировок требует глобального анализа поведения блокировок в программе. Недостаточно проверить ветки кода, которые получают несколько блокировок по отдельности; и метод `leftRight` и метод `rightLeft` представляют собой “разумные” способы получения двух блокировок, они просто не совместимы. Когда дело доходит до блокировки, левая рука должна знать, что делает правая рука.

10.1.2 Взаимоблокировки, вызванные динамическим порядком блокировок

Иногда не очевидно, что у вас достаточно контроля над порядком наложения блокировок, для предотвращения взаимоблокировки. Рассмотрим выглядящий безобидно код, представленный в листинге 10.2, который переводит средства с одного счета на другой. Он получает блокировки на обоих объектах `Account` перед выполнением перемещения, гарантируя, что балансы обновляются атомарно и без нарушения инвариантов, подобных утверждению “счет не может иметь отрицательный баланс”.

// Warning: deadlock-prone!




```

        throws InsufficientFundsException {
    class Helper {
        public void transfer() throws InsufficientFundsException {
            if (fromAcct.getBalance().compareTo(amount) < 0)
                throw new InsufficientFundsException();
            else {
                fromAcct.debit(amount);
                toAcct.credit(amount);
            }
        }
    }
    int fromHash = System.identityHashCode(fromAcct);
    int toHash = System.identityHashCode(toAcct);

    if (fromHash < toHash) {
        synchronized (fromAcct) {
            synchronized (toAcct) {
                new Helper().transfer();
            }
        }
    } else if (fromHash > toHash) {
        synchronized (toAcct) {
            synchronized (fromAcct) {
                new Helper().transfer();
            }
        }
    } else {
        synchronized (tieLock) {
            synchronized (fromAcct) {
                synchronized (toAcct) {
                    new Helper().transfer();
                }
            }
        }
    }
}

```

Листинг 10.3 Индуцирование упорядочивания блокировок для исключения возможности возникновения взаимоблокировки

В тех редких случаях, когда два объекта имеют один и тот же хэш-код, мы должны использовать произвольные средства упорядочения захвата блокировок, и это вновь приводит к возможности возникновения взаимоблокировки. Чтобы в этой ситуации избежать несогласованности в порядке захвата блокировок, используется третья “разрывающая связь” (*tie breaking*) блокировка. Захватывая блокировку разрывающую связь до получения любой из блокировок на объекте `Account`, мы гарантируем, что только один поток одновременно выполняет рискованную задачу получения двух блокировок в произвольном порядке, исключая возможность взаимоблокировки (если этот механизм используется согласованно). Если бы конфликты хэша были распространены, этот метод мог бы стать узким местом параллелизма (так же, как и наличие единственной блокировки), но в связи с тем,

что коллизии хэша, возвращаемого методом `System.identityHashCode` исчезающие редки, этот метод обеспечивает наличие последнего бита безопасности при небольших затратах.

Если класс `Account` имеет уникальный, неизменяемый, сопоставимый ключ, такой как номер учетной записи, индуцирование упорядочения наложения блокировки еще проще: упорядочивайте объекты по их ключу, тем самым устранив необходимость в использовании блокировки, разрывающей связь.

Вы можете решить, что мы завышаем риск возникновения взаимоблокировки, потому что блокировки обычно удерживаются на очень краткий промежуток времени, но взаимоблокировки являются серьезной проблемой в реальных системах. Производственное приложение может выполнять в день миллиарды циклов захвата-освобождения блокировок. Достаточно одной из них быть установленной неправильно, чтобы приложение попало в ситуацию взаимоблокировки, и даже тщательный режим нагружочного тестирования не может выявить все латентные взаимоблокировки¹⁰⁴. Класс `DemonstrateDeadlock` представленный в листинге 10.4¹⁰⁵, довольно быстро попадает в состояние взаимоблокировки на большинстве систем.

```
public class DemonstrateDeadlock {  
    private static final int NUM_THREADS = 20;  
    private static final int NUM_ACCOUNTS = 5;  
    private static final int NUM_ITERATIONS = 1000000;  
  
    public static void main(String[] args) {  
        final Random rnd = new Random();  
        final Account[] accounts = new Account[NUM_ACCOUNTS];  
  
        for (int i = 0; i < accounts.length; i++)  
            accounts[i] = new Account();  
  
        class TransferThread extends Thread {  
            public void run() {  
                for (int i=0; i<NUM_ITERATIONS; i++) {  
                    int fromAcct = rnd.nextInt(NUM_ACCOUNTS);  
                    int toAcct = rnd.nextInt(NUM_ACCOUNTS);  
                    DollarAmount amount =  
                        new DollarAmount(rnd.nextInt(1000));  
                    transferMoney(accounts[fromAcct],  
                                  accounts[toAcct], amount);  
                }  
            }  
        }  
        for (int i = 0; i < NUM_THREADS; i++)  
            new TransferThread().start();  
    }  
}
```

¹⁰⁴ По иронии судьбы, удержание блокировок в течение коротких периодов времени, что, как вы предполагаете, должно привести к уменьшению конкуренции между блокировками, увеличивает вероятность того, что тестирование не раскроет латентные риски возникновения взаимоблокировки.

¹⁰⁵ В целях упрощения, класс `DemonstrateDeadlock` игнорирует проблему отрицательного баланса на счёте.

```
}
```

Листинг 10.4 Управляющий цикл, индуцирующий взаимоблокировку при типичных условиях

10.1.3 Взаимоблокировки между взаимодействующими объектами

Захват нескольких блокировок далеко не всегда так очевиден, как в методах `LeftRightDeadlock` или `transferMoney`; две блокировки не обязательно должны быть захвачены одним и тем же методом. Рассмотрим взаимодействующие классы из листинга 10.5, которые могут быть использованы в приложении диспетчеризации машин такси. Класс `Taxi` представляет индивидуальное такси, с указанием местоположения и пункта назначения, класс `Dispatcher` представляет парк такси.

```
// Warning: deadlock-prone!
class Taxi {
    @GuardedBy("this") private Point location, destination;
    private final Dispatcher dispatcher;

    public Taxi(Dispatcher dispatcher) {
        this.dispatcher = dispatcher;
    }

    public synchronized Point getLocation() {
        return location;
    }

    public synchronized void setLocation(Point location) {
        this.location = location;
        if (location.equals(destination))
            dispatcher.notifyAvailable(this);
    }
}

class Dispatcher {
    @GuardedBy("this") private final Set<Taxi> taxis;
    @GuardedBy("this") private final Set<Taxi> availableTaxis;

    public Dispatcher() {
        taxis = new HashSet<Taxi>();
        availableTaxis = new HashSet<Taxi>();
    }

    public synchronized void notifyAvailable(Taxi taxi) {
        availableTaxis.add(taxi);
    }

    public synchronized Image getImage() {
        Image image = new Image();
```



```
        for (Taxi t : taxis)
            image.drawMarker(t.getLocation());
        return image;
    }
}
```

Листинг 10.5 Взаимоблокировка, вызванная порядком наложения блокировок между взаимодействующими объектами. Не делайте так.

Хотя ни один метод явно не захватывает две блокировки, объекты, вызывающие методы `setLocation` и `getImage` могут всё же могут попытаться захватить две одинаковых блокировки. Если поток вызывает метод `setLocation` в ответ на обновление, пришедшее от приемника GPS, он сначала обновляет местоположение такси, а затем проверяет, достигло ли оно места назначения. Если это так, он сообщает диспетчеру, что ему нужен новый пункт назначения. Так как оба метода - `setLocation` и `notifyAvailable` – объявлены как `synchronized`, поток, вызывающий метод `setLocation`, захватывает блокировку экземпляра `Taxi` и затем блокировку экземпляра `Dispatcher`. Аналогично, поток, вызывающий метод `getImage`, захватывает блокировку экземпляра `Dispatcher`, а затем блокировку на каждом экземпляре `Taxi` (по одному). Также как в классе `LeftRightDeadLock`, две блокировки будут захвачены двумя потоками в различном порядке, что приводит к риску возникновения взаимоблокировки.

В методах `LeftRightDeadlock` и `transferMoney` было достаточно легко обнаружить возможность возникновения взаимоблокировки, просто найдя методы, захватывающие две блокировки. Обнаружить возможность возникновения взаимоблокировки в классах `Taxi` и `Dispatcher` немного сложнее: предупреждающим знаком может служить то, что чужой метод вызывается, пока удерживается блокировка.

Вызов чужого метода с удерживаемой блокировкой напрашивается на проблемы с живучестью. Чужой метод может захватывать другие блокировки (рискуя вызвать взаимоблокировку) или заблокироваться на неожиданно долгое время, вынуждая останавливаться другие потоки, которым нужна блокировка, которую вы удерживаете.

10.1.4 Открытые вызовы

Конечно, классы `Taxi` и `Dispatcher` не знали, что каждый из них был половиной взаимоблокировки, ожидающей своего часа. И они не должны этого делать; вызов метода – это абстрактный барьер, предназначенный для защиты вас от знания деталей того, что происходит на другой стороне. Но, поскольку вы не знаете, что происходит на другой стороне вызова, вызов чужого метода в процессе удержания блокировки сложно проанализировать, и поэтому он довольно рискован.

Вызов метода без удерживаемых блокировок называется *открытым вызовом* (*open call*)[CPJ 2.4.1.3], а классы, использующие открытые вызовы, более корректны и компонуемы, чем классы, выполняющие вызовы с удерживаемыми блокировками. Использование открытых вызовов для того, что избежать возникновения взаимоблокировок, аналогично использованию инкапсуляции для обеспечения потокобезопасности: хотя можно, конечно, создать потокобезопасную программу без инкапсуляции, анализ потокобезопасности программы, которая

эффективно использует инкапсуляцию, намного проще, чем той, которая этого не делает. Аналогичным образом, анализ живучести программы, которая полагается исключительно на открытые вызовы, намного проще, чем той, которая этого не делает. Самоограничение использованием открытых вызовов позволяет значительно упростить определение веток кода, которые получают множество блокировок, и, следовательно, обеспечивает согласованный порядок получения блокировок¹⁰⁶.

Классы `Taxi` и `Dispatcher`, приведённые в листинге 10.5 можно легко отрефакторить с использованием открытых вызовов, и, таким образом, устраниТЬ риск возникновения взаимоблокировки. Это влечёт за собой сокращение блоков `synchronized` только для защиты операций затрагивающих совместно используемое состояние, как показано в листинге 10.6. Очень часто причиной проблем, подобных описанным в листинге 10.5, является использование синхронизированных методов вместо меньших, по размеру, синхронизированных блоков по причине компактного синтаксиса или простоты, а не потому, что весь метод должен быть защищён блокировкой. (В качестве бонуса, сокращение синхронизированного блока также может улучшить масштабируемость; за более подробными сведениями о размере синхронизируемых блоков см. раздел [11.4.1](#).)

```
@ThreadSafe
class Taxi {
    @GuardedBy("this") private Point location, destination;
    private final Dispatcher dispatcher; ...

    public synchronized Point getLocation() {
        return location;
    }

    public void setLocation(Point location) {
        boolean reachedDestination;
        synchronized (this) {
            this.location = location;
            reachedDestination = location.equals(destination);
        }
        if (reachedDestination)
            dispatcher.notifyAvailable(this);
    }
}

@ThreadSafe
class Dispatcher {
    @GuardedBy("this") private final Set<Taxi> taxis;
    @GuardedBy("this") private final Set<Taxi> availableTaxis;
    ...
    public synchronized void notifyAvailable(Taxi taxi) {
        availableTaxis.add(taxi);
    }
}
```

¹⁰⁶ Необходимость полагаться на открытые вызовы и тщательный порядок блокировок отражает фундаментальный беспорядок в компоновке синхронизированных объектов, а не в синхронизации компонуемых объектов.

```
}

public Image getImage() {
    Set<Taxi> copy;
    synchronized (this) {
        copy = new HashSet<Taxi>(taxis);
    }
    Image image = new Image();
    for (Taxi t : copy)
        image.drawMarker(t.getLocation());
    return image;
}
}
```

Листинг 10.6 Использование открытых вызовов для исключения взаимоблокировок между взаимодействующими объектами

Старайтесь в своей программе использовать открытые вызовы. Программы, которые полагаются на открытые вызовы, гораздо легче анализировать на отсутствие взаимоблокировок, чем те, которые допускают вызов чужих методов в процессе удержания блокировки.

Реструктуризация блока `synchronized`, с целью разрешения открытых вызовов, иногда может иметь нежелательные последствия, так как она берёт атомарную операцию и делает её не атомарной. Во многих случаях потеря атомарности вполне приемлема; нет никаких причин, по которым обновление местоположения такси и уведомление диспетчера о том, что оно готово получить новое место назначения, должны быть атомарными операциями. В других случаях потеря атомарности заметна, но семантические изменения все еще приемлемы. В версии, подверженной взаимоблокировкам, метод `getImage`, в момент вызова, создает полный снимок местоположений парка автомобилей такси; в переработанной версии, метод `getImage` получает местоположение каждого такси в немного различающиеся моменты времени.

Однако, в некоторых случаях, потеря атомарности является проблемой, и в этой ситуации вам придется использовать другой подход для достижения атомарности. Одним из таких подходов является структурирование параллельного объекта таким образом, чтобы только один поток мог выполнить ветку кода, следующую за открытым вызовом. Например, при завершении работы службы может потребоваться дождаться завершения текущих операций, а затем освободить ресурсы, используемые службой. Процесс удержания блокировки службы во время ожидания завершения операций, по самой своей сути, подвержен взаимоблокировке, но снятие блокировки со службы до завершения работы службы может позволить другим потокам начать новые операции. Решение состоит в том, чтобы удерживать блокировку достаточно долго для того, чтобы установить флаг состояния службы в “завершение работы”, и чтобы другие потоки, желающие начать новые операции - включая завершение работы службы - видели, что служба недоступна, и не предпринимали подобных попыток. Затем, после завершения выполнения открытого вызова, можно дождаться завершения работы, зная, что только поток завершения работы имеет доступ к состоянию службы. Таким образом, вместо того, чтобы использовать блокировку, для удержания других

потоков вне критической части кода, этот подход основан на построении таких протоколов, при которых другие потоки не будут пытаться в неё войти.

10.1.5 Взаимоблокировки ресурсов

Точно так же, как потоки могут попасть в ситуацию взаимоблокировки, когда каждый из них ожидает блокировку, удерживаемую другим потоком, и не освобождает свою, они также могут попасть в ситуацию взаимоблокировки при ожидании ресурсов.

Предположим, у вас есть два объединяющих ресурса, например пулы соединений для двух разных баз данных. Пулы ресурсов, как правило, реализуются с использованием семафоров (см. раздел [5.5.3](#)) для облегчения блокирования, когда пул пуст. Если задача требует подключения к обеим базам данных и оба ресурса не всегда запрашиваются в одном и том же порядке, поток *A* может удерживать подключение к базе *D₁*, при ожидании подключения к базе *D₂*, а поток *B* может удерживать подключение к базе *D₂*, при ожидании подключения к базе данных *D₁*. (Чем больше размер пулов, тем реже это происходит; если каждый пул имеет *N* соединений, для возникновения взаимоблокировки требуется *N* наборов циклически ожидающих потоков и множество неудачных моментов времени.)

Другой формой взаимоблокировки на основе ресурсов является *взаимоблокировка, вызванная голоданием потоков* (*thread-starvation deadlock*). Мы видели пример реализации этой угрозы в разделе 8.1.1, в котором задача, отправляющая другую задачу и ожидающая от неё результата, выполняется однопоточным экземпляром *Executor*. В этом случае, первая задача будет ожидать вечно, мгновенно останавливая эту задачу и все остальные, ожидающие выполнения в том же экземпляре *Executor*. Задачи, ожидающие результатов других задач, являются основным источником взаимоблокировки, вызванной голоданием потоков; ограниченные пулы и взаимозависимые задачи сочетаются плохо.

10.2 Предотвращение и диагностика взаимоблокировок

Программа, которая никогда не получает более одной блокировки за один раз, не может столкнуться с взаимоблокировкой, вызванной порядком захвата блокировок. Конечно, это не всегда практически целесообразно, но если вам это сойдет с рук, в результате будет намного меньше работы. Если вам необходимо захватить несколько блокировок, упорядочение блокировок должно быть частью дизайна программы: попробуйте свести к минимуму число потенциальных взаимодействий между блокировками, а также следуйте протоколу упорядочения блокировок и документируйте его для блокировок, которые могут быть захвачены совместно.

В программах, использующих детальную блокировку, выполняйте аудит кода на отсутствие взаимоблокировок, используя стратегию, состоящую из двух частей: сначала определите места, в которых можно получить несколько блокировок (попробуйте сделать собрать небольшой набор), а затем выполните глобальный анализ всех таких экземпляров, чтобы обеспечить согласованность порядка блокировок во всей программе. Использование открытых вызовов там, где это возможно, существенно упрощает проведение такого анализа. При отсутствии не открытых вызовов, найти экземпляры, в которых получено несколько блокировок,

довольно легко, либо путем ревью кода (*code review*), либо с помощью автоматического анализа байт-кода или исходного кода.

10.2.1 Блокировки, ограниченные по времени

Другим подходом, применяемым для обнаружения взаимоблокировок и восстановления при их возникновении, является использование ограниченной по времени версии метода `tryLock`, предоставляемой явными классами блокировок `Lock` (см. главу 13), вместо внутренней блокировки. В той ситуации, когда внутренние блокировки ожидают вечно, если не могут получить блокировку, явные блокировки позволяют указать время ожидания, по истечении которого метод `tryLock` возвращает ошибку. При использовании тайм-аута, значение которого превышает время, в течение которого вы ожидаете захватить блокировку, в случае, если что-то неожиданно пойдёт не так, вы сможете восстановить контроль. (В листинге 13.3 приведён альтернативный вариант реализации метода `transferMoney`, с использованием опрашиваемой версии `tryLock`, повторяющей попытки для избегания вероятного возникновения взаимоблокировки.)

Когда ограниченная по времени попытка захвата блокировки провалится, вам не обязательно знать, *почему это произошло*. Может быть, произошла взаимоблокировка; может быть, поток по ошибке вошел в бесконечный цикл, пока удерживал эту блокировку; или, может быть, просто какая-то активность работает намного медленнее, чем вы того ожидали. Но, по крайней мере, у вас есть возможность оставить запись о том, что попытка не удалась, логировать любую полезную информацию о том, что вы пытались сделать, и изящно перезапустить вычисление, вместо того, чтобы убить весь процесс.

Использование ограниченного по времени захвата блокировки, для захвата нескольких блокировок, может быть эффективным средством против возникновения взаимоблокировки, даже если блокировка по времени используется во всей программе не согласованно. Если время захвата блокировки истекло, вы можете освободить блокировки, отступить и подождать некоторое время, а затем повторить попытку, возможно, очистив условие взаимоблокировки и позволив программе восстановиться. (Этот подход работает только тогда, когда две блокировки захватываются вместе; если несколько блокировок захватываются во вложенных вызовах методов, вы не можете просто освободить внешнюю блокировку, даже если вы знаете о том, что удерживаете ее вы.)

10.2.2 Анализ взаимоблокировок с использованием дампов потоков

Несмотря на то, что предотвращение взаимоблокировок является, в основном, вашей проблемой, JVM может помочь идентифицировать их, когда они происходят, с помощью использования *дампа потоков*. Дамп потока содержит трассировку стека для каждого запущенного потока, аналогичную трассировке стека, сопровождающей исключение. Дампы потоков также включают в себя информацию о блокировании, например о том, какие блокировки удерживаются каждым потоком, в каких кадрах стека они были захвачены, и какую блокировку ожидает захватить заблокированный поток.¹⁰⁷ Перед созданием дампа потока, JVM просматривает граф ожидания (*is-waiting-for graph*) в поисках циклов, для

¹⁰⁷ Эта отладочная информация полезна даже в том случае, если у вас нет взаимоблокировок; периодический сброс дампов потоков позволяет наблюдать за поведением блокировки в программе.

определения наличия взаимоблокировок. Если среда JVM находит цикл, она включает информацию о взаимоблокировке, идентифицирующую, какие блокировки и потоки участвуют, где в программе нарушен захват блокировки.

Для инициирования сброса дампа потока, можно отправить процессу JVM сигнал SIGQUIT (`kill -3`) на платформах Unix или нажать комбинацию клавиш `Ctrl-\` в Unix или `Ctrl-Break` на платформах Windows. Многие IDE также могут запросить дамп потока.

Если вы используете явные классы блокировок `Lock`, вместо встроенной блокировки, среда Java 5.0 не поддерживает связывание информации класса `Lock` с дампом потока; явные блокировки вообще не отображаются в дампах потоков. Среда Java 6 включает поддержку дампа потока и обнаружение взаимоблокировки с явными блокировками класса `Lock`, но информация о том, где блокировки были захвачены, будет неизбежно менее точна, чем для встроенных блокировок. Встроенные блокировки связаны с кадром стека, в котором они были захвачены; явные блокировки `Lock` связаны только с захватившим их потоком.

В листинге 10.7 приведена часть дампа потока, взятая из рабочего приложения J2EE. Сбой, вызвавший взаимоблокировку, включает три компонента - приложение J2EE, контейнер J2EE и драйвер JDBC, каждый из них принадлежит разным поставщикам. (Имена были изменены, в целях защиты виновных.) Все трое были коммерческими продуктами, которые прошли через интенсивные циклы тестирования; каждый из них имел ошибку, которая сама по себе была безвредна, пока все они не начинали взаимодействовать друг с другом и не вызывали фатальный сбой сервера.

```
Found one Java-level deadlock:
=====
"ApplicationServerThread":
    waiting to lock monitor 0x080f0cdc (a MumbleDBConnection),
    which is held by "ApplicationServerThread"
"ApplicationServerThread":
    waiting to lock monitor 0x080f0ed4 (a MumbleDBCallableStatement),
    which is held by "ApplicationServerThread"

Java stack information for the threads listed above:
"ApplicationServerThread":
    at MumbleDBConnection.remove_statement
    - waiting to lock <0x650f7f30> (a MumbleDBConnection)
    at MumbleDBStatement.close
    - locked <0x6024ffb0> (a MumbleDBCallableStatement)
    ...
"ApplicationServerThread":
    at MumbleDBCallableStatement.sendBatch
    - waiting to lock <0x6024ffb0> (a MumbleDBCallableStatement)
    at MumbleDBConnection.commit
    - locked <0x650f7f30> (a MumbleDBConnection)
    ...
```

Листинг 10.7 Часть дампа потока, после наступления взаимоблокировки

Мы показали только ту часть дампа потока, что относится к идентификации взаимоблокировки. Среда JVM проделала для нас большую работу, в процессе

диагностики взаимоблокировки - какие блокировки вызывают проблему, какие потоки участвуют, какие другие блокировки они удерживают, и причиняются ли другим потокам косвенные неудобства. Один поток удерживает блокировку на объекте `MumbleDBConnection` и ожидает получения блокировки на объекте `MumbleDBCallableStatement`; другой удерживает блокировку на объекте `MumbleDBCallableStatement` и ожидает захвата блокировки на объекте `MumbleDBConnection`.

Драйвер JDBC, используемый здесь, имеет явную ошибку в порядке захвата блокировок: различные цепочки вызовов захватывают несколько блокировок в различном порядке, с использованием драйвера JDBC. Но эта проблема не проявилась бы, если бы не другая ошибка: несколько потоков пытались одновременно использовать один и тот же экземпляр `JDBC Connection`. Приложение работало совсем не так, как ожидалось - разработчики были удивлены, увидев один и то же экземпляр `Connection`, используемый одновременно двумя потоками. В спецификации JDBC нет явного требования о том, чтобы реализация `Connection` была потокобезопасной, и в обычной практике использование экземпляра `Connection` ограничивается одним потоком, как предполагалось и в этом случае. В свою очередь, поставщик попытался предоставить потокобезопасный драйвер JDBC, о чем свидетельствует синхронизация множества объектов JDBC в коде драйвера. К сожалению, поскольку поставщик не принял во внимание порядок захвата блокировок, драйвер был подвержен взаимоблокировке, но проблема обнаруживала себя только при взаимодействии драйвера, подверженного взаимоблокировке, и неправильном совместном использовании экземпляра `Connection` в приложении. Поскольку ни одна ошибка, в изоляции одна от другой, не была фатальной, обе сохранялись, несмотря на обширное тестирование.

10.3 Прочие угрозы живучести

В то время как взаимоблокировка представляет собой наиболее широко встречающуюся угрозу живучести, существует несколько других угроз живучести, с которыми вы можете столкнуться в параллельных программах, включая голодание (*starvation*), пропущенные сигналы (*missed signals*), и динамическую взаимоблокировку (*livelock*). (Пропущенные сигналы рассматриваются в разделе 14.2.3).

10.3.1 Голодание

Голодание возникает тогда, когда потоку постоянно отказывают в доступе к ресурсам, в которых он нуждается для достижения прогресса; наиболее подверженный голоданию ресурс - циклы ЦП. Голодание в приложениях Java может быть вызвано неправильным использованием приоритетов потоков. Оно также может быть вызвано выполнением непрерывающихся конструкций (бесконечные циклы или ожидания ресурсов, которые не завершаются) с удерживаемой блокировкой, в связи с чем, другие потоки, которым нужна эта блокировка, никогда не смогут ее захватить.

Приоритеты потоков, определенные в API потоков, представляют собой просто советы по планированию. В API потоков определено десять уровней приоритета, которые среда JVM, по своему усмотрению, может сопоставить приоритетам планирования операционной системы. Это сопоставление является специфичным для платформы, таким образом, два приоритета Java могут быть сопоставлены с

одним и тем же приоритетом ОС в одной системе, и другим приоритетам ОС в другой. Некоторые операционные системы имеют меньше чем десять уровней приоритета, и в этом случае несколько приоритетов Java сопоставляется с тем одним и тем же приоритетом ОС.

Планировщики операционной системы идут на многое, чтобы обеспечить справедливость планирования и живучесть, сверх того, что требуется спецификацией языка Java. В большинстве приложений Java все потоки приложений имеют одинаковый приоритет, `Thread.NORM_PRIORITY`. Механизм приоритетов потоков является довольно грубым инструментом, и не всегда очевидно, к какому эффекты приведёт изменение приоритетов; повышение приоритета потока может не оказаться никаких изменений, или может привести к тому, что один поток будет всегда планироваться в предпочтении к другим, тем самым провоцируя голодание.

Как правило, разумной политикой будет сопротивление искушению настроить приоритеты потоков. Как только вы начинаете изменять приоритеты, поведение вашего приложения станет специфичным для платформы, и вы введёте в программу риск возникновения голодания. Вы часто можете определить программу, которая пытается восстановиться после настройки приоритета или других проблем с отзывчивостью, по наличию вызовов `Thread.sleep` или `Thread.yield` в странных местах, применяемых для того, чтобы дать больше времени потокам с более низким приоритетом¹⁰⁸.

Избегайте соблазна использовать приоритеты потоков, так как они увеличивают зависимость от платформы и могут вызвать проблемы с живучестью. Большинство параллельных приложений могут использовать приоритет по умолчанию для всех потоков.

10.3.2 Плохая отзывчивость

Отступая на один шаг от голодания, следует плохая отзывчивость, не являющаяся редкостью в приложениях GUI, использующих фоновые потоки. В главе 9 был разработан фреймворк, применяемый для разгрузки долговременных задач в фоновые потоки, чтобы не происходило “заморозки” пользовательского интерфейса. Фоновые задачи с интенсивным использованием ЦП могут по-прежнему влиять на скорость отклика, поскольку они могут конкурировать с потоком событий за циклы ЦП. Это один из случаев, когда изменение приоритетов потоков имеет смысл; когда ресурсоемкие фоновые вычисления повлияют на скорость отклика. Если работа, выполняемая другими потоками, действительно является фоновой, снижение их приоритета может сделать задачи переднего плана более отзывчивыми.

Плохая отзывчивость также может быть вызвана плохим управлением блокировками. Если поток удерживает блокировку в течение длительного времени (возможно, при итерации большой коллекции и выполнении существенной работы для каждого элемента), другим потокам, которым требуется доступ к той же самой коллекции, возможно, придется ожидать очень долго.

¹⁰⁸ Семантика вызовов `Thread.yield` (и `Thread.sleep(0)`) не определена [JLS 17.9]; среда JVM свободна в реализации их как “пустых операций” (*no-ops*) или рассмотрении их как подсказок планирования. В частности, они не обязаны реализовывать семантику вызова `sleep(0)` в Unix системах - поместить текущий поток в конец очереди выполнения потоков с таким же приоритетом, уступая другим потокам того же приоритета - хотя некоторые реализации JVM реализуют вызов `yield` именно таким образом.

10.3.3 Динамическая взаимоблокировка

Динамическая взаимоблокировка (*livelock*) является формой проблемы с живучестью, при которой поток, фактически оставаясь не заблокированным, все же не может добиться прогресса выполнения, потому что продолжает повторять выполнение операции, которая всегда будет терпеть неудачу. Динамическая взаимоблокировка часто возникает в транзакционных приложениях обмена сообщениями, в которых инфраструктура обмена сообщениями откатывает транзакцию, если сообщение не может быть обработано успешно, и помещает сообщение в начало очереди. Если ошибка в обработчике сообщений, приводит к сбою для определенного типа сообщений, каждый раз, когда сообщение будет помещаться в очередь и передаваться неисправному обработчику, будет выполняться откат транзакции. Так как сообщение будет возвращаться в начало очереди, обработчик будет вызываться снова и снова, с тем же результатом. (Такую ситуацию иногда называют проблемой *отравленного сообщения* (*poison message*).) Поток обработки сообщений не блокируется, но он также никогда не будет достигать прогресса. Эта форма динамической взаимоблокировки часто берёт своё начало с чрезмерно ретивого кода восстановления ошибок, который ошибочно рассматривает неустранимую ошибку как восстанавливаемую.

Динамические взаимоблокировки также могут возникать, когда несколько взаимодействующих потоков изменяют свое состояние в ответ на изменение состояния другими потоками, таким образом, что ни один поток никогда не сможет добиться прогресса. Это похоже на то, что происходит, когда два чрезмерно вежливых человека идут в противоположных направлениях в коридоре: каждый отступает в другую сторону, и теперь они снова на пути друг у друга. Так что они оба отступают снова, и снова, и снова...

Решение для приведённого разнообразия динамических взаимоблокировок заключается в том, чтобы ввести некоторую случайную составляющую в механизм повтора. Например, когда две станции в сети ethernet пытаются передать пакет на совместно используемом носителе¹⁰⁹ в одно и то же время, пакеты сталкиваются. Станции обнаруживают коллизию, и каждая из них пытается снова, позже, передать свой пакет. Если каждая из них будет повторять попытку ровно через одну секунду, пакеты будут сталкиваться снова и снова, и ни один пакет никогда не будет отправлен, даже если доступна большая часть пропускной способности. Чтобы избежать этого, мы заставляем каждую станцию ожидать некоторое время, которое включает случайную составляющую. (Протокол ethernet также включает экспоненциальный откат после повторяющихся столкновений, уменьшая как перегрузку, так и риск повторного сбоя с несколькими сталкивающимися пакетами от разных станций.) Повторная попытка, как со случайными задержками, так и с откатами, может быть одинаково эффективной для предотвращения динамической взаимоблокировки в параллельных приложениях.

10.4 Итоги

Сбои живучести являются серьезной проблемой, потому что нет никакого более короткого способа восстановиться от них, чем прерывание работы приложения. Наиболее распространенной формой проблем с живучестью является взаимоблокировка, вызванная нарушением порядка захвата блокировок. Предотвращение взаимоблокировки, вызванной нарушением порядка захвата блокировок, начинается с момента проектирования: гарантируйте, что когда

¹⁰⁹ Имеется в виду сетевой кабель, wi-fi и т.д.

потоки захватывают несколько блокировок, они делают это в согласованном порядке. Лучший способ добиться этого - использовать открытые вызовы во всей программе. Это бы значительно уменьшило количество мест, где одновременно удерживаются несколько блокировок, и сделало бы более очевидным местонахождение таких мест.

Глава 11 Производительность и масштабируемость

Одной из основных причин использования потоков является повышение производительности¹¹⁰. Использование потоков может привести к улучшению использования ресурсов, упрощая использование приложениями доступной вычислительной мощности, а также повысить скорость отклика, позволяя приложениям начинать обработку новых задач немедленно, пока существующие задачи еще выполняются.

В этой главе рассматриваются подходы к анализу, мониторингу и повышению производительности параллельных программ. К сожалению, многие из подходов к повышению производительности также увеличивают сложность, тем самым увеличивая вероятность появления проблем с безопасностью и живучестью. Хуже того, некоторые методы, предназначенные для повышения производительности, на самом деле контрпродуктивны или обменивают одну проблему с производительностью на другую. В то время как более высокая производительность часто желательна - и повышение производительности может принести большое удовлетворение - безопасность всегда на первом месте. Сначала сделайте свою программу правильной, затем сделайте ее быстрой - и только в том случае, когда ваши требования к производительности и измерения скажут вам, что она должна быть быстрее. При проектировании параллельного приложения, выжимание последнего бита из производительности часто является наименьшей из проблем.

11.1 Размышления о производительности

Повышение производительности означает выполнение большего объема работы с использованием меньшего количества ресурсов. Значение термина "ресурсы" может варьироваться; для данной активности, определенным ресурсом обычно будет являться то, в чём активность нуждается в самое кратчайшее время, будь то циклами ЦП, памятью, пропускной способностью сети, пропускной способностью подсистемы ввода/вывода, запросами к базе данных, дисковым пространством или любым количеством других ресурсов. Когда производительность активности ограничена доступностью определенного ресурса, мы говорим, что активность *ограничена* этим ресурсом: ограничена CPU, базой данных и т. д.

Хотя целью может быть повышение производительности в целом, использование нескольких потоков всегда приводит к некоторым затратам на производительность по сравнению с однопоточным подходом. К ним относятся накладные расходы, связанные с координацией между потоками (блокировкой, сигнализацией и синхронизацией памяти), накладные расходы на переключение контекста, накладные расходы на создание и завершение работы потоков, а также накладные расходы на планирование. При эффективном использовании потоков эти затраты с лихвой компенсируются большей пропускной способностью, быстродействием или производительностью. С другой стороны, плохо

¹¹⁰ Некоторые могут утверждать, что это единственная причина, по которой мы миримся со сложностью потоков.

спроектированное параллельное приложение может работать даже хуже, чем аналогичное последовательное¹¹¹.

Используя параллелизм для повышения производительности, мы пытаемся сделать две вещи: более эффективно использовать имеющиеся у нас обрабатывающие ресурсы, и позволить нашей программе использовать дополнительные обрабатывающие ресурсы, если они станут доступны. С точки зрения мониторинга производительности это означает, что мы хотим, чтобы процессоры были максимально заняты выполнение работы. (Конечно, это не означает что надо “сжигать циклы” в бесполезных вычислениях; мы хотим, чтобы процессоры были заняты выполнением полезной работы.) Если программа связана с вычислениями, то мы можем увеличить ее производительность, добавив больше процессоров; если программа не может полностью загрузить работой доступные процессоры, добавление больше количества процессоров не поможет. Поточность предлагает средства для того, чтобы всегда держать процессоры “горячими”, это достигается за счёт декомпозиции приложения, так что всегда будет работа, которая может быть выполнена любым доступным процессором.

11.1.1 Производительность и масштабируемость

Производительность приложения можно измерить несколькими способами, такими как время обслуживания, задержка, пропускная способность, эффективность, масштабируемость или производительность. Некоторые из них (время обслуживания, задержка) являются показателями того, “как быстро” данная единица работы может быть обработана или подтверждена; другие (производительность, пропускная способность) являются показателями того, “сколько” работы может быть выполнено с заданным количеством вычислительных ресурсов.

Масштабируемость (Scalability) описывает возможность повышения пропускной способности или производительности при добавлении дополнительных вычислительных ресурсов (например, дополнительных CPU, памяти, хранилища или пропускной способности подсистемы ввода/вывода).

Проектирование и настройка параллельных приложений для обеспечения масштабируемости, может сильно отличаться от традиционной оптимизации производительности. При настройке производительности, цель обычно состоит в том, чтобы выполнить ту же работу с меньшими усилиями, например, повторно использовать ранее вычисленные результаты с помощью кэширования или заменить алгоритм со сложностью $O(n^2)$ на алгоритм со сложностью $O(N \log n)$. При настройке масштабируемости, вместо приведённого выше подхода, вы пытаетесь найти способы распараллеливания проблемы, чтобы использовать дополнительные обрабатывающие ресурсы для выполнения *дополнительной* работы с *большим* количеством ресурсов.

Эти два аспекта производительности - *как быстро* и *сколько* - полностью разделены, а иногда даже вступают в противоречие друг с другом. Для достижения

¹¹¹ Коллега рассказал забавный анекдот: он участвовал в тестировании дорогостоящего и сложного приложения, которое управляло своей работой через настраиваемый пул потоков. После того, как система была завершена, тестирование показало, что оптимальное количество потоков для пула было... 1. Это должно было быть очевидно с самого начала; целевая система была однопроцессорной системой, а приложение было почти полностью привязано к процессору.

более высокой масштабируемости или лучшего использования оборудования мы часто *увеличиваем* объем работы, выполняемой каждой отдельной задачей, например при разделении задач на несколько “конвейерных” подзадач. По иронии судьбы, многие приемы, повышающие производительность однопоточных программ, плохо влияют на масштабируемость (см. пример в разделе 11.4.4).

Знакомая трехуровневая модель приложения, в которой представление, бизнес-логика и хранилище разделены и могут обрабатываться различными системами, иллюстрирует, как повышение масштабируемости часто происходит за счет снижения производительности. Монолитное приложение, в котором переплетены уровни представления, уровень бизнес-логики и уровень хранилища, почти наверняка обеспечит лучшую производительность для первой единицы работы, чем хорошо продуманная многоуровневая реализация, распределенная по нескольким системам. Как так могло получиться? Монолитное приложение не будет иметь накладываемой сетью задержки, присущей передаче задач между уровнями, и ему не придется оплачивать затраты, присущие разделению вычислительного процесса на отдельные абстрактные слои (например, накладные расходы на помещение в очередь и на копирования данных).

Однако когда монолитная система достигнет, в операциях обработки, потолка производительности, у нас может возникнуть серьезная проблема: может случиться так, что будет непомерно сложно значительно поднять уровень производительности. Поэтому мы часто соглашаемся с затратами производительности на более длительное время выполнения или с потреблением больших вычислительных ресурсов, затрачиваемых на единицу работы, в результате чего, наше приложение может масштабироваться для обработки большей нагрузки, при добавлении большего количества ресурсов.

Из различных аспектов производительности, аспекты “сколько” – масштабируемости, пропускной способности и производительности – обычно имеют для серверных приложений большее значение, чем аспекты “как быстро”. (Для интерактивных приложений задержка имеет тенденцию быть более важным показателем, поэтому пользователи не нуждаются в индикации прогресса и не нуждаются в том, чтобы выяснить, что происходит.) В этой главе основное внимание уделяется масштабируемости, а не производительности при однопоточной обработке.

11.1.2 Оценка компромиссов производительности

Почти все инженерные решения предполагают некую форму компромисса. Использование более толстой стали в пролёте моста может увеличить его вместимость и безопасность, но также приведёт к увеличению цены конструкции. Хотя решения в области разработки программного обеспечения обычно не предполагают компромиссов между деньгами и риском для жизни человека, у нас часто меньше информации, с помощью которой можно прийти к правильному компромиссу. Например, алгоритм “быстрая сортировка” (*quicksort*) очень эффективен для больших наборов данных, но менее сложная “сортировка пузырьком” (*bubble sort*) на самом деле более эффективна для небольших наборов данных. Если вас попросят внедрить эффективную процедуру сортировки, вам нужно знать о размерах наборов данных, которые ей придется обрабатывать, а также метрики, указывающие, пытаетесь ли вы оптимизировать среднее время, худшее время или предсказуемость. К сожалению, эта информация часто не входит в требования, предъявляемые к автору библиотечной процедуры сортировки. Это

является одной из причин, в связи с которыми большинство оптимизаций преждевременны: они часто проводятся до того, как станет доступен четкий набор требований.

Избегайте преждевременной оптимизации. Сначала сделайте это правильно, затем сделайте это быстро - если это еще не достаточно быстро.

Иногда, при принятии инженерных решений, вам придется обменивать одну форму затрат на другую (например, время обработки и потребление памяти); иногда вы обмениваете затраты на безопасность. Безопасность не обязательно означает риск для человеческих жизней, как это было в примере с мостом. Множество оптимизаций производительности выполняется за счет удобочитаемости или сопровождаемости - чем более "умный" или неочевидный код, тем сложнее его понять и поддерживать. Иногда оптимизация влечет за собой компрометацию хороших принципов объектно-ориентированного проектирования, например, нарушение инкапсуляции; иногда она сопряжена с большим риском ошибок, поскольку более быстрые алгоритмы обычно сложнее. (Если вы не можете определить затраты или риски, вы, вероятно, не продумали всё достаточно тщательно для того, чтобы продолжать.)

Большинство решений касательно производительности включают в себя несколько переменных и очень ситуативны. Прежде чем решить, что один подход "быстрее" другого, задайте себе несколько вопросов:

- Что вы имеете в виду под "быстрее"?
- При каких условиях этот подход будет более быстрым? Под легкой или тяжелой нагрузкой? С большими или маленькими наборами данных? Можете ли вы подкрепить свой ответ измерениями?
- Как часто эти условия могут возникнуть в вашей ситуации? Можете ли вы подкрепить свой ответ измерениями?
- Может ли этот код использоваться в других ситуациях, когда условия могут отличаться?
- Какие скрытые затраты, такие как усложнение разработки или риск в обслуживании, вы готовы обменять на эту улучшенную производительность? Это хороший компромисс?

Эти соображения применимы к любым инженерным решениям, связанным с производительностью, но эта книга о параллелизме. Почему мы рекомендуем такой консервативный подход к оптимизации? *Стремление к производительности, вероятно, является самым большим источником ошибок в параллелизме.* Убеждение в том, что синхронизация была "слишком медленной", привело к появлению многих умных, но опасных идиом, применяемых для уменьшения синхронизации (таких как двойная проверка блокировки, обсуждаемая в разделе 16.2.4), и часто приводится в качестве оправдания за несоблюдение правил, касающихся синхронизации. Поскольку ошибки параллелизма являются одними сложнейших в плане отслеживания и устранения, все, что несет риск их введения в код, должно предприниматься с особой тщательностью.

Хуже того, когда вы торгуете безопасностью в угоду производительности, вы можете не получить ни того, ни другого. Особенно, когда дело доходит до параллелизма, интуитивные предположения многих разработчиков о том, в чём заключается проблема с производительностью или какой подход будет быстрее или более масштабируемым, часто неверна. Поэтому крайне важно, чтобы любое занятие по настройке производительности сопровождалось конкретными требованиями к производительности (чтобы вы знали, когда настраивать и когда останавливать настройку), а также программой измерений с использованием реалистичной конфигурации и профиля нагрузки. После проведения настройки вновь проведите измерения, чтобы убедиться, что вы достигли желаемых улучшений. Риски для безопасности и технического обслуживания, связанные с проведением множества оптимизаций, довольно велики - вы не хотите оплачивать эти расходы, если вам это не нужно - и вы определенно не хотите их оплачивать в том случае, если не получаете ожидаемой выгоды.

Измерьте, не догадывайтесь.

На рынке существуют сложные инструменты профилирования для измерения производительности и отслеживания узких мест производительности, но вам не нужно тратить много денег, чтобы понять, что делает ваша программа. Например, бесплатное приложение `perfbar` может вам помочь получить хорошее представление о том, насколько загружены процессоры, и так как ваша цель, как правило, заключается в том, чтобы держать процессоры загруженными, это очень хороший способ оценить, нужно ли вам проводить настройку производительности или насколько эффективной была ваша настройка.

11.2 Закон Амдала

Некоторые проблемы могут быть решены быстрее за счёт большего количества ресурсов – чем больше рабочей силы используется для уборки урожая, тем быстрее она может быть завершена. Другие задачи принципиально последовательны – никакое количество дополнительных работников не заставит урожай расти быстрее. Если одной из основных причин использования потоков является использование мощности нескольких процессоров, мы также должны убедиться, что проблема поддаётся параллельной декомпозиции и что наша программа эффективно использует доступный для распараллеливания потенциал.

Большинство параллельных программ имеют много общего с фермерством, состоя из смеси параллельных и последовательных частей. Закон Амдала описывает, насколько теоретически программа может быть ускорена при добавлении дополнительных вычислительных ресурсов, исходя из соотношения параллельных и последовательных компонентов. Если F представляет собой долю вычислений, которая должна выполняться последовательно, то закон Амдала гласит, что на машине с N процессорами, мы можем добиться ускорения не более:

$$\text{Speedup} \leq \frac{1}{F + \frac{1-F}{N}}$$

Когда значение N приближается к бесконечности, максимальное ускорение сходится к значению $1/F$, это означает, что программа, в которой пятьдесят процентов работы должно быть выполнено последовательно, может быть ускорена

только в два раза, независимо от того, сколько процессоров доступно, а программа, в которой только десять процентов работы должно быть выполнено последовательно, может быть ускорена не более чем в десять раз. Закон Амдала также позволяет количественно оценить эффективность затрат на последовательное выполнение. С десятью процессорами, программа с 10% последовательно выполняемого кода, может достигнуть максимального увеличения в 5.3 раза (при 53% использования ресурсов), а со 100 процессорами она сможет достичь увеличения в 9.2 раза (при 9% использовании ресурсов). Потребуется множество неэффективно используемых процессоров, но никогда не получится добиться коэффициента равного десяти.

На рисунке 11.1 иллюстрируется максимально возможная загрузка процессора при различных значениях степени последовательно выполняемого кода и количества используемых процессоров. (Использование определяется как ускорение, делённое на количество процессоров.) Очевидно, что по мере увеличения количества процессоров, даже небольшой процент последовательно выполняемого кода ограничивает то количество пропускной способности, которое можно увеличить за счет использования дополнительных вычислительных ресурсов.

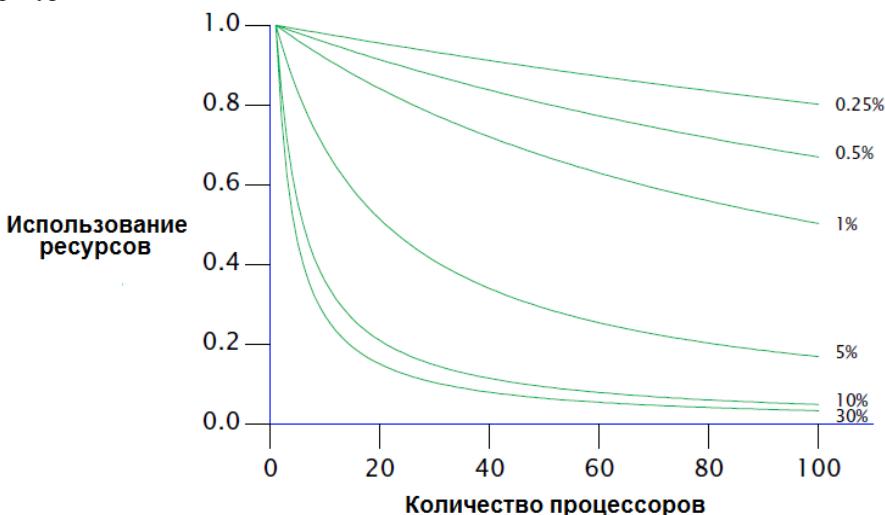


Рисунок 11.1 Максимальное использование ресурсов при различных размерах последовательно выполняемой части кода, согласно закону Амдала

В главе 6 рассматриваются вопросы определения логических границ, для разбиения приложений на задачи. Но для того, чтобы предсказать, какого рода ускорение возможно получить при запуске приложения в многопроцессорной системе, необходимо также определить предпосылки возникновения последовательно выполняемого кода в ваших задачах.

Представьте себе приложение, в котором N потоков выполняют метод `doWork` из листинга 11.1, извлекая задачи из общей рабочей очереди и обрабатывая их; предположим, что задачи не зависят от результатов или побочных эффектов, возникающих при выполнении других задач. Проигнорируем на мгновение то, как задачи попадают в очередь, насколько хорошо это приложение будет масштабироваться по мере добавления процессоров? На первый взгляд может показаться, что приложение полностью распараллеливается: задачи не ждут друг друга, и чем больше процессоров, тем больше задач может обрабатываться

одновременно. Однако имеется и последовательный компонент - извлечение задач из рабочей очереди. Рабочая очередь совместно используется всеми рабочими потоками, и для обеспечения ее целостности, в условиях параллельного доступа, потребуется некоторая синхронизация. Если для защиты состояния очереди используется блокировка, то в то время как один поток извлекает задачу из очереди, другие потоки, которым нужно извлечь из очереди свою следующую задачу, вынуждены ожидать - и именно в этом месте обработка задачи выполняется последовательным кодом

```
public class WorkerThread extends Thread {  
    private final BlockingQueue<Runnable> queue;  
  
    public WorkerThread(BlockingQueue<Runnable> queue) {  
        this.queue = queue;  
    }  
  
    public void run() {  
        while (true) {  
            try {  
                Runnable task = queue.take();  
                task.run();  
            } catch (InterruptedException e) {  
                break; /* Allow thread to exit */  
            }  
        }  
    }  
}
```

Листинг 11.1 Последовательный доступ к очереди задач

Время обработки одной задачи включает в себя не только время выполнения задачи, представленной экземпляром `Runnable`, но и время на извлечение задачи из совместно используемой рабочей очереди. Если рабочая очередь представляет собой реализацию `LinkedBlockingQueue`, операция извлечения задачи из очереди может блокироваться на меньшее время, чем при использовании синхронизированной реализации `LinkedList`, поскольку реализация `LinkedBlockingQueue` использует более масштабируемый алгоритм, но доступ к любой совместно используемой структуре данных фундаментально вводит в программу элемент сериализации.

В этом примере также игнорируется другой распространенный источник последовательно выполняемого кода: обработка результатов. Все полезные вычисления порождают какой-то результат или вызывают побочные эффекты - в противном случае они могут быть исключены как мертвый код. Поскольку интерфейс `Runnable` не предоставляет возможности явной обработки результатов, выполнение этих задач должно приводить к возникновению каких-то побочных эффектов, например, запись результатов выполнения в файл лога или их помещение в структуру данных. Файлы логов и контейнеры результатов обычно совместно используются несколькими рабочими потоками и поэтому также являются источником последовательно выполняемого кода. Если вместо этого каждый поток будет поддерживать свои собственные структуры данных для хранения результатов, которые будут объединяться после выполнения всех задач,

то окончательное слияние также будет являться источником последовательно выполняемого кода.

У всех параллельных приложений есть некоторые источники последовательно выполняемого кода; если вы думаете, что у вашего кода их нет, посмотрите снова.

11.2.1 Пример: скрытое последовательное выполнение в фреймворках

Чтобы увидеть, как последовательно выполняемый код может быть скрыт в структуре приложения, можно сравнить пропускную способность при добавлении потоков и определить различия в последовательном выполнении на основе наблюдаемых различий в масштабируемости. На рис. 11.2 демонстрируется результат выполнения простого приложения, в котором несколько потоков многократно удаляют элемент из общей очереди и обрабатывают его, аналогично примеру из листинга 11.1. Шаг обработки включает в себя только локальное, относительно потока, вычисление. Если поток обнаруживает, что очередь пуста, он помещает пакет новых элементов в очередь, чтобы другим потокам было что обрабатывать на следующей итерации. Доступ к совместно используемой очереди явным образом влечет за собой некоторую степень последовательного выполнения, но шаг обработки полностью распараллеливается, так как он не включает в себя совместно используемые данные.

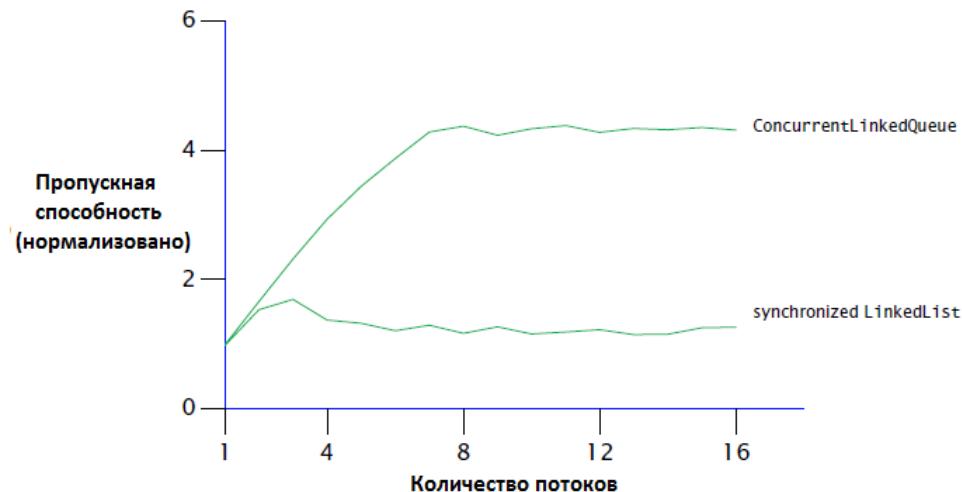


Рисунок 11.2 Сравнение реализаций очередей

Кривые на рис. 11.2 позволяют сравнить пропускную способность для двух потокобезопасных реализаций очереди: реализации на основе `LinkedList`, обернутой реализацией `synchronizedList`, и реализации `ConcurrentLinkedQueue`. Тесты проводились на 8-ядерной системе Sparc V880, под управлением ОС Solaris. Хотя каждый запуск представляет собой выполнение одинакового объема “работы”, мы видим, что простое изменение реализаций очередей может оказать

Пропускная способность реализации `ConcurrentLinkedQueue` продолжает улучшаться, пока количество потоков не достигнет количества процессоров, а затем остается почти постоянной. С другой стороны, пропускная способность

синхронизированной реализации `LinkedList` показывает некоторое улучшение до трех потоков, но затем падает по мере нарастания издержек синхронизации. К тому времени, когда количество потоков доберется до четырех или пяти, конкуренция станет настолько тяжела, что каждая попытка доступа к блокировке очереди будет оспариваться, и в объеме всей пропускной способности начнет доминировать переключение контекста.

Различие в пропускной способности происходит от различных степеней последовательно выполняемого кода, в двух разных реализациях очередей. Синхронизированная реализация `LinkedList` защищает все состояние очереди с помощью единственной блокировки, которая удерживается на протяжении вызовов методов `offer` или `remove`; реализация `ConcurrentLinkedQueue` использует сложный алгоритм неблокирующей очереди (см. раздел [15.4.2](#)), который использует атомарные ссылки для обновления отдельных указателей на ссылки. В одном случае - последовательно выполняется вставка или удаление; в другом случае - последовательно выполняется только обновление отдельных указателей.

11.2.2 Качественное применение закона Амдала

Закон Амдала количественно определяет возможное ускорение, когда доступно больше вычислительных ресурсов, если мы можем точно оценить долю последовательно выполняемого кода. Хотя прямое измерение доли последовательно выполняемого кода может быть затруднено, закон Амдала всё же может быть полезен и без такого измерения.

Поскольку наши ментальные модели находятся под влиянием окружающей нас среды, многие из нас привыкли думать, что многопроцессорная система, имеет два или четыре процессора, или, может быть (если у нас большой бюджет) несколько десятков таковых, потому что эта технология стала широкодоступна в последние годы. Но по мере того, как многоядерные процессоры станут популярнее, системы будут иметь сотни или даже тысячи процессоров¹¹². Алгоритмы, которые кажутся масштабируемыми в четырех процессорной системе, могут иметь скрытые узкие места (*bottlenecks*) в масштабируемости, которые еще не были обнаружены.

В процессе оценки алгоритма, размышление “с ограничениями” о том, что произойдет с сотнями или тысячами процессоров, может дать некоторое представление, где могут проявиться ограничения масштабирования. Например, в разделах 11.4.2 и 11.4.3 обсуждается два способа повышения детализации блокировок: разделение блокировок (разделение одной блокировки на две) и чередование блокировок (разделение одной блокировки на несколько блокировок). Глядя на них через призму закона Амдала, мы видим, что разделение блокировки на две части, не сильно приближает нас к использованию множества процессоров, но чередование блокировок кажется гораздо более перспективным, потому что размер набора полос (*stripe*) может быть увеличен по мере увеличения количества процессоров. (Конечно, оптимизация производительности всегда должна рассматриваться в свете фактических требований к производительности; в некоторых случаях, разделения блокировки на две части, для удовлетворения требований, может быть достаточно.)

¹¹² Обновление рынка: на момент написания этой статьи Sun поставляет недорогие серверные системы на основе 8-ядерного процессора Niagara, а Azul поставляет высокопроизводительные серверные системы (96, 192 и 384-ядра) на основе 24-ядерного процессора Vega.

11.3 Затраты, вводимые потоками

Однопоточные программы не требуют ни планирования, ни синхронизации и не нуждаются в использовании блокировок для сохранения согласованности структур данных. Планирование и координация взаимодействия между потоками требуют затрат производительности; для потоков, предлагающих повышение производительности, преимущества распараллеливания должны перевешивать затраты на обеспечение параллелизма.

11.3.1 Переключение контекста

Если основной поток является единственным потоком, выполняемым по расписанию, он почти никогда не будет планироваться. С другой стороны, если количество выполняемых потоков больше, чем количество доступных ЦП, в конечном итоге ОС будет упреждать один поток, чтобы другой мог использовать ЦП. Это вызывает *переключение контекста*, которое требует сохранения контекста текущего потока и восстановления контекста выполнения нового запланированного потока.

Переключение контекста операция не бесплатная; планирование потоков требует манипулирования совместно используемыми структурами данных в ОС и JVM. ОС и JVM используют те же процессоры, что и ваша программа; большее количество времени процессора, потраченного на JVM, и код ОС означает, что вашей программе будет доступно меньше времени. Но деятельность ОС и JVM - это не единственные затраты, возникающие при переключении контекста. Когда происходит переключение на новый поток, данные, в которых он нуждается, вряд ли будут доступны в локальном кэше процессора, таким образом, переключение контекста приведёт к шквалу промахов попадания в кэш, и потоки будут работать немного медленнее в том случае, если их выполнение будет запланировано впервые. Это одна из причин, по которой планировщики дают каждому выполняемому потоку определенный минимальный квант времени, даже когда многие другие потоки вынуждены ожидать: такой подход позволяет амортизировать затраты на переключение контекста и его последствия в течение более длительного времени непрерывного выполнения, улучшая общую пропускную способность (за счет некоторых затрат на отзывчивость).

Когда поток блокируется из-за ожидания конкурирующей блокировки, JVM обычно его приостанавливает и позволяет ему выключиться. Если потоки блокируются часто, они не смогут использовать весь запланированный для выполнения квант времени. Программа, использующая множество блокировок (блокирующий ввод/вывод, ожидание конкурирующих блокировок или ожидание переменных условий), использует больше переключений контекста, чем та, что ограничена только ЦП, что увеличивает затраты на планирование и приводит к снижению пропускной способности. (Неблокирующие алгоритмы также могут помочь в снижении количества переключений контекста; см. главу [15.](#))

Фактические затраты на переключение контекста варьируются от платформы к платформе, но хорошее эмпирическое правило заключается в том, что переключение контекста, для большинства текущих процессоров, по затратам эквивалентно 5,000-10,000 тактам процессора или нескольким микросекундам.

Команда `vmstat` в системах Unix и средство `perfmon` в системах Windows позволяют получить отчёт о количестве переключений контекста, и проценте времени, проведенном в ядре. Высокая загрузка ядра (более 10%) часто указывает

на интенсивное планирование, которое может быть вызвано блокировками при выполнении операций ввода/вывода или конкуренцией за захват блокировок.

11.3.2 Синхронизация памяти

Затраты производительности на синхронизацию берут своё начало из нескольких источников. Гарантии видимости, предоставляемые операторами `synchronized` и `volatile`, могут повлечь за собой использование специальных инструкций, называемых *барьерами памяти* (*memory barriers*), которые могут сбрасывать или аннулировать кэши, сбрасывать аппаратные буферы на запись и останавливать выполняющие конвейеры. Использование барьераов памяти также может иметь косвенные последствия для производительности, поскольку они препятствуют проведению оптимизаций, выполняемых компиляторами; большинство операций нельзя переупорядочить, когда используются барьеры памяти.

При оценке влияния синхронизации на производительность, важно различать *конкурентную* (*contended*) и *неконкурентную* (*uncontended*) синхронизацию. Механизм `synchronized` оптимизирован для неконкурентного случая (оператор `volatile` всегда неконкурентен), и на момент написания этих строк, затраты на выполнение “кратчайшего варианта” неконкурентной синхронизации, для большинства систем, колеблются в пределах от 20 до 250 тактов процессора. Хотя это, конечно, не ноль, эффект необходимой, незапланированной синхронизации редко оказывает значимое влияние на общую производительность приложения, в свою очередь, альтернатива включает в себя угрозу безопасности и потенциальное обречение себя (или вашего преемника), в будущем, на очень болезненный поиск ошибок.

Современные среды JVM могут снижать затраты на случайную синхронизацию, проводя оптимизацию блокировок, за которые, как доказано, никогда не будет конкуренции. Если объект блокировки доступен только текущему потоку, среде JVM разрешено оптимизировать захват блокировки, поскольку нет никакого способа, которым другой поток мог бы синхронизироваться на той же самой блокировке. Например, захват блокировки, как показано в листинге 11.2, всегда может быть устранен средой JVM.

```
synchronized (new Object()) {  
    // do something  
}
```



Листинг 11.2 Синхронизация, не оказывающая никакого эффекта. Не делайте так

Более сложные среды JVM могут использовать последовательный анализ для определения того, что ссылка на локальный объект никогда не публикуется в куче и поэтому является локальной для потока. В методе `getStoogeNames` из листинга 11.3, единственной ссылкой на список является локальная переменная `stooges`, а переменные, ограниченные стеком, автоматически являются локальными для потока. Наивное выполнение метода `getStoogeNames` захватывает и освобождает блокировку на экземпляре `Vector` четыре раза, по одному разу для каждого вызова `add` или `toString`. Тем не менее, умный компилятор среды выполнения может встроить эти вызовы, а затем увидеть, что переменная `stooges` и её внутреннее

состояние никогда не сбегают, и, следовательно, что все четыре захвата блокировок могут быть устраниены¹¹³.

```
public String getStoogeNames() {  
    List<String> stooges = new Vector<String>();  
    stooges.add("Moe");  
    stooges.add("Larry");  
    stooges.add("Curly");  
    return stooges.toString();  
}
```

Листинг 11.3 Кандидат на оптимизацию *lock elision*

Даже без обращения к последовательному анализу, компиляторы также могут выполнять укрупнение блокировок, путём слияния соседних блоков `synchronized` с использованием одной и той же блокировки. В случае метода `getStoogeNames`, среда JVM, выполняющая укрупнение блокировки, может объединить три вызова метода `add` и вызов метода `toString` в один блок захвата и освобождения блокировки, используя эвристику относительных затрат на синхронизацию по сравнению с инструкциями внутри блока `synchronized`.¹¹⁴ Это не только снижает затраты на синхронизацию, но и предоставляет оптимизатору гораздо больший блок для работы, что, вероятно, позволяет выполнять и другие оптимизации.

Не стоит чрезмерно беспокоиться о стоимости неконкурентной синхронизации. Базовый механизм уже достаточно быстр, и среды JVM могут выполнять дополнительные оптимизации, которые еще больше уменьшают или устраняют затраты. Вместо этого сосредоточьте усилия по оптимизации на тех областях, в которых фактически происходит конфликт блокировок.

Синхронизация в одном потоке также может повлиять на производительность других потоков. Синхронизация создает трафик в шине совместно используемой памяти; эта шина имеет ограниченную пропускную способность и совместно используется всеми процессорами. Если потоки вынуждены конкурировать за пропускную способность синхронизации, все потоки, использующие синхронизацию, будут страдать.¹¹⁵

11.3.3 Блокирование

Неконкурирующая синхронизация может быть полностью обработана в среде JVM (Bacon et al., 1998); для конкурирующей синхронизации может потребоваться вмешательство ОС, что приводит к увеличению затрат. В случае, если блокировка конкурирующая, проигравший поток должен быть заблокирован. Среда JVM может реализовать блокировку или через *spin-waiting* (многократная попытка захвата блокировки, пока захват не завершится успехом) или путем *приостановки* (*suspending*) заблокированного потока средствами операционной системы. Что

¹¹³ Эта оптимизация компилятора, называемая *lock elision*, выполняется в IBM JVM и ожидается в HotSpot с Java 7.

¹¹⁴ Умный динамический компилятор может выяснить, что этот метод всегда возвращает одну и ту же строку, и после первого выполнения перекомпилировать метод `getStoogeNames` для простого возврата значения, полученного при первом выполнении.

¹¹⁵ Этот аспект иногда используется в спорах против использования неблокирующих алгоритмов без отката, потому что при тяжелой конкуренции неблокирующие алгоритмы генерируют больше синхронизирующего трафика, чем основанные на блокировках. См. Главу 15.

является более эффективным, зависит от отношения между накладными расходами на переключение контекста и временем, в течение которого блокировка будет доступна; механизм spin-waiting предпочтительнее для короткого ожидания, а приостановка предпочтительнее для длительного ожидания. Некоторые среды JVM выбирают между обоими вариантами адаптивно, основывая своё решение на данных профилирования прошлых периодов ожидания, но большинство просто приостанавливают выполнение потоков, ожидающих блокировку.

Приостановка потока, из-за того, что он не может получить блокировку, или из-за того, что он заблокирован в состоянии ожидания или заблокирован при выполнении операций ввода/вывода, влечёт за собой два дополнительных переключения контекста и всю сопровождающую активность операционной системы и кэша: блокированный поток выключается, прежде чем истечёт выделенный ему квант времени, и затем, позже, переключается обратно, когда блокировка или другой необходимый ресурс становятся доступны. (Блокировка из-за конфликта блокировок также приводит к затратам потока, удерживающего блокировку: когда он освобождает блокировку, он должен затем попросить ОС возобновить выполнение заблокированного потока.)

11.4 Уменьшение конкуренции за блокировку

Мы видели, что последовательно выполняемый код ухудшает масштабируемость, а переключение контекста - производительность. Конкурирующая блокировка является причиной возникновения обоих случаев, поэтому уменьшение конкуренции может улучшить и производительность, и масштабируемость.

Доступ к ресурсам, защищаемым монопольной блокировкой, строго последователен - одновременно только один поток может получить к ним доступ. Конечно, мы используем блокировки по уважительным причинам, таким как предотвращение повреждения данных, но эта безопасность имеет свою цену. Постоянная конкуренция за захват блокировки ограничивает масштабируемость.

Основной угрозой масштабируемости в параллельных приложениях является монопольная блокировка ресурсов.

Вероятность возникновения конкуренции за блокировку зависит от двух факторов: частоты захвата блокировки и времени её удержания после захвата¹¹⁶. Если произведение этих факторов достаточно мало, то большинство попыток захвата блокировки не будет приводить к конкуренции, и конфликт блокировок не будет создавать значительных препятствий для масштабируемости. В том случае, если блокировка востребована достаточно сильно, потоки будут блокироваться в ожидании её; в крайнем случае, процессоры будут простаивать, даже если есть множество работы.

Существует три способа уменьшения конкуренции за блокировку:

- Сокращение времени, в течение которого удерживается блокировка;
- Уменьшение частоты, с которой захватываются блокировки; или
- Замена монопольных блокировок механизмами координации, обеспечивающими больший параллелизм.

¹¹⁶ Это следствие закона Литтла, являющегося результатом теории очередей, в которой говорится, что “среднее число клиентов в стабильной системе, равно их средней скорости прибытия, умноженной на их среднее время нахождения в системе”. (Маленький, 1961)

11.4.1 Сужение области действия блокировки (“Get in, get out”)

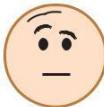
Эффективным подходом в снижении вероятности возникновения конкуренции, является удержание блокировки настолько краткое время, насколько это возможно. Этого можно добиться за счёт перемещения кода, не требующего блокировки в блоках `synchronized`, особенно дорогостоящих операций и потенциально блокирующих операций, таких как операции ввода/вывода.

Легко заметить, что слишком длительное удержание “горячей” блокировки может ограничить масштабируемость; мы видели пример демонстрирующий подобную ситуацию в классе `SynchronizedFactorizer`, приведённом в главе 2. Если операция удерживает блокировку в течение 2 миллисекунд и для каждой операции требуется эта блокировка, пропускная способность не сможет превысить 500 операций в секунду, независимо от количества доступных процессоров. Уменьшение времени удержания блокировки до 1 миллисекунды повышает предел ограниченной блокировкой пропускной способности до 1000 операций в секунду.¹¹⁷

В классе `AttributeStore` из листинга 11.4, демонстрируется пример, в котором блокировка удерживается дольше, чем это необходимо. Метод `userLocationMatches` ищет местоположение пользователя в экземпляре `Map` и использует соответствие регулярному выражению, чтобы определить, соответствует ли полученное значение предоставленному шаблону. Метод `userLocationMatches` в целом определён как `synchronized`, но единственная часть кода, которая действительно нуждается в блокировке, - это вызов `Map.get`.

```
@ThreadSafe
public class AttributeStore {
    @GuardedBy("this") private final Map<String, String>
        attributes = new HashMap<String, String>();

    public synchronized boolean userLocationMatches(String name,
                                                   String regexp) {
        String key = "users." + name + ".location";
        String location = attributes.get(key); if
            (location == null)
                return false;
        else
            return Pattern.matches(regexp, location);
    }
}
```



Листинг 11.4 Удержание блокировки дольше, чем это необходимо

В классе `BetterAttributeStore`, приведённом в листинге 11.5, переписан метод из класса `AttributeStore`, в целях значительного сокращения времени удержания блокировки. Первый шаг заключается в том, чтобы сформировать ключ - строку вида `users.name.location` - для экземпляра `Map`, связанного с местоположением пользователя. Это влечет за собой создание экземпляра объекта `StringBuilder`, добавление к нему нескольких строк и формирование результата в виде экземпляра

¹¹⁷ Фактически, это вычисление *занижает* затраты на удержание слишком длительных блокировок, потому что не учитывает издержек добавляемых переключением контекста, генерируемых увеличившейся конкуренцией за блокировку.

`String`. После извлечения расположения, регулярное выражение сопоставляется с результатирующей строкой местоположения. Поскольку построение строки-ключа и обработка регулярного выражения не имеют доступа к совместно используемому состоянию, они не нуждаются в выполнении с удерживаемой блокировкой. Класс `BetterAttributeStore` учитывает эти шаги вне блока `synchronized`, тем самым уменьшая время удержания блокировки.

```
@ThreadSafe
public class BetterAttributeStore {
    @GuardedBy("this") private final Map<String, String>
        attributes = new HashMap<String, String>();

    public boolean userLocationMatches(String name, String regexp) {
        String key = "users." + name + ".location";
        String location;
        synchronized (this) {
            location = attributes.get(key);
        }
        if (location == null)
            return false;
        else
            return Pattern.matches(regexp, location);
    }
}
```

Листинг 11.5 Сокращение времени удержания блокировки

Уменьшение области действия блокировки в методе `userLocationMatches` сокращает количество инструкций, выполняемых с удерживаемой блокировкой. Согласно закону Амдала, это приводит к устранению препятствий для масштабируемости, так как уменьшается объем последовательно выполняемого кода.

Поскольку класс `AttributeStore` имеет только одну переменную состояния, переменную `attributes`, мы можем в дальнейшем улучшить его, с использованием подхода заключающегося в *делегирования потокобезопасности* (*delegating thread safety*, раздел 4.3). С помощью замены переменной `attributes` потокобезопасной реализацией `Map` (`Hashtable`, `synchronizedMap`, или `ConcurrentHashMap`), класс `AttributeStore` может делегировать все свои обязательства по обеспечению потокобезопасности нижележащей потокобезопасной коллекции. Этот подход устраняет необходимость использования явной синхронизации в классе `AttributeStore`, уменьшает область и продолжительность действия блокировки, при получении доступа к экземпляру `Map`, и устраниет риск того, что те, кто будут в будущем сопровождать этот код, нарушают потокобезопасность, забыв получить соответствующую блокировку перед доступом к переменной `attributes`.

Хотя сжатие блоков `synchronized` может улучшить масштабируемость, блок `synchronized` может быть *слишком* маленьким - операции, которые должны быть атомарными (например, обновление нескольких переменных, участвующих в инварианте), должны содержаться в одном синхронизированном блоке. В связи с тем, что затраты на синхронизацию не равны нулю, разбиение одного блока `synchronized` на несколько блоков `synchronized` (при условии корректности) в какой-то момент становится контрпродуктивным, если смотреть с точки зрения

производительности.¹¹⁸ Идеальный баланс, конечно, зависит от платформы, но на практике имеет смысл беспокоиться о размере синхронизированного блока только тогда, когда за его пределы можно вынести “существенные” вычисления или блокирующие операции.

11.4.2 Уменьшение детализации блокировки

Другим способом уменьшить долю времени, в течение которого удерживается блокировка (и, следовательно, вероятность того, что возникнет конкуренция) заключается в том, чтобы потоки обращались к ней реже. Это может быть выполнено путем *разделения* (*lock splitting*) и *чредования блокировок* (*lock striping*), что приводит к использованию отдельных блокировок для защиты нескольких независимых переменных состояния, ранее защищенных одной блокировкой. Эти подходы позволяют уменьшить степень детализации, при которой происходит блокировка, что потенциально позволяет повысить масштабируемость, но использование большего количества блокировок также увеличивает риск возникновения взаимоблокировки.

В качестве мысленного эксперимента представьте, что произойдет, если для всего приложения будет существовать *только одна* блокировка, вместо отдельной блокировки для каждого объекта. Выполнение всех блоков *synchronized*, независимо от их блокировки, будет обрабатываться последовательным кодом. Поскольку множество потоков конкурирует за глобальную блокировку, вероятность того, что двум потокам одновременно понадобится блокировка, существенно возрастает, что приводит к большему количеству конфликтов. Поэтому, если бы запросы на захват блокировок были бы распределены по большему набору блокировок, было бы меньше конфликтов. Меньшее количество потоков блокировалось бы в ожидании доступности блокировок, что повысило бы масштабируемость.

Если блокировка защищает несколько *независимых* переменных состояния, можно улучшить масштабируемость, разделив ее на несколько отдельных блокировок, каждая из которых бы защищала разные переменные. В результате, обращение к каждой блокировке происходило бы реже.

Класс *ServerStatus* в листинге 11.6 демонстрируют часть интерфейса системы мониторинга сервера баз данных, поддерживающей набор пользователей, вошедших в систему, и набор запросов, выполняемых в данный момент. Когда пользователь входит в систему или выходит из неё, или начинается или завершается выполнение запроса, состояние объекта *ServerStatus* обновляется путем вызова соответствующего метода *add* или *remove*. Эти два типа информации полностью независимы; класс *ServerStatus* можно даже разделить на два отдельных класса, без потери функциональности.

```
@ThreadSafe
public class ServerStatus {
    @GuardedBy("this") public final Set<String> users;
    @GuardedBy("this") public final Set<String> queries;
    ...
}
```

¹¹⁸ Если JVM выполняет укрупнение блокировки, она может отменить, в некоторых случаях, разделение синхронизированных блоков на блоки меньшего размера.

```
public synchronized void addUser(String u) { users.add(u); }
public synchronized void addQuery(String q) { queries.add(q); }
public synchronized void removeUser(String u) {
    users.remove(u);
}
public synchronized void removeQuery(String q) {
    queries.remove(q);
}
}
```

Листинг 11.6 Кандидат на разделение блокировки

Вместо защиты как переменной `users`, так и переменной `queries` с помощью блокировки на классе `ServerStatus`, мы можем защитить каждую из них отдельной блокировкой, как показано в листинге 11.7. После разделения блокировки, каждая из новых, более детальных блокировок будет видеть меньше трафика блокировок, чем исходная, более грубая блокировка. (Делегирование обеспечения потокобезопасности переменных `users` и `queries` реализации `Set`, вместо использования явной синхронизации, будет неявно обеспечивать разделение блокировок, так как каждый из экземпляров `Set`, для защиты собственного состояния, будет использовать отдельную блокировку.)

```
@ThreadSafe
public class ServerStatus {
    @GuardedBy("users") public final Set<String> users;
    @GuardedBy("queries") public final Set<String> queries;
    ...
    public void addUser(String u) {
        synchronized (users) {
            users.add(u);
        }
    }

    public void addQuery(String q) {
        synchronized (queries) {
            queries.add(q);
        }
    }
    // remove methods similarly refactored to use split locks
}
```

Листинг 11.7 Класс `ServerStatus` отрефакторенный с использованием разделения блокировок

Разделение блокировки на две предлагает наибольшие возможности для улучшения, когда блокировка испытывает умеренную, но не тяжелую конкуренцию. Разделение блокировок, которые испытывают небольшую конкуренцию, на выходе дает небольшое улучшение производительности или пропускной способности, хотя также может быть увеличен порог загрузки, при котором производительность начинает ухудшаться из-за влияния конкуренции. Разделение блокировок, испытывающих умеренную конкуренцию, может фактически превратить их в неконкурентные блокировки, что является наиболее

желательным результатом, как для производительности, так и для масштабируемости.

11.4.3 Чередование блокировок

Разделение жёстко конкурентной блокировки на две, вероятнее всего приведет к получению двух жёстко конкурентных блокировок. Несмотря на то, что это приведет к небольшому улучшению масштабируемости, позволяя двум потокам выполнять одновременно, вместо выполнения по одному, это все еще не приведёт к значительному улучшению перспектив параллелизма в системе со многими процессорами. Пример разделения блокировки в классе `ServerStatus` не предполагает очевидной возможности для дальнейшего разделения блокировок.

Разделение блокировок иногда может быть расширено для секционирования блокировки на наборе независимых объектов переменного размера, и в этом случае называется *чертежом блокировок* (*lock striping*). Например, реализация `ConcurrentHashMap` использует массив из 16 блокировок, каждая из которых защищает 1/16 часть сегментов (*bucket*) хэша; сегмент N защищен блокировкой $N \bmod 16$. Предположим, что хэш-функция обеспечивает разумные характеристики распределения, а ключи доступны равномерно, это должно уменьшить спрос на любую заданную блокировку примерно в 16 раз. Именно этот метод позволяет классу `ConcurrentHashMap` поддерживать до 16 параллельных писателей. (Число блокировок может быть увеличено, чтобы обеспечить еще лучший параллелизм при интенсивном доступе, в системах с высоким числом процессоров, но число полос должно увеличиваться выше значения по умолчанию 16, только если у вас есть доказательства того, что параллельные записывающие операции генерируют конкуренцию в достаточной степени, чтобы гарантировать повышение предела.)

Один из недостатков чередования блокировок заключается в том, что блокировка коллекции для эксклюзивного доступа сложнее и затратнее, чем с использованием одиночной блокировки. Как правило, операцию можно выполнить, захватив не более одной блокировки, но иногда необходимо заблокировать всю коллекцию, например, когда классу `ConcurrentHashMap` необходимо расширить экземпляр `Map` и перехешировать значения в больший набор блоков. Обычно это осуществляется с помощью захвата всех блокировок в наборе полос.¹¹⁹

Класс `StripedMap` приведённый в листинге 11.8, иллюстрирует основанную, с помощью чередования блокировок, на хэше реализацию `Map`. Есть `N_LOCKS` блокировок, каждая защищает собственный поднабор сегментов. Большинство методов, подобных `get`, нуждаются в получении блокировки только для одного сегмента. Некоторым методам может требоваться захват всех блокировок, но, как и в случае реализации метода `clear`, может быть не обязательно захватывать их все одновременно¹²⁰.

```
@ThreadSafe
public class StripedMap {
```

¹¹⁹ Единственный способ получить произвольный набор внутренних блокировок - рекурсия.

¹²⁰ Используемый подход к очистке экземпляра `Map` не является атомарным, поэтому нет необходимости в ожидании момента времени, когда экземпляр `StripedMap` будет фактически пуст, если другие потоки параллельно добавляют элементы; чтобы сделать операцию атомарной, необходимо захватить все блокировки одновременно. Однако для параллельных коллекций, клиенты которых, как правило, не могут захватывать блокировку для эксклюзивного доступа, результат, возвращаемый такими методами, как `size` или `isEmpty`, может быть устаревшим к моменту возврата ими значения, так что такое поведение, хотя возможно это несколько удивительно, обычно является приемлемым.

```

// Synchronization policy: buckets[n] guarded by locks[n%N_LOCKS]
private static final int N_LOCKS = 16;
private final Node[] buckets;
private final Object[] locks;

private static class Node { ... }

public StripedMap(int numBuckets) {
    buckets = new Node[numBuckets];
    locks = new Object[N_LOCKS];
    for (int i = 0; i < N_LOCKS; i++)
        locks[i] = new Object();
}

private final int hash(Object key) {
    return Math.abs(key.hashCode() % buckets.length);
}

public Object get(Object key) {
    int hash = hash(key);
    synchronized (locks[hash % N_LOCKS]) {
        for (Node m = buckets[hash]; m != null; m = m.next)
            if (m.key.equals(key))
                return m.value;
    }
    return null;
}

public void clear() {
    for (int i = 0; i < buckets.length; i++) {
        synchronized (locks[i % N_LOCKS]) {
            buckets[i] = null;
        }
    }
}
...
}

```

Листинг 11.8 Основанная на хэше реализация Map с чередованием блокировок

11.4.4 Как избежать использования “горячих полей”

Разделение и чередование блокировок может улучшить масштабируемость, поскольку оно позволяет разным потокам работать с разными данными (или разными частями одной и той же структуры данных), не создавая помех друг другу. Программа, которая выиграла бы от разделения блокировок, чаще проявляет конкуренцию за захват блокировок, чем конкуренцию за доступ к *данным*, защищаемым этими блокировками. Если блокировка защищает две независимые переменные: *X* и *Y*, и поток *A* хочет захватить блокировку для доступа к переменной *X*, в то время как поток *B* хочет захватить блокировку для доступа к переменной *Y* (как было бы в том случае, если бы один поток вызывал метод

`addUser`, в то время как другой поток вызвал метод `addQuery` в классе `ServerStatus`), то эти два потока не будут конкурировать за любые данные, даже если они будут конкурировать за блокировку.

Детализацию блокировки нельзя уменьшить, если для каждой операции требуются переменные. Это еще одна область, где “сырая” (*raw*) производительность и масштабируемость часто расходятся друг с другом; распространённые оптимизации, такие как кэширование часто вычисляемых значений, могут ввести “горячие поля” (*hot fields*), что приводит к ограничению масштабируемости.

Если бы вы реализовывали класс `HashMap`, перед вами встал бы выбор того, каким образом метод `size` должен был бы вычислять количество записей в экземпляре `Map`. Простейший способ реализации метода - подсчитывать количество записей при каждом вызове. Распространённая оптимизация заключается в обновлении отдельного счетчика, по мере добавления или удаления записей; это немножко увеличивает затраты на вызовы методов `put` или `remove`, за счёт необходимости поддержания счетчика в актуальном состоянии, но снижает затраты на вызов метода `size` с $O(n)$ до $O(1)$.

Сохранение отдельного счетчика для ускорения операций, подобных `size` и `isEmpty`, отлично работает в однопоточной или полностью синхронизированной реализации, но значительно затрудняет улучшение масштабируемости реализации, поскольку каждая операция, изменяющая данные в экземпляре `Map`, также должна обновлять совместно используемый счетчик. Даже если вы используете чередование блокировок для цепочек хэшей, синхронизация доступа к счетчику вновь поднимает проблему масштабируемости эксклюзивной блокировки. То, что выглядело как оптимизация производительности - кэширование результатов операции `size` - превратилось в ответственность за масштабируемость. В этом случае, счетчик называется *горячим полем* (*hot field*), потому что каждая операция изменения данных должна получать к нему доступ.

Класс `ConcurrentHashMap` позволяет избежать этой проблемы, за счёт наличия метода `size`, перечисляющего размеры полос (*stripes*) и увеличивающего возвращаемое количество за счёт добавления элементов в каждой полосе, вместо поддержки глобального счётчика. Чтобы избежать перечисления каждого элемента, класс `ConcurrentHashMap` поддерживает, для каждой полосы, отдельное поле счётчика, также защищенное блокировкой полосы.¹²¹

11.4.5 Альтернативы монопольным блокировкам

Третий способ смягчения последствий при возникновении конкуренции за доступ к блокировкам заключается в отказе от использования монопольных блокировок в пользу более удобных, для обеспечения параллелизма, средств управления совместно используемым состоянием. К ним относятся параллельные коллекции, блокировки на чтение-запись, неизменяемые объекты и атомарные переменные.

Интерфейс `ReadWriteLock` (см. главу [13](#)) обеспечивает соблюдение дисциплины блокировки с несколькими читателями и одним писателем: несколько читателей могут параллельно получать доступ к совместно используемому ресурсу,

¹²¹ Если метод `size` вызывается часто, по сравнению с изменяющими данные операциями, чередующиеся структуры данных могут оптимизировать его, путем кэширования размера коллекции в переменной типа `volatile` всякий раз, когда вызывается метод `size` и аннулирования кэша (путём установки его значения в `-1`) всякий раз, когда коллекция изменяется. Если кэшированное значение неотрицательно при вызове метода `size`, оно является точным и может быть возвращено; в противном случае, перед возвратом оно пересчитывается.

пока никто из них не захочет его изменить, но писатели должны захватывать блокировку монопольно. Для структур данных, используемых в основном для чтения, интерфейс `ReadWriteLock` может предложить большую степень параллелизма, чем предлагает монопольная блокировка; в структурах данных, используемых только для чтения, неизменяемость может в целом устраниć необходимость в захвате блокировки.

Атомарные переменные (см. главу [15](#)) позволяют снизить затраты на обновление “горячих полей”, таких как статистические счётчики, генераторы последовательностей или ссылки на первые узлы в связанных структурах данных. (Мы использовали класс `AtomicLong` для сохранения счётчика попаданий в примере с сервером, приведённом в главе [2](#).) Классы атомарных переменных обеспечивают очень детальные (и, следовательно, более масштабируемые) атомарные операции над целыми числами или ссылками на объекты и реализуются с использованием примитивов параллелизма низкого уровня (таких как операция проверить-затем-поменять, *compare-and-swap*), предоставляемых большинством современных процессоров. Если ваш класс имеет небольшое количество горячих полей, которые не участвуют в инвариантах с другими переменными, их замена атомарными переменными может улучшить масштабируемость. (Изменение вашего алгоритма для использования меньшего количества горячих полей может улучшить масштабируемость еще больше - атомарные переменные уменьшают затраты на обновление горячих полей, но не устраняют их.)

11.4.6 Мониторинг использования CPU

При тестировании на масштабируемость, целью обычно является полное использование процессоров. Такие инструменты, как `vmstat` и `mpstat` в системах Unix или `perfmon` в системах Windows, могут рассказать вам о том, насколько “торячо” процессоры работают.

Если процессоры используются асимметрично (некоторые процессоры раскалены под нагрузкой, а другие нет), вашей первой целью должен стать поиск возросшего параллелизма в вашей программе. Асимметричное использование означает, что большая часть вычислений выполняется в небольшом наборе потоков, и приложение не может воспользоваться преимуществами, предоставляемыми дополнительными процессорами.

Если процессоры используются не полностью, необходимо выяснить, почему. Существует несколько вероятных причин:

Недостаточная нагрузка. Может оказаться, что тестируемое приложение просто не подвергается достаточной нагрузке. Это можно проверить с помощью увеличения нагрузки и измерения показателей использования, времени отклика или времени обслуживания. Создание достаточной нагрузки для насыщения приложения может потребовать значительных вычислительных ресурсов; проблема может заключаться в том, что клиентские системы, а не тестируемая система, работают на полную мощность.

Ограничения ввода/вывода. Можно определить, ограничено ли приложение дисковой подсистемой, с помощью утилит `iostat` или `perfmon`, а также ограничена ли пропускная способность, за счёт мониторинга уровня трафика в сети.

Внешние ограничения. Если приложение зависит от внешних служб, таких как база данных или веб-служба, узкое место может быть не в вашем коде. Это можно проверить с помощью профилировщика или инструментов администрирования базы данных, чтобы определить, сколько времени тратится на ожидание ответов от внешней службы.

Конфликт блокировок. Инструменты профилирования могут определить, сколько конфликтов блокировок возникает в приложении и какие блокировки являются “горячими”. Часто можно получить ту же информацию и без профилировщика, путем случайной выборки, снимая несколько дампов потоков и ища потоки, конкурирующие за блокировки. Если поток заблокирован в ожидании блокировки, соответствующий кадр стека в дампе потока указывает на “ожидание блокировки монитора...”. Блокировки, которые в основном являются неконкурентными, редко появляются в дампе потока; высоко конкурентные блокировки, почти всегда имеют хотя бы один поток, ожидающий возможности захвата, и таким образом часто появляются в дампах потоков.

Если приложение поддерживает достаточно высокую температуру процессоров, можно использовать инструменты мониторинга, чтобы сделать выводы, даёт ли выигрыш в производительности использование дополнительных процессоров. Программа, имеющая четыре потока, может быть в состоянии полностью загрузить 4-процессорную систему, но маловероятно, что мы увидим повышение производительности при перемещении программы в 8-процессорную систему, так как возникнет необходимость в ожидании выполняемых потоков, чтобы воспользоваться преимуществами, предоставляемыми дополнительными процессорами. (Вы также можете перенастроить программу, чтобы разделить ее рабочую нагрузку на большее количество потоков, например, настроить размер пула потоков.) Один из столбцов, публикуемых утилитой `vmstat`, - это число потоков, которые могли бы выполняться, но не выполняются в данный момент, поскольку CPU недоступен; если загрузка CPU высока и всегда есть выполняемые потоки, ожидающие CPU, ваше приложение, вероятно, выиграет от увеличения числа процессоров.

11.4.7 Просто скажите “нет” помещению объектов в пул

В ранних версиях JVM, операции выделения памяти под объект и сборка мусора были медленными¹²², но с тех пор их производительность существенно улучшилась. На самом деле, выделение памяти в Java теперь быстрее, чем вызов функции `malloc` в C: общее в ветках кода, выполняющих команду `new Object` в HotSpot 1.4.x и 5.0 составляет примерно десять машинных команд.

Чтобы обойти “медленные” жизненные циклы объектов, многие разработчики обращаются к пулу объектов, в котором объекты перерабатываются (*recycled*), вместо сборки мусора и выделения заново при необходимости. Даже принимая во внимание снижение издержек на сборку мусора, было показано, что помещение объектов в пул приводит к потерям быстродействия¹²³ во всех случаях, кроме

¹²² Как было со всем остальным - синхронизацией, графикой, запуском JVM, рефлексией - то же самое предсказуемо и для других первых версий экспериментальной технологии.

¹²³ В дополнение к потерям, выраженным в циклах CPU, помещение объектов в пул влечёт за собой ряд других проблем, среди которых проблема определения корректного размера пула объектов (если размер слишком мал - помещение объектов в пул не окажет никакого влияния, если размер слишком велик - будет оказываться давление на сборщик мусора и будет удерживаться память, которая могла бы быть использована на что-то другое, то есть более эффективно); риск того, что объект не будет надлежащим

самых дорогостоящих объектов (и к серьезным потерям для объектов с легким и средним весом) в однопоточных программах (Click, 2005).

В параллельных приложениях “тарифы” за размещение в пуле еще хуже. Когда потоки выделяют (*allocate*) новые объекты, межпоточная координация требуется в очень малой степени, так как координаторы (*allocators*), отвечающие за выделение памяти, обычно используют локальные для потоков блоки выделения (*allocation blocks*), для устранения большей части синхронизации в структурах данных кучи. Но если потоки вместо этого запрашивают получение объекта из пула, необходима некоторая синхронизация для координации доступа к структуре данных пула, что приводит к возникновению возможности блокировки потока. Поскольку блокировка потока из-за конфликта блокировок в сотни раз дороже, чем операция выделения памяти, даже небольшое количество конфликтов, связанных с пулем, приведёт к возникновению узкого места в масштабируемости. (Даже неконкурентная синхронизация обычно обходится дороже, чем выделение объекта.) Это еще один подход, предназначенный для оптимизации производительности, но превратившийся в угрозу масштабируемости. Помещение объектов в пул находит своё применение¹²⁴, но имеет достаточно ограниченную полезность, в качестве средства оптимизации производительности.

Выделение объектов обычно дешевле, чем синхронизация.

11.5 Пример: Сравнение производительности реализаций Map

Однопоточная производительность класса `ConcurrentHashMap` немного лучше, чем у синхронизированной реализации `HashMap`, но при параллельном использовании, она проявляет себя во всей красе. Реализация `ConcurrentHashMap` предполагает, что наиболее распространенной операцией является получение уже существующего значения, и поэтому оптимизирована для обеспечения максимальной производительности и параллелизма при выполнении успешных операций `get`.

Основным препятствием для обеспечения масштабируемости в синхронизированных реализациях `Map`, является наличие единственной блокировки для всего экземпляра `Map`, поэтому одновременно только один поток может получить доступ к экземпляру `Map`. С другой стороны, класс `ConcurrentHashMap` не накладывает блокировку в большинстве случаев успешных операций чтения и использует чередование блокировок для операций записи и тех немногих операций чтения, которые нуждаются в блокировке. В результате, несколько потоков могут одновременно получить доступ к экземпляру `Map`, без возникновения блокировки.

На рисунке 11.3 иллюстрируются различия в масштабируемости между несколькими реализациями интерфейса `Map`: `ConcurrentHashMap`, `ConcurrentSkipListMap`, и обёрнутые с помощью метода `synchronizedMap` реализации `HashMap` и `TreeMap`. Первые две реализации имеют потокобезопасный дизайн; последние две сделаны потокобезопасными с помощью синхронизирующей обёртки. При каждом запуске, N потоков одновременно

образом сбрасываться в его вновь выделенное состояние, вводит вероятность появления незначительных ошибок; риск того, что поток будет возвращать объект в пул, но всё равно будет продолжать использовать его; и что это приводит к большему количеству работы для сборщиков мусора поколений, поощряя шаблон формирования ссылок от старого к молодому (*old-to-young references*).

¹²⁴ В ограниченных средах, таких как J2ME или RTSJ, помещение объектов в пул может требоваться для эффективного управления памятью или для управления отзывчивостью.

выполняют цикл нагрузочных операций, в котором выбирается случайный ключ и производится попытка получения значения, соответствующего этому ключу. Если значение отсутствует, оно добавляется в экземпляр Map с вероятностью $p = 0.6$, и если оно присутствует, удаляется с вероятностью $p = 0.02$. Тесты проводились в пред релизной сборке Java 6 на 8-ядерном процессоре Sparc V880, и график отражает пропускную способность, нормализованную к случаю использования одного потока в классе `ConcurrentHashMap`. (Разрыв масштабируемости между параллельными и синхронизированными коллекциями даже больше чем в Java 5.0.)

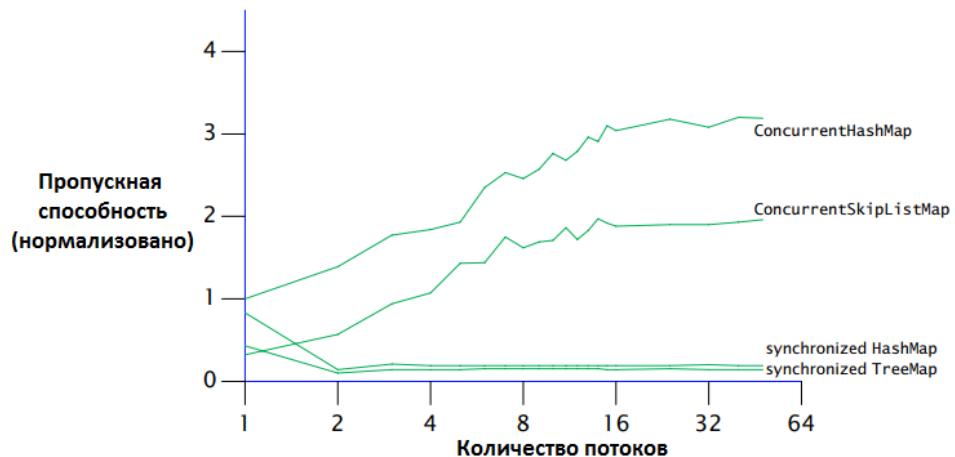


Рисунок 11.3. Сравнение производительности реализаций интерфейса Map

Данные для классов `ConcurrentHashMap` и `ConcurrentSkipListMap` показывают, что они хорошо масштабируются для большого числа потоков; пропускная способность продолжает улучшаться по мере добавления потоков. Хотя количество потоков, представленных на рис. 11.3 может показаться небольшим, эта тестовая программа создает большее количество конфликтов, в пересчёте на каждый поток, чем типичное приложение, так как она делает немного больше, чем заполнение экземпляра Map; реальная программа будет выполнять некоторую дополнительную, локальную для потока, работу в каждой итерации.

Цифры по синхронизированным коллекциям не столь обнадеживающие. Производительность для однопоточного случая сравнима с производительностью класса `ConcurrentHashMap`, но как только происходит переход нагрузки от “в основном неконкурентной” до “в основном конкурентной” - что происходит в приведённом в примере случае уже при двух потоках - производительность синхронизированных коллекций сильно страдает. Для кода, масштабируемость которого ограничивается конкуренцией за захват блокировок, это обычное явление. До тех пор, пока уровень конкуренции низок, во времени выполнения операции преобладает время фактически выполняемой работы, и пропускная способность может улучшаться по мере добавления потоков. Как только конкуренция становится существенной, во времени выполнения операции начинает преобладать время переключения контекста и задержки, вызванные затратами на планирование, и добавление большего количества потоков оказывает слабое влияние на пропускную способность.

11.6 Сокращение накладных расходов на переключение контекста

Множество задач включает операции, которые могут быть заблокированы; переход между состояниями выполнения (*running*) и блокировки (*blocked*) влечет за собой переключение контекста. Одним из источников блокировки в серверных приложениях является создание сообщений в логах, в процессе обработки запросов; чтобы проиллюстрировать, как можно повысить пропускную способность за счет уменьшения переключений контекста, мы проанализируем планируемое поведение двух подходов к логированию.

Большинство фреймворков логирования представляют собой тонкие обёртки вокруг вызова `println`; когда у вас есть что-то, предназначенное для логирования, просто пишите это прямо здесь и сейчас. Другой подход был продемонстрирован в классе `LogWriter` (см. раздел [7.2.1](#)): логирование выполняется в выделенном фоновом потоке, вместо потока, инициировавшего запрос. С точки зрения разработчика, оба подхода примерно равнозначны. Однако производительность может различаться в зависимости от объема логируемой информации, количества потоков, выполняющих ведение журнала, и других факторов, таких как затраты на переключение контекста¹²⁵.

Время обслуживания для операции ведения журнала включает в себя все вычисления, связанные с классами потоков ввода/вывода; если операция ввода/вывода блокируется, она также включает в себя время, в течение которого поток был заблокирован. Операционная система приостанавливает плановое выполнение заблокированного потока до завершения операции ввода/вывода и, возможно, немного дольше. Когда операция ввода/вывода завершается, другие потоки, вероятнее всего, активируются и им будет позволено завершить выполнение в отведённых им квантах планирования, и потоки, возможно, уже будут ожидать впереди очереди планирования – дальнейшее прибавка к времени обслуживания. В качестве альтернативы, если несколько потоков одновременно выполняют операцию логирования, может возникнуть конкуренция за блокировку выходного потока, и в этом случае результат будет таким же, как и при блокировке операции ввода/вывода – поток блокируется в ожидании блокировки и переключается. Встроенное логирование включает в себя операции ввода/вывода и блокировку, что может привести к увеличению переключения контекста и, следовательно, увеличению времени обслуживания.

Увеличение времени обслуживания запросов нежелательно по нескольким причинам. Во-первых, время обслуживания оказывает влияние на качество обслуживания: более длительное время обслуживания означает, что кто-то ждет результата дольше. Но что более важно, более длительное время обслуживания, в этом случае, означает, что происходит больше конфликтов за захват блокировок. Принцип “войти, выйти” (*get in, get out*), представленный в разделе [11.4.1](#), диктует нам, что мы должны освобождать блокировки как можно скорее, потому что чем дольше удерживается блокировка, тем больше вероятность того, что за возможность захвата блокировки возникнет конкуренция. Если поток блокируется

¹²⁵ Создание логгера, перемещающего операции ввода/вывода в другой поток, может повысить производительность, но также вносит ряд конструктивных сложностей, таких как прерывание (что произойдет, если поток, заблокированный в операции логирования, будет прерван?), гарантии обслуживания (гарантирует ли логгер, что сообщение, помещенное в очередь сообщений для логирования, будет записано до завершения работы службы?), политика насыщения (что произойдёт, если производители будут регистрировать сообщения быстрее, чем поток логгера сможет обработать их?), и жизненный цикл службы (каким образом мы завершаем работу логгера и каким образом мы сообщаем о состоянии службы производителям?).

на ожидании операций ввода/вывода, и в то же время удерживает другую блокировку, весьма вероятно, что другому потоку, также понадобится блокировка, удерживаемая первым потоком. Параллельные системы работают намного лучше, когда большинство операций захвата блокировок не конкуренты, поскольку захват блокировки на конкурентной основе означает большее количество переключений контекста. Стиль кодирования, который поощряет большее количество переключений контекста, в результате приводит к более низкой общей пропускной способности.

Перемещение операций ввода/вывода из обрабатывающего запросы потока может сократить среднее время обслуживания обработки запросов. Потоки, вызывающие метод `log`, больше не блокируются в ожидании блокировки выходного потока или в ожидании завершения операций ввода/вывода; им нужно только поместить сообщение в очередь, а затем вернуться к своей задаче. С другой стороны, мы ввели возможность возникновения конкуренции за доступ к очереди сообщений, но операция `put` имеет меньший вес, чем логирующий ввод/вывод (который может потребовать системных вызовов), и поэтому с меньшей вероятностью будет блокироваться при фактическом использовании (пока очередь не заполнена). Поскольку поток, инициировавший запрос, теперь с меньшей вероятностью будет блокироваться, вероятность переключения контекста в середине выполнения запроса снижается. То, что мы сделали, превратило сложную и неопределенную ветку исполнения кода, включающую в себя операции ввода/вывода и возможное возникновение конкуренции за блокировку, в достаточно прямолинейную ветку исполняемого кода.

В какой-то степени мы просто перемещаем работу, за счёт перемещения операций ввода/вывода в поток, в котором затраты на них не будут восприниматься пользователем (что само по себе может быть выигрышем). Но, перемещая *все* операции ввода/вывода логов в один поток, мы также исключаем вероятность возникновения конкуренции за выходной поток и, таким образом, устранием источник блокировки. Это повышает общую пропускную способность, поскольку на планирование, переключение контекста и управление блокировками расходуется меньше ресурсов.

Перемещение операций ввода/вывода из множества обрабатывающих запросы на логирование потоков в один логирующий поток, аналогично различию между бригадой с ведрами и группой лиц, борющихся с пожаром. С подходом "сто парней, бегающих с ведрами", у вас больше шансов на возникновение конкуренции за доступ к источнику с водой и к пламени (в результате чего к пламени будет доставлено меньшее количество воды), плюс большая неэффективность, потому что каждый работник постоянно переключает режимы (заполнение, бег, сброс, бег и т. д.). При Ведёрно-бригадном подходе, поток воды от источника к горящему зданию постоянный, затрачивается меньше энергии на транспортировку воды к огню, и каждый работник непрерывно сосредотачивается на выполнении одной работы. Точно так же, как прерывания разрушительны и снижают производительность людей, блокировка и переключение контекста разрушительны для потоков.

11.7 Итоги

Поскольку одной из наиболее распространенных причин использования потоков является желание использовать несколько процессоров, при обсуждении производительности параллельных приложений нас обычно больше беспокоит пропускная способность или масштабируемость, чем "сыре" время обслуживания.

Закон Амдала гласит, что масштабируемость приложения пропорциональна доле последовательно выполняемого кода. Поскольку основным источником последовательно выполняемого кода, в программах Java, является монопольная блокировка ресурсов, масштабируемость часто можно улучшить, затрачивая меньше времени на удержание блокировок, либо уменьшая степень детализации блокировок, уменьшая продолжительность удержания блокировок, либо заменяя монопольные блокировки не монопольными или неблокирующими альтернативами.

Глава 12 Тестирование параллельных программ

Параллельные программы используют принципы проектирования и шаблоны, подобные тем, что используются в последовательных программах. Разница в том, что параллельные программы, в отличие от последовательных, обладают некоторой степенью недетерминированности, что приводит к увеличению количества потенциальных взаимодействий и ситуаций, в которых происходят сбои, и это должно быть проанализировано и предусмотрено.

Аналогично, тестирование параллельных программ использует и расширяет идеи, подчерпнутые из тестирования последовательных программ. Те же самые методы, что применяются для проверки корректности и производительности в последовательных программах, могут быть применены к параллельным программам, но с параллельными программами пространство вариантов, которые могут пойти не так, как было запланировано, гораздо больше. Основная проблема при построении тестов для параллельных программ заключается в том, что потенциальные сбои могут быть, с большой вероятностью, редкими, а не детерминированными; тесты, раскрывающие такие сбои, должны быть более обширными и выполняться дольше, чем типичные последовательные тесты.

Большинство тестов параллельных классов подпадают под одну или обе классические категории: *безопасность* и *живучесть*. В главе 1 мы определили безопасность, как “ничего плохого никогда не случается”, а живость как “что-то хорошее когда-нибудь случается”.

Тесты на безопасность, которые проверяют, что поведение класса соответствует его спецификации, обычно принимают форму тестирования инвариантов. Например, в реализации связанного списка, которая кэширует размер списка при каждом его изменении, одним из тестов на безопасность будет сравнение кэшированного количества с фактическим количеством элементов в списке. В однопоточной программе это легко, так как содержимое списка не изменяется при тестировании его свойств. Но в параллельной программе такой тест, вероятно, будет чреват гонками, если вы не сможете наблюдать состояние поля подсчета и подсчитывать элементы в одной атомной операции. Это можно выполнить, заблокировав список для монопольного доступа, и применив какую-то функцию, предоставляемую реализацией, для создания “атомарного моментального снимка” (*atomic snapshot*) или используя “тестовые точки”, предоставляемые реализацией, которые позволяют тестам утверждать состояние инвариантов (*assert invariants*) или выполнять тестовый код атомарно.

В этой книге мы использовали временные диаграммы, чтобы изобразить “неудачные” взаимодействия, которые могут приводить к сбоям в неправильно построенных классах; тестовые программы пытаются подобрать достаточное пространство состояний, в котором, в конечном итоге, и случается неудача. К сожалению, тестовый код может вводить задержки или артефакты синхронизации, которые могут приводить к маскировке ошибок, которые, в противном случае, могли бы проявляться¹²⁶.

Свойства живучести представляют тестированию свои собственные вызовы. Тесты на живучесть включают в себя тесты на достижение прогресса и на отсутствие прогресса, которые трудно оценить количественно - как вы можете

¹²⁶ Ошибки, которые исчезают при добавлении отладки или тестового кода игриво называются Heisenbugs (гейзенбаг).

проверить, что метод блокируется, а не просто медленно работает? Аналогично, каким образом проверить, что алгоритм *не* попадает в состояние взаимоблокировки? Сколько вы должны ожидать, прежде чем объявить, что тест провален?

С тестами на живучесть связаны тесты производительности. Производительность можно измерить несколькими способами, в том числе:

Пропускная способность: скорость выполнения набора параллельных задач;

Отзывчивость: задержка между моментом поступления запроса и выполнением какого-либо действия (также называемая “временем ожидания”, *latency*); или

Масштабируемость: улучшение пропускной способности (или отсутствие такового) по мере увеличения доступности ресурсов (обычно ЦП).

12.1 Тестирование на корректность

Разработка модульных тестов для параллельного класса начинается с проведения того же анализа, что и для последовательного класса - определения инвариантов и постусловий, поддающихся механической проверке. Если вам повезёт, многие из них присутствуют в спецификации; в остальное время написание тестов - это приключение по итеративному разбору спецификации.

В качестве конкретной иллюстрации мы создадим набор тестовых случаев для ограниченного буфера. В листинге 12.1 показана реализация класса `BoundedBuffer`, использующая экземпляр `Semaphore` для реализации требуемого ограничения и блокировки.

```
@ThreadSafe
public class BoundedBuffer<E> {
    private final Semaphore availableItems, availableSpaces;
    @GuardedBy("this") private final E[] items;
    @GuardedBy("this") private int putPosition = 0, takePosition = 0;

    public BoundedBuffer(int capacity) {
        availableItems = new Semaphore(0);
        availableSpaces = new Semaphore(capacity);
        items = (E[]) new Object[capacity];
    }
    public boolean isEmpty() {
        return availableItems.availablePermits() == 0;
    }
    public boolean isFull() {
        return availableSpaces.availablePermits() == 0;
    }

    public void put(E x) throws InterruptedException {
        availableSpaces.acquire();
        doInsert(x);
        availableItems.release();
    }
```

```

public E take() throws InterruptedException {
    availableItems.acquire();
    E item = doExtract();
    availableSpaces.release();
    return item;
}

private synchronized void doInsert(E x) {
    int i = putPosition;
    items[i] = x;
    putPosition = (++i == items.length)? 0 : i;
}

private synchronized E doExtract() {
    int i = takePosition;
    E x = items[i];
    items[i] = null;
    takePosition = (++i == items.length)? 0 : i;
    return x;
}
}

```

Листинг 12.1 Ограниченнный буфер с использованием экземпляра Semaphore

Класс `BoundedBuffer` реализует очередь на основе массива фиксированной длины с блокирующими методами `put` и `take`, управляемыми парой подсчитывающих семафоров. Семафор `availableItems` представляет собой количество элементов, которое может быть *удалено* из буфера, и изначально равен нулю (так как буфер изначально пуст). Аналогично, семафор `availableSpaces` представляет собой количество элементов, которое может быть *добавлено* в буфер, и инициализируется размером буфера.

Операция `take` требует, чтобы сначала было получено разрешение от семафора `availableItems`. Эта операция выполняется немедленно, если буфер не пустой, в противном случае операция блокируется, до появления элементов в буфере. Как только разрешение получено, следующий элемент из буфера удаляется, и освобождается разрешение семафора `availableSpaces`¹²⁷. Операция `put` определяется обратным образом, так что при выходе из методов `put` или `take` сумма счетчиков обоих семафоров всегда равна значению границы. (На практике, если вам нужен ограниченный буфер, вы должны использовать классы `ArrayBlockingQueue` или `LinkedBlockingQueue`, а не писать своё собственное решение, но используемый здесь подход иллюстрирует, как вставки и удаления могут контролироваться и в других структурах данных.)

12.1.1 Простые модульные тесты

Самые основные модульные тесты для класса `BoundedBuffer` похожи на те, что мы используем в последовательном контексте - создаем ограниченный буфер, вызываем его методы и выполняем утверждение постусловий и инвариантов. Некоторые инварианты из тех, что быстрее всего приходят на ум, это то, что только что созданный буфер должен идентифицировать себя и как пустой, и как не полный. Аналогичен, но чуть более сложен, тест на безопасность при вставке N

¹²⁷ В подсчитывающем семафоре разрешения явно не представляются и не связаны с потоком-владельцем; операция `release` создает разрешение, а операция `acquire` использует его.

элементов в буфер с емкостью N (который должен успешно завершиться без блокировки), и проверить, что буфер распознает, что он полон (а не является пустым). Тестовые методы JUnit для этих свойств показаны в листинге 12.2.

```
class BoundedBufferTest extends TestCase {  
    void testIsEmptyWhenConstructed() {  
        BoundedBuffer<Integer> bb = new BoundedBuffer<Integer>(10);  
        assertTrue(bb.isEmpty());  
        assertFalse(bb.isFull());  
    }  
  
    void testIsFullAfterPuts() throws InterruptedException {  
        BoundedBuffer<Integer> bb = new BoundedBuffer<Integer>(10);  
        for (int i = 0; i < 10; i++)  
            bb.put(i);  
        assertTrue(bb.isFull());  
        assertFalse(bb.isEmpty());  
    }  
}
```

Листинг 12.2 Простой модульный тест для класса BoundedBuffer

Эти простые тестовые методы полностью последовательны. Включение набора последовательных тестов в ваш набор тестов часто полезно, так как они могут раскрывать ситуации, когда проблема *не* связана с проблемами параллелизма, перед началом поиска гонок данных.

12.1.2 Тестирование блокирующих операций

Для проверки основных свойств параллелизма требуется введение нескольких потоков. Большинство фреймворков для тестирования не очень дружелюбны по отношению к параллелизму: они редко включают средства для создания потоков или мониторинга их состояния, чтобы гарантировать, что не происходит непредвиденного завершения. Если вспомогательный поток, созданный тестовым случаем, натыкается на сбой, фреймворк обычно не знает, с каким тестом связан поток, таким образом, для того, чтобы соотнести информацию об успехе или неудаче, могут потребоваться некоторые дополнительные усилия, для предоставления данных, при возврате к главному, выполняющему тесты, потоку, в удобном для составления отчётов виде.

Для соответствия тестов пакету `java.util.concurrent`, очень важно, чтобы сбои четко соотносились с конкретным тестом. Поэтому, группой экспертов, разрабатывающей JSR 166, был создан базовый класс¹²⁸, предоставивший методы для передачи и формирования отчётов об ошибках, во время выполнения метода `tearDown`, следя соглашению, о том, что каждый тест должен ожидать, пока не завершатся все созданные им потоки. Скорее всего, вам нет нужды так глубоко погружаться в вопрос; основное требование заключается в том, чтобы было ясно, успешно ли были выполнены тесты, и что информация о сбое куда-то отправляется, для использования при диагностике проблемы.

Если предполагается, что метод блокируется при определенных условиях, то проверка такого поведения должна завершиться успешно только в том случае, если

¹²⁸ <http://gee.cs.oswego.edu/cgi-bin/viewcvs.cgi/jsr166/src/test/tck/JSR166TestCase.java>

поток не выполняется. Проверка того, что метод блокируется, подобна проверке того, что метод бросает исключение; если метод возвращает управление нормально, выполнение теста провалено.

Тестирование на предмет того, что метод блокируется, вводит дополнительные сложности: как только метод успешно блокируется, вы должны вынудить его каким-то образом разблокироваться. Очевидным способом сделать это, является использование прерывания - запустить блокирующую действие в отдельном потоке, подождать, пока поток не заблокируется, прервать его, а затем выдвинуть утверждение (*assert*), что блокирующая операция завершена. Конечно, это требует, чтобы ваши блокирующие методы реагировали на прерывания, раньше возвращая управление или бросая исключение `InterruptedException`.

Про часть “ждать, пока поток заблокирован” легче сказать, чем сделать; на практике вы должны принять произвольное решение о том, сколько времени может занять выполнение нескольких инструкций, и ждать дольше. Вы должны быть готовы увеличить это значение, если вы не правы (в этом случае вы увидите ложные сбои в процессе тестирования).

В листинге 12.3 демонстрируется подход к тестированию блокирующих операций. В нём создаётся поток “получателя”, который пытается взять элемент из пустого буфера. Если метод `take` выполняется успешно, то поток регистрирует сбой. Выполняющий тесты поток запускает поток получателя, ожидает долгое время, а затем прерывает его. Если поток получателя корректно заблокирован в операции `take`, будет брошено исключение `InterruptedException`, и блок `catch`, для этого исключения, будет рассматривать его как успешное прохождение теста, и позволит потоку вернуть управление. Затем, основной выполняющий тесты поток пытается выполнить метод `join` потока получателя и проверяет, что возврат управления из метода `join` произошёл успешно, путем вызова метода `Thread.isAlive`; если поток получателя среагировал на прерывание, выполнение метода `join` должно завершиться быстро.

```
void testTakeBlocksWhenEmpty() {
    final BoundedBuffer<Integer> bb = new BoundedBuffer<Integer>(10);
    Thread taker = new Thread() {
        public void run() {
            try {
                int unused = bb.take();
                fail(); // if we get here, it's an error
            } catch (InterruptedException success) { }
        }
    };
    try {
        taker.start();
        Thread.sleep(LOCKUP_DETECT_TIMEOUT);
        taker.interrupt();
        taker.join(LOCKUP_DETECT_TIMEOUT);
        assertFalse(taker.isAlive());
    } catch (Exception unexpected) {
        fail();
    }
}
```

Листинг 12.3 Тестирование блокировки и отзывчивости на прерывание

Ограниченнная по времени версия метода `join` гарантирует, что тест завершится, даже если выполнение метода `take` каким-то неожиданным образом застопорится. Этот тестовый метод проверяет несколько свойств метода `take` - не только то, что он блокируется, но и то, что при прерывании он бросает исключение `InterruptedException`. Это один из тех немногих случаев, при которых уместно явно использовать подкласс класса `Thread`, вместо использования экземпляра `Runnable` в пуле: для тестирования надлежащего завершения с использованием метода `join`. Подобный подход можно также использовать для проверки разблокировки потока получателя, после помещения элемента в очередь основным потоком.

Заманчиво использовать метод `Thread.getState`, для проверки того, что поток фактически заблокирован на условии ожидания, но этот подход не надежен. Нет такого требования, чтобы заблокированный поток всегда входил в состояние `WAITING` или `TIMED_WAITING`, так как JVM, вместо этого, может выбрать для реализации блокировки ожидание спина. Точно так же, в связи с тем, что разрешены ложные пробуждения (*spurious wakeups*) при выполнении методов `Object.wait` или `Condition.await` (см. главу [14](#)), поток, находящийся в состоянии `WAITING` или `TIMED_WAITING`, может временно перейти в состояние `RUNNABLE`, даже если условие, выполнение которого он ожидает, еще не истинно. Даже игнорируя эти параметры реализации, целевому потоку может потребоваться некоторое время для перехода в заблокированное состояние. *Результат, возвращаемый методом `Thread.getState`, не должен использоваться для управления параллелизмом, и имеет ограниченную полезность для тестирования - его основная утилита является источником отладочной информации.*

12.1.3 Тестирование безопасности

Тесты, приведённые в листингах 12.2 и 12.3, проверяют важные свойства ограниченного буфера, но вряд ли позволят выявить ошибки, связанные с гонками данных. Чтобы проверить, что параллельный класс работает правильно при непредсказуемом параллельном доступе, нам нужно настроить несколько потоков, выполняющих операции `put` и `take` в течение некоторого времени, а затем как-то проверить, что всё прошло так, как было запланировано.

Построение тестов для выявления ошибок безопасности в параллельных классах является проблемой “курицы и яйца”: тестовые программы сами по себе являются параллельными программами. Разработка хороших параллельных тестов может оказаться сложнее разработки классов, которые они тестируют.

Задача построения эффективных тестов на безопасность, для параллельных классов, заключается в определении легко проверяемых свойств, которые с высокой вероятностью завершатся неудачей, если что-то пойдет не так, и в то же время не позволят коду, выполняющему аудит отказов, искусственно ограничить параллелизм. Лучше всего, если проверка тестируемого свойства не требует синхронизации.

Один подход, который хорошо работает с классами, используемыми в шаблоне производитель-потребитель (подобным классу `BoundedBuffer`) заключается в том, чтобы проверить, что все, что помещается в очередь или буфер покидает её, и что больше ничего не происходит. Наивная реализация этого подхода будет помещать элемент в “теневой” (*shadow*) список, когда он помещается в очередь, удалять его

из списка, когда он удаляется из очереди, и выдвинет утверждение, что теневой список пуст, когда тест будет завершен. Но такой подход исказит расписание запуска тестовых потоков, поскольку для изменения теневого списка потребуется синхронизация и, возможно, блокировка.

Распространение этого подхода на ситуацию с несколькими производителями и потребителями требует использования функции, вычисляющей контрольную сумму, *нечувствительную* к порядку, в котором элементы объединяются, чтобы после завершения теста можно было объединить несколько контрольных сумм. В противном случае, синхронизация доступа к совместно используемому полю контрольной суммы может стать узким местом параллелизма или исказить время выполнения теста. (Любые коммутативные операции, такие как сложение или XOR, отвечают этим требованиям.)

Чтобы убедиться, что ваш тест действительно проверяет то, о чём вы думаете, важно, чтобы сами контрольные суммы не были угадываемыми компилятором. Было бы плохой идеей использовать последовательные целые числа в качестве тестовых данных, потому что тогда результат всегда будет одинаковым, и умный компилятор мог бы просто предварительно вычислить его.

Чтобы избежать этой проблемы, тестовые данные должны генерироваться случайным образом, но многие другие эффективные тесты компрометируются плохим выбором генератора случайных чисел (*random number generator, RNG*). Генерация случайных чисел может создавать связи между классами и артефактами синхронизации, так как большинство классов генераторов случайных чисел потокобезопасны и поэтому обеспечивают дополнительную синхронизацию¹²⁹. Предоставление каждому потоку собственного экземпляра *RNG*, позволяет использовать *RNG*, не являющиеся потокобезопасными.

Вместо использования *RNG* общего назначения лучше использовать простые псевдослучайные функции. Вам не нужна высококачественная степень случайности; все, что вам нужно, это достаточно случайное значение, чтобы обеспечить изменение чисел от запуска к запуску. Функция *xorShift* в листинге 12.4, (Marsaglia, 2003) среди самых дешевых функций генерирующих случайные числа среднего качества. Её запуск со значениями, основанными на результатах методов *hashCode* и *nanoTime*, делает суммы как неочевидными, так и почти всегда отличными, для каждого последующего запуска.

```
static int xorShift(int y) {
    y ^= (y << 6);
    y ^= (y >>> 21);
    y ^= (y << 7);
    return y;
}
```

Листинг 12.4 Генератор случайных чисел среднего качества, подходящий для тестирования

Класс *PutTakeTest*, представленный в листингах 12.5 и 12.6, запускает N потоков-производителей, которые генерируют элементы и помещают их в очередь, и N потоков-потребителей, которые извлекают элементы из очереди. Каждый поток обновляет контрольную сумму элементов по мере их поступления или изъятия, используя индивидуальную, для каждого потока, контрольную сумму, которая

¹²⁹ Множество тестов на производительность, без ведома разработчиков или пользователей, фактически представляет собой тесты того, насколько велико узкое место в параллелизме, введённое использованием *RNG*.

объединяется в конце выполнения теста, чтобы не добавлять большие синхронизации или конкуренции, чем требуется для тестирования буфера.

```
public class PutTakeTest {  
    private static final ExecutorService pool  
        = Executors.newCachedThreadPool();  
    private final AtomicInteger putSum = new AtomicInteger(0);  
    private final AtomicInteger takeSum = new AtomicInteger(0);  
    private final CyclicBarrier barrier;  
    private final BoundedBuffer<Integer> bb;  
    private final int nTrials, nPairs;  
  
    public static void main(String[] args) {  
        new PutTakeTest(10, 10, 100000).test(); // sample parameters  
        pool.shutdown();  
    }  
  
    PutTakeTest(int capacity, int npairs, int ntrials) {  
        this.bb = new BoundedBuffer<Integer>(capacity);  
        this.nTrials = ntrials;  
        this.nPairs = npairs;  
        this.barrier = new CyclicBarrier(npairs * 2 + 1);  
    }  
  
    void test() {  
        try {  
            for (int i = 0; i < nPairs; i++) {  
                pool.execute(new Producer());  
                pool.execute(new Consumer());  
            }  
            barrier.await(); // wait for all threads to be ready  
            barrier.await(); // wait for all threads to finish  
            assertEquals(putSum.get(), takeSum.get());  
        } catch (Exception e) {  
            throw new RuntimeException(e);  
        }  
    }  
}  
  
class Producer implements Runnable { /* Listing 12.6 */ }  
  
class Consumer implements Runnable { /* Listing 12.6 */ }  
}
```

Листинг 12.5 Тест на основе шаблона производитель-потребитель, для класса BoundedBuffer

В зависимости от платформы создание и запуск потока может быть умеренно тяжелой операцией. Если ваш поток является кратковременным, и вы запускаете несколько потоков в цикле, в худшем случае потоки выполняются последовательно, а не параллельно. Тот факт, что даже не в самом худшем случае первый поток имеет преимущество над другими, означает, что вы можете получить меньше наложений, чем ожидалось: первый поток выполняется сам по себе в течение некоторого времени, а затем первые два потока, в течение некоторого

времени, выполняются параллельно, и только со временем все потоки начинают работать параллельно. (То же самое происходит и в конце выполнения: потоки, которые получили фору, также завершаются раньше.)

```
/* inner classes of PutTakeTest (Listing 12.5) */
class Producer implements Runnable {
    public void run() {
        try {
            int seed = (this.hashCode() ^ (int)System.nanoTime());
            int sum = 0;
            barrier.await();
            for (int i = nTrials; i > 0; --i) {
                bb.put(seed);
                sum += seed;
                seed = xorShift(seed);
            }
            putSum.getAndAdd(sum);
            barrier.await();
        } catch (Exception e) {
            throw new RuntimeException(e);
        }
    }
}

class Consumer implements Runnable {
    public void run() {
        try {
            barrier.await();
            int sum = 0;
            for (int i = nTrials; i > 0; --i) {
                sum += bb.take();
            }
            takeSum.getAndAdd(sum);
            barrier.await();
        } catch (Exception e) {
            throw new RuntimeException(e);
        }
    }
}
```

Листинг 12.6 Классы производитель и потребитель, используемые в классе PutTakeTest

Мы представили подход для смягчения этой проблемы в разделе [5.5.1](#), используя один экземпляр `CountDownLatch` в качестве начального затвора, а другой - в качестве конечного. Другой способ получить тот же эффект заключается в том, чтобы использовать класс `CyclicBarrier`, инициализированный числом рабочих потоков плюс один, и имеющиеся рабочие потоки и тестовый драйвер, ожидающие у барьера в моменты начала и завершения выполнения. Это гарантирует, что все потоки запускаются и находятся в рабочем состоянии, до фактического начала выполнения работы. Класс `PutTakeTest` использует этот подход, чтобы координировать запуск и остановку рабочих потоков, потенциально создавая большее количество чередований параллелизма. Мы все еще не можем

гарантировать, что планировщик не будет обрабатывать каждый поток последовательно, до самого завершения, но сделав выполнение достаточно длительным, мы уменьшили степень, с которой планирование вносит искажение в наши результаты.

Последним трюком, использованным классом `PutTakeTest`, является использование детерминированного критерия завершения, так что никакой дополнительной координации между потоками, для выяснения, завершён ли тест, не требуется. Метод тестирования запускает ровно тоже количество производителей, как и потребителей, и каждый из них отправляет (`put`) или принимает (`take`) одинаковое количество элементов, поэтому общее количество добавленных и удаленных элементов одинаково.

Тесты, подобные классу `PutTakeTest`, как правило, хороши при поиске нарушений безопасности. Например, распространенная ошибка при реализации буферов, управляемых семафорами, заключается в том, что забывают о том, что код, фактически выполняющий вставку и извлечение, требует взаимного исключения (с помощью блока `synchronized` или с использованием класса `ReentrantLock`). Пример запуска класса `PutTakeTest` с версией класса `BoundedBuffer`, в которой опущена синхронизация методов `doInsert` и `doExtract`, довольно быстро терпит неудачу. Запуск класса `PutTakeTest` с несколькими десятками потоков, выполняющими итерацию несколько миллионов раз на буферах различной мощности, в различных системах, повышает нашу уверенность в отсутствии повреждения данных в методах `put` и `take`.

Тесты должны выполняться на многопроцессорных системах для потенциального увеличения разнообразия чередования. Однако наличие более чем нескольких процессоров не обязательно делает тесты более эффективными. Чтобы максимально увеличить вероятность обнаружения гонки данных, зависящей от момента времени, должно быть больше активных потоков, чем кол-во доступных ЦП, чтобы в любой момент времени некоторые потоки выполнялись, а некоторые отключались, что снижает предсказуемость взаимодействий между потоками.

Для тестов, выполняемых до тех пор, пока не будет выполнено фиксированное число операций, возможна ситуация, при которой тестовый случай никогда не завершится, если тестируемый код поймает исключение из-за возникновения ошибки. Самый распространенный способ это исправить заключается в том, чтобы тестовый фреймворк прерывал процесс тестирования, который не завершается в течение определенного промежутка времени; сколько времени ожидать завершения, следует определять опытным путем, а сбои должны быть проанализированы, чтобы убедиться, что проблема заключается не в том, что вы ждёте не достаточно долго. (Эта проблема не является уникальной для процесса тестирования параллельных классов; последовательные тесты также должны различать длительные и бесконечные циклы.)

12.1.4 Тестирование управления ресурсами

Тесты до сих пор были связаны с соблюдением классом его спецификации – контроль того, что он делает то, что он должен делать. Второстепенным аспектом тестирования является проверка того, что он *не* делает тех вещей, которые он *не*

должен делать, например, такие как утечка ресурсов. Любой объект, который содержит другие объекты или управляет ими, не должен продолжать удерживать ссылки на эти объекты дольше, чем необходимо. Такие утечки памяти не позволяют сборщикам мусора освобождать память (или потоки, дескрипторы файлов, сокеты, подключения к базе данных или другие ограниченные ресурсы) и могут привести к исчерпанию ресурсов и сбою приложения.

Проблемы управления ресурсами особенно важны для классов, подобных классу `BoundedBuffer` – в целом причина ограничения буфера заключается в предотвращении сбоя приложения из-за исчерпания ресурсов, когда производители слишком опережают потребителей. Ограничение заставляет чрезмерно продуктивных производителей блокироваться, а не продолжать создавать работу, которая будет потреблять все больше и больше памяти или других ресурсов.

Нежелательное удержание памяти можно легко протестировать с помощью инструментов проверки кучи, которые измеряют использование памяти приложением; это можно сделать с помощью различных коммерческих инструментов профилирования кучи, или с помощью инструментов с открытым исходным кодом. Метод `testLeak` в листинге 12.7 содержит маркеры для создания снимка кучи инструментом инспектирования кучи, который принудительно запускает сборщик мусора¹³⁰, а затем записывает информацию о размере кучи и использовании памяти.

```
class Big { double[] data = new double[100000]; }

void testLeak() throws InterruptedException {
    BoundedBuffer<Big> bb = new BoundedBuffer<Big>(CAPACITY);

    int heapSize1 = /* snapshot heap */;
    for (int i = 0; i < CAPACITY; i++)
        bb.put(new Big());
    for (int i = 0; i < CAPACITY; i++)
        bb.take();

    int heapSize2 = /* snapshot heap */;
    assertTrue(Math.abs(heapSize1 - heapSize2) < THRESHOLD);
}
```

Листинг 12.7 Тестирование на предмет утечек ресурсов

Метод `testLeak` вставляет несколько больших объектов в ограниченный буфер, а затем удаляет их; использование памяти в моментальном снимке кучи #2 должно быть примерно таким же, как в моментальном снимке кучи #1. С другой стороны, если метод `doExtract` забыл обнулить ссылку на возвращаемый элемент (`items[i]=null`), использование памяти в двух моментальных снимках, определенно, не будет одинаковым. (Это один из немногих случаев, когда необходимо явное обнуление; большую часть времени оно либо не полезно, либо фактически вредно [ЕJ пункт 5].)

¹³⁰ Технически, невозможно заставить сборщик мусора начать уборку; метод `System.gc` только предлагает JVM выполнить уборку, так как это может быть хорошим моментом времени для выполнения сборки мусора. Среда HotSpot может быть проинструктирована игнорировать вызов метода `System.gc`, с помощью параметров запуска командной `-XX:+DisableExplicitGC`.

12.1.5 Использование обратных вызовов

Обратные вызовы к клиентскому коду могут быть полезны при построении тестовых случаев; обратные вызовы часто выполняются в известных точках жизненного цикла объекта, которые предоставляют хорошие возможностями для утверждения инвариантов. Например, класс `ThreadPoolExecutor` выполняет вызовы задач - экземпляров `Runnable`, а также экземпляра `ThreadFactory`.

Тестирование пула потоков включает в себя тестирование ряда элементов политики выполнения: проверка того, что дополнительные потоки создаются только тогда, когда это ожидается, но не тогда, когда они не ожидаются; что простоявающие потоки пожинаются, когда это необходимо, и т.д. Создание комплексного набора тестов, охватывающего все возможности, является серьезной задачей, но многие аспекты могут быть достаточно просто протестированы индивидуально.

Мы можем инструментировать создание потока, путем использования собственной фабрики потока. Класс `TestingThreadFactory` представленный в листинге 12.8, поддерживает подсчет количества созданных потоков; затем тестовые случаи могут проверить количество потоков, созданных во время тестового запуска. Класс `TestingThreadFactory` может быть расширен, чтобы вернуть собственную реализацию класса `Thread`, который также записывает, когда поток завершается, так что тестовые случаи могут проверить, что потоки пожинаются в соответствии с политикой выполнения.

```
class TestingThreadFactory implements ThreadFactory {
    public final AtomicInteger numCreated = new AtomicInteger();
    private final ThreadFactory factory
        = Executors.defaultThreadFactory();

    public Thread newThread(Runnable r) {
        numCreated.incrementAndGet();
        return factory.newThread(r);
    }
}
```

Листинг 12.8 Фабрика потоков для тестирования класса `ThreadPoolExecutor`

Если корневой размер пула меньше максимального размера, размер пула потоков должен увеличиваться по мере увеличения спроса на выполнение. При отправке долговременных задач в пул потоков, количество выполняемых задач остается постоянным достаточно долго, так что можно сделать несколько утверждений, например, что пул расширяется должным образом, как показано в листинге 12.9.

```
public void testPoolExpansion() throws InterruptedException { int
MAX_SIZE = 10;
ExecutorService exec = Executors.newFixedThreadPool(MAX_SIZE);

for (int i = 0; i < 10 * MAX_SIZE; i++)
    exec.execute(new Runnable() {
        public void run() {
            try {
```

```

        Thread.sleep(Long.MAX_VALUE);
    } catch (InterruptedException e) {
        Thread.currentThread().interrupt();
    }
}
});
for (int i = 0;
    i < 20 && threadFactory.numCreated.get() < MAX_SIZE;
    i++)
    Thread.sleep(100);
assertEquals(threadFactory.numCreated.get(), MAX_SIZE);
exec.shutdownNow();
}

```

Листинг 12.9 Тестовый метод для проверки расширения пула потоков

12.1.6 Увеличение степени чередования

Поскольку многие из потенциальных сбоев в параллельном коде являются маловероятными событиями, тестирование на наличие ошибок параллелизма - это игра чисел, но есть некоторые вещи, которые вы можете сделать, чтобы улучшить свои шансы. Мы уже упоминали о том, что работа на многопроцессорных системах с меньшим количеством процессоров, чем количество активных потоков, может генерировать больше чередований, чем в случае однопроцессорной системы или одного потока со многими процессорами. Аналогичным образом, тестирование на множестве систем с разным количеством процессоров, операционных систем и процессорных архитектур может выявить проблемы, которые могут возникать не во всех системах.

Полезный трюк для увеличения числа чередований и, следовательно, более эффективного изучения пространства состояний ваших программ, заключается в использовании метода `Thread.yield`, для поощрения большего количества переключений контекста во время операций, которые обращаются к совместно используемому состоянию. (Эффективность этого метода зависит от платформы, так как JVM свободна в своём решении, обрабатывать вызов `Thread.yield` как по-определению [JLS 17.9]; использование короткого, но ненулевого вызова `sleep` было бы медленнее, но более надежным.) Метод, представленный в листинге 12.10, осуществляет передачу кредитов с одного счета на другой; между двумя операциями обновления, инварианты подобные “сумма всех счетов равна нулю” не проводятся¹³¹. Иногда, уступая в процессе выполнения операции, вы можете активировать чувствительные ко времени ошибки в коде, который не использует адекватную синхронизацию для доступа к состоянию. Неудобство добавления этих вызовов для тестирования и удаления их в продуктиве может быть уменьшено, путем добавления их с помощью инструментов аспектно-ориентированного программирования (AOP).

```

public synchronized void transferCredits(Account from,
                                         Account to,
                                         int amount) {

```

¹³¹ Имеется в виду банковская операция – проводка.

```
        from.setBalance(from.getBalance() - amount);
        if (random.nextInt(1000) > THRESHOLD)
            Thread.yield();
        to.setBalance(to.getBalance() + amount);
    }
```

Листинг 12.10 Использование метода `Thread.yield` для увеличения степени чередования

12.2 Тестирование производительности

Тесты производительности часто представляют собой расширенные версии тестов функциональности. На самом деле, почти всегда стоит включать в тесты производительности базовое тестирование функциональности, чтобы убедиться, что вы не тестируете производительность сломанного кода.

Хотя между тестами производительности и функциональными тестами определенно есть перекрытие, у них разное предназначение. Тесты производительности предназначены для измерения сквозных показателей производительности для репрезентативных вариантов использования. Выбор разумного набора сценариев использования не всегда прост; в идеале тесты должны отражать фактическое использование тестируемых объектов в приложении.

В некоторых случаях соответствующий сценарий тестирования очевиден. Ограниченные буферы почти всегда используются в дизайне производитель-потребитель, поэтому имеет смысл измерять пропускную способность производителей, передающих данные потребителям. Мы можем легко расширить класс `PutTakeTest`, чтобы превратить его в тест производительности для предложенного сценария.

Общая вторичная цель тестирования производительности заключается в эмпирическом выборе размеров, для различных граничных условий - количества потоков, ёмкости буфера и т. д. Хотя эти значения могут оказаться достаточно чувствительными к характеристикам платформы (таким как тип процессора или даже уровень пошагового выполнения процессора, количеству ЦП или объему памяти) и требовать динамической конфигурации, разумные варианты для этих значений достаточно часто хорошо работают в широком диапазоне систем.

12.2.1 Расширение класса `PutTakeTest` с добавлением учёта времени

Основное расширение, которое мы должны сделать в классе `PutTakeTest`, это добавить возможность измерения времени, затрачиваемого на выполнение. Вместо того, чтобы пытаться измерить время, затрачиваемое на выполнение одной операции, мы получаем более точное измерение, путем учёта времени, в целом затрачиваемого на выполнение и деления полученного значения на количество операций, чтобы рассчитать время, затрачиваемое на выполнение одной операции. Ранее мы уже использовали класс `CyclicBarrier` для запуска и остановки рабочих потоков, поэтому мы можем расширить его с помощью барьерного действия, измеряющего время начала и окончания выполнения, как показано в листинге 12.11.

```
public class BarrierTimer implements Runnable {
    private boolean started;
```

```
private long startTime, endTime;

public synchronized void run() {
    long t = System.nanoTime();
    if (!started) {
        started = true;
        startTime = t;
    } else
        endTime = t;
}
public synchronized void clear() {
    started = false;
}
public synchronized long getTime() {
    return endTime - startTime;
}
}
```

Листинг 12.11 Барьерный таймер

Мы можем изменить инициализацию барьера, чтобы использовать это барьерное действие, с помощью конструктора класса `CyclicBarrier`, принимающего в качестве параметра барьерное действие:

```
this.timer = new BarrierTimer();
this.barrier = new CyclicBarrier(npairs * 2 + 1, timer);
```

В листинге 12.12 представлен модифицированный тестовый метод, переписанный с использованием барьерного таймера. Запустив выполнение класса `TimedPutTakeTest`, мы можем выяснить несколько вещей. Одна из них – величина пропускной способности операции передачи при использовании шаблона производитель-потребитель, с различными комбинациями параметров; другая – величина масштабирования ограниченного буфера с различным числом потоков; третья – каким образом мы можем выбрать размер ограничения.

```
public void test() {
    try {
        timer.clear();
        for (int i = 0; i < nPairs; i++) {
            pool.execute(new Producer());
            pool.execute(new Consumer());
        }
        barrier.await();
        barrier.await();

        long nsPerItem = timer.getTime() / (nPairs * (long)nTrials);
        System.out.print("Throughput: " + nsPerItem + " ns/item");
        assertEquals(putSum.get(), takeSum.get());
    } catch (Exception e) {
        throw new RuntimeException(e);
    }
}
```

Листинг 12.12 Тестирование с использованием барьерного таймера

Для ответа на эти вопросы необходимо выполнить тест с различными комбинациями параметров, для этого нам понадобится основной тестовый драйвер, приведённый в листинге 12.13.

```
public static void main(String[] args) throws Exception {
    int tpt = 100000; // trials per thread
    for (int cap = 1; cap <= 1000; cap *= 10) {
        System.out.println("Capacity: " + cap);
        for (int pairs = 1; pairs <= 128; pairs *= 2) {
            TimedPutTakeTest t =
                new TimedPutTakeTest(cap, pairs, tpt);
            System.out.print("Pairs: " + pairs + "\t");
            t.test();
            System.out.print("\t");
            Thread.sleep(1000);
            t.test();
            System.out.println();
            Thread.sleep(1000);
        }
    }
    pool.shutdown();
}
```

Листинг 12.13 Программа драйвера для класса TimedPutTakeTest.

На рис. 12.1 приведён пример результатов, полученных на четырех ядерной машине, с емкостями буфера 1, 10, 100 и 1000 элементов. Мы сразу видим, что размер буфера равный единице приводит к очень низкой пропускной способности; это происходит от того, что каждый поток может выполнить только крошечный бит работы, тем самым едва увеличив прогресс, перед блокировкой и ожиданием другого потока. Увеличение размера буфера до десяти элементов, значительно увеличивает пропускную способность, но увеличение размера буфера выше значения в десять элементов, приводит к снижению пропускной способности.

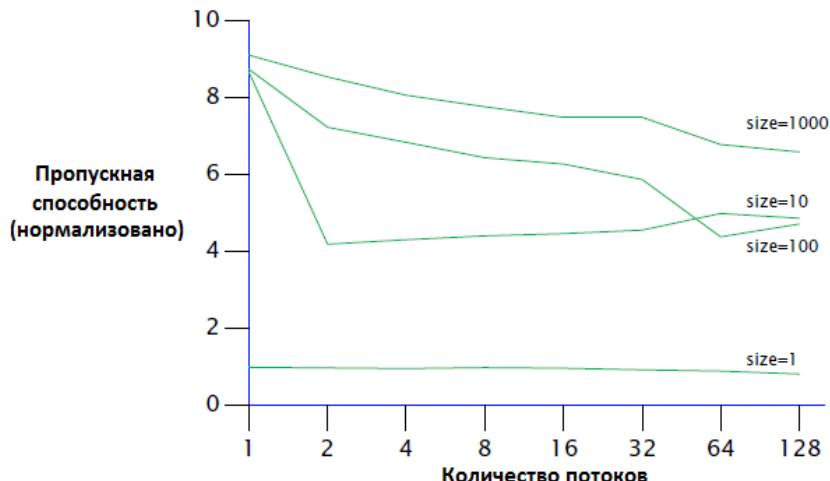


Рисунок 12.1 Запуск TimedPutTakeTest с различной ёмкостью буфера

На первый взгляд может несколько озадачивать, что добавление большего количества потоков достаточно слабо снижает производительность. Причину такого поведения трудно понять, рассматривая только данные, но легко понять, если во время выполнения теста также смотреть на измерения производительности процессора, например, с помощью утилиты perfbar: даже с множеством потоков производится не так уж и много вычислений, и большая их часть тратится на блокирование и разблокирование потоков. Таким образом, достаточно слабый CPU будет выполнять то же самое, без значительного ущерба производительности.

Однако будьте осторожны, делая выводы из этих данных, потому что вы всегда можете добавить большее количество потоков в программу, реализующую шаблон производитель-потребитель и использующую ограниченный буфер. Этот тест достаточно искусственен в том, что касается имитации *работы приложения*; производители почти не выполняют работу по созданию элемента, помещаемого в очередь, а потребители почти не выполняют никакой работы с получаемым элементом. Если рабочие потоки в реальном приложении, реализующем шаблон производитель-потребитель, выполняют некоторую нетривиальную работу по созданию и потреблению элементов (как это обычно и бывает), то эта слабина исчезнет, и последствия наличия слишком большого количества потоков могут стать очень заметными. Основная цель этого теста состоит в том, чтобы измерить, какие ограничения накладываются на общую пропускную способность при реализации шаблона производитель-потребитель с передачей данных с использованием ограниченного буфера.

12.2.2 Сравнение нескольких алгоритмов

Хотя реализация класса `BoundedBuffer` является достаточно надёжной и работающей достаточно хорошо, оказывается, что она не подходит ни для реализации `ArrayBlockingQueue`, ни для реализации `LinkedBlockingQueue` (это объясняет, почему этот алгоритм буфера не был выбран для включения в библиотеку классов). Алгоритмы в пакете `java.util.concurrent` были выбраны и настроены, частично с использованием тестов, подобных описанным выше, чтобы быть столь же эффективными, как те, превращение которых нам известно, но, в то же время, предлагая широчайший спектр функциональных возможностей.¹³² Главная причина, по которой класс `BoundedBuffer` работает плохо, заключается в том, что каждый из методов `put` и `take` имеет несколько операций, которые могут сталкиваться с конкуренцией - захват семафора, захват блокировки, освобождение семафора. Другие подходы к реализации имеют меньше точек соприкосновения, в которых может возникать конкуренция с другими потоками.

На рис. 12.2 приведена сравнительная пропускная способность на двухъядерной гиперпоточной машине, для всех трех классов с буфером на 256 элементов, с использованием варианта теста `TimedPutTakeTest`. Этот тест предполагает, что класс `LinkedBlockingQueue` масштабируется лучше, чем класс `ArrayBlockingQueue`. Сначала это может показаться странным: связанная очередь должна выделять память для ссылки на объект узла при каждой операции вставки и, следовательно, выполнять больше работы, чем очередь на основе массива. Однако, несмотря на то, что в ней выполняется больше операций выделения ресурсов и больше накладные расходы на сборку мусора, связанная очередь обеспечивает большую степень параллелизма при доступе методов `put` и `take`, чем

¹³² Вы можете превзойти их, если являетесь экспертом по параллелизму и можете отказаться от некоторой предоставляемой функциональности.

очередь на основе массива, потому что лучшие алгоритмы работы связанный очереди позволяют обновлять голову и хвост независимо друг от друга. Поскольку выделение обычно выполняется локально для потока, алгоритмы, которые могут уменьшать конкуренцию за счёт выполнения большего количества операций выделения памяти, обычно масштабируются лучше. (Это еще один пример, в котором традиционная настройка производительности, основанная на интуиции, противоречит тому, что действительно необходимо сделать для улучшения масштабируемости.)

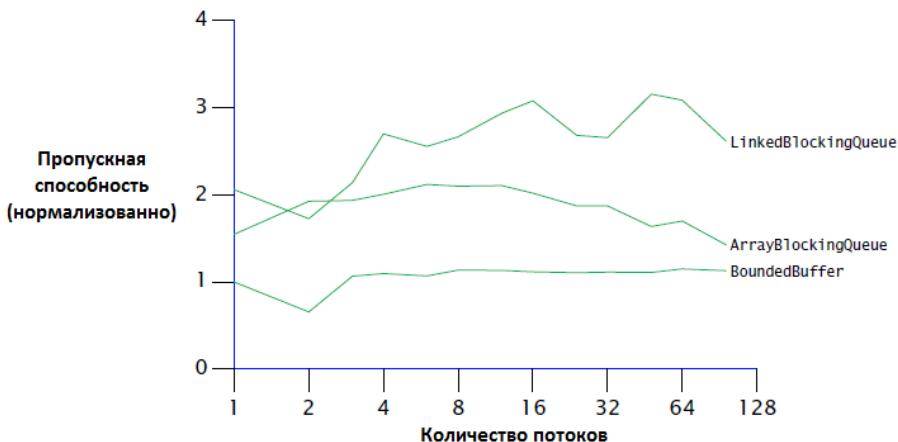


Рисунок 12.2 Сравнение реализаций блокирующих очередей

12.2.3 Измерение отзывчивости

До сих пор мы фокусировались на измерении пропускной способности, которая обычно является самой важной метрикой производительности для параллельных программ. Но иногда важнее знать, сколько времени может длиться выполнение отдельного действия, и в этом случае мы хотим измерить *отклонение* (*variance*) времени обслуживания. Иногда имеет смысл допустить более длительное среднее время обслуживания, если это позволяет нам получить меньшую дисперсию; предсказуемость также является ценной характеристикой производительности. Измерение дисперсии позволяет нам оценить ответы на вопросы, касающиеся качества обслуживания, подобные “Какой процент операций будет успешно выполнен за 100 миллисекунд?”

Гистограммы времени выполнения задачи обычно являются лучшим способом визуализации дисперсии времени обслуживания. Дисперсии лишь немного сложнее измерить, чем средние значения - вам нужно отслеживать время выполнения каждой задачи в дополнение к совокупному времени завершения. Так как детализация таймера может выступать в качестве фактора, оказывающего влияние на измерения времени отдельной задачи (отдельная задача может потребовать меньшее или близкое к наименьшему “циклу таймера” время на выполнение, которое будет вносить искажения в измерение длительности задачи), для того, чтобы избежать появления артефактов измерения, мы можем измерять время выполнения небольших пакетов операций *put* и *take*.

На рис. 12.3 показано время выполнения каждой задачи для варианта теста `TimedPutTakeTest` с размером буфера на 1000 элементов, в котором каждая из 256 параллельных задач выполняет итерацию только 1000 элементов для

несправедливых (зелёные полосы) и справедливых (красные полосы) семафоров. (В разделе 13.3 объясняется различие между справедливой и несправедливой организацией очередей блокировок и семафоров.) Время выполнения с несправедливыми семафорами находится в диапазоне от 1,04 до 8,714 мс, коэффициент больше восьмидесяти. Можно уменьшить этот диапазон, принудительно вводя больше справедливости в управление параллелизмом; это легко сделать в классе `BoundedBuffer`, инициализировав семафоры в справедливом режиме. Как показано на рис. 12.3, это позволяет значительно уменьшить отклонение (в текущий момент оно составляет от 38,194 до 38,207 мс), но, к сожалению, также значительно снижает пропускную способность. (Длительно выполняющийся тест с более типичными задачами, вероятно, покажет еще большее снижение пропускной способности.)

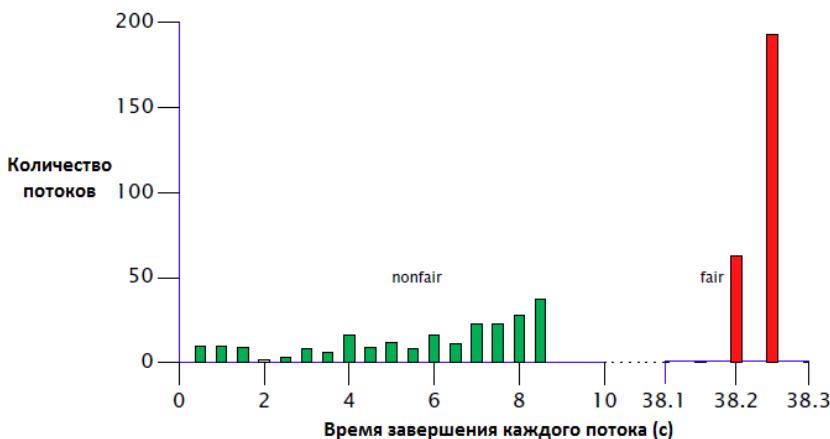


Рисунок 12.3 Гистограмма времени завершения потоков класса `TimedPutTakeTest` с несправедливыми (nonfair, по умолчанию) и справедливыми (fair) семафорами

Ранее мы уже видели, что очень малые размеры буфера приводят к тяжелому переключению контекста и плохой пропускной способности, даже в несправедливом режиме, потому что почти каждая операция включает в себя переключение контекста. В качестве показателя того, что затраты на справедливость является, в первую очередь, результатом блокировки потоков, мы можем перезапустить этот тест с размером буфера равным единице и увидеть, что несправедливые семафоры работают теперь, по времени, сравнимо со справедливыми семафорами. Рисунок 12.4 показывает, что в этом случае справедливость не приводит к ухудшению среднего значения или значительному улучшению дисперсии.

Таким образом, если потоки постоянно блокируются из-за жестких требований синхронизации, несправедливые семафоры обеспечивают гораздо лучшую пропускную способность, а справедливые семафоры обеспечивают меньшую дисперсию. Поскольку результаты в обоих режимах различаются очень сильно, класс `Semaphore` вынуждает своих клиентов принимать решение о том, для какого из этих двух факторов проводить оптимизацию.

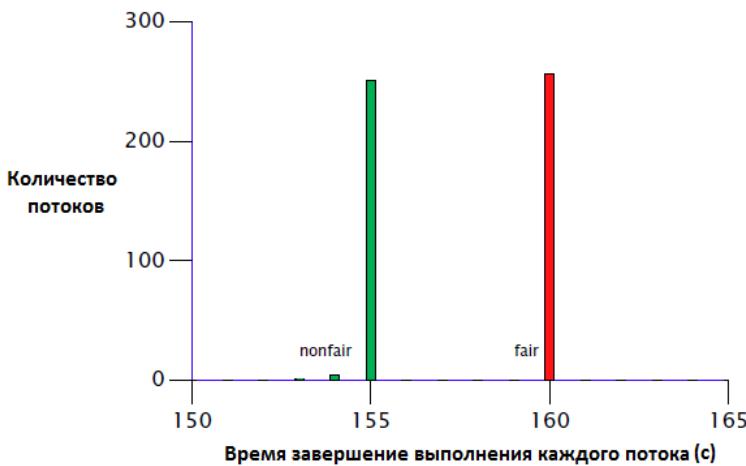


Рисунок 12.4 Гистограмма времени выполнения для теста TimedPutTakeTest, с буфером на один элемент

12.3 Как избежать ошибок тестирования производительности

Теоретически, разрабатывать тесты производительности просто - найдите типичный сценарий использования, напишите программу, которая выполняет этот сценарий множество раз, и зафиксируйте время. На практике, необходимо остерегаться ряда ошибок кодирования, которые не позволяют тестам производительности получить значимые результаты.

12.3.1 Сборка мусора

Время начала сборки мусора непредсказуемо, поэтому всегда существует вероятность того, что сборщик мусора запустится во время выполнения теста, осуществляющего измерения. Если тестовая программа выполняет N итераций и сборка мусора не запускается, но на итерации $N + 1$ сборка мусора запускается, небольшое отклонение в размере выполнения может оказаться большое (но ложное) влияние на измерение времени, затрачиваемого на каждую итерацию.

Существует две стратегии для предотвращения искажения результатов процессом сборки мусора. Во-первых, необходимо убедиться, что сборка мусора вообще не выполняется во время теста (можно вызвать JVM с параметром `-verbose:gc`, чтобы выяснить это); в качестве альтернативы можно убедиться, что сборщик мусора выполняется несколько раз во время выполнения, чтобы тестовая программа адекватно отражала затраты на текущее выделение и сборку мусора. Последняя стратегия часто является лучшим вариантом - она требует более длительного тестирования и, скорее всего, отражает реальную производительность.

Большинство приложений, реализующих шаблон производитель-потребитель, включают в себя достаточное количество выделения памяти и сборки мусора - производители выделяют память под новые объекты, которые используются и отбрасываются потребителями. Длительное по времени выполнение тестирования ограниченного буфера, достаточное для выполнения нескольких сборок мусора, порождает более точные результаты.

12.3.2 Динамическая компиляция

Написание и интерпретация тестов производительности для динамически скомпилированных языков, подобных Java, намного сложнее, чем для статически скомпилированных языков, таких как C или C++. HotSpot JVM (и другие современные JVM) использует комбинацию интерпретации байт-кода и динамической компиляции. При первой загрузке класса JVM выполняет его, интерпретируя байт-код. В какой-то момент, если метод выполняется достаточно часто, динамический компилятор убирает байт-код и преобразует его в машинный код; после завершения компиляции он переключается с интерпретации на прямое выполнение.

Момент начала компиляции непредсказуем. Тесты затрат времени должны выполняться только после компиляции всего кода; измерение скорости интерпретируемого кода не имеет значения, поскольку большинство программ выполняются достаточно долго, чтобы компилировались все часто выполняемые ветви кода. Возможность компилятора запускаться во время выполнения теста снимающего измерения, может двумя способамиказать влияние на результаты теста: компиляция потребляет ресурсы CPU, а измерение времени выполнения комбинации интерпретируемого и скомпилированного кода не является значимой метрикой производительности. На рисунке 12.5 показано, каким образом это может привести к искажению результатов. Три временные шкалы отражают выполнение одного и того же числа итераций: временная шкала **A** представляет собой всё интерпретируемое выполнение, **B** представляет собой компиляцию посреди процесса выполнения, а **C** представляет собой компиляцию на ранней стадии выполнения. Момент времени, в который выполняется компиляция, оказывает серьезное влияние на измеряемое время выполнения каждой операции.¹³³

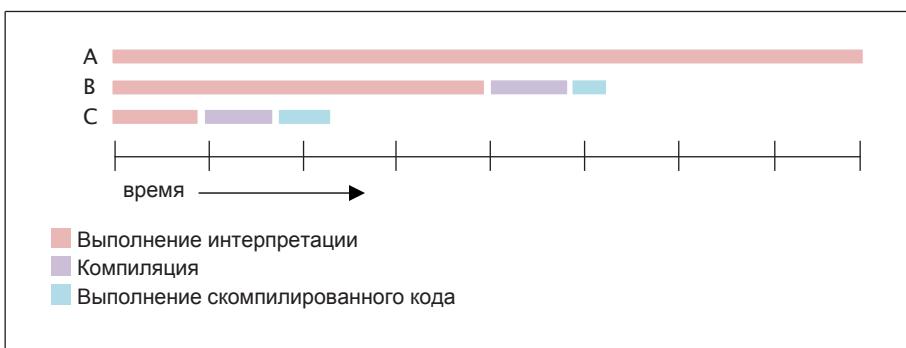


Рисунок 12.5 Необъективные результаты из-за динамической компиляции

Код также может быть декомпилирован (возвращён к интерпретируемому выполнению) и перекомпилирован по различным причинам, таким как загрузка класса, аннулирующего предположения, сделанные в предыдущих компиляциях, или в связи со сбором достаточного количества данных профилирования, для принятия решения о том, что ветка кода должна быть перекомпилирована с различными оптимизациями.

Одним из способов предотвращения оказания влияния компиляции на получаемые результаты является выполнение программы в течение длительного времени (по крайней мере, в течение нескольких минут), чтобы компиляция и интерпретируемое выполнение отнимали небольшую часть от общего времени

¹³³ Среда JVM может выбрать для выполнения компиляции поток приложения или фоновый поток; любой из вариантов может по своему оказывать влияние на результаты фиксации затрат времени.

выполнения. Другой подход заключается в использовании не измеряемого “прогрева”, при котором код выполняется достаточное количество времени, чтобы быть полностью скомпилированным к моменту фактического запуска синхронизации. В HotSpot, при запуске программы с параметром – “XX:+PrintCompilation”, при выполнении динамической компиляции выводится сообщение, так что вы можете убедиться, что запуск компиляции произошёл до, а не во время проведения измерений в тестовых запусках.

Для проверки методологии тестирования можно использовать выполнение одного и того же теста несколько раз, в одном и том же экземпляре JVM. Первая группа результатов должна быть отброшена, как используемая для “прогрева”; наблюдение несогласованных результатов в остальных группах предполагает, что в дальнейшем тест должен быть проверен, с целью определения, по какой причине результаты измерения затрат времени не повторяются.

Среда JVM использует различные фоновые потоки для выполнения сервисных задач. При измерении затрат времени в нескольких *несвязанных* вычислительно интенсивных активностей за один проход, рекомендуется размещать явные паузы между выполнением измерений, чтобы предоставить JVM возможность нагнать выполнение фоновых задач с минимальными помехами от задач, в которых происходят замеры. (Однако, при выполнении измерений затрат времени нескольких связанных активностей, таких как несколько запусков одного и того же теста, исключение фоновых задач JVM подобным образом, может привести к получению оптимистичных, но далёких от реальности результатов.)

12.3.3 Нереалистичная выборка веток выполнения кода

Компиляторы времени выполнения используют сведения, полученные при профилировании, для выполнения оптимизации компилируемого кода. Среде JVM разрешено использовать информацию, специфичную для времени выполнения, чтобы порождать лучший код, что означает, что компиляция метода *M* в одной программе может генерировать код, отличный от компиляции метода *M* в другой программе. В некоторых случаях JVM может выполнять оптимизацию на основе предположений, которые могут быть верны некоторое время, а затем возвращать их, аннулируя скомпилированный код, если они становятся неверными¹³⁴.

В результате, важно чтобы тестовые программы не только адекватно аппроксимировали шаблоны использования типичного приложения, но и аппроксимировали набор веток исполняемого кода, используемых таким приложением. В противном случае динамический компилятор может выполнить специальную оптимизацию, применяемую для чисто однопоточной тестовой программы, которая не может быть применена в реальных приложениях, включающих в себя, по крайней мере, случайный параллелизм. Поэтому тесты многопоточной производительности обычно следует смешивать с тестами однопоточной производительности, даже если требуется провести измерения только однопоточной производительности. (Эта проблема не возникает в классе `TimedPutTakeTest`, поскольку даже наименьший тестовый случай использует два потока.)

¹³⁴ Например, среда JVM может использовать *трансформацию мономорфного вызова* (*monomorphic call transformation*), для преобразования вызова виртуального метода в прямой вызов метода, если загруженные в данный момент классы не переопределяют этот метод, но скомпилированный код становится недействительным, если впоследствии загружается класс, переопределяющий тот же метод.

12.3.4 Нереалистичные предположения о степенях конкуренции

Параллельные приложения, как правило, чередуют два очень разных типа работы: доступ к совместно используемым данным, например, получение следующей задачи из совместно используемой рабочей очереди, и вычисления, локальные для потока (выполнение задачи, при условии, что задача сама по себе не имеет доступа к совместно используемым данным). В зависимости от относительных пропорций двух типов работы, приложение будет сталкиваться с различными уровнями конкуренции, и будет демонстрировать различную производительность и поведение масштабирования.

Если N потоков извлекают задачи из общей рабочей очереди и выполняют их, а задачи являются вычислительно-ресурсоемкими и длительно выполняющимися (и не имеют доступа к совместно используемым данным), конкуренция практически никогда возникать не будет; пропускная способность будет определяться доступностью ресурсов ЦП. С другой стороны, если задачи очень кратковременны, возникает множественная конкуренция за доступ к рабочей очереди, а пропускная способность будет определяться затратами на синхронизацию.

Чтобы в рамках исследования получить реалистичные результаты, параллельные тесты производительности должны пытаться аппроксимировать локальные для потока вычисления, выполняемые типичным приложением, в дополнение к параллельной координации. Если работа, выполняемая для каждой задачи в приложении, значительно отличается по характеру или объему от той, что выполняется тестовой программой, можно легко прийти к необоснованным выводам о том, в каком месте расположено “бутылочное горлышко”¹³⁵ производительности. Ранее, в разделе [11.5](#), мы видели, что для основанных на блокировках классов, таких как синхронизированные реализации Map, факт того, является ли доступ к блокировке по большей части конкурентным или в основном неконкурентным, может иметь очень сильное влияние на пропускную способность. Тесты в этом разделе не выполняют ничего, кроме “бомбардировки” экземпляра Map; даже в случае двух потоков, все попытки получить доступ к экземпляру Map оспариваются. В том случае, если бы приложение выполняло значительный объем локальных для потока вычислений, уровень конкуренции, при каждом доступе к совместно используемой структуре данных, мог бы быть достаточно низок, обеспечивая хороший уровень производительности.

В этой связи, класс `TimedPutTakeTest` может быть плохой моделью для некоторых приложений. Так как рабочие потоки делают не очень много, пропускная способность определяется издержками на координацию, и не обязательно, что такая ситуация имеет место во всех приложениях, которые обмениваются данными между производителями и потребителями через ограниченные буферы.

12.3.5 Устранение мёртвого кода

Одним из вызовов, возникающих в процессе написания хороших тестов (на любом языке) является то, что оптимизирующие компиляторы умеют выявлять и устранять мертвый код - код, который не влияет на результат. Поскольку бенчмарки часто ничего не вычисляют, они являются легкой мишенью для оптимизатора. В большинстве случаев хорошо, когда оптимизатор удаляет

¹³⁵ Бутылочное горлышко – узкое место, ограничивающее производительность.

мертвый код из программы, но для бенчмарка это большая проблема, потому что фактически вы выполняете меньше измерений исполняемого кода, чем ожидаете. Если вам повезёт, оптимизатор *целиком* удалит вашу программу, и тогда будет очевидно, что ваши данные фиктивные. Если вам не повезет, устранение мертвого кода просто ускорит вашу программу, что станет каким-то фактором, который *может* быть объяснён другими средствами.

Устранение мертвого кода также является проблемой в бенчмаркинге статически скомпилированных языков, но обнаружение того, что компилятор устранил хороший кусок вашего бенчмарка, намного проще, потому что вы можете посмотреть на машинный код и увидеть, что часть вашей программы отсутствует. В случае с динамически компилируемыми языками, получить доступ к этой информации не так просто.

Многие микро бенчмарки работают намного "лучше" при работе с параметром компилятора *HotSpot -server*, чем с параметром компилятора *-client*, не только потому, что серверный компилятор может порождать более эффективный код, но и потому, что он более искусен в оптимизации мертвого кода. К сожалению, устранение мертвого кода, сделавшего работу вашего бенчмарка весьма короткой, не приведёт к тем же последствиям и с выполняющимся кодом. Но вы все равно должны предпочитать параметр *-server* параметру *-client*, как для продуктива, так и в случае выполнения тестирования на многопроцессорных системах - вам просто нужно писать свои тесты так, чтобы они не были подвержены удалению мертвого кода.

Написание эффективных тестов производительности приводит к необходимости обманом заставлять оптимизатор не оптимизировать тесты производительности в качестве мертвого кода. Для этого необходимо, чтобы каждый вычисленный результат каким-то образом использовался вашей программой - таким образом, который не требует синхронизации или существенных вычислений.

В классе *PutTakeTest*, мы вычисляем контрольную сумму элементов, добавленных и удаленных из очереди, и объединяем эти контрольные суммы по всем потокам, но этот код все еще может быть оптимизирован, если мы фактически не используем полученное значение контрольной суммы. Нам это необходимо для проверки правильности алгоритма, но вы можете гарантировать, что значение используется, просто распечатав его. Однако следует избегать операций ввода/вывода во время выполнения теста, чтобы не оказывать влияние на измерение времени выполнения.

Дешевый трюк для предотвращения оптимизации вычисления, без внесения слишком больших накладных расходов, заключается в выполнении вычисления с помощью метода *hashCode*, принадлежащего полю некоторого производного объекта, сравнении его с произвольным значением, таким как текущее значение *System.nanoTime*, и печати бесполезного и игнорируемого сообщения, если оба значения совпадают:

```
if (foo.x.hashCode() == System.nanoTime())
    System.out.print(" ");
```

Сравнение редко оказывается успешным, и если это всё же произойдет, единственным эффектом будет вставка безвредного символа пробела в выходной

поток. (Метод `print` буферизует вывод до момента вызова метода `println`, поэтому в тех редких случаях, когда результаты вызовов методов `hashCode` и `System.nanoTime` оказываются равны, никаких операций ввода/вывода не происходит.)

Мало того, что каждый вычисляемый результат должен использоваться, но результаты также должны быть неочевидными. В противном случае интеллектуальный динамический оптимизирующий компилятор может заменить действия предварительно вычисленными результатами. Мы рассматривали это при построении класса `PutTakeTest`, но любая тестовая программа, на вход которой подаются статические данные, уязвима для этой оптимизации.

12.4 Комплементарные подходы к тестированию

Хотя мы хотели бы верить, что эффективная программа тестирования должна “выявить все ошибки”, но это нереалистичная цель. Агентство NASA посвящает больше своих инженерных ресурсов тестированию (по оценкам, они используют до 20 тестировщиков на каждого разработчика), чем может себе позволить любая коммерческая организация - и созданный код по-прежнему содержит дефекты. В сложных программах никакое количество тестов не может помочь выявить все ошибки кодирования.

Цель тестирования заключается не столько в том, чтобы *найти ошибки*, сколько в *увеличении уверенности* в том, что код работает так, как ожидалось. Поскольку нереально предположить, что вы сможете найти все ошибки, цель плана обеспечения качества (QA, *quality assurance*) должна заключаться в достижении максимально возможной уверенности, учитывая доступные ресурсы тестирования. В параллельной программе может произойти больше ошибок, чем в последовательной, и поэтому для достижения того же уровня уверенности требуется больше тестов. До сих пор мы сосредоточились в основном на методах построения эффективных модульных тестов и тестов производительности. Тестирование критически важно для создания уверенности в том, что параллельные классы ведут себя правильно, но это должно быть только одной из используемых методологий QA.

Различные методологии контроля качества более эффективны при поиске одних типов дефектов и менее эффективны при поиске других. Используя комплементарные методологии тестирования, такие как ревю кода и статический анализ, можно добиться большей уверенности, чем при использовании только какого-то одного подхода.

12.4.1 Ревю кода

Как бы ни были эффективны и важны модульные и стресс тесты для поиска ошибок параллелизма, они не смогут заменить тщательное ревю кода, осуществляющее множество людей. (С другой стороны, ревю кода также не заменяет тестирование.) Вы можете и должны проектировать тесты таким образом, чтобы максимизировать их шансы на обнаружение ошибок безопасности, и вы должны запускать их часто, но вы не должны пренебрегать тем, чтобы параллельный код тщательно проверялся кем-то, кроме его автора. Даже эксперты по параллелизму допускают ошибки; выделять время на то, чтобы кто-то другой провел ревю кода - почти всегда стоящее дело. Эксперты в параллельном программировании часто проявляют себя лучше в обнаружении тонких состояний

гонок, чем большинство тестовых программ. (Кроме того, особенности платформы, такие как детали реализации JVM или модели памяти процессора, могут препятствовать обнаружению ошибок в определенных конфигурациях оборудования или программного обеспечения.) Процесс ревю кода также имеет и другие преимущества; он не только позволяет находить ошибки, но и часто улучшает качество комментариев, описывающих детали реализации, тем самым уменьшая последующие затраты на поддержку и риски.

12.4.2 Инструменты для проведения статического анализа

На момент написания этой главы, *инструменты статического анализа (static analysis tools)* представляют собой быстро развивающееся и эффективное дополнение к формальному тестированию и анализу кода. Статический анализ кода представляет собой процесс анализа кода без его выполнения, и инструменты выполняющие аудит кода могут анализировать классы путём поиска экземпляров распространенных шаблонов ошибок (*bug patterns*). Инструменты статического анализа с открытым исходным кодом, например FindBugs¹³⁶, содержат детекторы для множества распространенных ошибок кодирования, многие из которых можно легко пропустить при тестировании или просмотре кода.

Инструменты статического анализа создают список предупреждений, которые необходимо проверять вручную, чтобы определить, указывают ли они на фактические ошибки. Исторически так сложилось, что инструменты подобные `lint` порождают множество ложных предупреждений, чтобы напугать разработчиков, но инструменты подобные FindBugs были настроены так, чтобы порождать намного меньше ложных предупреждений. Инструменты статического анализа все еще несколько примитивны (особенно в случае с интеграцией с инструментами разработки и подстройкой под жизненный цикл), но они уже достаточно эффективны, чтобы быть ценным дополнением для процесса тестирования.

На момент написания этой главы, инструмент FindBugs включает в себя детекторы для следующих шаблонов ошибок, связанных с параллелизмом, и их список всё время пополняется:

Несогласованная синхронизация. Многие объекты следуют политике синхронизации, защищая все переменные с помощью встроенной блокировки объекта. Если к полю обращаются часто, но, при этом, не всегда удерживая блокировку на `this`, это может указывать на то, что последовательное соблюдение политики синхронизации не выполняется. Инструменты, осуществляющие анализ, вынуждены угадывать политику синхронизации, поскольку классы Java не имеют формальных спецификаций параллелизма. В будущем, если аннотации, подобные `@GuardedBy`, будут стандартизированы, инструменты, выполняющие аудит кода, смогут интерпретировать аннотации, вместо того, чтобы гадать о взаимосвязи между переменными и блокировками, что приведёт к повышению качества анализа кода.

Вызов метода `Thread.run`. Класс `Thread` реализует интерфейс `Runnable`, и поэтому обладает методом `run`. Однако вызов метода `Thread.run` напрямую, практически всегда является ошибкой; по обыкновению, программист вызывает метод `Thread.start`.

¹³⁶ <http://findbugs.sourceforge.net>

Неосвобождённая блокировка. В отличие от внутренних блокировок, явные блокировки (см. главу [13](#)) автоматически не освобождаются, когда управление покидает область, в которой они были захвачены. Стандартная идиома заключается в том, что освобождать блокировку в блоке `finally`; в противном случае, при возбуждении исключения `Exception` блокировка может остаться неосвобождённой.

Пустой блок `synchronized`. Несмотря на то, что пустые блоки `synchronized` прописаны в семантике модели памяти Java, они часто используются некорректно, и обычно существует лучшее решение той проблемы, которую разработчик пытался решить.

Блокировка с двойной проверкой. Блокировка с двойной проверкой представляет собой неустойчивую идиому, применяемую для уменьшения издержек синхронизации, возникающих при отложенной инициализации (см. раздел [16.2.4](#)), которая приводит к чтению совместно используемого изменяемого поля без соответствующей синхронизации.

Запуск потока из конструктора. Запуск потока из конструктора может привести к риску возникновения проблем с подклассами и может позволить ссылке на `this` сбежать из конструктора.

Ошибки уведомления. Методы `notify` и `notifyAll` указывают на то, что состояние объекта могло измениться таким образом, чтобы разблокировать потоки, ожидающие в ассоциированной с условием очереди. Эти методы следует вызывать только тогда, когда изменяется состояние, ассоциированное с очередью условий (*condition queue*). Вызов методов `notify` или `notifyAll` из блока `synchronized` без внесения изменений в состояние, вероятнее всего, будет ошибкой. (См. главу [14](#).)

Ошибки ожидания условия. Когда происходит ожидание на очереди условий, методы `Object.wait` или `Condition.await` должны вызываться в цикле, с удержанием соответствующей блокировки, после проверки некоторого предиката состояния (см. главу [14](#)). Вызов методов `Object.wait` или `Condition.await` без удержания блокировки, не в цикле, или без проверки предиката состояния, почти наверняка является ошибкой.

Неправильное использование Lock и Condition. Использование интерфейса `Lock` в качестве аргумента блокировки в блоке `synchronized`, может быть опечаткой, как и вызов метода `Condition.wait`, а не `await` (хотя последний, вероятнее всего, будет перехвачен во время тестирования, так как он бросит исключение `IllegalMonitorStateException` при первом же вызове).

Засыпание или ожидание в процессе удержания блокировки. Вызов метода `Thread.sleep` с удерживаемой блокировкой может препятствовать другим потокам в продвижении прогресса выполнения в течение длительного времени и, следовательно, потенциально является серьезной угрозой живучести. Вызов `Object.wait` или `Condition.await` с двумя удерживаемыми блокировками представляет собой аналогичную угрозу.

Оборачиваемые циклы. Код, который не выполняет ничего, кроме оборачиваемости (занят ожиданием), выполняет проверку состояния поля на соответствие ожидаемому значению и может расходовать процессорное время и, если поле не является *volatile*, завершение выполнения кода не гарантируется. Ожидание на защелках и условиях часто являются лучшим подходом, в ожидании перехода состояния.

12.4.3 Аспектно-ориентированные подходы к тестированию

На момент написания этой главы, методы аспектно-ориентированного программирования (AOP) имели ограниченную применимость к параллелизму, поскольку большинство популярных средств AOP еще не поддерживают срезы (*pointcuts*) в точках синхронизации. Однако AOP может применяться для утверждения инвариантов или соответствия некоторых аспектов политикам синхронизации. Например, в (Laddad, 2003) приведен пример использования аспекта для того, чтобы обернуть все вызовы не потокобезопасных методов Swing с утверждением, что вызов происходит в потоке событий. Поскольку внесения изменений в код не требуется, применение этого метода достаточно просто и может раскрывать тонкие ошибки, связанные с публикацией и ограничением потока.

12.4.4 Профилировщики и инструменты мониторинга

Большинство коммерческих инструментов для профилирования поддерживают потоки. Они различаются по набору функций и эффективности, но часто могут помочь составить представление о том, что делает ваша программа (хотя средства профилирования обычно навязчивы и могут существенно повлиять на синхронизацию и поведение программы). Большинство из них предлагает дисплей с графиками для каждого потока, с различными цветами для различных состояний потока (выполняется, заблокирован, ожидает захвата блокировки, заблокирован в ожидании завершения операций ввода/вывода и т.д.). Такой дисплей может показать, насколько эффективно ваша программа использует доступные ресурсы процессора, и если она работает плохо, где искать причину. (Многие профилировщики также заявляют о возможностях по идентификации того, какие блокировки вызывают конкуренцию, но на практике эти функции часто являются более грубым инструментом, чем требуется для осуществления анализа поведения блокировок программы.)

Встроенный агент JMX также предлагает некоторые ограниченные возможности по мониторингу поведения потоков. Класс `ThreadInfo` включает в себя текущее состояние потока и, если поток заблокирован, блокировку или очередь условий, на которой он блокируется. Если включена функция “мониторинг конкуренции потоков” (по умолчанию она отключена из-за влияния на производительность), класс `ThreadInfo` также включает информацию о количестве времени, в течение которого поток был заблокирован в ожидании захвата блокировки или ожидания уведомления, а также совокупное время ожидания.

12.5 Итоги

Тестирование параллельных программ на корректность может оказаться чрезвычайно сложной задачей, поскольку многие из возможных режимов сбоя параллельных программ являются маловероятными событиями, чувствительными

ко времени, нагрузке и другим трудновоспроизводимым условиям. Кроме того, предназначенная для тестирования инфраструктура может ввести дополнительные ограничения времени или синхронизации, которые могут привести к маскировке проблем параллелизма в тестируемом коде. Тестирование параллельных программ на производительность может быть не менее сложной задачей; Java-программы сложнее тестировать, чем программы, написанные на статически скомпилированных языках, подобных C, поскольку на измерения времени может влиять динамическая компиляция, сборка мусора и адаптивная оптимизация.

Чтобы иметь наилучшие шансы найти скрытые ошибки до их проявления в продуктиве, объединяйте традиционные методы тестирования (стараясь избежать ошибок, описанных здесь) с ревю кода и средствами автоматизированного анализа. Каждый из этих методов позволяет выявить проблемы, которые могут быть пропущены другими.

Часть IV Дополнительные темы

Глава 13 Явные блокировки

До Java 5.0 единственными механизмами координации доступа к совместно используемым данным были `synchronized` и `volatile`. В Java 5.0 была добавлена другая опция: класс `ReentrantLock`. Вопреки тому, что некоторые писали, класс `ReentrantLock` не является заменой внутренней блокировки, а скорее альтернативой с расширенными функциями, когда внутренняя блокировка оказывается слишком ограниченной.

13.1 Интерфейсы `Lock` и `ReentrantLock`

Интерфейс `Lock`, приведённый в листинге 13.1, определяет несколько абстрактных блокирующих операций. Отличаясь от внутренней блокировки, интерфейс `Lock` предлагает возможность выбора из безусловного, опрашиваемого, ограниченного по времени и прерываемого захвата блокировки, и все операции блокировки и разблокировки выполняются явно. Реализации интерфейса `Lock` должны предоставлять ту же семантику видимости памяти, что и внутренние блокировки, но могут различаться в семантике блокировки, алгоритмах планирования, гарантиях упорядочивания, и характеристиках производительности. (Метод `Lock.newCondition` рассматривается в главе [14](#)).

```
public interface Lock {  
    void lock();  
    void lockInterruptibly() throws InterruptedException;  
    boolean tryLock();  
    boolean tryLock(long timeout, TimeUnit unit)  
        throws InterruptedException;  
    void unlock();  
    Condition newCondition();  
}
```

Листинг 13.1 Интерфейс `Lock`

Класс `ReentrantLock` реализует интерфейс `Lock`, обеспечивая те же гарантии взаимного исключения и видимости памяти, что и блок `synchronized`. Захват блокировки с помощью класса `ReentrantLock` имеет ту же семантику памяти, что и вход в блок `synchronized`, а освобождение блокировки, захваченной с помощью класса `ReentrantLock`, имеет ту же семантику памяти, что и выход из блока `synchronized`. (Видимость памяти рассматривается в разделе [3.1](#) и в Главе [16](#).) И, как и блок `synchronized`, класс `ReentrantLock` предлагает реентерабельную¹³⁷ семантику блокировки (см. раздел [2.3.2](#)). Класс `ReentrantLock` поддерживает все режимы блокировки, определенные в интерфейсе `Lock`, при этом обеспечивая большую гибкость при работе с недоступностью блокировки, чем блок `synchronized`.

Зачем создавать новый механизм блокировки, который так похож на внутреннюю блокировку? Внутренняя блокировка отлично работает в большинстве ситуаций, но имеет некоторые функциональные ограничения - невозможно прервать поток, ожидающий захвата блокировки, или попытаться захватить блокировку, не желая ожидать ее вечно. Внутренние блокировки также должны

¹³⁷ Реентерабельность – повторная входимость.

освобождаться в том же блоке кода, в котором они были захвачены; это упрощает кодирование и обеспечивает прекрасное взаимодействие с обработкой исключений, но делает невозможными блокировки с не-блочкой структурой. Нет причин отказываться от использования блока `synchronized`, но в некоторых случаях более гибкий механизм блокировки обеспечивает лучшую живучесть или производительность.

В листинге 13.2 приведена каноническая форма использования блокировки. Эта идиома несколько сложнее, чем использование встроенных блокировок: блокировка должна быть освобождена в блоке `finally`. В противном случае блокировка может быть никогда не освобождена, если защищаемый код бросит исключение. При использовании блокировки необходимо также учитывать, что происходит, если исключение бросается из блока `try`; если существует вероятность того, что объект может остаться в несогласованном состоянии, может потребоваться использование дополнительных блоков `try-catch` или `try-finally`. (Всегда следует учитывать влияние исключений при использовании любой формы блокировки, включая внутреннюю блокировку.)

```
Lock lock = new ReentrantLock();
...
lock.lock();
try {
    // update object state
    // catch exceptions and restore invariants if necessary
} finally {
    lock.unlock();
}
```

Листинг 13.2 Защита состояния объекта с использованием класса `ReentrantLock`

Пренебрежение использованием блока `finally`, в целях освобождения экземпляра `Lock`, представляет собой бомбу с часовым механизмом. Когда блокировка станет не нужна, вам будет трудно отследить её происхождение, поскольку не будет никаких записей о том, где и когда экземпляр `Lock` должен был быть освобожден. Это одна из причин, из-за которой не следует использовать класс `ReentrantLock` в качестве полной замены блока `synchronized`: это более “опасно”, потому что блокировка автоматически не освобождается, когда управление покидает защищаемый блок. До тех пор, пока вы помните о том, что необходимо освобождать блокировку в блоке `finally`, всё не так сложно, но всё равно существует вероятность того, что вы забудете это сделать.¹³⁸

13.1.1 Опрашиваемый и ограниченный по времени захват блокировки

Опрашиваемые и ограниченные по времени режимы захвата блокировок, предоставляемые методом `tryLock`, позволяют выполнять более сложное восстановление после возникновения ошибок, чем при безусловном захвате. При использовании внутренних блокировок, взаимоблокировка неустранима - единственный способ восстановления заключается в перезапуске приложения, а единственная защита - создавать программу таким образом, чтобы был невозможен

¹³⁸ Инструмент FindBugs имеет детектор “неосвобождённой блокировки”, определяющий случай, когда блокировка не освобождается во всех ветках кода вне блока, в котором она была захвачена.

несогласованный порядок захвата блокировок. Ограниченные по времени и опрашиваемые блокировки могут предложить другой вариант: вероятностное предотвращение взаимоблокировки.

Использование ограниченных по времени или опрашиваемых блокировок (`tryLock`) позволяет восстановить управление в том случае, если не удаётся захватить все необходимые блокировки, освободить те, которые были захвачены, и повторить попытку (или, по крайней мере, логировать информацию о сбое и выполнить что-то еще). В листинге 13.3 показан альтернативный способ устранения динамической взаимоблокировки, вызванной порядком захвата блокировок, из раздела [10.1.2](#): используйте метод `tryLock` для попытки захвата обеих блокировок, выполните откат и повторите попытку, если обе блокировки одновременно захватить невозможно. Время сна имеет фиксированный и случайный компоненты, для уменьшения вероятности возникновения динамической взаимоблокировки. Если блокировки не удаётся получить в течение указанного времени, метод `transferMoney` возвращает статус, указывающий на то, что произошёл сбой, чтобы операция могла завершиться с ошибкой. (См. [[CPJ 2.5.1.2](#)] и [[CPJ 2.5.1.3](#)] для дополнительных примеров использования опрашиваемых блокировок, для избежания взаимоблокировок.)

```
public boolean transferMoney(Account fromAcct,
                             Account toAcct,
                             DollarAmount amount,
                             long timeout,
                             TimeUnit unit)
        throws InsufficientFundsException, InterruptedException {
    long fixedDelay = getFixedDelayComponentNanos(timeout, unit);
    long randMod = getRandomDelayModulusNanos(timeout, unit); long
    stopTime = System.nanoTime() + unit.toNanos(timeout);

    while (true) {
        if (fromAcct.lock.tryLock()) {
            try {
                if (toAcct.lock.tryLock()) {
                    try {
                        if (fromAcct.getBalance().compareTo(amount)
                            < 0)
                            throw new InsufficientFundsException();
                        else {
                            fromAcct.debit(amount);
                            toAcct.credit(amount);
                            return true;
                        }
                    } finally {
                        toAcct.lock.unlock();
                    }
                }
            } finally {
                fromAcct.lock.unlock();
            }
        }
    }
}
```

```
        if (System.nanoTime() > stopTime)
            return false;
        NANoseconds.sleep(fixedDelay + rnd.nextLong() % randMod);
    }
}
```

Листинг 13.3 Избегание возникновения взаимоблокировок с помощью метода `tryLock`

Блокировки, ограниченные по времени, также полезны при реализации активностей, управляющих бюджетом времени (см. раздел [6.3.7](#)). Когда действие с определённым бюджетом времени вызывает блокирующий метод, оно может предоставить время ожидания, соответствующее оставшемуся времени в бюджете. Это позволяет завершать действия раньше, если они не смогут предоставить результат в течение указанного промежутка времени. При использовании внутренних блокировок, отменить захват блокировки после запуска операции уже невозможно, поэтому внутренние блокировки ставят под угрозу возможность реализации действий, запланированных по времени.

В примере портала путешествий, приведённом в листинге 6.17 (в разделе [6.3.8](#)), создаётся отдельная задача для каждой компании по прокату автомобилей, от которой запрашивались предложения. Запрос ставки, вероятнее всего, включает в себя какой-то механизм запроса по сети, например, запрос к веб-службе. Однако, для получения предложения может также потребоваться эксклюзивный доступ к дефицитному ресурсу, подобному прямой линии связи с компанией.

Ранее, в разделе [9.5](#), мы видели один способов обеспечить сериализованный доступ к ресурсу: однопоточный исполнитель. Другой подход заключается в использовании эксклюзивной блокировки для ограничения доступа к ресурсу. Код, приведенный в листинге 13.4, пытается отправить сообщение по совместно используемой линии связи, защищённой экземпляром `Lock`, но корректно завершает работу, если не может выполнить это в рамках своего бюджета времени. Ограниченнная по времени версия метода `tryLock` практически включает эксклюзивную блокировку в такую ограниченную по времени активность.

```
public boolean trySendOnSharedLine(String message,
                                    long timeout, TimeUnit unit)
                                    throws InterruptedException {
    long nanosToLock = unit.toNanos(timeout)
                           - estimatedNanosToSend(message);
    if (!lock.tryLock(nanosToLock, NANoseconds))
        return false;
    try {
        return sendOnSharedLine(message);
    } finally {
        lock.unlock();
    }
}
```

Листинг 13.4 Блокировка с ограниченным бюджетом времени

13.1.2 Прерываемый захват блокировки

Также как ограниченный по времени захват блокировки позволяет эксклюзивным блокировкам использоваться в активностях, ограниченных по времени, также и прерываемый захват блокировки позволяет блокировкам использоваться в

отменяемых активностях. В разделе [7.1.6](#) было определено несколько механизмов, подобных захвату внутренней блокировки, не реагирующей на прерывание. Эти не прерываемые механизмы блокировки усложняют реализацию отменяемых задач. Метод `lockInterruptibly` позволяет вам попытаться захватить блокировку, оставаясь отзывчивым к прерыванию, и его включение в интерфейс `Lock` позволяет избежать создания другой категории не прерываемых механизмов блокировки.

Каноническая структура прерываемого захвата блокировки немного сложнее, чем обычный захват блокировки, поскольку необходимы два блока `try`. (Если прерываемый захват блокировки может бросить исключение `InterruptedException`, работает стандартная идиома захвата блокировки `try-finally`). В листинге 13.5 метод `lockInterruptibly` используется для реализации метода `sendOnSharedLine` из листинга 13.4, так что мы можем вызвать его из задачи, которую можно отменить. Ограниченный по времени метод `tryLock` также отзывчив к прерыванию, и поэтому может быть использован, когда вам необходим и ограниченный по времени, и прерываемый захват блокировки.

```
public boolean sendOnSharedLine(String message)
    throws InterruptedException {
    lock.lockInterruptibly();
    try {
        return cancellableSendOnSharedLine(message);
    } finally {
        lock.unlock();
    }
}

private boolean cancellableSendOnSharedLine(String message)
    throws InterruptedException { ... }
```

Листинг 13.5 Прерываемый захват блокировки

13.1.3 Блокировка с не-блочкой структурой

В случае с внутренними блокировками, пары захват-освобождение имеют блочную структуру - блокировка всегда освобождается в том же базовом блоке, в котором она была захвачена, независимо от того, каким образом управление покинет блок. Автоматическое освобождение блокировки упрощает анализ и предотвращает возможные ошибки кодирования, но иногда требуется более гибкая дисциплина блокировки.

Ранее, в главе [11](#), мы видели, как уменьшение детализации блокировок может повысить масштабируемость. Чертеживание блокировок позволяет различным цепочкам хэшей, в коллекции на основе хэшей, использовать разные блокировки. Мы можем применить аналогичный принцип, чтобы уменьшить степень детализации блокировки в связанном списке, используя отдельную блокировку для *каждой ссылки на узел*, позволяя различным потокам работать независимо в разных частях списка. Блокировка для данного узла защищает указатели ссылок и данные, хранящиеся в этом узле, поэтому при обходе или изменении списка мы должны удерживать блокировку на одном узле, пока не получим блокировку на следующем узле; только после этого можно отпустить блокировку на первом узле. Пример такого подхода, называемый *блокировкой рука об руку (hand-over-hand)* или *связью блокировок (lock coupling)*, приведен в [CPJ 2.5.1.4].

13.2 Вопросы производительности

Когда класс `ReentrantLock` был добавлен в Java 5.0, он предлагал гораздо лучшую производительность, чем внутренние блокировки. Для примитивов синхронизации, производительность при конкуренции является ключом к масштабируемости: если на управление блокировками и планирование расходуется больше ресурсов, для приложения, соответственно, доступно меньше. Лучшая реализация блокировки осуществляет меньше системных вызовов, вынуждает выполнять меньше переключений контекста и инициирует меньше трафика синхронизации доступа к памяти в совместно используемойшине доступа к памяти - перечень операций, которые требуют много времени на выполнение и отвлекают вычислительные ресурсы программы.

В Java 6 используется улучшенный алгоритм управления внутренними блокировками, подобный тому, который используется в классе `ReentrantLock`, что приводит к значительному сокращению разрыва в масштабируемости. На рис. 13.1 показана разница в производительности между встроенными блокировками и классом `ReentrantLock` в Java 5.0 и в предварительной сборке Java 6, на четырехядерной системе Opteron запущенной на ОС Solaris. Кривые демонстрируют “ускорение” класса `ReentrantLock` по сравнению с внутренней блокировкой на одной и той же версией JVM. На кривой Java 5.0 класс `ReentrantLock` предлагает значительно лучшую пропускную способность, но на кривой Java 6 оба значения довольно близки¹³⁹. Тестовая программа та же, что использовалась в разделе 11.5, но на этот раз сравнивается пропускная способность `HashMap`, защищенная встроенной блокировкой и классом `ReentrantLock`.

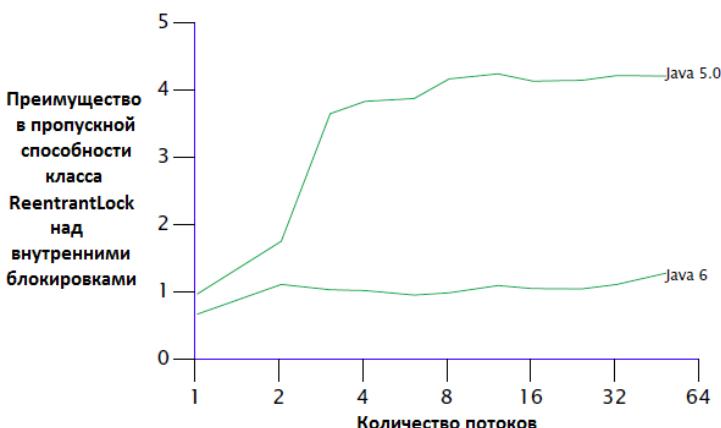


Рисунок 13.1 Сравнение производительности внутренней блокировки и класса `ReentrantLock`, в Java 5 и Java 6

В Java 5.0 производительность внутренней блокировки резко падает при переходе от выполнения одного потока (без конкуренции) к выполнению нескольких потоков; производительность класса `ReentrantLock` страдает намного меньше, демонстрируя его лучшую масштабируемость. Но в Java 6 совсем другая история - внутренние блокировки больше не разваливаются из-за конкуренции, и оба масштаба довольно похожи.

Графики подобные тому, что приводится на рисунке 13.1, напоминают нам, что утверждения вида “ X быстрее Y ” в лучшем случае недолговечны.

¹³⁹ Хотя график этого и не отражает, разница в масштабируемости между Java 5.0 и Java 6 действительно связана с улучшением внутренней блокировки, а не с регрессией производительности в классе `ReentrantLock`.

Производительность и масштабируемость зависят от таких факторов, определяемых платформой, как ЦП, количество процессоров, размер кэша и характеристики JVM, которые со временем могут изменяться¹⁴⁰.

Производительность представляет собой изменяющийся показатель; вчерашний бенчмарк, показывающий, что X быстрее, чем Y, возможно сегодня уже устарел.

13.3 Справедливость

Конструктор класса `ReentrantLock` предлагает выбор из двух вариантов справедливости: создать *несправедливую* (*nonfair*) блокировку (по умолчанию) или *справедливую* (*fair*) блокировку. Потоки захватывают справедливую блокировку в том порядке, в котором они ее запросили, в то время как несправедливая блокировка разрешает *баржирование* (*barging*)¹⁴¹: потоки, запрашивающие блокировку, могут опережать очередь ожидающих потоков, если блокировка оказывается доступной при запросе. (Класс `Semaphore` также предлагает выбор из честного и нечестного порядка захвата блокировок.) “Несправедливые” экземпляры класса `ReentrantLock` не стараются изо всех сил способствовать баржированию - они просто не препятствуют потоку в баржировании, если он появляется в нужное время. При справедливой блокировке, вновь запрашивающие потоки помещаются в очередь, если блокировка удерживается другим потоком или, если потоки помещены в очередь на ожидание блокировки; с несправедливой блокировкой, поток помещается в очередь только в том случае, если блокировка удерживается в текущий момент¹⁴².

Разве мы не хотим, чтобы все блокировки были справедливыми? В конце концов, справедливость - это хорошо, а несправедливость - плохо, не так ли? (Просто спросите своих детей.) Однако, когда дело доходит до блокировки, справедливость приводит к значительным затратам производительности, из-за накладных расходов на приостановку и возобновление потоков. На практике гарантия статистической справедливости - это обещание, что заблокированный поток в *конечном счёте* захватит блокировку - часто достаточно хороша и затраты на её предоставление значительно дешевле. Некоторые алгоритмы полагаются на справедливую очередь для обеспечения своей корректности, но обычно так не делается. В большинстве случаев, преимущества несправедливой блокировки перевешивают преимущества справедливого помещения в очередь.

На рис. 13.2 демонстрируется запуск еще одного теста производительности реализации `Map`, на этот раз сравнивается производительность экземпляра `HashMap`, обёрнутого со справедливой и не справедливой реализацией `ReentrantLock`, на четырёх ядерном процессоре Opteron под управлением ОС Solaris, с результатом выведенным по логарифмической шкале¹⁴³. Штраф за использование

¹⁴⁰ Когда мы начинали работу над этой книгу, класс `ReentrantLock` казался последним словом в масштабируемости блокировок. Менее чем год спустя, внутренняя блокировка предоставляет хорошую производительность за свои деньги. Производительность - это не просто изменяющийся показатель, она может быть быстро изменяющимся показателем.

¹⁴¹ Баржирование – перевозка на барже, в данном контексте перенос вне очереди.

¹⁴² Опрашиваемый метод `tryLock` всегда баржируется, даже для справедливой блокировки.

¹⁴³ График для `ConcurrentHashMap` довольно волнистый в области между четырьмя и восемью потоками. Эти колебания почти наверняка возникают из-за измерительного шума, который может быть введен случайными взаимодействиями с хэш-кодами элементов, планированием потоков, изменением размера `Map`, сборкой мусора или другими эффектами системы памяти или ОС, решающей запустить некоторую периодическую задачу по уборке во время выполнения тестового случая. Реальность такова, что в тестах производительности существуют всевозможные вариации, о контроле которых обычно

справедливости составляет почти два порядка. *Не платите за справедливость, если она вам не нужна.*

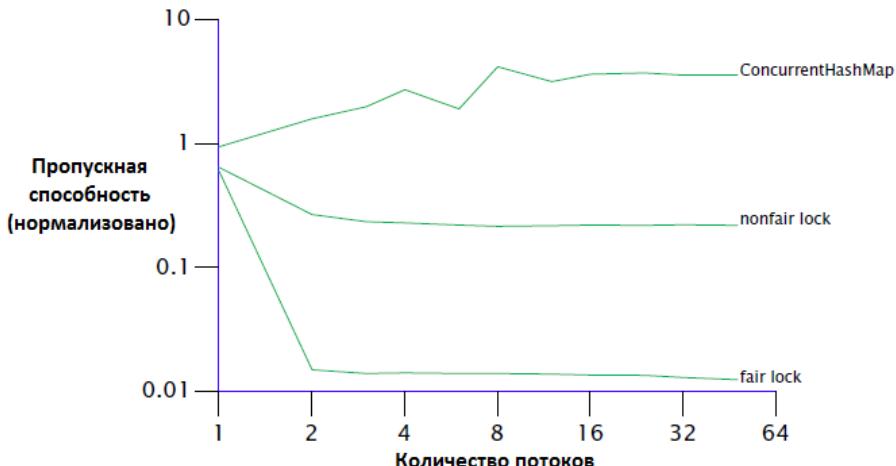


Рисунок 13.2 Сравнение производительности справедливой и несправедливой блокировок

Одной из причин, из-за которой баржированные блокировки работают намного лучше, чем справедливые блокировки, при интенсивной конкуренции, является то, что между возобновлением приостановленного потока и его фактическим запуском может быть значительная задержка. Предположим, поток **A** удерживает блокировку, а поток **B** запрашивает эту блокировку. Так как блокировка занята, поток **B** приостанавливается (*suspended*). Когда поток **A** освобождает блокировку, выполнение потока **B** возобновляется (*resumed*), чтобы он мог повторить попытку. Тем временем, однако, если поток **C** запрашивает блокировку, есть хороший шанс того, что **C** сможет захватить блокировку, использовать ее и отпустить её до того, как поток **B** закончит просыпаться. В этом случае выигрывают все: поток **B** захватывает блокировку не позже, чем получил бы её в ином случае, поток **C** захватывает блокировку намного раньше, а пропускная способность улучшается.

Справедливые блокировки, как правило, работают лучше всего, когда они удерживаются в течение относительно длительного времени или когда среднее время между запросами на захват блокировки относительно велико. В этих случаях условие, при котором баржирование обеспечивает преимущество в пропускной способности - когда блокировка освобождена, но поток в настоящее время просыпается, чтобы заявить на неё права - удерживается с меньшей вероятностью.

Как и в настройке класса `ReentrantLock` по умолчанию, встроенная блокировка не даёт никаких детерминированных гарантий справедливости, но статистические гарантии справедливости большинства реализаций блокировки достаточно хороши практически для всех ситуаций. Спецификация языка не требует, чтобы JVM реализовывала встроенные блокировки справедливо, и никакие производственные JVM этого не делают. Класс `ReentrantLock` не угнетает справедливость блокировки до новых минимумов - он только делает явным то, что присутствовало всё время.

беспокоиться не стоит. Мы не пытались искусственно очистить наши графики, потому что реальные измерения производительности также полны шума.

13.4 Выбор между synchronized и ReentrantLock

Класс `ReentrantLock` предоставляет ту же семантику блокировки и памяти, что и внутренняя блокировка, а также дополнительные функции, такие как ожидание ограниченной по времени блокировки, прерываемое ожидание блокировки, справедливость и возможность реализации не блочней структурированной блокировки. Производительность класса `ReentrantLock`, по всей видимости, доминирует над внутренней блокировкой, немного выигрывая в Java 6 и значительно в Java 5.0. Так почему бы не объявить устаревшим блок `synchronized` и не рекомендовать всему новому параллельному коду использовать класс `ReentrantLock`? Некоторые авторы фактически и предложили это, рассматривая блок `synchronized` как “унаследованную” конструкцию. Но это *зашло слишком далеко*.

Внутренние блокировки по-прежнему имеют значительные преимущества по сравнению с явными блокировками. Нотация знакома и компактна, и многие существующие программы уже используют встроенную блокировку - и смешивание этих двух видов блокировок может привести к путанице и седлать код подверженным ошибкам. Класс `ReentrantLock`, безусловно, более опасный в использовании инструмент, чем синхронизация; если вы забыли обернуть вызов метода `unlock` блоком `finally`, ваш код, вероятно, будет работать должным образом, но фактически вы создали бомбу замедленного действия, которая может повредить невинным наблюдателям. Приберегите класс `ReentrantLock` для ситуаций, в которых вам нужно что-то обеспечиваемое возможностями класса `ReentrantLock`, чего не делает встроенная блокировка.

Класс `ReentrantLock` представляет собой продвинутый инструмент, применяемый в тех ситуациях, в которых применение внутренних блокировок не практично. Используйте его, если вам нужны расширенные функции: ограничение по времени, опрос, или прерывание захвата блокировки, справедливое помещение в очередь, или не блочная структура блокировки. В ином случае, отдавайте предпочтение блоку `synchronized`.

В Java 5.0 встроенная блокировка имеет еще одно преимущество по сравнению с классом `ReentrantLock`: дампы потоков показывают кадры, в которых видно какими вызовами какие блокировки были захвачены и могут обнаруживать и идентифицировать потоки, попавшие в состояние взаимоблокировки. Среда JVM ничего не знает о том, какими потоками удерживается экземпляр `ReentrantLock` и поэтому не может помочь в отладке проблем с потоками, использующими класс `ReentrantLock`. Это несоответствие устранено в Java 6 путем предоставления интерфейса управления и мониторинга, с помощью которого блокировки могут регистрироваться, позволяя информацию об экземпляре блокировки `ReentrantLock` появляться в дампах потоков и, вследствие этого, в других интерфейсах управления и отладки. Доступность этой информации для отладки является существенным, хотя в основном и времененным, преимуществом блока `synchronized`; информация о блокировке в дампах потоков спасла многих программистов от полного ужаса. Не блочная структура класса `ReentrantLock` по-прежнему означает, что захват блокировок не может быть связан с определёнными кадрами стека, как это происходит в случае с внутренними блокировками.

Улучшения производительности в будущем, вероятно, предпочтут блок `synchronized` классу `ReentrantLock`. Поскольку блокировка `synchronized`

встроена в JVM, она может выполнять такие оптимизации, как устранение блокировки для ограниченных потоком объектов блокировки и укрупнение блокировок для устранения синхронизации со встроенными блокировками (см. раздел [11.3.2](#)); выполнение таких операций с блокировками на основе библиотек кажется гораздо менее вероятным. Если вы не планируете выполнять разворачивание в среде Java 5.0 в обозримом будущем, и у вас есть очевидная потребность в преимуществах масштабируемости, предоставляемых классом `ReentrantLock` на этой платформе, будет плохой идеей предпочесть класс `ReentrantLock` блоку `synchronized` по соображениям производительности.

13.5 Блокировки на чтение-запись

Класс `ReentrantLock` реализует стандартную взаимно-исключающую блокировку: экземпляр `ReentrantLock` может одновременно удерживаться не более чем одним потоком. Но взаимное исключение часто является более строгой дисциплиной блокировки, чем необходимо для сохранения целостности данных, и, таким образом, приводит к большему ограничению параллелизма, чем это необходимо. Взаимное исключение представляет собой консервативную стратегию блокировки, которая предотвращает перекрытие операций `запись/запись` и `запись/чтение`, но также предотвращает перекрытие операций `чтение /чтение`. Во многих случаях, структуры данных предназначены “в основном для чтения” - они изменяются и иногда модифицируются, но большинство обращений к ним происходит только для чтения данных. В этих случаях было бы неплохо ослабить требования к блокировке, чтобы позволить нескольким читателям одновременно получать доступ к структуре данных. До тех пор, пока каждому потоку гарантировано актуальное представление данных, и никакой другой поток не изменяет данные представления, в то время как читатели просматривают его, не будет никаких проблем. Это то, что позволяют получить блокировки на чтение-запись: ресурс может быть доступен одновременно либо нескольким читателям, либо одному писателю, но не обоим сразу.

Интерфейс `ReadWriteLock`, приведённый в листинге 13.6, предоставляет два объекта `Lock` - один для чтения и один для записи. Для чтения данных, защищаемых блокировкой `ReadWriteLock`, вы должны сначала захватить блокировку на чтение, а для изменения данных, защищаемых блокировкой `ReadWriteLock` вы должны сначала захватить блокировку на запись. Хотя может показаться, что существует две отдельные блокировки, блокировка на чтение и блокировка на запись - это просто разные представления интегрированного объекта блокировки на чтение-запись.

```
public interface ReadWriteLock {  
    Lock readLock();  
    Lock writeLock();  
}
```

Листинг 13.6 Интерфейс `ReadWriteLock`

Стратегия блокировки, реализуемая блокировками на чтение-запись, позволяет получать одновременный доступ нескольким читателям, но только одному писателю. Подобно интерфейсу `Lock`, интерфейс `ReadWriteLock` допускает множество реализаций, которые могут отличаться производительностью, гарантиями планировщика, преференциями при захвате, справедливостью или семантикой блокировки.

Блокировки на чтение-запись представляют собой оптимизацию производительности и были разработаны для обеспечения большего параллелизма в определенных ситуациях. На практике блокировки на чтение-запись, применяемые в многопроцессорных системах, могут повысить производительность часто используемых структур данных, предназначенных в основном для чтения; в других условиях они работают несколько хуже, чем монопольные блокировки, из-за большей сложности реализации. Являются ли они улучшением, лучше всего определяется в любой конкретной ситуации, с помощью выполнения профилирования; поскольку блокировка `ReadWriteLock` использует экземпляры `Lock` для чтения и записи заблокированными частями, относительно легко поменять блокировку типа чтение-запись на монопольную, если профилирование покажет, что блокировка на чтение-запись не даёт преимуществ.

Взаимодействие между блокировками на чтение и на запись допускает ряд возможных реализаций. Некоторым из вариантов реализации для `ReadWriteLock` являются:

Преимущества при освобождении. Когда писатель освобождает блокировку на запись, и оба, и читатель и писатель помещены в очередь, кому из них следует отдать предпочтение - читателям, писателям или тому, кто спросил первым?

Баржирование читателя. Если блокировка удерживается читателями, но есть ожидающие писатели, должен ли вновь прибывающим читателям быть предоставлен немедленный доступ, или они должны ожидать в очереди позади писателей? Предоставление читателям возможности баржирования перед писателями повышает параллелизм, но создает риск "голодания" для писателей.

Реентерабельность. Являются ли блокировки чтения и записи реентерабельными?

Понижение. Если поток удерживает блокировку на запись, может ли он захватить блокировку на чтение без освобождения блокировки на запись? Это позволило бы писателю "понизиться" до блокировки чтения, не позволяя при этом другим писателям изменять защищенный ресурс в то же самое время.

Повышение. Можно ли повысить блокировку на чтение до блокировки на запись в предпочтении другим ожидающим читателям или писателям? Большинство реализаций блокировки на чтение-запись не поддерживают повышение, поскольку без явной операции повышения оно подвержено взаимоблокировкам. (Если два читателя попытаются одновременно выполнить повышение до блокировки на запись, никто из них не освободит блокировку на чтение.)

Класс `ReentrantReadWriteLock` предоставляет реентерабельную семантику блокировки для обеих блокировок. Как и класс `ReentrantLock`, класс `ReentrantReadWriteLock` может быть построен несправедливым (по умолчанию) или справедливым. В случае справедливой блокировки предпочтение отдается потоку, который ожидал дольше; если блокировка удерживается читателями и поток запрашивает блокировку на запись, читателям не будет позволяться захватывать блокировку на чтение, пока писатель не будет обслужен и не освободит блокировку на запись. В случае несправедливой блокировки, порядок в котором потокам предоставляется доступ, не определен. Понижение от писателя до

читателя допускается; повышение с читателя до писателя нет (попытка выполнить это приведёт к взаимоблокировке).

Как и в классе `ReentrantLock`, блокировка на запись в классе `ReentrantReadWriteLock` имеет уникального владельца и может быть освобождена только тем потоком, который ее захватил. В Java 5.0 блокировка на чтение ведет себя скорее как семафор, чем как блокировка, поддерживая только количество активных читателей, но не их идентичности. Это поведение было изменено в Java 6, чтобы отслеживать, каким потокам были предоставлены блокировки на чтение.¹⁴⁴

Блокировки на чтение-запись могут улучшать параллелизм в тех случаях, когда блокировки обычнодерживаются в течение умеренно длительного времени, и большинство операций не изменяют защищённые ресурсы. Класс `ReadWriteMap` приведённый в листинге 13.7, использует класс `ReentrantReadWriteLock`, чтобы обернуть экземпляр `Map` таким образом, чтобы он мог безопасно использоваться несколькими читателями и при этом предотвращать конфликты типа чтение/запись или запись/запись.¹⁴⁵ На самом деле производительность класса `ConcurrentHashMap` настолько высока, что вы, вероятнее всего, будете использовать его вместо приведённого подхода, если всё, что вам нужно, это параллельная основанная на хэше реализация `Map`, но этот метод будет полезен, если вы хотите предоставить больший параллелизм при доступе к альтернативной реализации `Map`, такой как `LinkedHashMap`.

```
public class ReadWriteMap<K,V> {
    private final Map<K,V> map;
    private final ReadWriteLock lock = new ReentrantReadWriteLock();
    private final Lock r = lock.readLock(); private final Lock w =
lock.writeLock();

    public ReadWriteMap(Map<K,V> map) {
        this.map = map;
    }

    public V put(K key, V value) {
        w.lock();
        try {
            return map.put(key, value);
        } finally {
            w.unlock();
        }
    }
    // Do the same for remove(), putAll(), clear()

    public V get(Object key) {
        r.lock();
        try {
            return map.get(key);
        } finally {

```

¹⁴⁴ Одной из причин этого изменения является то, что в Java 5.0 реализация блокировки не может различать поток, запрашивающий блокировку чтения в первый раз от повторного запроса блокировки, что делает справедливые блокировки на чтение-запись подверженными взаимоблокировкам.

¹⁴⁵ Класс `ReadWriteMap` не реализует интерфейс `Map`, потому что реализация методов представления, таких как `entrySet` и `values`, довольно сложна, и, как правило, “простых” методов обычно достаточно.

```

        r.unlock();
    }
}

// Do the same for other read-only Map methods
}

```

Листинг 13.7 Обёртывание реализации Map с помощью блокировки на чтение-запись

На рисунок 13.3 приведено сравнение пропускной способности класса `ArrayList`, обернутого с помощью блокировки `ReentrantLock` и с помощью блокировки `ReadWriteLock`, выполненное в четырех-ядерной системе Opteron под управлением ОС Solaris. Тестовая программа, используемая здесь, похожа на тест производительности реализации `Map`, который мы использовали на протяжении всей книги - каждая операция случайным образом выбирает значение и ищет его в коллекции, а небольшой процент от общего количества операций изменяет содержимое коллекции.

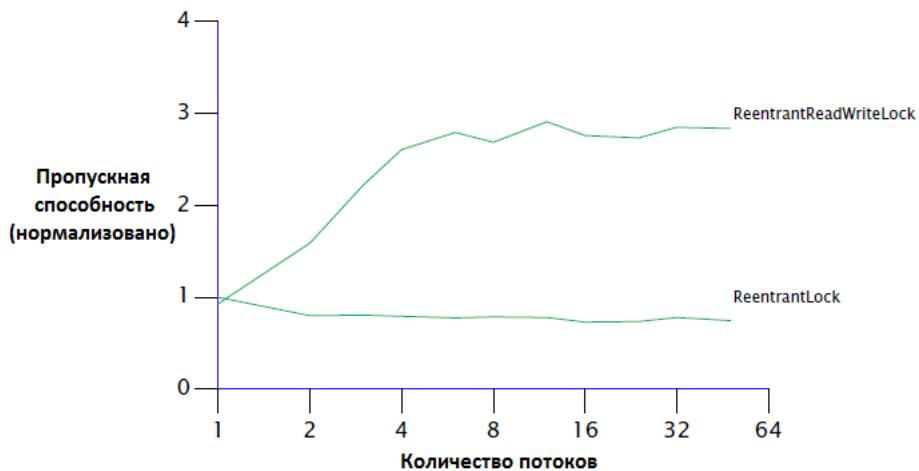


Рисунок 13.3 Производительность блокировок на чтение-запись

13.6 Итоги

Явные реализации `Lock` предлагают расширенный набор функций по сравнению со встроенной блокировкой, включая большую гибкость при работе с недоступностью блокировки и больший контроль над поведением очереди. Но класс `ReentrantLock` не является полной заменой блока `synchronized`; используйте его только тогда, когда вам нужны те функции, которых у блока `synchronized` нет.

Блокировки на чтение-запись позволяют нескольким читателям параллельно обращаться к защищаемому объекту, обеспечивая возможность улучшения масштабируемости при доступе к структурам данных, использующим преимущественно чтение.

Глава 14 Разработка собственных синхронизаторов

Библиотеки классов включают в себя несколько *зависящих от состояния классов* (*state-dependent classes*), – то есть, имеющих операции, зависящие от *предусловий определяемых состоянием* (*state-based preconditions*) - таких как `FutureTask`, `Semaphore` и `BlockingQueue`. Например, нельзя удалить элемент из пустой очереди или получить результат еще не завершенной задачи; перед выполнением этих операций необходимо дождаться, когда очередь перейдет в состояние “не пусто” или задача перейдет в состояние “завершено”.

Самый простой способ создать класс, зависящий от состояния, это построить его поверх существующего библиотечного класса, зависящего от состояния; мы поступили аналогичным образом с классом `ValueLatch` из раздела [8.5.1](#), используя класс `CountDownLatch` для обеспечения требуемого поведения блокировки. Но если библиотечные классы не предоставляют необходимых функциональных возможностей, можно также создать собственные синхронизаторы, используя низкоуровневые механизмы, предоставляемые языком и библиотеками, включая внутренние *очереди условий* (*condition queues*), явные объекты `Condition` и фреймворк `AbstractQueuedSynchronizer`. В этой главе рассматриваются различные варианты реализации зависимости от состояния, а также правила использования механизмов зависимости от состояния, предусмотренных платформой.

14.1 Управление зависимостью от состояния

В однопоточной программе, если предусловие на основе состояния (например, “пул соединений непустой”) при вызове метода не выполняется, оно никогда не станет истинным. Поэтому классы в последовательных программах могут быть закодированы на сбой, если их предварительные условия не выполняются. Но в параллельной программе условия, основанные на состоянии, могут изменяться под влиянием других потоков: пул, который был пуст несколько инструкций назад, может стать непустым, потому что другой поток возвратил элемент. Методы, зависящие от состояния параллельных объектов, иногда могут выйти из строя, когда их предусловия не выполняются, но часто существует лучшая альтернатива: дождаться момента, когда предусловие станет истинным.

Операции, зависящие от состояния и блокирующиеся до тех пор, пока операция не сможет продолжить выполнение, более удобны и менее подвержены ошибкам, в отличие от тех, что просто падают с ошибкой. Встроенный механизм очереди условий позволяет потокам блокироваться до тех пор, пока объект не войдёт в состояние, которое позволяет продолжить прогресс выполнения и пробуждать заблокированные потоки, когда они будут в состоянии в дальнейшем продолжить прогресс выполнения. Мы детально рассмотрим очереди условий в разделе [14.2](#), но с целью мотивации в оценке эффективности механизма ожидания условий, мы сначала покажем, как зависимость от состояния может быть (с большими сложностями) решена с помощью опроса и сна.

Блокирующееся зависящее от состояния действие, принимает форму, приведённую в листинге 14.1. Схема блокировки несколько необычна тем, что блокировка освобождается и вновь захватывается посреди выполнения операции.

Переменные состояния, составляющие предварительное условие, должны быть защищены блокировкой объекта, чтобы оставаться константными во время проверки предусловия. Но если предусловие не выполняется, блокировка должна быть снята, чтобы другой поток мог изменить состояние объекта, иначе предусловие никогда не станет истинным. Блокировка должна быть захвачена перед тем, как вновь проверять выполнение предусловия.

```
acquire lock on object state
while (precondition does not hold) {
    release lock
    wait until precondition might hold
    optionally fail if interrupted or timeout expires
    reacquire lock
}
perform action
release lock
```

Листинг 14.1 Структура блокирующихся зависимых от состояния действий

Ограничные буферы, подобные классу `ArrayBlockingQueue`, как правило, используются в дизайнах производитель-потребитель. Ограниченный буфер предоставляет операции *поместить* и *взять*, каждая из которых обладает предусловиями: нельзя взять элемент из пустого буфера или поместить элемент в полный буфер. Операции, зависящие от состояния, могут сталкиваться со сбоями предусловий, бросая при этом исключение или возвращая статус ошибки (перекладывая решение проблемы на вызывающего), или блокируясь до тех пор, пока объект не перейдет в нужное состояние.

Мы собираемся разработать несколько реализаций ограниченного буфера, которые используют различные подходы к обработке сбоев предусловий. Каждая из реализаций расширяет класс `BaseBoundedBuffer`, представленный в листинге 14.2, реализующий классический колцевой буфер на основе массива, в котором переменные состояния буфера (`buf`, `head`, `tail`, и `count`) защищены внутренней блокировкой буфера. Класс предоставляет синхронизированные методы `doPut` и `doTake`, которые используются подклассами для реализации операций `put` и `take`; базовое состояние скрыто от подклассов.

```
@ThreadSafe
public abstract class BaseBoundedBuffer<V> {
    @GuardedBy("this") private final V[] buf;
    @GuardedBy("this") private int tail;
    @GuardedBy("this") private int head;
    @GuardedBy("this") private int count;

    protected BaseBoundedBuffer(int capacity) {
        this.buf = (V[]) new Object[capacity];
    }

    protected synchronized final void doPut(V v) {
        buf[tail] = v;
        if (++tail == buf.length)
            tail = 0;
        ++count;
    }
```

```

}

protected synchronized final V doTake() {
    V v = buf[head];
    buf[head] = null;
    if (++head == buf.length)
        head = 0;
    --count;
    return v;
}

public synchronized final boolean isFull() {
    return count == buf.length;
}

public synchronized final boolean isEmpty() {
    return count == 0;
}
}

```

Листинг 14.2 Базовый класс для реализаций ограниченного буфера

14.1.1 Пример: распространение сбоев предусловий на вызывающий код

Класс GrumpyBoundedBuffer приведённый в листинге 14.3, представляет собой первую, и довольно грубую, попытку реализации ограниченного буфера. Методы `put` и `take` объявлены как `synchronized`, для обеспечения монопольного доступа к состоянию буфера, так как оба, при осуществлении доступа к буферу, используют логику проверить-затем-выполнить.

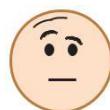
```

@ThreadSafe
public class GrumpyBoundedBuffer<V> extends BaseBoundedBuffer<V> {
    public GrumpyBoundedBuffer(int size) { super(size); }

    public synchronized void put(V v) throws BufferFullException {
        if (isFull())
            throw new BufferFullException();
        doPut(v);
    }

    public synchronized V take() throws BufferEmptyException {
        if (isEmpty())
            throw new BufferEmptyException();
        return doTake();
    }
}

```



Листинг 14.3 Ограниченный буфер, прерывающий свою работу, когда предварительные условия не выполняются

Хотя этот подход достаточно прост в реализации, он вызывает раздражение. Исключения должны возникать при исключительных условиях [ЕJ пункт 39]. Условие “Буфер заполнен” является исключительным условием для ограниченного буфера не больше, чем “красный” является исключительным условием в качестве сигнала трафику движения. Упрощение реализации буфера (принуждение вызывающего объекта к управлению зависимостью от состояния) с лихвой компенсируется существенным усложнением его использования, так как теперь вызывающий объект должен быть готов перехватывать исключения и, возможно, повторно выполнять каждую операцию буфера.¹⁴⁶ Хорошо структурированный вызов метода `take` показан в листинге 14.4 - не очень красиво, особенно если методы `put` и `take` вызываются по всей программе.

```
while (true) {
    try {
        V item = buffer.take();
        // use item
        break;
    } catch (BufferEmptyException e) {
        Thread.sleep(SLEEP_GRANULARITY);
    }
}
```

Листинг 14.4 Клиентская логика, используемая для обращения к классу `GrumpyBoundedBuffer`

Один из вариантов такого подхода заключается в том, чтобы возвращать значение ошибки, когда буфер находится в несогласованном состоянии. Это незначительное улучшение, по своей сути, заключается в том, что не происходит злоупотребления механизмом исключений, за счёт бросания исключения, которое в действительности означает “извините, попробуйте еще раз”, но это не решает фундаментальную проблему: вызывающие объекты должны самостоятельно разбираться со сбоями в предусловиях.¹⁴⁷

Клиентский код, приведённый в листинге 14.4, не является единственным способом реализации логики повтора. Вызывающий объект мог бы немедленно повторить вызов метода `take`, без засыпания - подход, известный как *напряжённое ожидание* (*busy waiting*) или *ожидание с проверкой* (*spin waiting*). Такое ожидание может отнимать довольно много процессорного времени, если состояние буфера не изменяется в течение некоторого времени. С другой стороны, если вызывающий объект решает уснуть, чтобы не потреблять так много времени CPU, он может легко “проспать”, если состояние буфера изменится вскоре после вызова метода `sleep`. Так что клиентскому коду остается выбор между плохим использованием CPU из-за прокручивания вызовов и плохой отзывчивостью из-за засыпания. (В каждой итерации, где-то между напряжённым ожиданием и сном, будет располагаться вызов метода `Thread.yield`, представляющий собой подсказку планировщику о том, что в текущий момент было бы разумно позволить работать

¹⁴⁶ Передача зависимого состояния обратно вызывающему объекту, также приводит к невозможности выполнения таких действий, как сохранение порядка FIFO; вынуждая вызывающий объект повторить попытку, вы теряете информацию о том, кто обратился первым.

¹⁴⁷ Интерфейс `Queue` предлагает оба варианта - метод `poll` возвращает значение `null`, если очередь пуста, а метод `remove` бросает исключение – но интерфейс `Queue` не предназначен для использования в дизайне производитель-потребитель. Интерфейс `BlockingQueue`, операции которого блокируются до тех пор, пока очередь, для продолжения работы, не перейдёт в согласованное состояние, является лучшим выбором, когда производители и потребители будут выполняться параллельно.

другому потоку. Если вы ожидаете выполнения другого потока, чтобы сделать что-то, что-то могло бы произойти быстрее в том случае, если бы вы уступали процессор, а не потребляли весь выделенный при планировании квант времени.)

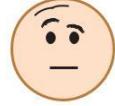
14.1.2 Пример: грубая блокировка путем опроса и сна

Класс `SleepyBoundedBuffer` представленный в листинге 14.5, пытается избавить вызывающие объекты от неудобств, вызванных реализацией логики повтора при каждом вызове, путём инкапсуляции повторяющихся вызовов грубого механизма “опросить и уснуть” внутрь операций `put` и `take`. Если буфер пуст, метод `take` засыпает до момента, пока другой поток не поместит некоторые данные в буфер; если буфер полон, метод `put` засыпает до момента, пока другой поток не освободит место, удалив некоторые данные. Этот подход инкапсулирует управление предусловиями и упрощает использование буфера – определенно, шаг в верном направлении.

```
@ThreadSafe
public class SleepyBoundedBuffer<V> extends BaseBoundedBuffer<V> {
    public SleepyBoundedBuffer(int size) { super(size); }

    public void put(V v) throws InterruptedException {
        while (true) {
            synchronized (this) {
                if (!isFull())
                    doPut(v);
                return;
            }
        }
        Thread.sleep(SLEEP_GRANULARITY);
    }

    public V take() throws InterruptedException {
        while (true) {
            synchronized (this) {
                if (!isEmpty())
                    return doTake();
            }
        }
        Thread.sleep(SLEEP_GRANULARITY);
    }
}
```



Листинг 14.5 Ограниченнный буфер, использующий грубую блокировку

Реализация класса `SleepyBoundedBuffer` сложнее, чем предыдущая попытка.¹⁴⁸ Код буфера должен проверять соответствующее условие состояния (*state condition*) удерживая блокировку буфера, поскольку переменные, представляющие состояние, защищены блокировкой буфера. Если проверка завершается неудачей, выполняющийся поток засыпает на некоторое время, перед этим освободив

¹⁴⁸ Мы избавим вас от деталей, касающихся других пяти реализаций ограниченного буфера, особенно `SneezyBoundedBuffer`.

блокировку, чтобы другие потоки могли получить доступ к буферу.¹⁴⁹ Когда поток проснется (*wakes up*), он повторно захватит блокировку и попытается снова, чередуя сон и проверку условия состояния до тех пор, пока выполнение операции не будет продолжено.

С точки зрения вызывающего объекта, это работает отлично - если операция может быть выполнена немедленно, так и будет сделано, а иначе она будет заблокирована - и абоненту не нужно иметь дело с механикой сбоя и повтором. Выбор детализации сна - это компромисс между отзывчивостью и использованием CPU; чем меньше степень детализации сна, тем больше отзывчивость, но и тем больше потребление ресурсов процессора. На рис. 14.1 показано, каким образом детализация сна может влиять на скорость отклика: может быть задержка между моментом времени, когда пространство буфера становится доступным и моментом времени, когда поток просыпается и снова выполняет проверку.

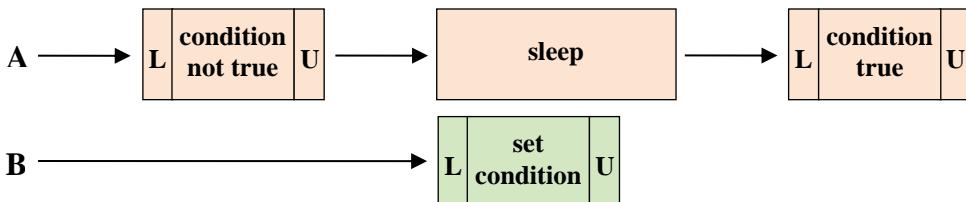


Рисунок 14.1 Поток проспал, потому что условие приняло значение `true` после того, как он захотел уснуть

Класс `SleepyBoundedBuffer` также порождает другое требование для вызывающего объекта - работа с исключением `InterruptedException`. Когда метод блокируется на ожидании выполнения условия, вежливое действие заключается в предоставлении механизма отмены (см. главу 7). Подобно большинству библиотечных блокирующих методов с хорошим поведением, класс `SleepyBoundedBuffer` поддерживает отмену через прерывание, возвращая управление раньше и бросая исключение `InterruptedException`, если выполнение было прервано.

Эти попытки синтезировать блокирующую операцию из опроса и сна были довольно болезненными. Было бы неплохо иметь способ приостановки потока, но гарантировать, что он быстро пробудится, как только определенное условие (например, буфер больше не полон) становится истинным. Это именно то, чем занимаются очереди условий.

14.1.3 Очереди условий на освобождение

Очереди условий похожи на звонок “тост готов” в вашем тостере. Если вы слышите звонок, вы будете немедленно уведомлены, когда ваш тост будет готов и сможете бросить то, чем занимались (или нет, может быть, вы хотите в первую очередь закончить чтение газеты) и получить тост. Если вы не слышите звонок (возможно, вы вышли на улицу, чтобы получить газету), вы можете пропустить уведомление, но по возвращении на кухню вы можете наблюдать за состоянием тостера и либо получить тост, если он готов, либо вновь начать слушать звонок, если это не так.

¹⁴⁹ Как правило, для потока является плохой идеей осуществление перехода в спящий режим или блокирование иным образом совместно с удерживаемой блокировкой, но в данном случае всё еще хуже, потому что желаемое условие (буфер полон/пуст) никогда не сможет стать истинным, если блокировка не будет освобождена!

Очередь условий (*condition queue*) получила свое имя, поскольку она предоставила группе потоков - называемой *набором ожидания* (*wait set*) - способ ожидания выполнения определенного условия. В отличие от обычных очередей, в которых элементы являются единицами данных, элементами очереди условий являются потоки, ожидающие выполнения условия.

Так же, как каждый объект Java может выступать в качестве блокировки, каждый объект может также выступать в качестве очереди условий, и методы `wait`, `notify`, и `notifyAll` класса `Object` составляют API внутренних очередей условий. Внутренняя блокировка объекта и его внутренняя очередь условий связаны: чтобы вызвать любой из методов очереди условий объекта *X*, необходимо удерживать блокировку на объекте *X*. Это связано с тем, что механизм ожидания условий на основе состояний тесно связан с механизмом сохранения согласованности состояний: не получится дождаться условия, пока не удастся проверить состояние, и нельзя освободить другой поток из состояния ожидания по условию, пока не удастся изменить состояние.

Метод `Object.wait` атомарно освобождает блокировку и просит ОС приостановить выполнение текущего потока, позволяя другим потокам захватить блокировку и, следовательно, изменить состояние объекта. После пробуждения он повторно захватывает блокировку перед возвращением. Интуитивно, вызов `wait` означает “Я хочу пойти спать, но разбудите меня, когда произойдет что-то интересное”, а вызов методов уведомления означает “что-то интересное произошло”.

Представленный в листинге 14.6 класс `BoundedBuffer`, реализует ограниченный буфер с помощью методов `wait` и `notifyAll`. Этот вариант проще, чем в случае со “спящей” версией, и является более эффективным (просыпаться приходится реже, если состояние буфера не изменяется) и более отзывчивым (когда происходит интересующее изменение состояния, просыпаться получается быстро). Это значительное улучшение, но обратите внимание, что введение очередей условий не привело к внесению изменений в семантику, по сравнению со “спящей” версией. Это просто оптимизация в нескольких измерениях: эффективности использования ЦП, издержек переключения контекста и отзывчивости. Очереди условий не позволяют делать то, чего вы не можете сделать с помощью опроса и сна¹⁵⁰, но они позволяют сделать всё намного проще и эффективнее выражать и управлять зависимостью от состояния.

```
@ThreadSafe
public class BoundedBuffer<V> extends BaseBoundedBuffer<V> {
    // CONDITION PREDICATE: not-full (!isFull())
    // CONDITION PREDICATE: not-empty (!isEmpty())

    public BoundedBuffer(int size) { super(size); }

    // BLOCKS-UNTIL: not-full
    public synchronized void put(V v) throws InterruptedException {
        while (isFull())
            wait();
        doPut(v);
    }
}
```

¹⁵⁰ Это не совсем верно; справедливая очередь условий может гарантировать относительный порядок, в котором потоки освобождаются из набора ожидания. Внутренние очереди условий, подобно внутренним блокировкам, не предлагают справедливого помещения в очередь; явные реализации интерфейса `Condition` предлагают выбор из справедливого и несправедливого помещения в очередь.

```

        notifyAll();
    }

    // BLOCKS-UNTIL: not-empty
    public synchronized V take() throws InterruptedException {
        while (isEmpty())
            wait();
        V v = doTake();
        notifyAll();
        return v;
    }
}

```

Листинг 14.6 Ограниченный буфер с использованием очередей условий

Класс `BoundedBuffer`, наконец-то, достаточно хорош в использовании - он прост в использовании и разумно управляет зависимостью от состояния.¹⁵¹ Производственная версия должна также включать синхронизированные версии методов `put` и `take`, чтобы блокировка операций могла прерываться по таймауту, если они не смогут завершиться в течение выделенного бюджета времени. Синхронизированная версия метода `Object.wait` упрощает реализацию.

14.2 Использование очередей условий

Очереди условий упрощают создание эффективных и отзывчивых классов, зависящих от состояния, но их всё ещё довольно просто использовать некорректно; существует множество правил относительно их корректного использования, которые не применяются компилятором или платформой. (Это одна из причин для того, чтобы строить ваши классы, в тех случаях, когда это возможно, на основе классов, подобных `LinkedBlockingQueue`, `CountDownLatch`, `Semaphore`, и `FutureTask`; если вам это удастся, это значительно упростит дело.)

14.2.1 Предикат условия

Ключом к правильному использованию очередей условий является определение *предикатов условий* (*condition predicates*), которые может ожидать объект. Это предикат условия приводит к путанице вокруг методов `wait` и `notify`, потому что у него нет создания экземпляра в API, и ни спецификация языка, ни реализация JVM не гарантирует его правильного использования. Фактически, о нём напрямую не упоминается ни в спецификации языка, ни в Javadoc. Но без него, ожидание по условию работать не будет.

Предикат условия представляет собой предусловие, которое в первую очередь делает операцию зависимой от состояния. В ограниченном буфере метод `take` может продолжить выполнение только в том случае, если буфер не пуст; в противном случае он должен ожидать. Для метода `take` предикатом условия является утверждение “буфер не пуст”, которое метод `take` должен проверить перед продолжением выполнения. Точно так же предикатом условия для метода `put` будет “буфер не полон”. Предикаты условий представляют собой выражения, построенные из переменных состояния класса; проверка экземпляра класса `BaseBoundedBuffer` на соответствие утверждению “буфер не пуст” выполняется

¹⁵¹ Класс `ConditionBoundedBuffer` из раздела 14.3 еще лучше: он более эффективен, потому что может использовать одно уведомление вместо вызова метода `notifyAll`.

путём сравнения значения поля `count` со значением ноль, а проверка утверждения “буфер не полон” выполняется путём сравнения значения поля `count` с размером буфера.

Документируйте предикаты условий, связанные с очередью условий, и ожидающими на них операциями.

Условие ожидания включает в себя важное трехстороннее отношение: блокировку, метод `wait` и предикат условия. Предикат условия включает в себя переменные состояния, а переменные состояния защищены блокировкой, поэтому перед проверкой предиката условия, мы должны удерживать эту блокировку. Объект блокировки и объект очереди условий (объект, на котором вызываются методы `wait` и `notify`) должны быть одним и тем же объектом.

В классе `BoundedBuffer` состояние буфера защищено блокировкой буфера, а объект буфера используется в качестве очереди условий. Метод `take` захватывает блокировку буфера, а затем проверяет предикат условия (буфер не пуст). Если буфер действительно не пуст, он удаляет первый элемент, в отношении которого может это выполнить, потому что он по-прежнему удерживает блокировку, защищающую состояние буфера.

Если предикат условия не выполняется (буфер пуст), метод `take` должен ожидать, пока другой поток не поместит объект в буфер. Это делается с помощью вызова метода `wait` на внутренней очереди условий буфера, что требует удержания блокировки на объекте очереди условий. Так как в продуманном дизайне это было реализовано, метод `take` уже удерживает ту блокировку, которая ему нужна для проверки предиката условия (и если бы предикат условия выполнялся, то метод изменил бы состояние буфера в той же атомарной операции). Метод `wait` освобождает блокировку, блокирует текущий поток и ожидает истечения указанного времени ожидания, прерывания потока или пробуждения потока поступившим уведомлением. После того, как поток просыпается, метод `wait`, перед возвратом управления, повторно захватывает блокировку. Поток просыпается из метода `wait` без получения особого приоритета в повторном захвате блокировки; он конкурирует за блокировку, подобно любому другому потоку, пытаясь войти в блок `synchronized`.

Каждый вызов метода `wait` неявно связан с определенным *предикатом условия*. При вызове метода `wait` относительно определенного предиката условия, вызывающий объект должен уже удерживать блокировку, связанную с очередью условий, и эта блокировка также должна защищать переменные состояния, из которых состоит предикат условия.

14.2.2 Раннее пробуждение

Как будто бы трехстороннее отношение между блокировкой, предикатом условия и очередью условий не было достаточно сложным, так ещё и возврат управления из метода `wait` не обязательно означает, что предикат условия, ожидаемый потоком, возвращает `true`.

Единственная внутренняя очередь условий может использоваться более чем одним предикатом условия. Когда ваш поток пробуждается из-за того, что кто-то вызвал метод `notifyAll`, это не означает, что предикат условия, которого вы ожидали, теперь истинен. (Это похоже на то, как если бы ваш тостер и кофеварка

имели один и тот же звонок; когда он звонит, вам все равно нужно посмотреть, какое устройство подало сигнал.)¹⁵² Кроме того, методу `wait` разрешается даже “ложно” возвращать управление - не в ответ на вызов метода `notify` каким-либо потоком.¹⁵³

Когда управление повторно входит в код, вызвавший метод `wait`, метод повторно захватывает блокировку, связанную с очередью условий. Теперь предикат условия истинен? Возможно. Это могло быть истиной на момент уведомления потока вызовом метода `notifyAll`, но могло бы опять стать ложью, когда вы повторно захватите блокировку. Другие потоки могли захватить блокировку и изменить состояние объекта, в промежутке между моментом пробуждения потока и моментом повторного захвата блокировки. Или, может быть, условие вообще не было истинным с момента, как вы вызвали метод `wait`. Вы не знаете, почему другой поток вызвал методы `notify` или `notifyAll`; возможно, это произошло потому, что другой предикат условия, связанный с той же очередью условий, стал истинным. Наличие множество предикатов условий для каждой очереди состояния довольно распространённое явление - класс `BoundedBuffer` использует одну и ту же очередь условий в таких предикатах, как “не полон”, так и “не пуст”.¹⁵⁴

В связи со всеми этими причинами, когда вы просыпаетесь из метода `wait`, вы должны вновь проверить предикат условия и вернуться к ожиданию (или упасть со сбоем), если условие ещё не истинно. Так как вы можете несколько раз просыпаться без выполнения предиката условия (значение предиката не равно `true`), вы должны всегда вызывать метод `wait` в цикле, проверяя предикат условия в каждой итерации. В листинге 14.7 приведена каноническая форма для ожидания по условию.

```
void stateDependentMethod() throws InterruptedException {
    // condition predicate must be guarded by lock
    synchronized(lock) {
        while (!conditionPredicate())
            lock.wait();
        // object is now in desired state
    }
}
```

Листинг 14.7 Каноническая форма методов, зависимых от состояния

При использовании ожидания по условию (`Object.wait` или `Condition.await`):

- Всегда иметь предикат условия - некоторая проверка состояния объекта, которая должен выполняться перед продолжением;
- Всегда проверяйте предикат условия перед вызовом метода `wait` и после возврата из метода `wait`;

¹⁵² Эта ситуация, на самом деле, прекрасно описывает кухню Тима; так много устройств подает звуковой сигнал, что, когда вы слышите его, вы должны проверить тостер, микроволновую печь, кофе-машину и несколько других устройств, чтобы определить причину подачи сигнала.

¹⁵³ Если и дальше развивать аналогию с завтраком, это похоже на тостер со свободным соединением, которое заставляет звонок включаться, когда тост готов, но также иногда и в том случае, когда он ещё не готов.

¹⁵⁴ Фактически возможна ситуация, когда потоки могут ожидать выполнения предикатов “не полон” и “не пуст” одновременно! Это может произойти, когда несколько производителей / потребителей достигнут пределов емкости буфера.

- Убедитесь, что переменные состояния, составляющие предикат условия, защищены блокировкой, связанной с очередью условий;
- Удерживайте блокировку, связанную с очередью условий, при вызове методов `wait`, `notify` или `notifyAll`; и
- Не освобождайте блокировку после проверки предиката условия, но до выполнения действий с ним.

14.2.3 Пропущенные сигналы

В главе 10 обсуждались проблемы живучести, такие как взаимоблокировка и динамическая блокировка. Другой формой проблем с живучестью являются *пропущенные сигналы* (*missed signals*). Пропущенный сигнал возникает тогда, когда поток должен дождаться определенного условия, которое уже истинно, но не может проверить предикат условия перед началом ожидания. Таким образом, поток ожидает уведомления о событии, которое уже произошло. Это все равно, что начать готовить тост, выйти за газетой, пока вы находитесь снаружи, сработает звонок, а затем сесть за кухонный стол в ожидании теста. Вы можете ожидать длительное время - потенциально бесконечно.¹⁵⁵ В отличие от мармелада для тостов, уведомление не является “липким” (*sticky*) - если поток **A** уведомляет очередь условий, а поток **B** впоследствии ожидает в этой же очереди условий, поток **B** просыпается *не* сразу - для пробуждения потока **B** требуется другое уведомление. Пропущенные сигналы являются результатом ошибок кодирования, подобных тем, предупреждение о которых было приведено в списке выше, например, провал при попытке выполнения проверки предиката условия, перед вызовом метода `wait`. Если вы структурируете свое состояние ожидания, как показано в листинге 14.7, у вас не будет проблем с пропущенными сигналами.

14.2.4 Уведомление

До сих пор мы описывали только половину того, что происходит в ожидании по условию: ожидание. Другая половина - уведомление. В ограниченном буфере метод `take` блокируется, если вызывается в тот момент, когда буфер пуст. Для того чтобы *разблокировать* метод `take`, когда буфер становится непустым, мы должны гарантировать, что в каждой ветке выполнения кода, в которой буфер может стать непустым, выполняется уведомление. В классе `BoundedBuffer` есть только одно такое место – после метода `put`. Поэтому метод `put`, после успешного добавления объекта в буфер, вызывает метод `notifyAll`. Подобным образом, метод `take` вызывает метод `notifyAll` после удаления элемента, чтобы указать, что буфер больше не полон, в случае, если какие-либо потоки ожидают выполнения условия “буфер не полон”.

Всякий раз, ожидая выполнения условия, убедитесь, что кто-то выполнит уведомление, когда предикат условия станет истинным.

Существует два метода уведомления в API очереди условий - `notify` и `notifyAll`. Для вызова любого из них необходимо удерживать блокировку, связанную с объектом очереди условий. Вызов метода `notify` приводит к тому, что

¹⁵⁵ Чтобы выйти из этого ожидания, кто-то другой должен был бы сделать тост, но это только ухудшит ситуацию; когда раздастся звонок, у вас возникнут разногласия по поводу владения тестом.

JVM выбирает один поток, ожидающий пробуждения на очереди условий; вызов `notifyAll` пробуждает все потоки, ожидающие на очереди условий. Поскольку при вызове `notify` или `notifyAll` необходимо удерживать блокировку объекта очереди условий, а ожидающие потоки не могут вернуться из режима ожидания без повторного захвата блокировки, уведомляющий поток должен быстро освободить блокировку, чтобы гарантировать, что ожидающие потоки будут разблокированы как можно скорее.

Класс `BoundedBuffer` представляет собой хорошую иллюстрацию того, почему следует отдавать предпочтение методу `notifyAll`, вместо метода `notify`, в большинстве случаев. Очередь условий используется для двух различных предикатов условий: “не полон” и “не пуст”. Предположим, поток **A** ожидает в очереди условий предиката P_A , а поток **B** ожидает в той же очереди условий предиката P_B . Теперь предположим, что условие предиката P_B становится истинным, а поток **C** выполняет только одно уведомление: среда JVM разбудит один поток по своему выбору. Если будет выбран поток **A**, он проснется, увидит, что предикат P_A еще не истинен, и вернется к ожиданию. Между тем, поток **B**, который смог бы продолжить выполнение, не просыпается. Это не совсем пропущенный сигнал - это скорее “захваченный сигнал” (*hijacked signal*), но проблема одна и та же: поток ожидает сигнала, который уже произошел (или должен был произойти).

Единственный вызов метода `notify` может использоваться вместо вызова метода `notifyAll` только в том случае, когда выполняются оба следующих условия:

Единообразное ожидание. Только один предикат условия связан с очередью условий, и каждый из поток выполняет одну и ту же логику при возврате управления из метода `wait`; и

Один вход, один выход. Уведомление на условной переменной позволяет продолжить выполнение не более чем одного потока.

Класс `BoundedBuffer` встречается с требованием “один вход, один выход”, но не встречается с требованием единообразного ожидания, потому что ожидающие потоки, могут ожидать как условия “не полон”, так и условия “не пуст”. Защелка “стартового затвора”, подобная той, что используется в классе `TestHarness` из раздела [5.5.1](#), в которой одно событие освобождает набор потоков, не соответствует требованию “один вход, один выход”, поскольку открытие начального затвора позволяет продолжить выполнение множества потоков.

Большинство классов не отвечают этим требованиям, поэтому использование метода `notifyAll` вместо одиночного метода `notify` будет являться мудрым решением. Хотя это может быть и неэффективно, при использовании метода `notifyAll` вместо метода `notify`, гораздо проще обеспечить правильное поведение классов.

Это “устоявшееся мнение” (*prevailing wisdom*) заставляет некоторых людей чувствовать себя некомфортно, и не зря. Использование метода `notifyAll`, когда только один поток может добиться прогресса, неэффективно - иногда немного, иногда очень весомо. Если десять потоков ожидают на очереди условий, вызов `notifyAll` приводит к тому, что каждый из них просыпается и конкурирует за блокировку; затем большинство из них, или все, снова переходят в спящий режим. Это означает, что выполняется множество переключений контекста и множество

конкурирующих захватов блокировки для каждого события, которое позволяет (возможно) одному потоку продолжить выполнение. (В худшем случае использование метода `notifyAll` приводит к $O(n^2)$ пробуждениям, где n пробуждений будет достаточно.) Это еще одна ситуация, когда проблемы производительности поддерживаются одним подходом, а проблемы безопасности - другим.

Уведомление, выполняемое методами `put` и `take` класса `BoundedBuffer`, является консервативным: уведомление выполняется каждый раз, когда объект помещается в буфер или удаляется из него. Это можно оптимизировать, наблюдая, что поток может быть освобожден из ожидания, только если буфер переходит от пустого к не пустому или от полного к не полному, и уведомляя, только если методы `put` или `take` оказывают влияние на эти переходы состояния. Это называется *уведомлением по условию* (*conditional notification*). В то время как уведомление по условию может улучшить производительность, его сложно использовать правильно (оно также усложняет реализацию подклассов) и поэтому должно использоваться аккуратно. В листинге 14.8 приведён пример использования уведомления по условию в методе `BoundedBuffer.put`.

```
public synchronized void put(V v) throws InterruptedException {
    while (isFull())
        wait();
    boolean wasEmpty = isEmpty();
    doPut(v);
    if (wasEmpty)
        notifyAll();
}
```

Листинг 14.8 Использование уведомления по условию в методе `BoundedBuffer.put`

Одиночное уведомление и условное уведомление являются оптимизациями. Как всегда, при использовании этих оптимизаций, следуйте принципу “Сначала сделайте это правильно, а затем сделайте это быстро - если это еще не достаточно быстро”; легко ввести странные сбои живучести, из-за неправильного применения уведомлений.

14.2.5 Пример: класс затвора

Защёлка “стартового затвора” в классе `TestHarness`, приведённом в разделе [5.5.1](#), была построена с начальным значением счётчика равным единице, создав, таким образом, *бинарную защёлку* (*binary latch*): одну с двумя состояниями, начальным состоянием и конечным состоянием. Защёлка препятствует потокам в преодоления стартового затвора до тех пор, пока она не будет раскрыта, и в этот момент все потоки смогут пройти. Несмотря на то, что механизм защёлки часто является именно тем, что необходимо, иногда является недостатком, что затвор, построенный таким образом, не может быть повторно закрыт после открытия.

Достаточно просто разработать повторно закрываемый класс `ThreadGate`, используя ожидание по условию, как показано в листинге 14.9. Класс `ThreadGate` позволяет затвору открываться и закрываться, предоставляя метод `await`, который блокируется, пока затвор открыт. Метод `open` использует метод `notifyAll`, так как семантика этого класса не позволяет выполнить проверку “один вход, один выход” для единственного уведомления.

```
@ThreadSafe
public class ThreadGate {
    // CONDITION-PREDICATE: opened-since(n) (isOpen || generation>n)
    @GuardedBy("this") private boolean isOpen; @GuardedBy("this")
    private int generation;

    public synchronized void close() {
        isOpen = false;
    }

    public synchronized void open() {
        ++generation;
        isOpen = true;
        notifyAll();
    }

    // BLOCKS-UNTIL: opened-since(generation on entry)
    public synchronized void await() throws InterruptedException {
        int arrivalGeneration = generation;
        while (!isOpen && arrivalGeneration == generation)
            wait();
    }
}
```

Листинг 14.9 Повторно закрываемый затвор, использующий методы `wait` и `notifyAll`

Предикат условия, используемый методом `await`, сложнее простой проверки метода `isOpen`. Это необходимо, потому что если N потоков ожидают у затвора в то время, когда он открыт, им всем должно быть позволено продолжить выполнение. Но, если бы затвор открывался и закрывался в быстрой последовательности, не все потоки могли бы быть освобождены, если бы метод `await` проверял только метод `isOpen`: к тому времени, когда все потоки получают уведомление, повторно захватывают блокировку и выходят из метода `wait`, затвор, возможно, уже вновь закрылся. Поэтому класс `ThreadGate` использует несколько более сложный предикат условия: каждый раз, когда затвор открывается, счетчик “генерации” увеличивается, и поток может пройти метод `await`, если затвор открыт сейчас, или если затвор открылся с тех пор, как этот поток прибыл к нему.

Так как класс `ThreadGate` поддерживает только ожидание открытия затвора, он выполняет уведомление только в методе `open`; для поддержки обеих операций “ожидать открытия” и “ожидать закрытия”, ему придётся выполнять уведомления и в методе `open`, и в методе `close`. Это является иллюстрацией того, почему зависящие от состояния классы могут быть хрупкими в поддержке - добавление новой, зависящей от состояния, операции может потребовать изменения множества веток выполнения кода, которые изменяют состояние объекта, так чтобы соответствующие уведомления могли быть выполнены.

14.2.6 Вопросы безопасности подклассов

Использование одиночного уведомления или уведомления по условию вводит ограничения, которые могут усложнить реализацию подклассов [CPJ 3.3.3.3]. Если

вы в целом хотите поддерживать подклассы, вам необходимо структурировать свой класс так, чтобы подклассы могли добавлять соответствующее уведомление от имени базового класса, если подкласс выполняет это таким образом, который нарушает одно из требований одиночного уведомления или уведомления по условию.

Класс, зависящий от состояния, должен либо полностью раскрывать (и документировать) свои протоколы ожидания и уведомления для подклассов, либо вообще не включать в них подклассы. (Это расширение принципа “дизайн и документирование для наследования, или же его запрет” [ЕJ пункт 15].) По крайней мере, при разработке класса, зависящего от состояния, в целях наследования требуется предоставить доступ к очередям условий и блокировкам, а также документировать предикаты условий и политику синхронизации; также может потребоваться предоставить переменные состояния, лежащие в основе. (Самое худшее, что может сделать зависящий от состояния класс, - это предоставить своё состояние подклассам, но не документировать свои протоколы ожидания и уведомления; это похоже на класс, предоставляющий свои переменные состояния, но не документирующий свои инварианты.)

Одним из вариантов сделать это, является запрет на создание подклассов, либо объявив класс как `final`, либо скрывая от подклассов очереди условий, блокировки и переменные состояния. В противном случае, если подкласс делает что-то для дискредитации способа, с помощью которого базовый класс использует метод `notify`, он должен быть в состоянии возместить ущерб. Рассмотрим неограниченный блокирующий стек, в котором операция изъятия (`pop`) блокируется, если стек пуст, но операция вложить (`push`) всегда может быть выполнена. Это удовлетворяет требованию для одиночного уведомления. Если этот класс использует одиночное уведомление, а подкласс добавляет блокирующий метод “последовательно изъять два элемента”, то теперь существует два класса ожиданий: ожидающие появления одного элемента и ожидающие появления двух. Но если базовый класс предоставляет очередь условий и документирует свои протоколы по её использования, подкласс может переопределить метод вложения, для выполнения метода `notifyAll`, таким образом, восстанавливая безопасность.

14.2.7 Инкапсулирование очередей условий

Обычно лучше инкапсулировать очередь условий, чтобы она была недоступна вне иерархии классов, в которой она используется. В противном случае у вызывающих объектов может возникнуть соблазн подумать, что они понимают ваши протоколы ожидания и уведомления и могут использовать их в не соответствующей вашему дизайну манере. (Невозможно принудительно применить требование о единообразном ожидании для одиночного уведомления, если объект очереди условий недоступен для кода, который вы не контролируете; если чужой код ошибочно ожидает на очереди условий, это может привести к нарушению протокола уведомлений и вызвать захват сигнала.)

К сожалению, этот совет - инкапсулировать объекты, используемые в качестве очередей условий - не согласуется с наиболее распространенным шаблоном проектирования для потокобезопасных классов, в котором внутренняя блокировка объекта используется для защиты его состояния. Класс `BoundedBuffer` иллюстрирует эту распространенную идиому, в нём объект буфера сам по себе является и блокировкой и очередью условий. Однако класс `BoundedBuffer` может быть легко реструктурирован для использования приватного объекта блокировки и

очереди условий; единственное различие будет заключаться в том, что он больше не будет поддерживать любую форму блокировки на стороне клиента.

14.2.8 Протоколы входа и выхода

Веллингс (Wellings, 2004) характеризует надлежащее использование методов `wait` и `notify` с точки зрения *протоколов входа (entry)* и *выхода (exit)*. Для каждой операции, зависящей от состояния, и для каждой операции, изменяющей состояние, от которого зависит другая операция, необходимо определить и задокументировать протоколы входа и выхода. Протокол входа - это предикат условия операции; протокол выхода включает в себя проверку всех переменных состояния, которые были изменены операцией, чтобы узнать, не вызвали ли они выполнение какого-либо другого предиката условия, и если это так, то отправить уведомление связанной очереди условий.

Класс `AbstractQueuedSynchronizer`, на основе которого построено большинство зависимых от состояния классов из пакета `java.util.concurrent` (см. раздел [14.4](#)), использует концепцию протокола выхода. Вместо того чтобы позволить классам синхронизатора выполнять своё собственное уведомление, он вместо этого требует, чтобы методы синхронизатора возвращали значение, указывающее, могут ли его действия разблокировать один или несколько ожидающих потоков. Это явное требование API затрудняет возникновение ситуации, при которой можно “забыть” выполнить уведомление о некоторых переходах состояний.

14.3 Явные объекты условия

Как мы видели в главе [13](#), явные блокировки могут быть полезны в тех случаях, когда встроенные блокировки не обладают достаточной гибкостью. Подобно тому, как интерфейс `Lock` является обобщением внутренних блокировок, интерфейс `Condition` (см. листинг 14.10) является обобщением внутренних очередей условий.

```
public interface Condition {  
    void await() throws InterruptedException;  
    boolean await(long time, TimeUnit unit)  
        throws InterruptedException;  
    long awaitNanos(long nanosTimeout) throws InterruptedException;  
    void awaitUninterruptibly();  
    boolean awaitUntil(Date deadline) throws InterruptedException;  
  
    void signal();  
    void signalAll();  
}
```

Листинг 14.10 Интерфейс `Condition`

Внутренние очереди условий имеют несколько недостатков. Каждая внутренняя блокировка может быть связана только с одной очередью условий, что означает, что в классах подобных классу `BoundedBuffer`, множество потоков может ожидать на одной очереди условий различных предикатов условий, а наиболее распространенный шаблон блокировки включает предоставление объекта очереди

условий. Оба этих фактора делают невозможным применение единообразного требования к ожидающим, в отношении использования метода `notify`. Если вы хотите написать параллельный объект с несколькими предикатами условий или хотите осуществлять больший контроль над видимостью очереди условий, явные классы `Lock` и `Condition` предлагают более гибкую альтернативу внутренним блокировкам и очередям условий.

Экземпляр `Condition` связан с одиночным экземпляром `Lock`, так же как очередь условий связана с одиночной внутренней блокировкой; для создания экземпляра `Condition`, вызовите метод `Lock.newCondition` у связанной блокировки. Подобно тому, как интерфейс `Lock` предлагает более богатый набор функций, чем внутренняя блокировка, интерфейс `Condition` предлагает более богатый набор функций, чем внутренние очереди условий: несколько наборов ожиданий (*wait set*) для каждой блокировки, прерываемые и непрерываемые ожидания по условию, ожидания, ограниченные по сроку, и выбор из справедливого и несправедливого помещения в очередь.

В отличие от встроенных очередей условий, можно иметь столько объектов `Condition` для каждого объекта `Lock`, сколько вам необходимо. Объекты `Condition` наследуют настройки справедливости от связанного с ними экземпляра `Lock`; в случае использования справедливых блокировок, потоки освобождаются из метода `Condition.await` в порядке FIFO.

Предупреждение об опасности: эквивалентом методов `wait`, `notify` и `notifyAll` для объектов `Condition` являются методы `await`, `signal` и `signalAll`. Однако класс `Condition`¹⁵⁶ расширяет класс `Object`, что означает, что он также имеет методы `wait` и `notify`. Обязательно используйте правильные версии методов - `await` и `signal` - вместо наследуемых!

В листинге 14.11 показана еще одна реализация ограниченного буфера, на этот раз использующая два экземпляра `Condition`, - `notFull` и `notEmpty` - для явного представления предикатов условий “не полон” и “не пуст”. Когда метод `take` блокируется из-за того, что буфер пуст, он ожидает на условии `notEmpty`, и метод `put` разблокирует любые потоки, заблокированные в методе `take`, отправляя сигнал условию `notEmpty`.

```
@ThreadSafe
public class ConditionBoundedBuffer<T> {
    protected final Lock lock = new ReentrantLock();
    // CONDITION PREDICATE: notFull (count < items.length)
    private final Condition notFull = lock.newCondition();
    // CONDITION PREDICATE: notEmpty (count > 0)
    private final Condition notEmpty = lock.newCondition();
    @GuardedBy("lock")
    private final T[] items = (T[]) new Object[BUFFER_SIZE];
    @GuardedBy("lock") private int tail, head, count;

    // BLOCKS-UNTIL: notFull
```

¹⁵⁶ Подразумевается конкретная реализация интерфейса `Condition`.

```

public void put(T x) throws InterruptedException {
    lock.lock();
    try {
        while (count == items.length)
            notFull.await();
        items[tail] = x;
        if (++tail == items.length)
            tail = 0;
        ++count;
        notEmpty.signal();
    } finally {
        lock.unlock();
    }
}

// BLOCKS-UNTIL: notEmpty
public T take() throws InterruptedException {
    lock.lock();
    try {
        while (count == 0)
            notEmpty.await();
        T x = items[head];
        items[head] = null;
        if (++head == items.length)
            head = 0;
        --count;
        notFull.signal();
        return x;
    } finally {
        lock.unlock();
    }
}
}

```

Листинг 14.11 Ограниченный буфер, использующий явные условные переменные

Поведение класса `ConditionBoundedBuffer` аналогично поведению класса `BoundedBuffer`, но использование им очередей условий удобнее для чтения - легче анализировать класс, использующий множество экземпляров `Condition`, чем тот, что использует единственную внутреннюю очередь условий с несколькими предикатами условий. Разделив два предиката условия на отдельные наборы ожидания, интерфейс `Condition` упрощает выполнение требований для одиночного уведомления. Использование более эффективного метода `signal`, вместо метода `signalAll`, уменьшает количество переключений контекста и количество захватов блокировок, вызванных каждой операцией буфера.

Как и во внутренних блокировках и в очередях условий, трехсторонняя связь между блокировкой, предикатом условия и условной переменной также должна выполняться при использовании явных реализаций `Lock` и `Condition`. Переменные, участвующие в предикате условия, должны быть защищены

экземпляром блокировки `Lock`, и блокировка `Lock` должна удерживаться при тестировании предиката условия и при вызове методов `await` и `signal`.¹⁵⁷

Выбирайте между использованием явных экземпляров `Condition` и внутренних очередей условий так же, как если бы вы выбирали между классом `ReentrantLock` и блоком `synchronized`: используйте экземпляр `Condition`, если вам нужны его расширенные функции, такие как справедливая очередь или несколько наборов ожидания на каждую блокировку, в противном случае, отдайте предпочтение использованию внутренних очередей условий. (Если вы уже используете класс `ReentrantLock`, потому что вам нужны его расширенные функции, выбор уже сделан.)

14.4 Анатомия синхронизатора

Интерфейсы классов `ReentrantLock` и `Semaphore` имеют много общего. Оба класса выступают в качестве “затворов”, позволяя одновременно только ограниченному числу потоков продолжать выполнение; потоки прибывают к затвору и допускаются через него (методы `lock` или `acquire` успешно возвращают управление), ожидают (методы `lock` или `acquire` блокируются), или возвращаются назад (методы `tryLock` или `tryAcquire` возвращают `false`, указывая, что блокировка или разрешение не были получены за отведённое на ожидание время). Далее, оба класса позволяют выполнять прерываемые, непрерываемые и ограниченные по времени попытки захвата, и оба позволяют выбрать справедливое или несправедливое помещение в очередь ожидания потоков.

Учитывая эту общность, вы можете подумать, что класс `Semaphore` был реализован поверх класса `ReentrantLock`, или, возможно, класс `ReentrantLock` был реализован как класс `Semaphore` с одним разрешением. Это было бы вполне практически; это распространённое упражнение, заключающееся в том, чтобы доказать, что подсчитывающий семафор может быть реализован с помощью блокировки (как в случае класса `SemaphoreOnLock`, приведённого в листинге 14.12) и что блокировка может быть реализована с помощью подсчитывающего семафора.

```
// Not really how java.util.concurrent.Semaphore is implemented
@ThreadSafe
public class SemaphoreOnLock {
    private final Lock lock = new ReentrantLock();
    // CONDITION PREDICATE: permitsAvailable (permits > 0)
    private final Condition permitsAvailable = lock.newCondition();
    @GuardedBy("lock") private int permits;

    SemaphoreOnLock(int initialPermits) {
        lock.lock();
        try {
            permits = initialPermits;
        } finally {
            lock.unlock();
        }
    }
}
```

¹⁵⁷ Класс `ReentrantLock` требует, чтобы при вызове методов `signal` или `signalAll` удерживалась блокировка, но реализациям `Lock` разрешается создавать экземпляры `Condition`, которые не имеют этого требования.

```

// BLOCKS-UNTIL: permitsAvailable
public void acquire() throws InterruptedException {
    lock.lock();
    try {
        while (permits <= 0)
            permitsAvailable.await();
        --permits;
    } finally {
        lock.unlock();
    }
}

public void release() {
    lock.lock();
    try {
        ++permits;
        permitsAvailable.signal();
    } finally {
        lock.unlock();
    }
}
}

```

Листинг 14.12 Подсчитывающий семафор, реализованный с помощью интерфейса Lock

На самом деле, они оба реализованы с использованием общего базового класса `AbstractQueuedSynchronizer` (AQS) - как, впрочем, и многие другие синхронизаторы. Класс AQS представляет собой фреймворк для построения блокировок и синхронизаторов, и с его использованием может быть легко и эффективно построен на удивление обширный ряд синхронизаторов. Не только классы `ReentrantLock` и `Semaphore` построены с использованием класса AQS, но также и `CountDownLatch`, `ReentrantReadWriteLock`, `SynchronousQueue`¹⁵⁸, и `FutureTask`.

Класс AQS обрабатывает множество деталей реализации синхронизатора, например, таких как очередь FIFO ожидающих потоков. Конкретные синхронизаторы могут определять гибкие критерии того, должно ли потоку быть позволено пройти или от него требуется ожидать.

Использование класса AQS для создания синхронизаторов имеет ряд преимуществ. Это не только приводит к значительному сокращению усилий, необходимых для реализации, но и не приводит к необходимости расплачиваться за возникновение нескольких точек конкуренции, в отличие от случая построения одного синхронизатора поверх другого. В классе `SemaphoreOnLock` получение разрешения может быть заблокировано в двух местах - один раз на блокировке, защищающей состояние семафора, а затем вновь, в случае если разрешение недоступно. Синхронизаторы, построенные с помощью класса AQS, имеют только одну точку, в которой они могут быть заблокированы, приводя, таким образом, к снижению издержек на переключение контекста и улучшению пропускной способности. Класс AQS был разработан для обеспечения масштабируемости, и

¹⁵⁸ В Java 6 класс `SynchronousQueue` на основе AQS был заменён на (более масштабируемую) неблокирующую версию.

все синхронизаторы в пакете `java.util.concurrent`, что были построены с его использованием, получают преимущества от этого.

14.5 Класс `AbstractQueuedSynchronizer`

Большинство разработчиков, вероятно, никогда не будут использовать класс AQS напрямую; стандартный набор синхронизаторов охватывает довольно широкий диапазон ситуаций. Но, видение того, как реализованы стандартные синхронизаторы, может помочь в прояснении того, как они работают.

Основные операции, которые выполняет синхронизатор, основанный на классе AQS, это несколько вариантов захвата и освобождения. Операция захвата зависит от состояния и всегда может быть заблокирована. В случае блокировки или семафора смысл операции захвата прост - захват блокировки или разрешение - и вызывающему объекту, возможно, придется ожидать, пока синхронизатор находится в том состоянии, в котором это может произойти. С классом `CountDownLatch`, захват означает “ожидать, пока защелка достигнет своего терминального состояния”, а в случае класса `FutureTask`, это означает “ожидать, пока задача не будет завершена”. Освобождение является не блокирующей операцией; освобождение может позволить потокам, заблокированным в процессе захвата, продолжить своё выполнение.

Чтобы класс зависел от состояния, он должен иметь некоторое состояние. Класс AQS берет на себя задачу управления некоторыми состояниями для класса синхронизатора: он управляет информацией о состоянии как единым целым, которым можно управлять с помощью защищенных методов `getState`, `setState` и `compareAndSetState`. Этот механизм можно использовать для представления произвольного состояния; например, класс `ReentrantLock` использует его для отображения того, сколько раз поток-владелец захватывал блокировку, класс `Semaphore` использует его для представления количества оставшихся разрешений, а класс `FutureTask` использует его для отображения состояния задачи (не начата, выполняется, завершена, отменена). Синхронизаторы также могут сами управлять дополнительными переменными состояния; например, класс `ReentrantLock` отслеживает текущего владельца блокировки, чтобы различать реентерабельные и конкурирующие запросы на захват блокировки.

Захват и освобождение блокировки принимает в классе AQS форму, приведённую в листинге 14.13. В зависимости от синхронизатора, захват может быть *исключительным (exclusive)*, как в случае класса `ReentrantLock`, или *неисключительным (nonexclusive)*, как в случае классов `Semaphore` и `CountDownLatch`. Операция захвата состоит из двух частей. Во-первых, синхронизатор решает, позволяет ли текущее состояние выполнить захват; если это так, потоку позволяет продолжить выполнение, а если нет, операция захвата блокируется или происходит сбой. Это решение определяется семантикой синхронизатора; например, захват блокировки может быть успешен, если блокировка не удерживается, и захват защелки может быть успешен, если защелка находится в терминальном состоянии.

```
boolean acquire() throws InterruptedException {
    while (state does not permit acquire) {
        if (blocking acquisition requested) {
```

```

        enqueue current thread if not already queued
        block current thread
    }
    else
        return failure
}
possibly update synchronization state
dequeue thread if it was queued
return success
}

void release() {
    update synchronization state
    if (new state may permit a blocked thread to acquire)
        unblock one or more queued threads
}

```

Листинг 14.13 Каноническая форма захвата и освобождения в классе AQS

Вторая часть включает в себя, возможно, обновление состояния синхронизатора; один поток, захвативший синхронизатор, может оказывать влияние на то, смогут ли другие потоки также захватить его. Например, захват блокировки изменяет состояние блокировки с “не удерживается” на “удерживается”, а получение разрешения от семафора уменьшает количество оставшихся разрешений. С другой стороны, захват защёлки одним потоком не оказывает влияние на возможность её захвата другими потоками, поэтому захват защёлки не изменяет её состояние.

Синхронизатор, поддерживающий эксклюзивный захват, должен реализовывать защищенные методы `tryAcquire`, `tryRelease`, и `isHeldExclusively`, а тот, что поддерживает совместный захват, должен реализовывать методы `tryAcquireShared` и `tryReleaseShared`. Методы `acquire`, `acquireShared`, `release` и `releaseShared` класса AQS вызывают формы `try` этих методов в подклассе синхронизатора, чтобы определить, может ли операция продолжиться. Подкласс синхронизатора может использовать методы `getState`, `setState` и `compareAndSetState` для проверки и обновления состояния, в соответствии с семантикой захвата и освобождения и сообщить базовому классу через возврат статуса, была ли успешна попытка захвата или освобождения синхронизатора. Например, возврат отрицательного значения из класса `tryAcquireShared` свидетельствует о крахе захвата; возврат нулевое значения указывает на то, что синхронизатор был захвачен эксклюзивно; и возврат положительного значения указывает на то, что синхронизатор был захвачен не эксклюзивно. Методы `tryRelease` и `tryReleaseShared` должны возвращать значение `true`, если при освобождении, могут быть разблокированы потоки, пытающиеся захватить синхронизатор.

Для упрощения реализации блокировок, поддерживающих очереди условий (например, класс `ReentrantLock`), класс AQS также предоставляет механизмы для построения условных переменных, связанных с синхронизаторами.

14.5.1 Простая защёлка

Класс `OneShotLatch` из листинга 14.14, представляет собой бинарную защёлку, реализованную с использованием класса `AQS`. Он имеет два публичных метода, `await` и `signal`, которые соответствуют захвату и освобождению. Изначально, защёлка закрыта; любой поток, вызвавший метод `await`, блокируется до тех пор, пока защёлка не будет открыта. Как только защёлка открывается с помощью вызова метода `signal`, ожидающие потоки освобождаются, и потоки, последовательно прибывшие к защёлке, продолжают своё выполнение.

```
@ThreadSafe
public class OneShotLatch {
    private final Sync sync = new Sync();

    public void signal() { sync.releaseShared(0); }

    public void await() throws InterruptedException {
        sync.acquireSharedInterruptibly(0);
    }

    private class Sync extends AbstractQueuedSynchronizer {
        protected int tryAcquireShared(int ignored) {
            // Succeed if latch is open (state == 1), else fail
            return (getState() == 1) ? 1 : -1;
        }

        protected boolean tryReleaseShared(int ignored) {
            setState(1); // Latch is now open
            return true; // Other threads may now be able to acquire
        }
    }
}
```

Листинг 14.14 Бинарная защёлка, реализованная с помощью класса `AbstractQueuedSynchronizer`

В классе `OneShotLatch`, состояние защёлки определяется состоянием класса `AQS` – закрыто (*ноль*) или открыто (*один*). Метод `await` вызывает метод `acquireSharedInterruptibly` класса `AQS`, который в свою очередь консультируется с методом `tryAcquireShared` класса `OneShotLatch`. Реализация метода `tryAcquireShared` должна возвращать значение, указывающее, может ли быть выполнен захват защёлки или нет. Если защёлка уже была ранее открыта, метод `tryAcquireShared` возвращает значение “успешно”, позволяя потоку продолжить выполнение; в противном случае, он возвращает значение, указывающее на то, что попытка захвата защёлки была неудачной. Метод `acquireSharedInterruptibly` интерпретирует неудачу при захвате защёлки в качестве признака того, что поток должен быть помещен в очередь ожидающих потоков. Аналогично, метод `signal` вызывает метод `releaseShared`, который, в свою очередь, консультируется с методом `tryReleaseShared`. Реализация метода `tryReleaseShared` безоговорочно устанавливает состояние защёлки в “открыто” и указывает (через своё возвращаемое значение), что синхронизатор находится в полностью освобождённом состоянии. Это приводит к тому, что класс `AQS`

позволяет всем ожидающим потокам попытаться повторно захватить синхронизатор, и захват защёлки теперь будет выполнен успешно, потому что метод `tryAcquireShared` возвращает значение “успешно”.

Класс `OneShotLatch` представляет собой полнофункциональный, удобный и производительный синхронизатор, реализованный всего в двадцати или около того строках кода. Конечно, ему не хватает некоторых полезных возможностей, - таких как ограниченного по времени захвата или возможности проверки состояния защёлки - но их также легко реализовать, поскольку класс `AQS` предоставляет ограниченные по времени версии методов захвата и методы-утилиты для распространённых операций проверок.

Класс `OneShotLatch` мог бы быть реализован путем расширения класса `AQS`, а не делегирования ему, но это нежелательно по некоторым причинам [ЕJ пункт 14]. Расширение разрушит простотой (состоящий из двух методов) интерфейс класса `OneShotLatch`, и, несмотря на то, что открытые методы `AQS` не позволят вызывающим объектам нарушить состояние защёлки, вызывающие объекты могут легко использовать их неправильно. Ни один из синхронизаторов пакета `java.util.concurrent` не расширяет класс `AQS` напрямую - все они делегируют выполнение приватным внутренним подклассам `AQS`.

14.6 *AQS* в классах-синхронизаторах из пакета `java.util.concurrent`

Множество блокирующих классов из пакета `java.util.concurrent`, таких как `ReentrantLock`, `Semaphore`, `ReentrantReadWriteLock`, `CountDownLatch`, `SynchronousQueue` и `FutureTask`, построено с использованием класса `AQS`. Не вдаваясь в детали слишком глубоко (исходный код является частью поставки JDK¹⁵⁹), давайте кратко рассмотрим, как каждый из этих классов использует `AQS`.

14.6.1 Класс `ReentrantLock`

Класс `ReentrantLock` поддерживает только эксклюзивный захват блокировки, поэтому он реализует методы `tryAcquire`, `tryRelease`, и `isHeldExclusively`; реализация метода `tryAcquire` для несправедливой версии, приведена в листинге 14.15. Класс `ReentrantLock` использует состояние синхронизации для хранения счетчика захватов блокировки и поддерживает переменную `owner`, предназначенную для хранения идентификатора потока-владельца, который изменяется только тогда, когда текущий поток только что захватил блокировку или только собирается ее освободить.¹⁶⁰ В методе `tryRelease` проверяется поле `owner`, чтобы гарантировать, что текущий поток завладеет блокировкой прежде, чем позволит выполниться методу `unlock`; в методе `tryAcquire` это поле используется для различия между реентерабельным и конкурентным захватом.

¹⁵⁹ Или с меньшим количеством лицензионных ограничений <http://gee.cs.oswego.edu/dl/concurrency-interest>.

¹⁶⁰ Поскольку защищенные (protected) методы манипулирования состоянием, и при чтении, и при записи, имеют семантику памяти `volatile`, и класс `ReentrantLock` аккуратно считывает значение поля `owner` только после вызова метода `getState`, и записывает в него значение только перед вызовом метода `setState`, класс `ReentrantLock` может переключаться на семантику памяти синхронизированного состояния, и таким образом избежать дальнейшей синхронизации - см. раздел 16.1.4.

```
protected boolean tryAcquire(int ignored) {
    final Thread current = Thread.currentThread();
    int c = getState();
    if (c == 0) {
        if (compareAndSetState(0, 1)) {
            owner = current;
            return true;
        }
    } else if (current == owner) {
        setState(c+1);
        return true;
    }
    return false;
}
```

Листинг 14.15 Несправедливая реализация метода `tryAcquire` в классе `ReentrantLock`

Когда поток пытается захватить блокировку, метод `tryAcquire` сначала справляется о состоянии блокировки. Если она не удерживается, метод пытается обновить состояние блокировки, чтобы обозначить, что блокировка удерживается. В связи с тем, что состояние могло измениться из-за того, что сведения о нём были получены несколько инструкций назад, метод `tryAcquire` пытается использовать метод `compareAndSetState` для автоматического обновления состояния, чтобы указать, что блокировка удерживается и убедится, что оно не изменилось с момента прошлого наблюдения. (См. описание метода `compareAndSetState` в разделе 15.3.) Если состояние блокировки указывает на то, что она уже удерживается и текущий поток является владельцем блокировки - счетчик захватов увеличивается; если текущий поток не является владельцем блокировки, попытка захвата блокировки завершается неудачей.

Класс `ReentrantLock` также использует преимущества встроенной поддержки классом `AQS` множества условных переменных и наборов ожидания. Метод `Lock.newCondition` возвращает новый экземпляр `ConditionObject`, представляющий собой внутренний класс `AQS`.

14.6.2 Классы `Semaphore` и `CountDownLatch`

Класс `Semaphore` использует синхронизированное состояние класса `AQS` для хранения количества доступных разрешений. Метод `tryAcquireShared` (см. листинг 14.16) сначала вычисляет количество оставшихся разрешений, а если их недостаточно, возвращает значение, указывающее, что захват завершился неудачно. Если осталось достаточное количество разрешений, он пытается атомарно уменьшить количество разрешений с помощью метода `compareAndSetState`. Если операция выполнена успешно (это означает, что количество разрешений не изменилось с момента последнего просмотра), он возвращает значение, указывающее, что захват произведён успешно. Возвращаемое значение также содержит код успешности других попыток совместного захвата, в этом случае другие ожидающие потоки также будут разблокированы.

```
protected int tryAcquireShared(int acquires) {
    while (true) {
```

```

        int available = getState();
        int remaining = available - acquires;
        if (remaining < 0)
            || compareAndSetState(available, remaining))
        return remaining;
    }
}

protected boolean tryReleaseShared(int releases) {
    while (true) {
        int p = getState();
        if (compareAndSetState(p, p + releases))
            return true;
    }
}

```

Листинг 14.16 Методы `tryAcquireShared` и `tryReleaseShared` класса `Semaphore`.

Цикл `while` завершается тогда, когда нет достаточного количества разрешений или когда метод `tryAcquireShared` может атомарно обновить количество разрешений для соответствия количеству захватов. В то время как любой вызов метода `compareAndSetState` может завершиться неудачей из-за конкуренции с другими потоками (см. раздел 15.3), вынуждая его повторяться, один из двух критериев завершения станет истинным в пределах разумного количества повторных попыток. Аналогичным образом, метод `tryReleaseShared` увеличивает количество разрешений, потенциально разблокируя ожидающие потоки и повторяет попытки до успешного обновления. Возвращаемое методом `tryReleaseShared` значение указывает на то, могли ли другие потоки быть разблокированы освобождением.

Класс `CountDownLatch` использует класс `AQS` аналогично классу `Semaphore`: синхронизированное состояние содержит текущее количество. Метод `countDown` вызывает метод `release`, который приводит к уменьшению счетчика и разблокирует ожидающие потоки, если счетчик достигает нуля; метод `await` вызывает метод `acquire`, который немедленно возвращает управление, если счетчик достигает нуля, в противном случае блокируется.

14.6.3 Класс `FutureTask`

На первый взгляд класс `FutureTask` даже не выглядит как синхронизатор. Но метод `Future.get` имеет семантику, очень похожую на семантику защелки – если произошло какое-либо событие (завершение или отмена задачи, представленной экземпляром `FutureTask`), то потоки могут продолжать работу, иначе они помещаются в очередь до тех пор, пока это событие не произойдет.

Класс `FutureTask` использует синхронизированное состояние класса `AQS` для хранения состояния задачи – выполняется (*running*), завершено (*completed*) или отменено (*cancelled*). Он также поддерживает дополнительные переменные состояния для хранения результата вычисления или вызванного им исключения. Кроме того, он поддерживает ссылку на поток, в котором выполняется вычисление (если он в настоящее время находится в состоянии выполнения), чтобы его можно было прервать, если задача будет отменена.

14.6.4 Класс ReentrantReadWriteLock

Интерфейс `ReadWriteLock` предполагает наличие двух блокировок - блокировки на чтение и блокировки на запись, но в реализации класса `ReentrantReadWriteLock` на основе `AQS`, блокировкой на чтение и на запись управляет единственный подкласс `AQS`. Класс `ReentrantReadWriteLock` использует 16 бит состояния для подсчёта количества блокировок на запись, а остальные 16 бит для подсчёта количества блокировок на чтение. Операции блокировок на чтение совместно используют методы захвата и освобождения; операции блокировки на запись используют методы захвата и освобождения эксклюзивно.

Внутри себя класс `AQS` поддерживает очередь ожидающих потоков, отслеживая, запрашивал ли поток эксклюзивный или совместный доступ. В случае класса `ReentrantReadWriteLock`, когда блокировка становится доступной, если поток, находящийся в начале очереди, ищет доступ на запись, он получит его, и если поток, находящийся в начале очереди, ищет доступ на чтение, все ожидающие в очереди потоки будут получать его, вплоть до первого запроса записи.¹⁶¹

14.7 Итоги

Если вам нужно реализовать класс, зависимый от состояния, - тот, чьи методы должны блокироваться, если предусловие на основе состояния не выполняется - лучшая стратегия, как правило, заключается в реализации на основе существующего библиотечного класса, подобного `Semaphore`, `BlockingQueue` или `CountDownLatch`, как в случае класса `ValueLatch` (из раздела [8.5.1](#)). Однако иногда существующие библиотечные классы не обеспечивают достаточной базы; в этих случаях можно создать собственные синхронизаторы, используя внутренние очереди условий, явные объекты `Condition` или класс `AbstractQueuedSynchronizer`. Внутренние очереди условий тесно связаны с внутренней блокировкой, так как механизм управления зависимостью состояний обязательно связан с механизмом обеспечения согласованности состояний. Точно так же, явные экземпляры `Condition` тесно связаны с явными экземплярами `Lock` и предлагают расширенный набор функций по сравнению с внутренними очередями условий, включая несколько наборов ожидания для каждой блокировки, прерываемые или непрерывные условия ожидания, справедливое или несправедливое помещение в очередь и ожидание с ограничением по времени.

¹⁶¹ Этот механизм не позволяет выбирать политику предпочтения операциям чтения или записи, как это делают некоторые реализации блокировки на чтение-запись. Для этого или очередь ожидания класса `AQS` должна была бы быть отличной от очереди `FIFO`, или были бы необходимы две очереди. Однако, такая строгая политика упорядочивания редко нужна на практике; если несправедливая версия класса `ReentrantReadWriteLock` не предлагает приемлемой живучести, справедливая версия обычно предоставляет удовлетворительный порядок и гарантирует отсутствие голодания для читателей и писателей.

Chapter 15 Атомарные переменные и неблокирующая синхронизация

Многие из классов пакета `java.util.concurrent`, такие как `Semaphore` и `ConcurrentLinkedQueue`, обеспечивают лучшую производительность и масштабируемость по сравнению с альтернативами, использующими блок `synchronized`. В этой главе мы рассмотрим основной источник повышения производительности: атомарные переменные и неблокирующую синхронизацию.

Большая часть недавних исследований в области параллельных алгоритмов была сосредоточена на *неблокирующих алгоритмах*, которые используют для обеспечения целостности данных, при параллельном доступе, низкоуровневые атомарные машинные инструкции, подобные *сравнить-и-обменять* (*compare-and-swap*), а не блокировки. Неблокирующие алгоритмы широко используются в операционных системах и среде JVM для планирования потоков и процессов, сборки мусора, а также для реализации блокировок и других параллельных структур данных.

Неблокирующие алгоритмы значительно сложнее в разработке и реализации, чем основанные на блокировках альтернативы, но они могут предложить значительные преимущества в масштабируемости и живучести. Они координируются на более тонком уровне детализации и могут значительно сократить затраты на планирование, поскольку они не блокируются, когда несколько потоков конкурируют за одни и те же данные. Кроме того, они невосприимчивы к взаимоблокировкам и другим проблемам живучести. В алгоритмах, основанных на блокировках, другие потоки не могут прогрессировать, если поток переходит в спящий режим или вращается, пока удерживается блокировка, тогда как неблокирующие алгоритмы непроницаемы для сбоев отдельных потоков. Начиная с Java 5.0, можно построить эффективные неблокирующие алгоритмы в Java, используя *классы атомарных переменных*, такие как `AtomicInteger` и `AtomicReference`.

Атомарные переменные также могут быть использованы в качестве “лучших `volatile` переменных”, даже если вы не разрабатываете неблокирующие алгоритмы. Атомарные переменные предлагают ту же семантику памяти, что и переменные `volatile`, но с дополнительной поддержкой атомарных обновлений, что делает их идеально подходящими для счетчиков, генераторов последовательностей и сбора статистики, предлагая лучшую масштабируемость, чем альтернативы на основе блокировок.

15.1 Недостатки блокировки

Координация доступа к совместно используемому состоянию с помощью согласованного протокола блокировки гарантирует, что любой поток, удерживающий блокировку, защищающую набор переменных, имеет монопольный доступ к этим переменным и что любые изменения, внесенные в эти переменные, видны другим потокам, которые впоследствии захватят блокировку.

Современные среды JVM могут оптимизировать захват и освобождение незапланированных блокировок довольно эффективно, но если несколько потоков запрашивают блокировку одновременно, JVM обращается за помощью к операционной системе. Если он дойдет до этого момента, какой-то неудачливый

поток будет приостановлен, и его выполнение будет возобновлено позже¹⁶². При возобновлении этого потока, может потребоваться дождаться завершения квантов планирования другими потоками, прежде чем он фактически будет запланирован к выполнению. Приостановка и возобновление потока приводит к большим издержкам и обычно влечет за собой длительное прерывание. Для классов на основе блокировок с хорошо детализированными операциями (подобных синхронизированным классам коллекций, в которых большинство методов содержат только несколько операций), отношение издержек планирования к полезной работе может быть довольно высоким, когда блокировка часто оспаривается.

Переменные `volatile` - это более легкий механизм синхронизации, чем блокировка, поскольку они не приводят к переключению контекста или планированию потоков. Однако `volatile` переменные имеют некоторые ограничения по сравнению с блокировкой: хотя они предоставляют аналогичные гарантии видимости, их нельзя использовать для построения атомарных составных действий. Это означает, что переменные `volatile` нельзя использовать, когда одна переменная зависит от другой, или когда новое значение переменной зависит от ее старого значения. Это ограничивает спектр ситуаций, в которых допустимо применение переменных `volatile`, так как они не могут быть использованы для надежной реализации обычных инструментов, подобных счетчикам или мьютексам.¹⁶³

Например, в то время как операция инкремента (`++i`) может выглядеть как атомарная операция, на самом деле это три различные операции - получить текущее значение переменной, добавить к ней единицу, а затем записать обновленное значение обратно. Чтобы не потерять обновление, вся операция чтение-изменение-запись должна быть атомарной. До сих пор мы видели только один способ, которым можно это сделать – с помощью блокировки, как в классе `Counter` из раздела [4.1](#).

Класс `Counter` потокобезопасен, и при наличии небольшой конкуренции, или в её отсутствии, работает просто отлично. Но в условиях конфликта производительность страдает из-за накладных расходов на переключение контекста и задержек, вызванных планированием. Когда блокировки удерживаются на краткое время, приостановка выполнения потока является суровым наказанием за то, что блокировка была запрошена в неудачный момент времени.

Блокировка имеет ещё несколько недостатков. Когда поток ожидает блокировку, он не может делать ничего другого. Если поток, удерживающий блокировку, задерживается (из-за ошибки страницы, задержки, вызванной планированием или тому подобного), то ни один поток, которому нужна эта блокировка, не сможет прогрессировать. Это может быть серьезной проблемой, если заблокированный поток является высокоприоритетным потоком, но поток, удерживающий блокировку, является низкоприоритетным потоком - угроза производительности, известная как *инверсия приоритета* (*priority inversion*). Даже при том, что поток с более высоким приоритетом должен иметь приоритет, он вынужден ожидать, пока блокировка не будет освобождена, и это эффективно понижает его приоритет до потока с более низким приоритетом. Если поток, удерживающий блокировку, постоянно блокируется (из-за бесконечного цикла,

¹⁶² Умная среда JVM не обязательно должна приостанавливать поток, если он конкурирует за блокировку; она может использовать данные профилирования, чтобы принять адаптивное решение о приостановке или блокировке на основе прокручивания, в зависимости от того, как долго блокировка удерживалась во время предыдущих захватов.

¹⁶³ Теоретически возможно, хотя и совершенно непрактично, использовать семантику `volatile` для построения мьютексов и других синхронизаторов; см. (Raynal, 1986).

взаимоблокировки, динамической блокировки или другого сбоя живучести), любые потоки, ожидающие этой блокировки, никогда не смогут добиться прогресса.

Даже игнорируя эти угрозы, блокировка сама по себе является тяжеловесным механизмом для хорошо детализированных операций, подобных операции увеличения счётчика. Было бы неплохо иметь более детализированный подход для управления конкуренцией между потоками – что-то вроде *volatile* переменных, но предоставляющих возможность атомарных обновлений. К счастью, современные процессоры предлагают нам именно такой механизм.

15.2 Аппаратная поддержка параллелизма

Эксклюзивная блокировка является *пессимистичным* подходом - она предполагает худшее (если вы не запрете дверь, гремлины войдут и переставят ваши вещи) и не позволит продолжить, пока вы можете гарантировать, захватив соответствующую блокировку, что другие потоки не будут оказывать влияние.

Для хорошо детализированных операций существует альтернативный подход, который часто является более эффективным - *оптимистичный* подход, при котором вы продолжаете обновление, надеясь, что вы можете завершить его без постороннего вмешательства. Этот подход основан на *обнаружении коллизий* (*collision detection*) путём определения наличия помех со стороны других взаимодействующих сторон, во время выполнения обновления, в случае наличия таковых - операция завершается неудачно и может быть повторена (или нет). Оптимистичный подход подобен старой поговорке “легче получить прощение, чем разрешение”, где “легче” здесь означает “эффективнее”.

Процессоры, предназначенные для работы в многопроцессорных системах, предоставляют специальные инструкции для управления параллельным доступом к совместно используемым переменным. Процессоры ранних лет имели атомарные инструкции проверить-и-установить (*test-and-set*), получить-и-увеличить (*fetch-and-increment*) или обменять (*swap*), достаточные для реализации мьютексов, которые, в свою очередь, могли использоваться для реализации более сложных параллельных объектов. Сегодня почти каждый современный процессор имеет атомарную инструкцию прочитать-изменить-записать (*read-modify-write*), такую как *сравнить-и-обменять* (*compare-and-swap*) или загрузка-с-пометкой /попытка-записи (*load-linked/store-conditional*). Операционные системы и среда JVM используют эти инструкции для реализации блокировок и параллельных структур данных, но до Java 5.0 они не были непосредственно доступны классам Java.

15.2.1 Сравнить и обменять

Подход, используемый большинством процессорных архитектур, включая IA32 и Sparc, заключается в реализации инструкции *сравнить-и-обменять* (CAS). (Другие процессоры, такие как PowerPC, реализуют ту же функциональность с помощью пары инструкций: *загрузка-с-пометкой* и *попытка-записи*.) CAS имеет три операнда - ячейку памяти V , с которой следует работать, ожидаемое старое значение A и новое значение B . CAS атомарно обновляет V до нового значения B , но только если значение в V соответствует ожидаемому старому значению A ; в противном случае операция ничего не делает. В любом случае операция возвращает значение, находящееся в текущий момент в V . (Вариант, называемый *сравнить-и-обменять*, вместо этого возвращает результат, указывающий на то, была ли операция выполнена успешно.) CAS означает “Я думаю, что V должно иметь значение A ; если это так, поставьте там B , иначе не меняйте значение, но

скажите мне, что я ошибался⁷. Операция CAS является оптимистичным подходом - она продолжает обновление в надежде на успех и может обнаружить сбой, если другой поток обновил переменную с момента ее последнего получения. Класс SimulatedCAS приведённый в листинге 15.1 иллюстрирует семантику (но не реализацию или производительность) операции CAS.

```
@ThreadSafe
public class SimulatedCAS {
    @GuardedBy("this") private int value;

    public synchronized int get() { return value; }

    public synchronized int compareAndSwap(int expectedValue,
                                           int newValue) {
        int oldValue = value;
        if (oldValue == expectedValue)
            value = newValue;
        return oldValue;
    }

    public synchronized boolean compareAndSet(int expectedValue,
                                              int newValue) {
        return (expectedValue
                == compareAndSwap(expectedValue, newValue));
    }
}
```

Листинг 15.1 Симуляция операции CAS

Когда несколько потоков пытаются обновить одну и ту же переменную с помощью CAS одновременно, один выигрывает и обновляет значение переменной, а остальные проигрывают. Но проигравшие не наказываются отстранением, как это могло бы быть, если бы они не смогли захватить блокировку; вместо этого им говорят, что они не выиграли гонку на этот раз, но могут попробовать еще раз. Поскольку поток, который проигрывает в операции CAS, не блокируется, он может решить, хочет ли он повторить попытку, выполнить какое-либо другое действие по восстановлению или не делать ничего.¹⁶⁴ Эта гибкость позволяет исключить множество из угроз живучести, связанных с блокировками (хотя в нетипичных случаях может привести к введению риска возникновения динамической взаимоблокировки (*livelock*) - см. раздел 10.3.3).

Типичный шаблон использования CAS заключается в том, чтобы сначала прочитать значение из *V*, вывести новое значение *B* из *A*, а затем использовать операцию CAS для атомарного изменения значения *V* с *A* на *B*, до того момента, когда на значение *V* перестанут оказывать влияние другие потоки. Операция CAS решает проблему реализации атомарных последовательностей прочитать-изменить-записать без использования блокировки, потому что она может обнаруживать внесение помех другими потоками.

¹⁶⁴ Бездействие может быть вполне разумным ответом на провал операции CAS; в некоторых неблокирующих алгоритмах, подобных алгоритму связанной очереди из раздела 15.4.2, провал операции CAS означает, что кто-то уже выполнил работу, которую вы планировали выполнить.

15.2.2 Неблокирующий счётчик

Класс CasCounter из листинга 15.2, реализует потокобезопасный счетчик с помощью операции CAS. Операция инкремента следует канонической форме - извлечь старое значение, преобразовать его в новое значение (добавив единицу) и использовать операцию CAS для установки нового значения. Если операция CAS проваливается, операция немедленно повторяется. Осуществление повторной попытки, как правило, является разумной стратегией, хотя в случаях жёсткой конкуренции может быть желательным подождать некоторое время или откатится, прежде чем повторять попытку, с целью избежания возникновения динамической взаимоблокировки.

```
@ThreadSafe
public class CasCounter {
    private SimulatedCAS value;

    public int getValue() {
        return value.get();
    }

    public int increment() {
        int v;
        do {
            v = value.get();
        }
        while (v != value.compareAndSwap(v, v + 1));
        return v + 1;
    }
}
```

Листинг 15.2 Неблокирующий счётчик с использованием операции CAS

Класс CasCounter не блокируется, хотя, возможно, операцию придется повторить несколько¹⁶⁵ раз, если другие потоки будут обновлять значение счетчика в то же самое время. (На практике, если все, что вам нужно, это счетчик или генератор последовательности, просто используйте классы AtomicInteger или AtomicLong, которые предоставляют методы для атомарного увеличения и другие арифметические методы.)

На первый взгляд счетчик на основе CAS выглядит так, как будто он должен работать хуже, чем счетчик на основе блокировки; он состоит из большего количества операций и поток управления в нём более сложен, и он зависит от кажущейся сложной операции CAS. Но на самом деле счетчики на основе CAS значительно превосходят счетчики на основе блокировок, если существует даже небольшое количество конфликтов, и часто даже при отсутствии конфликтов. Для быстрого неконкурентного захвата блокировки, требуется, по крайней мере, одна операция CAS плюс другие связанные с блокировкой служебные действия, так что в лучшем случае для счетчика на основе блокировки выполняется больше работы, чем в нормальной ситуации для счетчика на основе операции CAS. Так как

¹⁶⁵ Теоретически, операция может повторяться произвольно большое количество раз, если другие потоки продолжают выигрывать гонку в операции CAS; на практике, голодание такого рода случается редко.

операция CAS преуспевает большую часть времени (при условии низкой и умеренной конкуренции), железо будет корректно предсказывать возникновение неявной ветви в цикле `while`, минимизируя издержки, введённые более сложной логикой управления.

Синтаксис языка для блокировки может быть компактным, но работа, выполняемая JVM и ОС для управления блокировками, таковой не является. Блокировка влечёт за собой пересечение относительно сложных веток кода в JVM и может повлечь за собой блокировку на уровне ОС, приостановку потока и переключение контекста. В лучшем случае блокировка требует, по крайней мере, выполнения одной операции CAS, поэтому использование блокировок выводит операции CAS из поля зрения, но не сохраняет фактическую стоимость выполнения. С другой стороны, выполнение операции CAS из программы не включает в себя код JVM, системные вызовы или планирование действий. То, что выглядит как более длинная ветка выполнения кода на уровне приложения, на самом деле является гораздо более короткой веткой кода, когда учитывается активность JVM и ОС. Основным недостатком операции CAS является то, что она вынуждает вызывающий объект сталкиваться с конкуренцией (повторение попытки, откат или отказ от выполнения операции), в то время как блокировки сталкиваются с конкуренцией за блокирование автоматически, пока блокировка доступна.¹⁶⁶

Производительность операции CAS широко варьируется между различными процессорами. В системе с одним процессором операция CAS обычно занимает порядка нескольких тактовых циклов, так как синхронизация между процессорами не требуется. На момент написания этой статьи, затраты на неконкурентную операцию CAS, в системе с несколькими ЦП, составляли примерно от десяти до 150 циклов; производительность CAS является быстро движущейся целью и варьируется не только между архитектурами, но даже между версиями одного и того же процессора. Конкурирующие силы, скорее всего, приведут к дальнейшему повышению эффективности CAS в течение следующих нескольких лет. Хорошее эмпирическое правило заключается в том, что затраты на “короткий путь” для неконкурентного захвата и освобождения блокировки на большинстве процессоров примерно в два раза превышает стоимость операции CAS.

15.2.3 Поддержка операций CAS в среде JVM

Итак, каким образом Java-код убеждает процессор выполнить операции CAS от его имени? До Java 5.0 не было способа сделать это короче, чем написать нативный (native) код. В Java 5.0 была добавлена низкоуровневая поддержка для предоставления операций CAS для типов `int`, `long` и объектных ссылок, и JVM компилирует их в наиболее эффективные средства, предоставляемые железом. На платформах, поддерживающих CAS, среда выполнения помещает их в соответствующие машинные инструкции; в худшем случае, если CAS-подобная инструкция недоступна, JVM использует спин-блокировку¹⁶⁷. Эта низкоуровневая поддержка JVM используется классами атомарных переменных (`AtomicXxx` в пакете `java.util.concurrent.atomic`) для обеспечения эффективной работы CAS с числовыми и ссылочными типами; эти классы атомарных переменных

¹⁶⁶ На самом деле, самым большим недостатком CAS является сложность корректного построения окружающих алгоритмов.

¹⁶⁷ Блокировка с повторами или блокировка с вращением.

используются, прямо или косвенно, для реализации большинства классов в пакете `java.util.concurrent`.

15.3 Классы атомарных переменных

Атомарные переменные являются более детализированными и облегченными, чем блокировки, и критически важны для реализации высокопроизводительного параллельного кода в многопроцессорных системах. Атомарные переменные ограничивают область конкуренции одной переменной; это настолько детализировано, насколько возможно получить (даже предполагая, что ваш алгоритм может быть реализован с использованием такой тонкой детализации). Быстрый (неконкурентный) путь для обновления атомарной переменной не медленнее, чем быстрый путь для захвата блокировки, и обычно быстрее; медленный путь определенно быстрее, чем медленный путь в случае использования блокировок, поскольку он не включает приостановку и перепланирование потоков. С алгоритмами, основанными на атомарных переменных вместо блокировок, потоки с большей вероятностью смогут продолжить работу без задержек и затратить меньшее время на восстановление, если столкнутся с конкуренцией.

Классы атомарных переменных предоставляют обобщение `volatile` переменных для поддержки атомарных условных операций прочитать-изменить-записать. Класс `AtomicInteger` представляет значение `int` и предоставляет методы `get` и `set` с той же семантикой памяти, как при чтении и записи `volatile int`. Он также предоставляет атомарный метод `compareAndSet` (который в случае успешного выполнения обладает теми же эффектами памяти, как при чтении и записи `volatile` переменной) и, для удобства, атомарные методы добавления (`add`), увеличения (`increment`) и уменьшения (`decrement`). Класс `AtomicInteger` имеет поверхностное сходство с расширенным классом `Counter`, но предлагает значительно лучшую масштабируемость при конкуренции, поскольку он может напрямую использовать аппаратную поддержку параллелизма железом.

Существует двенадцать классов атомарных переменных, разделенных на четыре группы: скаляры (*scalars*), средства обновления полей (*field updaters*), массивы и составные переменные. Наиболее часто используемыми атомарными переменными являются скаляры: `AtomicInteger`, `AtomicLong`, `AtomicBoolean`, и `AtomicReference`. Все поддерживают операции CAS; `Integer` и `Long` версии также поддерживают арифметику. (Для имитации атомарных переменных других примитивных типов можно привести `short` или `byte` значения к типу `int` и из него, и использовать методы `floatToIntBits` или `doubleToLongBits` для чисел с плавающей запятой.)

Классы атомарных массивов (доступные в `Integer`, `Long` и `Reference` версиях) - это массивы, элементы которых можно обновлять атомарно. Классы атомарных массивов предоставляют `volatile` семантику доступа к элементам массива, функцию, недоступную для обычных массивов - `volatile` массив обладает `volatile` семантикой только для ссылок на массив, а не для содержащихся в нём элементов. (Другие типы атомарных переменных обсуждаются в разделах [15.4.3](#) и [15.4.4](#).)

В то время как атомарные скалярные классы расширяют класс `Number`, они не расширяют примитивные классы-обёртки, такие как `Integer` или `Long`. Фактически, они не могут: примитивные классы-обёртки неизменяемы, тогда как классы атомарных переменных изменяемы. Классы атомарных переменных также

не переопределяют методы `hashCode` или `equals`; каждый экземпляр индивидуален. Как и большинство изменяемых объектов, они не являются хорошими кандидатами на роль ключей в коллекциях, основанных на хэшах.

15.3.1 Atomics as “better volatiles”

В разделе [3.4.2](#) мы использовали `volatile` ссылку на неизменяемый объект для атомарного обновления нескольких переменных состояния. Этот пример опирался на операцию проверить-затем-выполнить, но в этом конкретном случае гонка данных была безвредной, потому что нам было все равно, если мы иногда теряли обновление. В большинстве других ситуаций такая проверка не была бы безвредной и могла бы поставить под угрозу целостность данных. Например, класс `NumberRange` из раздела [4.3.3](#), не может быть безопасно реализован с помощью `volatile` ссылки на неизменяемый объект холдера для верхней и нижней границ, а также с использованием атомарных целых чисел для хранения границ. Поскольку инвариант ограничивает два числа, и они не могут быть обновлены одновременно, пока инвариант сохраняется, класс диапазона чисел, использующий `volatile` ссылки или несколько атомарных целых чисел, столкнётся с небезопасной последовательностью операций проверить-затем-выполнить.

Мы можем скомбинировать подход, принятый в классе `OneValueCache` с атомарными ссылками, чтобы исключить возможность возникновения условий гонки, *атомарно* обновляя ссылку на неизменяемый объект, содержащий значения нижней и верхней границ. Класс `CasNumberRange` из листинга 15.3 использует класс `AtomicReference` для хранения состояния экземпляра `IntPair`; с помощью метода `compareAndSet` он может обновлять верхнюю или нижнюю границы без возникновения условий гонки, в отличие от класса `NumberRange`.

```
public class CasNumberRange {  
    @Immutable  
    private static class IntPair {  
        final int lower; // Invariant: lower <= upper  
        final int upper;  
        ...  
    }  
    private final AtomicReference<IntPair> values =  
        new AtomicReference<IntPair>(new IntPair(0, 0));  
  
    public int getLower() { return values.get().lower; }  
    public int getUpper() { return values.get().upper; }  
  
    public void setLower(int i) {  
        while (true) {  
            IntPair oldv = values.get();  
            if (i > oldv.upper)  
                throw new IllegalArgumentException(  
                    "Can't set lower to " + i + " > upper");  
            IntPair newv = new IntPair(i, oldv.upper);  
            if (values.compareAndSet(oldv, newv))  
                return;  
        }  
    }  
}
```

```
// similarly for setUpper  
}
```

Листинг 15.3 Сохранение инвариантов из множества переменных с использованием CAS

15.3.2 Производительность: блокировки против атомарных переменных

Чтобы продемонстрировать различия в масштабируемости между блокировками и атомарными переменными, мы построили бенчмарк, сравнивающий несколько реализаций генератора псевдослучайных чисел (PRNG). В PRNG следующее "случайное" число является детерминированной функцией предыдущего числа, поэтому PRNG должен запоминать предыдущее число как часть своего состояния.

В листингах 15.4 и 15.5 приведены две реализации потокобезопасного PRNG, одна из которых использует класс `ReentrantLock`, а другая использует класс `AtomicInteger`.

```
@ThreadSafe  
public class ReentrantLockPseudoRandom extends PseudoRandom {  
    private final Lock lock = new ReentrantLock(false);  
    private int seed;  
  
    ReentrantLockPseudoRandom(int seed) {  
        this.seed = seed;  
    }  
  
    public int nextInt(int n) {  
        lock.lock();  
        try {  
            int s = seed;  
            seed = calculateNext(s);  
            int remainder = s % n;  
            return remainder > 0 ? remainder : remainder + n;  
        } finally {  
            lock.unlock();  
        }  
    }  
}
```

Листинг 15.4 Генератор случайных чисел с использованием класса `ReentrantLock`

Тестовый драйвер каждый раз вызывается повторно; каждая итерация генерирует случайное число (которое извлекает и изменяет совместно используемое состояние переменной `seed`), а также выполняет ряд итераций “занят выполнением работы”, которые работают строго с локальными для потока данными. Такое поведение позволяет имитировать типичные операции, которые включают в себя некоторую часть работы с совместно используемым состоянием и некоторую часть работы с состоянием, локальным для потока.

```
@ThreadSafe  
public class AtomicPseudoRandom extends PseudoRandom {  
    private AtomicInteger seed;
```

```

AtomicPseudoRandom(int seed) { this.seed
    = new AtomicInteger(seed);
}

public int nextInt(int n) {
    while (true) {
        int s = seed.get();
        int nextSeed = calculateNext(s);
        if (seed.compareAndSet(s, nextSeed)) {
            int remainder = s % n;
            return remainder > 0 ? remainder : remainder + n;
        }
    }
}

```

Листинг 15.5 Генератор случайных чисел с использованием класса AtomicInteger.

На рисунках 15.1 и 15.2 показана пропускная способность с низким и умеренным уровнями моделируемой работы в каждой итерации. При низком уровне локальных для потока вычислений, блокировка или атомарная переменная испытывают большую конкуренцию; при большем количестве локальных для потока вычислений, блокировка или атомарная переменная испытывают меньшую конкуренцию, так как каждый поток реже обращается к совместно используемым данным.

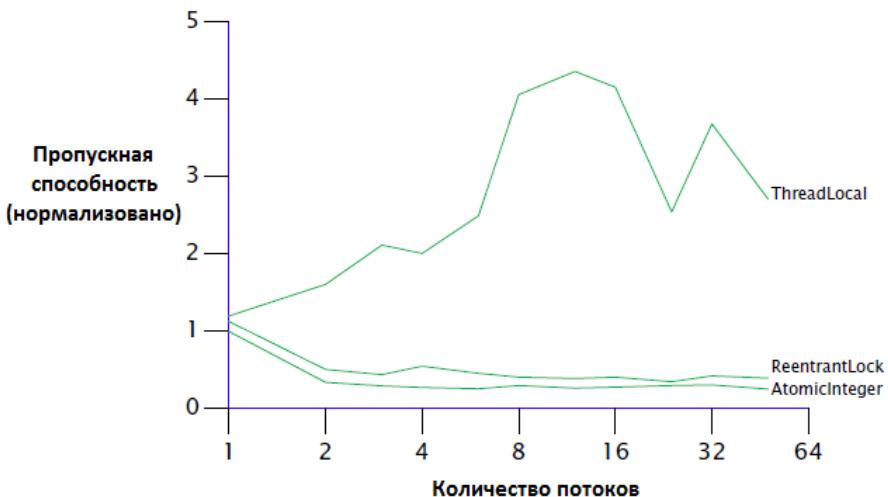


Рисунок 15.1 Производительность классов Lock и AtomicInteger при высокой конкуренции

Как показывают эти графики, при высоких уровнях конкуренции блокировка имеет тенденцию превосходить атомарные переменные, но при более реалистичных уровнях конкуренции, атомарные переменные превосходят блокировки.¹⁶⁸ Это вызвано тем, что блокировка реагирует на конкуренцию путем

¹⁶⁸ То же самое справедливо и в других областях: светофоры обеспечить лучшую производительность для высокого трафика, но которые обеспечивают лучшую пропускную способность при низкой интенсивности движения; конкурентная схема используемая в Ethernet сетях работает лучше при низкой

приостановки потоков, уменьшая, таким образом, использование ЦП и трафик синхронизации в совместно используемой шине памяти. (Это похоже на то, как блокировка производителя в дизайне производитель-потребитель снижает нагрузку на потребителей и тем самым позволяет им догнать производителя.) С другой стороны, в случае использования атомарных переменных управление конфликтами перекладывается на вызывающий класс. Подобно большинству алгоритмов на основе CAS, класс `AtomicPseudoRandom` реагирует на конкуренцию, немедленно повторяя попытку, что обычно является правильным подходом, но в среде с высокой конкуренцией просто приводит к большему количеству конфликтов.

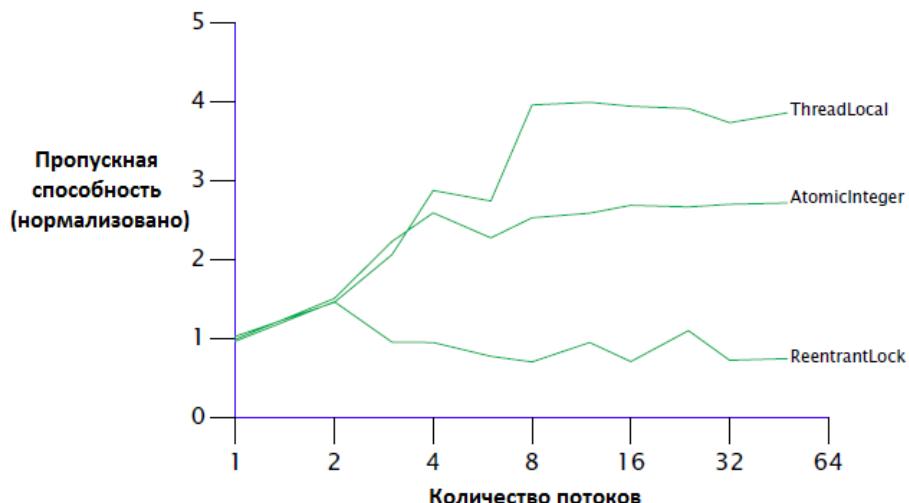


Рисунок 15.2 Производительность классов Lock и `AtomicInteger` при умеренной конкуренции

Прежде чем мы осудим класс `AtomicPseudoRandom` как плохо написанный или атомарные переменные как худший выбор по сравнению с блокировками, мы должны понять, что уровень конкуренции, приведённый на рисунке 15.1 нереально высок: никакая реальная программа не занимается только тем, что борется за блокировку или атомарную переменную. На практике атомики (*atomics*), как правило, масштабируется лучше, чем блокировки, потому что атомики более эффективно справляются с типичными уровнями конкуренции.

Изменение производительности между блокировками и атомиками при разных уровнях конкуренции, иллюстрирует сильные и слабые стороны каждого из них. При низкой и умеренной конкуренции, атомики обеспечивает лучшую масштабируемость; при высокой конкуренции, блокировки обеспечивают лучшее предотвращение конфликтов. (Алгоритмы на основе CAS, в системе с одним процессором, также превосходят основанные на блокировке алгоритмы, так как CAS всегда преуспевает, за исключением маловероятного случая, что поток будет вытеснен в середине операции чтение-изменение-запись.)

Рисунки 15.1 и 15.2 включают в себя третью кривую; реализацию класса `PseudoRandom`, использующего класс `ThreadLocal` для вывода состояния PRNG. Такой подход к реализации изменяет поведение класса - каждый поток видит свою собственную частную последовательность псевдослучайных чисел вместо всех потоков, совместно использующих одну последовательность, - но иллюстрирует, что часто дешевле вообще не использовать состояние, если этого можно избежать.

интенсивности траффика, но схема с передачей маркера, используемая в кольцевых сетях с маркером лучше подходит при высокой интенсивности траффика.

Мы можем улучшить масштабируемость, более эффективно справляясь с конкуренцией, но истинная масштабируемость достигается только за счет полного устранения конкуренции.

15.4 Неблокирующие алгоритмы

Алгоритмы, основанные на блокировках, подвержены риску возникновения ряда сбоев живучести. Если поток, удерживающий блокировку, задерживается из-за блокировки ввода/вывода, ошибки страницы или по другой причине, возможна ситуация, при которой поток не будет прогрессировать. Алгоритм называется *неблокирующими* (*nonblocking*), если сбой или приостановка какого-либо потока не может привести к сбою или приостановке другого потока; алгоритм называется *свободным от блокировок* (*lock-free*), если при выполнении каждого шага *какой-либо* поток сможет прогрессировать. Алгоритмы, использующие операции CAS исключительно для координации между потоками, при правильном построении могут быть как неблокирующими, так и свободными от блокировок. Конкуренция в операциях CAS всегда завершается успешно, и если несколько потоков конкурируют за операцию CAS, один всегда выигрывает и, следовательно, прогрессирует. Неблокирующие алгоритмы также невосприимчивы к взаимоблокировке или инверсии приоритета (хотя они могут сталкиваться с голоданием или динамической взаимоблокировкой, потому что могут выполнять повторные попытки). До сих пор мы видели один неблокирующий алгоритм: класс `CasCounter`. Хорошие неблокирующие алгоритмы известны для многих распространенных структур данных, включая стеки, очереди, приоритетные очереди и хэш-таблицы - хотя проектирование новых алгоритмов представляет собой задачу, которую лучше оставить экспертам.

15.4.1 Неблокирующий стек

Неблокирующие алгоритмы значительно сложнее, чем их эквиваленты на основе блокировок. Ключом к созданию неблокирующих алгоритмов является выяснение того, как ограничить область атомарных изменений одной переменной при сохранении согласованности данных. В классах связанных коллекций, подобных очередям, иногда можно избежать выражения изменяемых состояний в виде изменений отдельных ссылок и использования класса `AtomicReference` для представления каждой ссылки, которая должна обновляться атомарно.

Стек представляют собой простейшую связанную структуру данных: каждый элемент ссылается только на один другой элемент, а на каждый элемент ссылается только одна объектная ссылка. В класс `ConcurrentStack` из листинга 15.6 демонстрируется, как построить стек с использованием атомарных ссылок. Стек представляет собой связанный список элементов `Node`, вершина которого хранится в переменной `top`, каждый элемент списка содержит значение и ссылку на следующий элемент. Метод `push` готовит новую ссылку на элемент, следующее поле которого ссылается на текущую вершину стека, и затем использует операцию CAS, чтобы попытаться установить его качестве вершины стека. Если тот же самый узел все еще находится на вершине стека, как и в тот момент, когда мы начали, операция CAS выполняется успешно; если верхний узел изменился (потому что другой поток добавил или удалил элементы, во время начала выполнения операции), операция CAS завершается сбоем и метод `push` обновляет информацию о новом элементе на основе текущего состояния стека и

повторяет операцию вновь. В любом случае, стек все еще находится в согласованном состоянии, после выполнения операции CAS.

```
@ThreadSafe
public class ConcurrentStack <E> {
    AtomicReference<Node<E>> top = new AtomicReference<Node<E>>();

    public void push(E item) {
        Node<E> newHead = new Node<E>(item);
        Node<E> oldHead;
        do {
            oldHead = top.get();
            newHead.next = oldHead;
        } while (!top.compareAndSet(oldHead, newHead));
    }

    public E pop() {
        Node<E> oldHead;
        Node<E> newHead;
        do {
            oldHead = top.get();
            if (oldHead == null)
                return null;
            newHead = oldHead.next;
        } while (!top.compareAndSet(oldHead, newHead));
        return oldHead.item;
    }

    private static class Node <E> {
        public final E item;
        public Node<E> next;

        public Node(E item) {
            this.item = item;
        }
    }
}
```

Листинг 15.6 Неблокирующий стек использующий алгоритм Трейбера (Treiber, 1986)

Классы CasCounter и ConcurrentStack иллюстрируют характеристики всех неблокирующих алгоритмов: некоторая работа выполняется спекулятивно и, возможно, её придется переделывать. В классе ConcurrentStack, при создании экземпляра Node, представляющего собой новый элемент, мы надеемся, что значение ссылки next, к моменту установки в стеке, будет по-прежнему корректным, но готовы повторить попытку в случае возникновения конкуренции.

Неблокирующие алгоритмы, подобные классу ConcurrentStack, выводят свою потокобезопасность из того факта, что, подобно блокировке, метод compareAndSet предоставляет гарантии как атомарности, так и видимости. Когда поток изменяет состояние стека, он делает это с помощью метода compareAndSet, который обладает теми же эффектами памяти при записи, что и volatile. Когда поток

проверяет стек, он делает это, вызывая метод `get` на том же экземпляре `AtomicReference`, который обладает теми же эффектами памяти при чтении, что и `volatile`. Таким образом, любые изменения, сделанные одним потоком, безопасно публикуются в любом другом потоке, который проверяет состояние списка. И список модифицируется методом `compareAndSet`, который атомарно обновляет ссылку на элемент `top` или завершает своё выполнение сбоем, если обнаруживает влияние другого потока.

15.4.2 Неблокирующий связанный список

Два неблокирующих алгоритма, которые мы видели до сих пор, счетчик и стек, иллюстрируют основную схему использования операций CAS для спекулятивного обновления значения, с повторной попыткой, если выполнить обновление не удается. Хитрость построения неблокирующих алгоритмов заключается в ограничении области атомарных изменений одной переменной. Со счетчиками это тривиально, и со стеком всё достаточно прямолинейно, но для более сложных структур данных, подобных очередям, хэш-таблицам или деревьям, всё может быть намного сложнее.

Связанная очередь сложнее, чем стек, поскольку она должна поддерживать быстрый доступ как к голове (*head*), так и к хвосту (*tail*). Для этого она сохраняет отдельные указатели на голову и хвост. Два указателя ссылаются на узел в хвосте: `next`, указывающий на текущий последний элемент и указатель хвоста. Для успешной вставки нового элемента, оба указателя должны обновляться атомарно. На первый взгляд, это невозможно сделать с помощью атомарных переменных; для обновления двух указателей требуются отдельные операции CAS, и если первая завершается успешно, но второй это сделать не удаётся, очередь остается в несогласованном состоянии. И, даже если обе операции выполняются успешно, другой поток может попытаться получить доступ к очереди между первой и второй операциями. Построение неблокирующего алгоритма для связанный очереди требует конкретного плана для обеих ситуаций.

Для разработки этого плана нам необходимо несколько трюков. Во-первых, убедитесь, что структура данных всегда находится в согласованном состоянии, даже в процессе многоэтапного обновления. Таким образом, если поток **A** находится в процессе обновления, когда поток **B** выходит на сцену, потоку **B** могут сообщить, что операция была частично завершена, и он будет знать, что нельзя пытаться немедленно применить свое собственное обновление. Затем поток **B** может ожидать (многократно проверяя состояние очереди) до тех пор, пока поток **A** не завершит своё выполнение, чтобы они не мешали друг другу.

В то время как этого трюка самого по себе будет достаточно, чтобы позволить потокам получать доступ к структуре данных “по очереди” без её повреждения, если у одного потока произошёл сбой в середине выполнения обновления, ни один поток не сможет получить доступ к очереди. Чтобы сделать алгоритм неблокирующим, мы должны гарантировать, что сбой потока не помешает другим потокам добиться прогресса. Таким образом, второй трюк состоит в том, чтобы убедиться, что если поток **B** прибывает, чтобы получить доступ к структуре данных в середине операции обновления потока **A**, в структуре данных для потока **B** воплощено достаточно информации, чтобы завершить операцию обновления потока **A**. Если поток **B** “помогает” потоку **A**, завершая операцию **A**, поток **B** может продолжить свою собственную операцию, не дожидаясь потока **A**. Когда поток **A**

завершит свою работу, он обнаружит, что поток *B* уже выполнил эту работу за него.

Класс `LinkedQueue` из листинга 15.7, демонстрирует часть неблокирующего алгоритма Майкла-Скотта (Michael and Scott, 1996), касающуюся вставки в связанную очередь, который используется в классе `ConcurrentLinkedQueue`. Как и во многих других алгоритмах очередей, пустая очередь состоит из “дозорного” (*sentinel*) или “фиктивного” (*dummy*) узла, а указатели на голову и хвост инициализируются ссылками на дозорный узел. Указатель хвоста всегда ссылается на дозорный узел (если очередь пуста), последний элемент в очереди или (в случае, если операция находится в процессе обновления) предпоследний элемент.

```
@ThreadSafe
public class LinkedQueue <E> {
    private static class Node <E> {
        final E item;
        final AtomicReference<Node<E>> next;

        public Node(E item, Node<E> next) {
            this.item = item;
            this.next = new AtomicReference<Node<E>>(next);
        }
    }

    private final Node<E> dummy = new Node<E>(null, null);
    private final AtomicReference<Node<E>> head
        = new AtomicReference<Node<E>>(dummy);
    private final AtomicReference<Node<E>> tail
        = new AtomicReference<Node<E>>(dummy);

    public boolean put(E item) {
        Node<E> newNode = new Node<E>(item, null);
        while (true) {
            Node<E> curTail = tail.get();
            Node<E> tailNext = curTail.next.get();
            if (curTail == tail.get()) {
                if (tailNext != null) { A
                    // Queue in intermediate state, advance tail
                    tail.compareAndSet(curTail, tailNext); B
                } else {
                    // In quiescent state, try inserting new node
                    if (curTail.next.compareAndSet(null, newNode)) { C
                        // Insertion succeeded, try advancing tail
                        tail.compareAndSet(curTail, newNode);
                        return true;
                    }
                }
            }
        }
    }
}
```

```
    }  
}
```

Листинг 15.7 Фрагмент неблокирующего алгоритма Michael-Scott, касающийся вставки в очередь

На рисунке 15.3 иллюстрируется очередь с двумя элементами в нормальном, или *спокойном* (*quiescent*) состоянии.

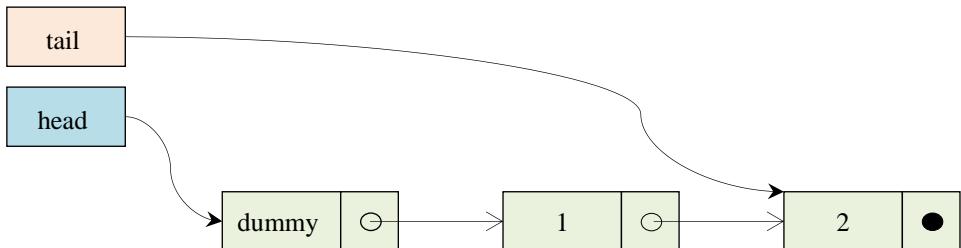


Рисунок 15.3 Очередь с двумя элементами в спокойном состоянии

Вставка нового элемента предполагает обновление двух указателей. Первый связывает новый узел с концом списка, обновляя указатель *next* на текущий последний элемент; второй перемещает указатель хвоста в обход, чтобы указать на новый последний элемент. Между этими двумя операциями, очередь находится в *промежуточном* (*intermediate*) состоянии, показанном на рис. 15.4.

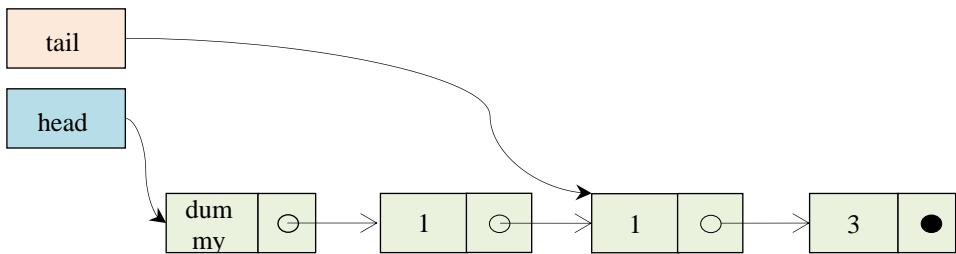


Рисунок 15.4 Очередь в промежуточном состоянии в процессе вставки

После второго обновления очередь снова находится в состоянии покоя, показанном на рисунке 15.5.

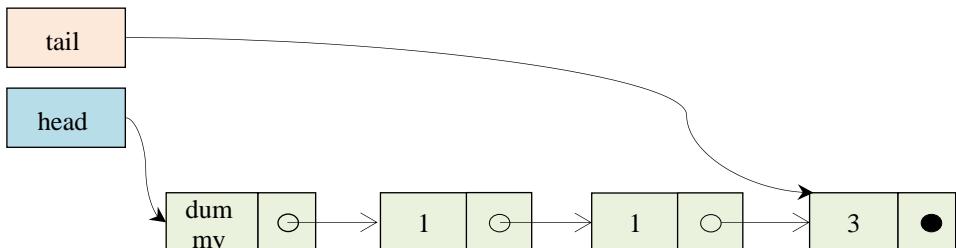


Рисунок 15.5 Очередь вновь в состоянии покоя, после завершения вставки

Ключевое наблюдение, которое включает оба из требуемых трюков, состоит в том, что, если очередь находится в состоянии покоя, связующее поле элемента *next*, на который указывает хвост, содержит `null`, и если это промежуточное состояние, значение поля `tail.next` не `null`. Таким образом, любой поток может сразу определить состояние очереди, проверяя значение поля `tail.next`. Кроме того, если очередь находится в промежуточном состоянии, ее можно восстановить

в состояние покоя, переместив указатель хвоста вперед на один узел, таким образом, завершив операцию для любого потока, находящегося в процессе вставки элемента.¹⁶⁹

Метод `LinkedQueue.put` сначала проверяет, находится ли очередь в промежуточном состоянии, прежде чем вставлять новый элемент (шаг **A**). Если это так, то какой-то другой поток уже находится в процессе вставки элемента (между шагами **C** и **D**). Вместо того чтобы ожидать завершения этого потока, текущий поток помогает ему, завершая операцию для него, продвигая указатель хвоста (шаг **B**). Затем он повторяет эту проверку в том случае, если другой поток начал вставлять новый элемент, продвигая указатель хвоста, пока он не обнаружит очередь в состоянии покоя, чтобы он мог начать свою собственную вставку.

Операция CAS на шаге **C**, в котором новый узел связывается хвостом очереди, может завершиться сбоем, если в одно и то же время два потока пытаются вставить элемент. В этом случае ничего страшного не происходит: никаких изменений не было сделано, и текущий поток может просто повторно загрузить указатель на хвост и повторить попытку. После успешного завершения шага **C**, вставка считается вступившей в силу; второй операцией CAS (стадия **D**) является “очистка”, поскольку она может быть выполнена либо вставляющим элемент потоком, либо любым другим потоком. Если на шаге **D** происходит сбой, вставляющий поток так или иначе возвращается, вместо того, чтобы повторно выполнить операцию CAS, потому что в повторной попытке нет необходимости – другой поток уже выполнил задание на шаге **B**! Это работает потому, что прежде чем какой-либо поток попытается связать новый узел с очередью, он сначала проверяет, нуждается ли очередь в очистке, проверяя, что значение поля `tail.next` не `null`. Если это так, он сначала перемещает указатель хвоста (возможно, несколько раз), до тех пор, пока очередь не будет находиться в состоянии покоя.

15.4.3 Обновления атомарного поля

В листинге 15.7 показан алгоритм, используемый классом `ConcurrentLinkedQueue`, но фактическая реализация немного отличается. Вместо представления каждого экземпляра `Node` с помощью атомарной ссылки, класс `ConcurrentLinkedQueue` использует `volatile` ссылку и обновляет ее с помощью основанного на `reflection` класса `AtomicReferenceFieldUpdater`, как показано в листинге 15.8.

```
private class Node<E> {
    private final E item;
    private volatile Node<E> next;

    public Node(E item) {
        this.item = item;
    }
}

private static AtomicReferenceFieldUpdater<Node, Node> nextUpdater
```

¹⁶⁹ Полный отчет о корректности этого алгоритма см. (Michael and Scott, 1996) или (Herlihy and Shavit, 2006).

```
= AtomicReferenceFieldUpdater.newUpdater(  
    Node.class, Node.class, "next");
```

Листинг 15.8 Использование обновлений атомарных полей в классе ConcurrentLinkedQueue

Классы-апдейтеры (*updater classes*) атомарных полей (доступные в `Integer`, `Long` и `Reference` версиях) представляют основанное на reflection “представление” существующего поля `volatile` так, чтобы операция CAS могла использоваться на существующих `volatile` полях. Классы-апдейтеры не имеют конструкторов; для их создания необходимо вызвать фабричный метод `newUpdater`, указав класс и имя поля. Классы-апдейтеры полей не привязаны к определенному экземпляру; их можно использовать для обновления целевого поля в любом экземпляре целевого класса. Гарантии атомарности для классы-апдейтеров слабее, чем для обычных атомарных классов, потому что вы не можете гарантировать, что базовые поля не будут изменены непосредственно - метод `compareAndSet` и арифметические методы гарантируют атомарность только по отношению к другим потокам, использующим методы апдейтеров для обновления атомарных полей.

В классе `ConcurrentLinkedQueue`, для обновления поля `next` экземпляра `Node` используется метод `compareAndSet`, вызываемый из метода `nextUpdater`. Этот, в некоторой степени, окольный подход используется исключительно по соображениям производительности. Для часто выделяемых (*allocated*) короткоживущих объектов, подобных связываемым узлам в очереди, исключение создания экземпляра `AtomicReference` для каждого экземпляра `Node` приводит к достаточно значительному снижению затрат на операции вставки. Тем не менее, практически во всех ситуациях обычные атомарные переменные работают просто отлично - только в нескольких случаях будут требоваться апдейтеры атомарных полей. (Апдейтеры атомарных полей также полезны тогда, когда вы хотите выполнить атомарные обновления, сохраняя сериализованную форму существующего класса.)

15.4.4 Проблема АВА

Проблема АВА является аномалией, которая может возникнуть из-за наивного использования алгоритмов сравни-и-обменять, в которых узлы могут быть переработаны¹⁷⁰ (в основном в средах без сборки мусора). Операция CAS фактически спрашивает “Значение *V* все еще *A*?”, и продолжает обновление, если это так. В большинстве ситуаций, включая примеры, представленные в этой главе, этого вполне достаточно. Однако иногда мы хотим спросить: “Изменялось ли значение *V* с тех пор, как я наблюдал, что оно было *A*?”. Для некоторых алгоритмов изменение значения *V* с *A* на *B*, а затем обратно на *A* по-прежнему считается изменением, которое требует от нас повторения некоторых алгоритмических шагов.

Проблема АВА может возникнуть в алгоритмах, которые выполняют своё собственное управление памятью для объектов связываемых узлов. В этом случае недостаточно того, что голова списка все еще ссылается на ранее наблюдаемый узел, чтобы подразумевать, что содержимое списка не изменилось. Если вы не можете избежать проблемы АВА, разрешив сборщику мусора управлять ссылочными узлами, существует относительно простое решение: вместо обновления значения ссылки, обновите пару значений, ссылку и номер версии.

¹⁷⁰ Каким-то образом «очищены» и снова пущены в работу, вместо уничтожения старого экземпляра и создания нового.

Даже если значение изменяется с *A* на *B* и обратно на *A*, номера версий будут отличаться. Класс `AtomicStampedReference` (и его кузен `AtomicMarkableReference`) обеспечивают условное атомарное обновление пары переменных. Класс `AtomicStampedReference` обновляет объект пары “ссылка - целое число”, позволяя “версионированным” ссылкам стать нечувствительными¹⁷¹ к проблеме АВА. Аналогично, класс `AtomicMarkableReference` обновляет пару “ссылка на объект - булево значение”, которая используется некоторыми алгоритмами, чтобы позволить узлу оставаться в списке, когда он помечен как удаленный.¹⁷²

15.5 Итого

Неблокирующие алгоритмы поддерживают потокобезопасность, используя вместо блокировок низкоуровневые примитивы параллелизма, такие как *сравнить-и-поменять*. Эти низкоуровневые примитивы предоставляются через классы атомарных переменных, которые также могут использоваться как “лучшие *volatile* переменные”, обеспечивая атомарные операции обновления для целых чисел и ссылок на объекты.

Неблокирующие алгоритмы трудно разработать и реализовать, но они могут предложить лучшую масштабируемость в типичных условиях и большую устойчивость к сбоям живучести. Многие достижения в области параллельной производительности от одной версии JVM к другой, связаны с использованием неблокирующих алгоритмов как в JVM, так и в библиотеках платформы.

¹⁷¹ На практике, в любом случае; теоретически счетчик можно обернуть.

¹⁷² Многие процессоры предоставляют операцию CAS двойной ширины (CAS2 или CASX), которая может работать с парой “указатель - целое”, что делает выполнение этого операции достаточно эффективным. Начиная с Java 6, класс `AtomicStampedReference` не использует операцию CAS двойной ширины даже на платформах, которые его поддерживают. (Двойная операция CAS отличается от DCAS, который работает с двумя несвязанными расположениями памяти; на момент написания статьи ни один текущий процессор не реализует операцию DCAS.)

Глава 16 Модель памяти Java

На протяжении всей книги мы в основном избегали низкоуровневых деталей модели памяти Java (JMM) и вместо этого сосредоточились на вопросах проектирования более высокого уровня, таких как безопасная публикация, спецификация и соблюдение политик синхронизации. Они порождают свою безопасность от JMM, и вы можете упростить себе фактическое использование этих механизмов, когда поймёте, *почему* они так работают. Эта глава приподнимает занавес, чтобы выявить низкоуровневые требования и гарантии модели памяти Java и идеи, лежащие в основе некоторых правил проектирования более высокого уровня, предлагаемых в этой книге.

16.1 Что такое модель памяти и зачем она мне нужна?

Предположим, что один поток присваивает значение переменной `aVariable`:

```
aVariable = 3;
```

Модель памяти озадачивается вопросом, “при каких условиях поток, который читает переменную `aVariable`, увидит значение 3?”? Это может прозвучать как глупый вопрос, но в отсутствии синхронизации, существует ряд причин, по которым поток может не сразу - или вообще никогда - увидеть результаты операции в другом потоке. Компиляторы могут генерировать инструкции в порядке, отличном от “очевидного”, прописанного в исходном коде, или хранить переменные в регистрах, а не в оперативной памяти; процессоры могут выполнять инструкции параллельно или не по порядку; кэши могут различаться порядком, в котором данные записываются в переменные или фиксируются в основной памяти; а значения, хранящиеся в локальных кэшах процессора, могут быть не видны другим процессорам. Эти факторы могут стать препятствием для того, чтобы поток видел самое последнее значение переменной и может привести к тому, что операции с памятью в других потоках будут выглядеть беспорядочными - если вы не используете адекватную синхронизацию.

В однопоточной среде, все эти трюки, играющие в нашей программе роль окружения, скрыты от нас, предназначены для ускорения выполнения и не оказывают никакого иного влияния. Спецификация языка Java требует, чтобы среда JVM поддерживала *в потоке семантику последовательного выполнения* (*within-thread as-if-serial semantics*): до тех пор, пока программа получает тот же самый результат, как если бы она выполнялась в среде со строго последовательным порядком выполнения, все эти игры допустимы. И это тоже хорошо, потому что, в последние годы, эти перестановки влекут за собой значительное улучшение производительности при выполнении вычислений. Конечно, более высокие тактовые частоты способствовали повышению производительности, но также увеличился параллелизм - конвейерные суперскалярные вычислительные модули, динамическое планирование инструкций, спекулятивное выполнение и сложные многоуровневые кэши памяти. По мере того как процессоры становятся все более изощренными, так же есть и компиляторы, переупорядочивающие инструкции для облегчения оптимального выполнения и использующие сложные глобальные алгоритмы распределения регистров. И поскольку производители процессоров

переходят на многоядерные процессоры, во многом из-за того, что увеличение тактовых частот становится экономически всё дороже, аппаратный параллелизм будет только увеличиваться.

В многопоточной среде, иллюзия последовательность не может поддерживаться без существенных затрат производительности. Поскольку большинство потоков в параллельном приложении ограничено по времени, и каждый из них “занят своим делом”, чрезмерная координация между потоками только замедлит работу приложения, без получения реальной выгоды. Только в том случае, когда несколько потоков совместно используют одни и те же данные, необходимо координировать их действия, и JVM в этом отношении полагается на программу, чтобы та определяла, когда это будет происходить, с помощью синхронизации.

Спецификация JMM определяет минимальные гарантии, которым должна следовать среда JVM, когда записывает значения в переменные, которые становятся видимыми для других потоков. Она был разработана, чтобы сбалансировать требование предсказуемости и простоты разработки программ с реалиями реализации высокопроизводительных JVM на широком спектре популярных процессорных архитектур. Некоторые аспекты JMM сначала могут вызвать беспокойство, если вы не знакомы с приемами, используемыми современными процессорами и компиляторами, для “выжимания” дополнительной производительности из вашей программы.

16.1.1 Основа моделей памяти

В мультипроцессорной архитектуре с совместно используемой памятью, каждый процессор имеет собственный кэш, который периодически согласовывается с основной памятью. Процессорные архитектуры обеспечивают различную степень согласованности кэша (*cache coherence*); некоторые предоставляют минимальные гарантии, позволяющие разным процессорам видеть различные значения для одного и того же участка памяти, практически в любое время. Операционная система, компилятор и среда выполнения (а иногда и программа) должны компенсировать разницу между возможностями, предоставляемыми оборудованием, и тем, что требуется для обеспечения потокобезопасности.

Обеспечение возможности того, чтобы каждый процессор знал, чем занят другой процессор, всегда обходится дорого. Большую часть времени эта информация не нужна, поэтому процессоры ослабляют гарантии согласованности памяти, с целью повышения производительности. Архитектура модели памяти сообщает программам, какие гарантии они могут ожидать от системы памяти, и определяет специальные инструкции, необходимые (называемые барьерами памяти (*memory barriers*) или ограждениями (*fences*)) для получения дополнительных гарантий координации памяти, необходимых при совместном использовании данных. Чтобы оградить разработчика Java от различий между моделями памяти в различных архитектурах, Java предоставляет свою собственную модель памяти, а JVM имеет дело с различиями между JMM и нежеллежащей моделью памяти платформы, путём вставки барьеров памяти в соответствующих местах.

Одним из удобных способов мысленно представить себе модель выполнения программы – это принять, что существует единый порядок, в котором в программе выполняются все операции, независимо от того, на каком процессоре они выполняются, и принять, что каждое чтение переменной будет видеть последнюю запись этой переменной, в порядке выполнения любым процессором. Эта счастливая, хотя и не реалистичная, модель называется согласованной последовательностью (*sequential consistency*). Разработчики программного

обеспечения часто ошибочно предполагают согласованную последовательность, но ни один современный мультипроцессор не предлагает согласованной последовательности, как, в прочем, и JMM. Классическая модель последовательных вычислений, модель фон Неймана, является лишь расплывчатым приближением того, как ведут себя современные мультипроцессоры.

Суть в том, что современные мультипроцессоры с совместно используемой памятью (а также компиляторы) могут делать некоторые удивительные вещи, когда данные совместно используются разными потоками, за исключением ситуации, при которой вы указали им не использовать барьеры памяти. К счастью, программам Java нет необходимости указывать расположение барьеров памяти; они должны только определить момент, когда выполняется обращение к совместно используемому состоянию, посредством надлежащего использования синхронизации.

16.1.2 Переупорядочивание

При описании условий гонки и сбоев атомарности в главе 2, мы использовали диаграммы взаимодействия, изображающие “неудачный момент времени”, в который планировщик чередовал операции, что приводило к появлению неправильных результатов в недостаточно синхронизированных программах. Усугубляя проблему, спецификация JMM может позволить действиям казаться выполняемыми в различном порядке, с точки зрения разных потоков, что делает рассуждения о порядке в отсутствие синхронизации еще более сложными. Различные причины, по которым операции могут задерживаться или выполняться не по порядку, можно сгруппировать в общую категорию *переупорядочивание* (*reordering*).

Класс `PossibleReordering` из листинга 16.1 демонстрирует, как трудно рассуждать о поведении даже самых простых параллельных программ, если они синхронизированы не правильно. Довольно легко себе представить, как класс `PossibleReordering` может печатать (1, 0) или (0, 1), или (1, 1): поток **A** может завершить выполнение до запуска потока **B**, поток **B** может завершить выполнение, до запуска потока **A**, или их действия могут чередоваться. Но, как ни странно, класс `PossibleReordering` также может напечатать (0, 0)! Действия в каждом потоке не зависят друг от друга и, соответственно, могут выполняться не по порядку. (Даже если они выполняются по порядку, задержка времени, с использованием которой кэши сбрасываются в основную память, может привести к тому, что с точки зрения потока **B** присваивания в потоке **A** будут выполняться в обратном порядке.)

```
public class PossibleReordering {
    static int x = 0, y = 0;
    static int a = 0, b = 0;

    public static void main(String[] args)
        throws InterruptedException {
        Thread one = new Thread(new Runnable() {
            public void run() {
                a = 1;
                x = b;
            }
        });
    }
}
```



```

Thread other = new Thread(new Runnable() {
    public void run() {
        b = 1;
        y = a;
    }
});
one.start(); other.start(); one.join();
other.join(); System.out.println("(" + x +
", " + y + ")");
}
}

```

Листинг 16.1 Недостаточно синхронизированная программа, которая может вернуть неожиданные результаты.

На рис. 16.1 показано возможное чередование с переупорядочением, приводящее к печати (0, 0).

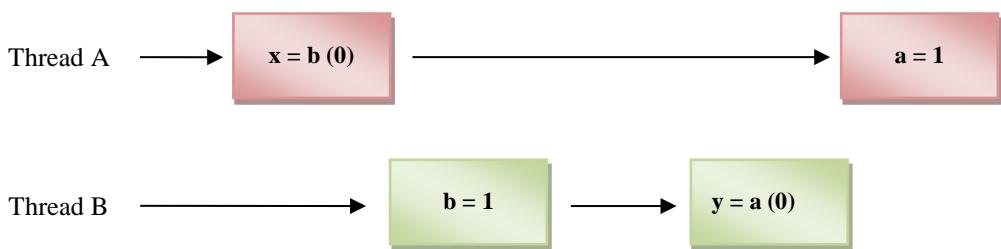


Рис. 16.1 Чередование отражает переупорядочивание в классе PossibleReordering

Класс `PossibleReordering` представляет собой тривиальную программу, но всё ещё на удивление сложно перечислить ее возможные результаты. Изменение порядка на уровне памяти может привести к неожиданному поведению программы. Рассуждать о порядке в отсутствие синхронизации непомерно сложно; гораздо проще убедиться, что ваша программа использует синхронизацию надлежащим образом. Синхронизация препятствует компилятору, среде выполнения и железу в переупорядочивании операций с памятью таким образом, чтобы нарушить гарантии видимости, предоставляемые JMM.¹⁷³

16.1.3 О модели памяти Java менее чем в 500 словах

Модель памяти Java определяется в терминах *действий*, которые включают чтение и запись переменных, блокировку и разблокировку мониторов, запуск и присоединение к потокам. Спецификация JMM определяет частичное упорядочение¹⁷⁴, называемое *happens-before*, на всех действиях программы. Чтобы гарантировать, что поток, выполняющий действие *B*, сможет видеть результаты действия *A* (независимо от того, происходят ли действия *A* и *B* в разных потоках), между действиями *A* и *B* должна существовать связь *happens-before*. В отсутствие

¹⁷³ В большинстве популярных процессорных архитектур модель памяти достаточно сильна, таким образом, затраты на `volatile` чтение стоят в одном ряду с затратами на не `volatile` чтение.

¹⁷⁴ Частичное упорядочение \prec - это отношение на множестве, являющееся антисимметричным, рефлексивным и транзитивным, но для любых двух элементов X и Y , не обязательно должно выполняться, что $X \prec Y$ или $Y \prec X$. Мы используем частичное упорядочивание каждый день, для выражения предпочтений; мы можем иметь предпочтения от суши до чизбургеров и от Моцарта до Малера, но у нас не всегда имеются чёткие предпочтения между чизбургерами и Моцартом.

упорядочивающей связи *happens-before* между двумя операциями, среда JVM может свободно переупорядочивать их по своему усмотрению.

Гонка данных случается, когда переменная считывается более чем одним потоком и записывается, по крайней мере, одним потоком, но операции чтения и записи не упорядочены связью *happens-before*. Корректно синхронизированная программа - это программа без гонки данных; правильно синхронизированные программы демонстрируют последовательную согласованность, что означает, что все действия в программе происходят в фиксированном, глобальном порядке.

Правила для *happens-before*:

Правило порядка программы. Каждое действие потока связывается через отношение *happens-before* с каждым действием в том потоке, что придет после, согласно порядку программы.

Правило блокировка монитора. Разблокировка блокировки на мониторе связана отношением *happens-before* с каждой последующей блокировкой на той же самой блокировке монитора.¹⁷⁵

Правило Volatile переменной. Запись поля `volatile` связана отношением *happens-before* со всеми последующими чтениями того же самого поля.¹⁷⁶

Правило запуска потока. Вызов метода `Thread.start` потока связан отношением *happens-before* с каждым действием в запущенном потоке.

Правило завершения потока. Любое действие в потоке связано отношением *happens-before* с любым другим потоком, определившим, что поток завершен, а также с успешным возвратом из метода `Thread.join` или из метода `Thread.isAlive`, вернувшим `false`.

Правило прерывания. Поток, вызвавший метод `interrupt` другого потока, связан отношением *happens-before* с прерванным потоком, определившим прерывание (а также бросившим исключение `InterruptedException`, или вызвавшим методы `isInterrupted` или `interrupted`).

Правило финализации. Завершение конструктора некоторого объекта связано отношением *happens-before* с запуском финализатора этого объекта.

Транзитивность. Если поток **A** связан отношением *happens-before* с потоком **B**, и поток **B** связан отношением *happens-before* с потоком **C**, тогда поток **A** связан отношением *happens-before* с потоком **C**.

¹⁷⁵ Блокировки и разблокировки явных объектов `Lock` имеют ту же семантику памяти, что и внутренние блокировки.

¹⁷⁶ Операции чтения и записи атомарных переменных имеют ту же семантику памяти, что и переменные `volatile`.

Хотя это лишь частично упорядоченные действия - захват и освобождение блокировки, чтение и запись `volatile` переменной - они полностью упорядочены. Это позволяет описать отношение *happens-before* в терминах “последовательных” захватов блокировок и чтений `volatile` переменных.

На рис. 16.2 показано отношение *happens-before*, когда два потока синхронизируются с помощью общей блокировки. Все действия потока **A** упорядочены согласно правилу порядка программы, также как и все действия потока **B**. Поскольку поток **A** освобождает блокировку **M** и поток **B** впоследствии захватывает блокировку **M**, все действия в потоке **A**, до освобождения блокировки, упорядочены до всех действий в потоке **B**, после захвата блокировки. Когда два потока синхронизируются на *разных* блокировках, мы не можем ничего сказать о порядке выполнения действий между ними - между действиями в двух потоках не существует отношения *happens-before*.

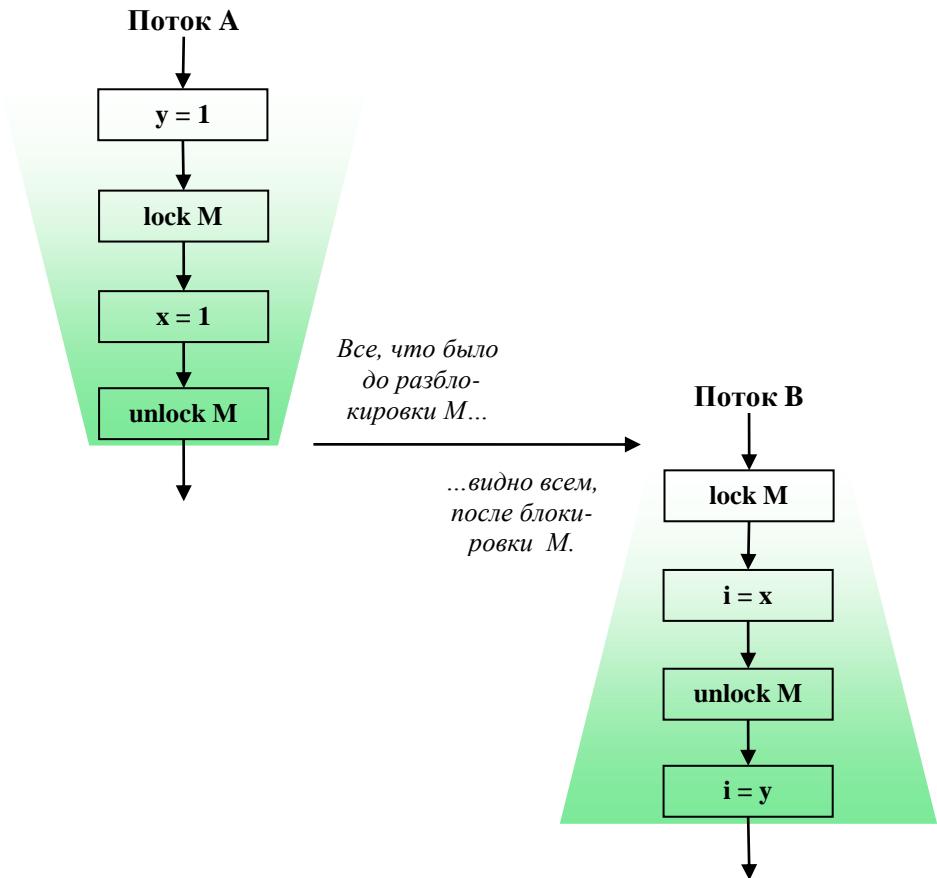


Рисунок 16.2 Иллюстрация отношения *happens-before* в модели памяти Java

16.1.4 Комбинирование в синхронизации

В связи с прочностью упорядочивания отношения *happens-before*, вы иногда можете комбинировать свойства видимости существующей синхронизации. Это влечет за собой объединение правила порядка программы для отношения *happens-before* с одним из других правил порядка (обычно, с правилом блокировки монитора или правилом переменной `volatile`), чтобы упорядочить доступ к

переменной, в ином случае не защищаемой блокировкой. Этот подход очень чувствителен к порядку, в котором происходят обращения, и поэтому довольно хрупок; это продвинутый подход, который должен быть зарезервирован для выжимания последней капли производительности из наиболее критичных для производительности классов, подобных классу `ReentrantLock`.

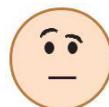
Реализация защищенных методов класса `AbstractQueuedSynchronizer` классом `FutureTask` иллюстрирует комбинирование. Класс AQS поддерживает некоторое целое число, характеризующее состояние синхронизатора, которое класс `FutureTask` использует для хранения состояния задачи: выполняется (*running*), завершено (*completed*) или отменено (*cancelled*). Но класс `FutureTask` также поддерживает дополнительные переменные, такие как результат вычисления. Когда один поток вызывает метод `set`, чтобы сохранить результат, и другой поток вызывает метод `get`, чтобы получить его, обоим лучше быть упорядоченными через отношение *happens-before*. Этого можно добиться, объявив ссылку на результат как `volatile`, но можно использовать и существующую синхронизацию для достижения того же результата при меньших затратах.

Класс `FutureTask` разрабатывался с особой тщательностью, чтобы гарантировать, что успешный вызов метода `tryReleaseShared` всегда связан отношением *happens-before* с последующим вызовом метода `tryAcquireShared`; метод `tryReleaseShared` всегда выполняет запись в `volatile` переменную, значение которой считывается методом `tryAcquireShared`. В листинге 16.2 приведены методы `innerSet` и `innerGet`, вызывающиеся при сохранении или извлечении результата; поскольку метод `innerSet` выполняет запись в переменную `result` перед вызовом метода `releaseShared` (который, в свою очередь, вызывает метод `tryReleaseShared`) и метод `innerGet` выполняет чтение из переменной `result` после вызова метода `acquireShared` (который, в свою очередь, вызывает метод `tryAcquireShared`), правило порядка программы сочетается с правилом `volatile` переменной, чтобы убедиться, что выполнение записи в переменную `result` в методе `innerSet` связано отношением *happens-before* с чтением из переменной `result` в методе `innerGet`.

```
// Inner class of FutureTask
private final class Sync extends AbstractQueuedSynchronizer {
    private static final int RUNNING = 1, RAN = 2, CANCELLED = 4;
    private V result;
    private Exception exception;

    void innerSet(V v) {
        while (true) {
            int s = getState();
            if (ranOrCancelled(s))
                return;
            if (compareAndSetState(s, RAN))
                break;
        }
        result = v;
        releaseShared(0);
        done();
    }

    V innerGet() {
        int s = getState();
        if (ranOrCancelled(s))
            return null;
        if (s == RAN)
            return result;
        if (exception != null)
            throw exception;
        return null;
    }
}
```



```

    V innerGet() throws InterruptedException, ExecutionException {
        acquireSharedInterruptibly(0);
        if (getState() == CANCELLED)
            throw new CancellationException();
        if (exception != null)
            throw new ExecutionException(exception);
        return result;
    }
}

```

Листинг 16.2 Внутренний класс класса FutureTask, иллюстрирующий комбинированную синхронизацию

Мы называем такой подход “комбинированным”, потому что он использует существующий порядок, определяемый отношением *happens-before*, которое было создано по какой-то другой причине, для обеспечения видимости объекта *X*, вместо создания упорядочивающего отношения *happens-before* специально для публикации объекта *X*.

Вид комбинирования, используемого в классе FutureTask довольно хрупок и не должен применяться по случаю. Тем не менее, в некоторых случаях применение комбинирования вполне разумно, например, когда класс фиксирует упорядочивающее отношение *happens-before* между методами, как часть своей спецификации. Например, безопасная публикация с использованием класса BlockingQueue предстает собой форму комбинирования. Один поток помещает объект в очередь, а другой поток впоследствии извлекает его – такая ситуация представляет собой безопасную публикацию, поскольку реализацией класса BlockingQueue обеспечивается достаточная внутренняя синхронизация, чтобы гарантировать, что операция помещения в очередь связана отношением *happens-before* с операцией извлечения объекта из очереди.

Другие отношения упорядочивания *happens-before*, обеспечиваемые библиотекой классов:

- Помещение элемента в потокобезопасную коллекцию связано отношением *happens-before* с другим потоком, извлекающим этот элемент из коллекции;
- Уменьшение счётчика в экземпляре CountDownLatch связано отношением *happens-before* с потоком, возвращаемым из метода await этой защёлки;
- Освобождение разрешения экземпляра Semaphore связано отношением *happens-before* с захватом разрешения того же самого экземпляра Semaphore;
- Действия, выполняемые задачей, представленной экземпляром Future, связаны отношением *happens-before* с другим потоком, успешно возвращённым из метода Future.get;
- Отправка экземпляров Runnable или Callable экземпляру Executor связана отношением *happens-before* с началом выполнения задачи; и
- Поток, прибывший к барьеру CyclicBarrier или двустороннему барьеру Exchanger, связан отношением *happens-before* с другими потоками, освободившимися у того же барьера или точки обмена. Если класс CyclicBarrier использует действие барьера, потоки, прибывающие к барьеру,

связаны отношением *happens-before* с действием барьера, которое, в свою очередь, связано отношением *happens-before* с потоком, освободившимся у барьера.

16.2 Публикация

Ранее, в главе 3, мы рассматривали, как объект может быть опубликован безопасно или неправильно. Описанные там методы безопасной публикации основаны на гарантиях, предоставляемых JMM; риски неправильной публикации являются следствием отсутствия упорядочивающего отношения *happens-before* между публикацией совместного используемого объекта и доступом к нему из другого потока.

16.2.1 Небезопасная публикация

Возможность переупорядочивания при отсутствии связи *happens-before* объясняет, почему публикация объекта без надлежащей синхронизации может позволить другому потоку увидеть *частично созданный объект* (см. раздел 3.5). Инициализация нового объекта включает в себя запись в переменные - поля нового объекта. Аналогично, публикация ссылки включает в себя запись в другую переменную - ссылку на новый объект. Если вы не гарантируете, что при публикации совместно используемой ссылки *happens-before*, другой поток загрузит эту совместно используемую ссылку, то запись ссылки на новый объект может быть переупорядочена (с точки зрения потока, потребляющего объект) с записями в его поля. В этом случае, другой поток может видеть обновленное значение для ссылки на объект, но также видеть устаревшие значения для некоторых или всех состояний этого объекта - частично созданный объект.

Небезопасная публикация может произойти в результате неправильной отложенной инициализации, как показано в листинге 16.3. На первый взгляд, единственной проблемой здесь кажется состояние гонки, описанное в разделе 2.2. При определенных обстоятельствах, например, когда все экземпляры `Resource` идентичны, вы можете игнорировать их (наряду с возможной неэффективностью создания экземпляра `Resource` более одного раза). К сожалению, даже если эти дефекты игнорируются, небезопасная отложенная инициализация все равно небезопасна, поскольку другой поток может наблюдать ссылку на частично созданный экземпляр `Resource`.

```
@NotThreadSafe
public class UnsafeLazyInitialization {
    private static Resource resource;

    public static Resource getInstance() {
        if (resource == null)
            resource = new Resource(); // unsafe publication
        return resource;
    }
}
```



Листинг 16.3 Небезопасная отложенная инициализация. Не делайте так.

Предположим, что поток A первым вызывает метод `getInstance`. Он видит, что объект `resource` имеет значение `null`, создает новый экземпляр `Resource` и

устанавливает объекту `resource` ссылку на него. Когда поток **B** позже вызывает метод `getInstance`, он мог бы видеть, что ресурс уже имеет ненулевое значение и просто использует уже построенный ресурс. Поначалу это может показаться безобидным, но отношения порядка *happens-before* между записью экземпляра `resource` потоком **A** и чтением экземпляра `resource` потоком **B** не существует. Гонка данных использовалась для публикации объекта, и поэтому потоку **B** не гарантируется, что он увидит корректное состояние объекта `Resource`.

Конструктор экземпляра `Resource` изменяет поля только что выделенного экземпляра `Resource` со значений по умолчанию (записанных конструктором `Object`) на исходные значения. Так как ни один поток не использовал синхронизацию, поток **B** мог видеть действия **A** в порядке отличном от того, в котором **A** фактически выполнял их. Таким образом, даже если поток **A** инициализировал экземпляр `Resource` перед установкой `resource` для ссылки на него, поток **B** мог видеть запись в ресурс как происходящую перед записями в поля ресурса. Таким образом, поток **B** может видеть частично созданный ресурс, который вполне может находиться в недопустимом состоянии - и состояние которого может неожиданно измениться позже.

За исключением неизменяемых объектов, небезопасно использовать объект, инициализированный другим потоком, если только публикация не используется потребляющим потоком с отношением *happens-before*.

16.2.2 Безопасная публикация

Идиомы безопасной публикации, описанные в главе 3, гарантируют, что опубликованный объект виден другим потокам, поскольку они гарантируют, что при публикации с отношением *happens-before*, потребляющий поток загрузит ссылку на опубликованный объект. Если поток **A** помещает объект *X* в экземпляр `BlockingQueue` (и иные потоки впоследствии не изменяют его), и поток **B** извлекает его из очереди, поток **B** гарантированно увидит объект *X*, когда поток **A** оставит его. Это вызвано тем, что реализации `BlockingQueue` имеют достаточную внутреннюю синхронизацию, чтобы гарантировать, что метод `put` связан отношением *happens-before* с методом `take`. Аналогичным образом, с помощью совместно используемой переменной, защищаемой блокировкой или совместно используемой переменной `volatile`, гарантируется, что чтение и запись переменной упорядочено отношением *happens-before*.

Отношение *happens-before* фактически даёт более сильные гарантии видимости и упорядоченности, чем можно добиться безопасной публикацией. Когда объект *X* безопасно публикуется от потока **A** к потоку **B**, безопасная публикация гарантирует видимость состояния объекта *X*, но не состояние других переменных потока **A**, которых возможно, коснулись. Но если поток **A** помещает объект *X* в очередь, связанную отношением *happens-before* с потоком **B**, получающим объект *X* из этой очереди, поток **B** не только видит объект *X* в состоянии, в котором его оставил поток **A** (при условии, что объект *X* не был впоследствии изменен потоком **A** или еще кем-либо), но поток **B** также видит всё, что делал поток **A** до передачи (опять же, с той же оговоркой).¹⁷⁷

Почему мы так сильно сосредоточились на аннотации `@GuardedBy` и безопасной публикации, когда спецификация JMM уже предоставляет нам более мощное

¹⁷⁷ Спецификация JMM гарантирует, что поток **B** видит значение, по крайней мере, таким же актуальным, как и значение, записанное **A**; последующие записи могут быть или не быть видимыми.

отношение *happens-before*? Мысление с точки зрения передачи права собственности на объект и публикации лучше вписывается в большинство программных проектов, чем мышление с точки зрения видимости отдельных записей в память. Упорядочивающее отношение *happens-before* работает на уровне отдельных обращений к памяти - это своего рода “язык ассемблера параллелизма”. Безопасная публикация работает на уровне, более близком к дизайну вашей программы.

16.2.3 Идиомы безопасной инициализации

Иногда имеет смысл отложить инициализацию дорогостоящих объектов до тех пор, пока они действительно не понадобятся, но ранее мы уже сталкивались с тем, что неправильное использование отложенной инициализации может привести к проблемам. Класс `UnsafeLazyInitialization` может быть исправлен путем объявления метода `getResource` как `synchronized`, как показано в листинге 16.4. Поскольку ветка кода, выполняющаяся в методе `getInstance`, является довольно короткой (проверка и прогнозируемая ветвь), и если метод `getInstance` вызывается множеством потоков достаточно редко, уровень конкуренции за блокировку `SafeLazyInitialization` будет низким, так что этот подход обеспечит адекватную производительность.

```
@ThreadSafe
public class SafeLazyInitialization {
    private static Resource resource;

    public synchronized static Resource getInstance() {
        if (resource == null)
            resource = new Resource();
        return resource;
    }
}
```

Листинг 16.4 Потокобезопасная ленивая (отложенная) инициализация

Обработка статических полей с инициализаторами (или полей, значения которых инициализируется в статическом блоке инициализации [JPL 2.2.1 и 2.5.3]) несколько специфична и предлагает дополнительные гарантии потокобезопасности. Статические инициализаторы запускаются JVM во время инициализации класса, после загрузки класса, но до того, как класс будет использоваться любым потоком. Поскольку JVM захватывает блокировку во время инициализации [JLS 12.4.2], и эта блокировка захватывается каждым потоком, по крайней мере, один раз, чтобы гарантировать, что класс был загружен, записи в память, сделанные во время статической инициализации, автоматически видны всем потокам. Таким образом, статически инициализированные объекты не требуют явной синхронизации ни во время построения, ни при ссылке. Однако это относится только к состоянию *as-constructed* - если объект является изменяемым, синхронизация по-прежнему требуется как для чтения, так и для записи, чтобы сделать последующие изменения видимыми и избежать повреждения данных.

Использование ранней (*eager*) инициализации, показанной в листинге 16.5, исключает затраты на синхронизацию при каждом вызове метода `getInstance` класса `SafeLazyInitialization`. Этот метод может быть объединен с отложенной

загрузкой класса средой JVM, для создания метода отложенной инициализации, который не требует синхронизации в общей ветви выполнения кода.

```
@ThreadSafe
public class EagerInitialization {
    private static Resource resource = new Resource();

    public static Resource getResource() { return resource; }
}
```

Листинг 16.5 Ранняя инициализация.

Идиома *ленивой инициализации классом холдером* [EJ Item 48], приведённым в листинге 16.6, использует класс, единственной задачей которого является инициализация экземпляра `Resource`. Среда JVM откладывает инициализацию класса `ResourceHolder` до его фактического использования [JLS 12.4.1], и поскольку объект `Resource` инициализируется статическим инициализатором, дополнительная синхронизация не требуется. Первый же вызов метода `getResource` любым потоком, вынуждает класс `ResourceHolder` загрузиться и инициализироваться, в это время происходит инициализация объекта `Resource` через статический инициализатор.

```
@ThreadSafe
public class ResourceFactory {
    private static class ResourceHolder {
        public static Resource resource = new Resource();
    }

    public static Resource getResource() {
        return ResourceHolder.resource;
    }
}
```

Листинг 16.6 Идиома ленивой инициализации классом холдером.

16.2.4 Блокировка с двойной проверкой

Ни одна книга по параллелизму не будет полной без обсуждения печально известного анти паттерна “блокировка с двойной проверкой” (*double-checked locking, DCL*), приведённого в листинге 16.7. В очень ранних версиях JVM, синхронизация, даже в отсутствии конкуренции, приводила к значительным затратам производительности. В результате было изобретено множество умных (или, по крайней мере, выглядящих таковыми) трюков, для уменьшения влияния синхронизации - некоторые хорошие, некоторые плохие и некоторые уродливые. Анти паттерн DCL попадает в категорию “уродливых”.

```
@NotThreadSafe
public class DoubleCheckedLocking {
    private static Resource resource;

    public static Resource getInstance() {
        if (resource == null) {
```



```

        synchronized (DoubleCheckedLocking.class) {
            if (resource == null)
                resource = new Resource();
        }
    }
    return resource;
}

```

Листинг 16.7 Анти паттерн “блокировка с двойной проверкой”. Не делайте так.

Опять же, поскольку производительность ранних версий JVM оставляла желать лучшего, ленивая инициализация часто использовалась, чтобы избежать потенциально ненужных дорогостоящих операций или для сокращения затрат времени на запуск приложения. Правильно написанный метод отложенной инициализации требует синхронизации. Но в то время синхронизация была медленной и, что более важно, не совсем понятной: аспекты исключения были поняты достаточно хорошо, а аспекты видимости - нет.

Анти паттерн DCL намеревался предложить лучшее из обоих миров - ленивую инициализацию, без оплаты штрафа за синхронизацию в общей ветке кода. Сначала нужно было проверить без синхронизации, нужна ли инициализация, и, если ссылка на ресурс не равна `null`, использовать ее. В противном случае, выполнить синхронизацию и проверить еще раз, инициализирован ли ресурс, гарантируя таким образом, что только один поток фактически инициализирует совместно используемый экземпляр `Resource`. Общая ветка кода - получение ссылки на уже созданный ресурс - не использует синхронизацию. И вот в чем проблема: как описано в разделе [16.2.1](#), поток может увидеть частично построенный экземпляр `Resource`.

Реальная проблема с DCL заключается в предположении, что худшее, что может произойти при чтении ссылки на совместно используемый объект без синхронизации, это ошибочно увидеть устаревшее значение (в этом случае `null`); в этом случае идиома DCL компенсирует риск, повторяя попытку с блокировкой. Но наихудший случай на самом деле значительно хуже - можно увидеть текущее значение ссылки, но устаревшие значения для состояния объекта, что означает, что объект может быть замечен в недопустимом или неправильном состоянии.

Последующие изменения в JMM (Java 5.0 и более поздние версии) позволили DCL работать, если `resource` объявлен как `volatile`, и влияние этого на производительность мало, так как `volatile` чтения не намного затратнее не `volatile` чтений. Тем не менее, полезность этой идиомы в значительной степени уже в прошлом – предпосылки к её возникновению (медленная синхронизация в отсутствии конкуренции, медленный запуск среды JVM), больше не играют роли, что делает ее менее эффективной в качестве оптимизации. Идиома ленивой инициализации с помощью холдера предлагает те же преимущества и легче в понимании.

16.3 Безопасность инициализации

Гарантия безопасности инициализации позволяет потокам безопасно совместно использовать, без синхронизации, правильно сконструированные неизменяемые объекты, независимо от способа их публикации, даже если они опубликованы с использованием гонки данных. (Это означает, что класс

`UnsafeLazyInitialization` фактически безопасен, если экземпляр `Resource` неизменяем.)

Без безопасной инициализации, кажущиеся неизменяемыми объекты, такие как `String`, могут изменить свое значение, если синхронизация не используется ни публикующим потоком, ни потребляющим. Архитектура безопасности основана на неизменности экземпляра `String`; отсутствие безопасности инициализации может создать уязвимости безопасности, позволяющие вредоносному коду обходить проверки безопасности.

Безопасность инициализации гарантирует, что для *правильно построенных объектов* все потоки будут видеть правильные значения `final` полей

установленных конструктором, независимо от способа публикации объекта.

Кроме того, любые переменные, которые могут быть *достигнуты (reached)* через `final` поле корректно построенного объекта (например, элементы `final` массива или содержимое экземпляра `HashMap`, на которое ссылается поле `final`), также будут гарантированно видны другим потокам.¹⁷⁸

Для объектов с `final` полями, безопасность инициализации запрещает переупорядочивание любой части конструкции с начальной загрузкой ссылки на этот объект. Все записи в `final` поля, сделанные конструктором, а также в любые переменные, доступные через эти поля, “замораживаются” по завершении конструктора, и любой поток, который получает ссылку на этот объект, гарантированно увидит значение, по крайней мере, такое же актуальное, как и замороженное значение. Записи, инициализирующие переменные доступные через `final` поля, не переупорядочиваются с операциями, следующими за замораживанием после построения.

Безопасность инициализации означает, что класс `SafeStates` из листинга 16.8 может быть безопасно опубликован даже через небезопасную отложенную инициализацию или хранение ссылки на экземпляр `SafeStates` в открытом статическом поле без синхронизации, даже если класс не использует синхронизацию и полагается на не потокобезопасный класс `HashSet`.

```
@ThreadSafe
public class SafeStates {
    private final Map<String, String> states;

    public SafeStates() {
        states = new HashMap<String, String>();
        states.put("alaska", "AK");
        states.put("alabama", "AL");
        ...
        states.put("wyoming", "WY");
    }

    public String getAbbreviation(String s) {
        return states.get(s);
    }
}
```

¹⁷⁸ Это относится только к тем объектам, которые доступны только через `final` поля создаваемого объекта.

}

Листинг 16.8 Безопасность инициализации неизменяемых объектов

Однако ряд небольших изменений в классе `SafeStates` отнимет у него потокобезопасность. Если бы переменная `states` не была объявлена как `final`, или если бы любой метод, отличный от конструктора, изменил содержимое экземпляра, безопасность инициализации была бы недостаточно сильной для безопасного доступа к классу `SafeStates` без синхронизации. Если бы класс `SafeStates` имел другие, не `final`, поля, другие потоки могли бы по-прежнему видеть неправильные значения этих полей. И разрешение объекту сбежать во время строительства, отменяет гарантию безопасности инициализации.

Безопасность инициализации гарантирует видимость только тех значений, которые доступны через поля `final` на момент завершения работы конструктора. Для значений, доступных через поля, не являющиеся `final`, или значений, которые могут изменяться после построения, для гарантии видимости необходимо использовать синхронизацию.

16.3 Итоги

Модель памяти Java определяет, когда действия с памятью одного потока гарантированно будут видимы другому. Специфика включает в себя обеспечение того, чтобы операции были упорядочены частичным порядком, называемым *happens-before*, который задается на уровне отдельных операций памяти и синхронизации. При отсутствии достаточной синхронизации, при доступе потоков к совместно используемым данным, могут происходить очень странные вещи. Однако правила более высокого уровня, предлагаемые в главах 2 и 3, такие как аннотация `@GuardedBy` и безопасная публикация, могут использоваться для обеспечения потокобезопасности, без необходимости касаться низкоуровневых деталей отношения *happens-before*.

Приложение А Описание аннотаций

We've used annotations such as `@GuardedBy` and `@ThreadSafe` to show how thread-safety promises and synchronization policies can be documented. This appendix documents these annotations; their source code can be downloaded from this book's website. (There are, of course, additional thread-safety promises and implementation details that should be documented but that are not captured by this minimal set of annotations.)

A.1 Аннотации уровня классов

Мы используем три аннотации уровня класса, предназначенные для описания обещанного уровня потокобезопасности: `@Immutable`, `@ThreadSafe` и `@NotThreadSafe`. Аннотация `@Immutable` означает, конечно, что класс является неизменяемым и включает в себя аннотацию `@ThreadSafe`. Аннотация `@NotThreadSafe` опциональна - если класс не аннотирован как потокобезопасный, его следует считать не потокобезопасным, но если вы хотите сделать это ещё более очевидным, используйте аннотацию `@NotThreadSafe`.

Эти аннотации относительно ненавязчивы и полезны как пользователям, так и сопровождающим. Пользователи могут сразу увидеть, является ли класс потокобезопасным, а сопровождающие могут сразу увидеть, должны ли быть предварительно предоставлены гарантии потокобезопасности. Аннотации также полезны для третьей группы: инструменты. Статические инструменты анализа кода могут быть в состоянии проверить, что код соответствует контракту, указанному аннотацией, такому как проверка, что класс, аннотированный `@Immutable` фактически, является неизменяемым.

A.2 Аннотации уровня полей и методов

Приведенные выше аннотации уровня класса являются частью общедоступной документации по классу. Другие аспекты стратегии обеспечения потокобезопасности класса предназначены исключительно для разработчиков и не являются частью общедоступной документации.

Классы, использующие блокировки, должны документировать, какие переменные состояния с помощью каких блокировок защищаются и какие блокировки используются для защиты этих переменных. Распространённым источником непреднамеренного внесения не потокобезопасности является ситуация, когда класс согласованно использует блокировку для защиты своего состояния, но позже подвергается изменению, с целью добавления новых переменных состояния, адекватно не защищённых блокировкой, или новых методов, не корректно использующих блокировку, для защиты существующих переменных состояния. Документирование информации о том, какие переменные защищаются какими блокировками, может помочь предотвратить оба типа пропусков.

Аннотация `@GuardedBy(lock)` документирует, что поле или метод должны быть доступны только с определенной блокировкой. Аргумент `lock` определяет блокировку, которая должна удерживаться при доступе к аннотируемому полю или методу. Возможные значения `lock`:

- Аннотация `@GuardedBy("this")`, означает внутреннюю блокировку вмещающего объекта (объектом которого является метод или поле);
- Аннотация `@GuardedBy("fieldName")`, означает блокировку, связанную с объектом, на который ссылается именованное поле, внутреннюю блокировку (для полей, незывающихся на экземпляр `Lock`) или явный экземпляр `Lock` (для полей,зывающихся на экземпляр `Lock`);
- Аннотация `@GuardedBy("ClassName.fieldName")`, подобна аннотации `@GuardedBy("fieldName")`, но ссылка на объект блокировки, удерживается в статическом поле другого класса;
- Аннотация `@GuardedBy("methodName()")` означает, что объект блокировки, возвращается при вызове именованного метода;
- Аннотация `@GuardedBy("ClassName.class")`, литерально означает объект `Class` именованного класса.

Использование аннотации `@GuardedBy` для идентификации каждой переменной состояния, которая нуждается в блокировке и того, какая блокировка что защищает, может помочь в обслуживании и проверке кода, а также помочь инструментам автоматического анализа выявить потенциальные ошибки в потокобезопасности.

Библиография

Ken Arnold, James Gosling, and David Holmes. *The Java Programming Language*, Fourth Edition. Addison–Wesley, 2005.

David F. Bacon, Ravi B. Konuru, Chet Murthy, and Mauricio J. Serrano. Thin Locks: Featherweight Synchronization for Java. In *SIGPLAN Conference on Programming Language Design and Implementation*, pages 258–268, 1998. URL <http://citeseer.ist.psu.edu/bacon98thin.html>.

Joshua Bloch. *Effective Java Programming Language Guide*. Addison–Wesley, 2001.

Joshua Bloch and Neal Gafter. *Java Puzzlers*. Addison–Wesley, 2005.

Hans Boehm. Destructors, Finalizers, and Synchronization. In *POPL '03: Proceedings of the 30th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 262–272. ACM Press, 2003. URL <http://doi.acm.org/10.1145/604131.604153>.

Hans Boehm. Finalization, Threads, and the Java Memory Model. JavaOne presentation, 2005. URL <http://developers.sun.com/learning/javaoneonline/2005/coreplatform/TS-3281.pdf>.

Joseph Bowbeer. The Last Word in Swing Threads, 2005. URL <http://java.sun.com/products/jfc/tsc/articles/threads/threads3.html>.

Cliff Click. Performance Myths Exposed. JavaOne presentation, 2003.

Cliff Click. Performance Myths Revisited. JavaOne presentation, 2005. URL <http://developers.sun.com/learning/javaoneonline/2005/coreplatform/TS-3268.pdf>.

Martin Fowler. Presentation Model, 2005. URL <http://www.martinfowler.com/eaaDev/PresentationModel.html>.

Erich Gamma, Richard Helm, Ralph Johnson, and John Vlissides. *Design Patterns*. Addison–Wesley, 1995.

Martin Gardner. The fantastic combinations of John Conway's new solitaire game 'Life'. *Scientific American*, October 1970.

James Gosling, Bill Joy, Guy Steele, and Gilad Bracha. *The Java Language Specification*, Third Edition. Addison–Wesley, 2005.

Tim Harris and Keir Fraser. Language Support for Lightweight Transactions. In *OOPSLA '03: Proceedings of the 18th Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications*, pages 388–402. ACM Press, 2003. URL <http://doi.acm.org/10.1145/949305.949340>.

Tim Harris, Simon Marlow, Simon Peyton-Jones, and Maurice Herlihy. Composable Memory Transactions. In *PPoPP '05: Proceedings of the Tenth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 48–60. ACM Press, 2005. URL <http://doi.acm.org/10.1145/1065944.1065952>.

Maurice Herlihy. Wait-Free Synchronization. *ACM Transactions on Programming Languages and Systems*, 13(1):124–149, 1991. URL <http://doi.acm.org/10.1145/114005.102808>.

Maurice Herlihy and Nir Shavit. *Multiprocessor Synchronization and Concurrent Data Structures*. Morgan-Kaufman, 2006.

C. A. R. Hoare. Monitors: An Operating System Structuring Concept. *Communications of the ACM*, 17(10):549–557, 1974. URL <http://doi.acm.org/10.1145/355620.361161>.

David Hovemeyer and William Pugh. Finding Bugs is Easy. *SIGPLAN Notices*, 39(12):92–106, 2004. URL <http://doi.acm.org/10.1145/1052883.1052895>.

Ramnivas Laddad. *AspectJ in Action*. Manning, 2003.

Doug Lea. *Concurrent Programming in Java*, Second Edition. Addison–Wesley, 2000.

Doug Lea. JSR-133 Cookbook for Compiler Writers. URL <http://gee.cs.oswego.edu/dl/jmm/cookbook.html>.

J. D. C. Little. A proof of the Queueing Formula $L = \lambda W$ ". *Operations Research*, 9:383–387, 1961.

Jeremy Manson, William Pugh, and Sarita V. Adve. The Java Memory Model. In *POPL '05: Proceedings of the 32nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 378–391. ACM Press, 2005. URL <http://doi.acm.org/10.1145/1040305.1040336>.

George Marsaglia. XorShift RNGs. *Journal of Statistical Software*, 8(13), 2003. URL <http://www.jstatsoft.org/v08/i14>.

Maged M. Michael and Michael L. Scott. Simple, Fast, and Practical Non-Blocking and Blocking Concurrent Queue Algorithms. In *Symposium on Principles of Distributed Computing*, pages 267–275, 1996. URL <http://citeseer.ist.psu.edu/michael96simple.html>.

Mark Moir and Nir Shavit. *Concurrent Data Structures*, In *Handbook of Data Structures and Applications*, chapter 47. CRC Press, 2004.

William Pugh and Jeremy Manson. Java Memory Model and Thread Specification, 2004. URL <http://www.cs.umd.edu/~pugh/java/memoryModel/jsr133.pdf>.

M. Raynal. *Algorithms for Mutual Exclusion*. MIT Press, 1986.

William N. Scherer, Doug Lea, and Michael L. Scott. Scalable Synchronous Queues. In *11th ACM SIGPLAN Symposium on Principles and Practices of Parallel Programming (PPoPP)*, 2006.

R. K. Treiber. Systems Programming: Coping with Parallelism. Technical Report RJ 5118, IBM Almaden Research Center, April 1986.

Andrew Wellings. *Concurrent and Real-Time Programming in Java*. John Wiley & Sons, 20