

Содержание

Предисловие	13
Глава 1. Введение.....	19
1.1. Обоснование.....	19
1.2. Абстракции распределенного программирования	21
1.2.1. Естественно распределенные системы.....	22
1.2.2. Искусственно распределенные системы	25
1.3. Сквозной аргумент	26
1.4. Программные компоненты	27
1.4.1. Композиционная модель.....	27
1.4.2. Программный интерфейс	30
1.4.3. Модули.....	32
1.5. Классы алгоритмов	36
1.6. Практикум	37
1.6.1. Модуль Print	37
1.6.2. Модуль BoundedPrint	39
1.6.3. Объединение модулей.....	43
1.6.4. Эксперимент.....	44
1.7. Примечания к главе.....	45
Глава 2. Базовые абстракции.....	46
2.1. Распределенные вычисления	47
2.1.1. Процессы и сообщения.....	47
2.1.2. Автоматы и шаги.....	47
2.1.3. Безопасность и живучесть.....	49
2.2. Абстракции процессов.....	51
2.2.1. Отказы процессов.....	51
2.2.2. Аварии	51
2.2.3. Пропуск данных.....	53
2.2.4. Аварии с восстановлением	53
2.2.5. Перехват сообщений.....	56
2.2.6. Произвольное поведение	56
2.3. Криптографические абстракции.....	58
2.3.1. Хэш-функции.....	58
2.3.2. Коды аутентификации сообщений.....	58
2.3.3. Цифровые подписи	59
2.4. Абстракции связи	60
2.4.1. Отказы связи.....	61

2.4.2. Связи с приемлемыми потерями	63
2.4.3. Связи с повторами.....	63
2.4.4. Идеальные связи.....	65
2.4.5. Идеальные связи с журналированием	67
2.4.6. Идеальная связь с аутентификацией.....	70
2.4.7. Об абстракциях связей	72
2.5. Предположения о синхронности	74
2.5.1. Асинхронные системы	74
2.5.2. Синхронные системы.....	75
2.5.3. Частичная синхронность	77
2.6. Абстракция времени.....	78
2.6.1. Обнаружение отказов	78
2.6.2. Идеальное обнаружение отказа.....	79
2.6.3. Выбор лидера	82
2.6.4. Эвентуально идеальный датчик отказов.....	84
2.6.5. Эвентуальный выбор лидера.....	87
2.6.6. Византийский выбор лидера	91
2.7. Модели распределенных систем.....	94
2.7.1. Комбинации абстракций.....	95
2.7.2. Подготовка.....	96
2.7.3. Кворумы.....	97
2.7.4. Измерение стоимости алгоритмов	98
2.8. Упражнения.....	99
2.9. Решения	100
2.10. Практикум.....	103
2.10.1. Передаваемое событие.....	103
2.10.2. Сообщение	104
2.10.3. Связь «точка-точка» с приемлемыми потерями	105
2.10.4. Идеальная связь «точка-точка»	105
2.10.5. Идеальный датчик отказов	106
2.11. Примечания к главе.....	108
Глава 3. Надежная рассылка.....	111
3.1. Обоснование.....	111
3.1.1. Клиент-серверные вычисления	111
3.1.2. Системы со множеством участников.....	112
3.2. Негарантированная рассылка.....	113
3.2.1. Спецификация.....	113
3.2.2. Алгоритм простой рассылки.....	114
3.3. Обычная надежная рассылка	115
3.3.1. Спецификация.....	115
3.3.2. Алгоритм для модели с остановкой после отказов: ленивая надежная рассылка.....	116

3.3.3. Алгоритм для модели без уведомления об отказах: жадная надежная рассылка.....	118
3.4. Унифицированная надежная рассылка	120
3.4.1. Спецификация.....	120
3.4.2. Алгоритм для модели с остановкой после отказов: унифицированная надежная рассылка с подтверждением всеми	121
3.4.3. Алгоритм для модели без уведомления об отказах: унифицированная надежная рассылка с подтверждением большинством	124
3.5. Рассылка с повторами.....	125
3.5.1. Спецификация.....	125
3.5.2. Алгоритм для модели с восстановлением после отказов: простая рассылка с повторами	126
3.6. Негарантированная рассылка с журналированием.....	127
3.6.1. Обзор.....	127
3.6.2. Спецификация.....	128
3.6.3. Алгоритм для модели с восстановлением после отказов: простая рассылка с журналированием	129
3.7. Унифицированная надежная рассылка с журналированием.....	130
3.7.1. Спецификация.....	130
3.7.2. Fail-Recovery Algorithm: Logged Majority-Ack Uniform Reliable Broadcast	131
3.8. Вероятностная рассылка.....	132
3.8.1. Масштабируемость надежной рассылки.....	133
3.8.2. Распространение эпидемии	134
3.8.3. Спецификация.....	134
3.8.4. Рандомизированный алгоритм: жадная вероятностная рассылка	135
3.8.5. Рандомизированный алгоритм: ленивая вероятностная рассылка	138
3.9. Рассылка в порядке FIFO и причинная рассылка.....	142
3.9.1. Обзор.....	142
3.9.2. Спецификация порядка FIFO	143
3.9.3. Алгоритм для модели без уведомления об отказах: рассылка с использованием последовательных номеров	143
3.9.4. Спецификация причинно-следственного порядка.....	144
3.9.5. Алгоритм для модели без уведомления об отказах: рассылка в причинно-следственном порядке без ожидания	146
3.9.6. Алгоритм для модели с остановкой после отказов: сборка мусора причинно-следственного прошлого.....	148
3.9.7. Алгоритм для модели без уведомления об отказах: причинно-следственная рассылка с ожиданием	150
3.10. Византийская согласованная рассылка	152
3.10.1. Обоснование	153
3.10.2. Спецификация	154
3.10.3. Алгоритм для модели с произвольным поведением: эхо-рассылка с аутентификацией	155

3.10.4. Алгоритм для модели с произвольным поведением: эхо-рассылка с подписанными сообщениями	158
3.11. Византийская надежная рассылка	160
3.11.1. Спецификация	160
3.11.2. Алгоритм для модели с произвольным поведением: двойная эхо-рассылка с аутентификацией	161
3.12. Византийские широковещательные каналы	164
3.12.1. Спецификации	164
3.12.2. Алгоритм для модели с произвольным поведением: Византийский согласованный канал	166
3.12.3. Алгоритм для модели с произвольным поведением: Византийский надежный канал	167
3.13. Упражнения.....	167
3.14. Решения.....	169
3.15. Практикум.....	178
3.15.1. Простая рассылка.....	178
3.15.2. Ленивая надежная рассылка	181
3.15.3. Унифицированная надежная рассылка с подтверждением всеми.....	184
3.15.4. Унифицированная надежная рассылка с подтверждением большинством.....	187
3.15.5. Вероятностная надежная рассылка	188
3.15.6. Рассылка в причинно-следственном порядке без ожидания.....	191
3.15.7. Рассылка в причинно-следственном порядке без ожидания и с уборкой мусора.....	197
3.15.8. Рассылка в причинно-следственном порядке с ожиданием	204
3.16. Примечания к главе.....	208

Глава 4. Разделяемая память210

4.1. Введение	211
4.1.1. Разделяемое хранилище в распределенной системе	211
4.1.2. Регистр	211
4.1.3. Завершенность и предшествование	214
4.2. Обычный регистр (1, N).....	216
4.2.1. Спецификация.....	216
4.2.2. Алгоритм для модели с остановкой после отказов: обычный регистр «Чтение в одном – запись во всех»	217
4.2.3. Алгоритм для модели без уведомления об отказах: обычный регистр с голосованием большинством	220
4.3. Атомарный регистр (1, N).....	223
4.3.1. Спецификация.....	223
4.3.2. Преобразование обычного регистра (1, N) в атомарный регистр (1, N)	226

4.3.3. Алгоритм атомарного регистра для модели с остановкой после отказов: чтение с записью – запись во всех	231
4.3.4. Алгоритм атомарного регистра для модели без уведомления об отказах: чтение с записью – запись в большинстве	233
4.4. Атомарный регистр (N, N)	235
4.4.1. Несколько пишущих процессов	235
4.4.2. Спецификация.....	236
4.4.3. Преобразование атомарного регистра ($1, N$) в атомарный регистр (N, N).....	237
4.4.4. Алгоритм для модели с остановкой после отказов: чтение с записью – запись с проверкой во всех	241
4.4.5. Алгоритм для модели без уведомления об отказах: чтение с записью – запись с проверкой в большинстве	244
4.5. Обычный регистр ($1, N$) с журналированием	247
4.5.1. Предшествование в модели с восстановлением после отказов	247
4.5.2. Спецификация.....	248
4.5.3. Алгоритм для модели с восстановлением после отказов: голосование большинством с журналированием.....	250
4.6. Византийский безопасный регистр ($1, N$)	253
4.6.1. Спецификация.....	254
4.6.2. Алгоритм для модели с произвольным поведением после отказов: Византийский маскирующий кворум	255
4.7. Византийский обычный регистр ($1, N$)	257
4.7.1. Спецификация.....	258
4.7.2. Алгоритм для модели с произвольным поведением: Аутентификация данных Византийским кворумом.....	258
4.7.3. Алгоритм для модели с произвольным поведением после отказов: Византийский кворум с двухфазной записью.....	261
4.8. Византийский атомарный регистр ($1, N$).....	268
4.8.1. Спецификация.....	268
4.8.2. Алгоритм для модели с произвольным поведением после отказов: Византийский кворум со слушателями.....	269
4.9. Упражнения.....	273
4.10. Решения.....	274
4.11. Практикум	280
4.11.1. Обычный регистр ($1, N$).....	280
4.11.2. Атомарный регистр ($1, N$)	284
4.11.3. Атомарный регистр (N, N).....	288
4.12. Примечания к главе.....	292

Глава 5. Консенсусы.....296

5.1. Обычный консенсус	297
5.1.1. Спецификация.....	297

5.1.2. Алгоритм для модели с остановкой после отказов: консенсус заполнением	298
5.1.3. Алгоритм для модели с остановкой после отказов: иерархический консенсус	302
5.2. Унифицированный консенсус	305
5.2.1. Спецификация.....	305
5.2.2. Алгоритм для модели с остановкой после отказов: унифицированный консенсус с заполнением	306
5.2.3. Алгоритм для модели с остановкой после отказов: унифицированный иерархический консенсус	308
5.3. Унифицированный консенсус для модели с уведомлением об отказах.....	311
5.3.1. Обзор.....	311
5.3.2. Смена эпох	312
5.3.3. Консенсус эпохи.....	316
5.3.4. Алгоритм для модели с уведомлением об отказах: Консенсус с лидером	322
5.4. Консенсус с журналированием	325
5.4.1. Спецификация.....	325
5.4.2. Смена эпох с журналированием.....	327
5.4.3. Консенсус эпохи с журналированием	328
5.4.4. Алгоритм для модели с восстановлением после отказов: Консенсус с лидером и журналированием	332
5.5. Рандомизированный консенсус	333
5.5.1. Спецификация.....	334
5.5.2. Общая монета	335
5.5.3. Рандомизированный алгоритм для модели без уведомления об отказах: рандомизированный двоичный консенсус	336
5.5.4. Рандомизированный алгоритм для модели без уведомления об отказах: рандомизированный консенсус с большой областью выбора.....	341
5.6. Византийский консенсус	343
5.6.1. Спецификация.....	343
5.6.2. Византийская смена эпох	346
5.6.3. Византийский консенсус эпохи.....	348
5.6.4. Алгоритм для модели с уведомлением об отказах и произвольным поведением после отказов: Византийский консенсус с лидером	360
5.7. Византийский рандомизированный консенсус	362
5.7.1. Спецификация.....	362
5.7.2. Рандомизированный алгоритм для модели с произвольным поведением после отказов: Византийский рандомизированный двоичный консенсус	362
5.8. Упражнения.....	367
5.9. Решения.....	369
5.10. Практикум.....	380

5.10.1. Протокол обычного консенсуса заполнением.....	380
5.10.2. Протокол обычного иерархического консенсуса	385
5.10.3. Унифицированный консенсус заполнением.....	389
5.10.4. Унифицированный иерархический консенсус	392
5.11. Примечания к главе.....	396

Глава 6. Варианты консенсуса400

6.1. Рассылка в едином порядке.....	400
6.1.1. Обзор.....	400
6.1.2. Спецификации.....	402
6.1.3. Алгоритм для модели без уведомления об отказах: рассылка в едином порядке на основе консенсуса	404
6.2. Византийская рассылка в едином порядке	407
6.2.1. Обзор.....	407
6.2.2. Спецификация.....	408
6.2.3. Алгоритм для модели с уведомлением об отказах и произвольным поведением после отказов: Византийская рассылка с ротацией отправителя	408
6.3. Обрывающаяся надежная рассылка	413
6.3.1. Обзор.....	413
6.3.2. Спецификация.....	414
6.3.3. Алгоритм для модели с остановкой после отказов: унифицированная обрывающаяся надежная рассылка на основе консенсуса	414
6.4. Быстрый консенсус.....	417
6.4.1. Обзор.....	417
6.4.2. Спецификация.....	418
6.4.3. Алгоритм для модели без уведомления об отказах: От унифицированного консенсуса к унифицированному быстрому консенсусу.....	419
6.5. Быстрый Византийский консенсус.....	422
6.5.1. Обзор.....	422
6.5.2. Спецификация.....	422
6.5.3. Алгоритм для модели с произвольным поведением: от Византийского консенсуса к быстрому Византийскому консенсусу	423
6.6. Неблокирующее атомарное подтверждение	425
6.6.1. Обзор.....	425
6.6.2. Спецификация.....	427
6.6.3. Алгоритм для модели с остановкой после отказов: неблокирующее атомарное подтверждение на основе консенсуса	427
6.7. Членство в группах.....	430
6.7.1. Обзор.....	430
6.7.2. Спецификация.....	431

6.7.3. Алгоритм для модели с остановкой после отказов: членство в группе на основе консенсуса	432
6.8. Синхронизация взаимодействий с представлением	435
6.8.1. Обзор	435
6.8.2. Спецификация	436
6.8.3. Алгоритм для модели с остановкой после отказов: синхронизация взаимодействий с представлением на основе TRB	438
6.8.4. Алгоритм для модели с остановкой после отказов: унифицированная синхронизация взаимодействий с представлением	443
6.9. Упражнения	447
6.10. Решения	449
6.11. Практикум	464
6.11.1. Унифицированная рассылка в едином порядке	464
6.11.2. Неблокирующее атомарное подтверждение на основе консенсуса	471
6.11.3. Членство в группе на основе консенсуса	474
6.11.4. Синхронизация с представлением на основе TRB	478
6.12. Примечания к главе	485

Глава 7. Заключительные замечания488

7.1. Реализация в Arria	488
7.2. Дополнительные реализации	489
7.3. Для дальнейшего чтения	491

Использованная литература493

Список модулей501

Список алгоритмов503

Предметный указатель506

Предисловие

Эта книга содержит введение в абстракции распределенного программирования и знакомит с фундаментальными алгоритмами и их реализациями в нескольких распределенных окружениях. Перед читателем будут раскрыты важные проблемы распределенных вычислений и основные алгоритмические приемы их решения. На подробных примерах читатель сможет понять, как с помощью этих приемов конструировать распределенные приложения. Главной темой книги является обеспечение устойчивости к случайным и преднамеренным воздействиям в распределенных системах, которые могут быть обусловлены задержками в сети, ошибками и даже злонамеренными атаками.

Содержание

В современных вычислениях программы нередко объединяют *несколько процессов*. Процесс в данном контексте – это простая абстракция, которая может представлять физический или виртуальный компьютер, процессор в компьютере или конкретный поток выполнения в многозадачной системе. Основная проблема таких распределенных программ состоит в том, чтобы заставить все процессы *вместе работать* над решением *общей* задачи. При этом в каждом процессе все еще приходится решать традиционные алгоритмические проблемы. Распределенные окружения от единичных компьютеров до вычислительных центров и даже глобальных систем, доступных круглосуточно, влекут за собой дополнительные сложности: как добиться надежного решения задач даже в случае аварийного завершения некоторых процессов, потери связи с некоторыми процессами и злонамеренного нападения на некоторые процессы? Распределенные алгоритмы должны быть безотказными, надежными, безопасными и иметь предсказуемое поведение даже при наличии в распределенном окружении отрицательно влияющих факторов.

Если бы не требовалось сотрудничество процессов, распределенные программы просто представляли бы собой множество независимых программ, каждая из которых выполняется в отдельном процессе, и мы были бы лишены преимуществ многозадачности в распределенных системах. Потребность в обеспечении сотрудничества обнажила множество сложных проблем, рассматривающихся в этой книге, которые необходимо было решить, чтобы сделать распределенные вычисления реальностью. Данная книга не просто знакомит читателей с описаниями этих проблем, но также демонстрирует пути их решения в разных контекстах.

Неудивительно, что распределенное программирование можно существенно упростить, если заключить сложности надежного сотрудничества в определенные *абстракции*. Изолируя сложные алгоритмические проблемы, такие абстракции распределенного программирования служат мостом, соединяющим уровни сетевых взаимодействий, которые не дают никаких гарантий надежности, и распределенные прикладные уровни, которые часто требуют высоконадежных примитивов.

Книга знакомит с разными абстракциями распределенного программирования и описывает алгоритмы для их реализации. В некотором смысле мы даем прикладному программисту библиотеку спецификаций абстрактных интерфейсов, а системному программисту – библиотеку алгоритмов для реализации этих спецификаций.

Работая над этой книгой, мы потратили много времени, чтобы сформировать коллекцию упражнений и их решений. Мы настоятельно рекомендуем читателям проработать эти упражнения. Никакие знания и навыки не могут быть приобретены пассивным путем. Это особенно верно в области распределенных вычислений, где человеческий разум часто стремится следовать за привлекательными, но ошибочными представлениями. Книга также включает решения для всех упражнений, чтобы подчеркнуть наше намерение сделать их неотъемлемой частью содержания книги. Многие упражнения просты и могут решаться коллективно. Другие – сложнее и требуют больше времени. Как правило, такие упражнения решаются индивидуально.

Представление

Эта книга получилась вполне самодостаточной. Это стало возможным, потому что сфера распределенных алгоритмов достигла определенного уровня зрелости, когда отвлекающие детали можно вынести за рамки обсуждения распределенных алгоритмов. К числу таких деталей можно отнести работу сетей связи, разного рода отказы, а также реализации криптографических примитивов; все они подробно рассматриваются в других работах. Безусловно, элементарные знания алгоритмов, логики предикатов, языков программирования, особенностей работы сетей, принципов обеспечения безопасности и операционных систем будут полезны при чтении книги. Но мы считаем, что большинство наших абстракций и алгоритмов будет понятно даже тем, кто обладает минимальными знаниями об этих понятиях.

Книга следует поэтапному подходу и писалась в первую очередь как учебник для высших учебных заведений. Она знакомит читателей с основными элементами распределенных вычислений на наглядных примерах и выстраивает сложные абстракции распределенного программирования постепенно, от простого к сложному. Всякий раз, разрабатывая алгоритм для реализации той или иной абстракции, мы сначала рассматриваем простую модель распределенной системы, а затем адаптируем этот алгоритм для более сложных моделей. Иными словами, мы сначала разрабатываем алгоритм, опираясь на упрощенные предположения об устройстве распределенного окружения, а затем обсуждаем, как избавиться от этих упрощений.

Мы старались сбалансировать простоту и наглядность, с одной стороны, и строгость – с другой. Местами строгости уделяется меньше внимания, и это не всегда делалось специально. Основной упор в этой книге делается на абстракциях и алгоритмах, а не на вычислимости и сложности. На самом деле в этой книге не обсуждаются теоремы. Доказательства приводятся лишь с целью разъяснения алгоритмов: сами по себе они не являются формальными доказательствами.

Организация

Эта книга содержит шесть глав в двух частях. Первая часть строит общую платформу.

В главе 1 мы *обоснуем* необходимость абстракций распределенного программирования, обсудив некоторые приложения, где используются эти абстракции. В этой главе также будут представлены блочная форма записи и псевдокод, используемые для записи алгоритмов в книге.

В главе 2 будут представлены разного рода *допущения и предположения* о распределенной среде. С этой целью мы введем в обиход семейство моделей распределенных систем. В основном модели описывают низкоуровневые абстракции, на которых будут строиться абстракции более высокого уровня. К числу таких абстракций относятся абстракции процесса и связей. Эту главу можно рассматривать как справочник для других глав.

Остальные четыре главы составляют вторую часть книги. Каждая глава посвящена одной проблеме, содержащей широкий класс абстракций, и разным алгоритмам их реализации. Мы будем двигаться от простых абстракций к сложным.

В главе 3 мы введем абстракции обмена данными для распределенного программирования. Они описывают *широковещательную рассылку* сообщений группам процессов и предлагают разные гарантии надежности доставки сообщений процессам. Например, мы обсудим, как гарантировать, что сообщение, доставленное одному процессу, будет доставлено всем другим процессам, даже в случае аварийного завершения процесса-отправителя.

В главе 4 мы обсудим абстракции *разделяемой памяти* как простые формы распределенных хранилищ объектов, доступных для чтения и записи. Это могут быть файлы в распределенной системе хранения или регистры в памяти многопроцессорного компьютера. Мы охватим методы чтения и записи данных клиентами, такие, что значения, хранимые множеством процессов, позднее могут быть получены, даже если некоторые из процессов потерпели аварию, стерли свои значения или записали ошибочные данные.

В главе 5 мы обратимся к абстракциям *консенсуса*, с помощью которых множество процессов может приходиться к общему решению, опираясь на исходные данные, изначально предложенные процессами. Они могут приходиться к тому же решению, несмотря на отказы в отдельных процессах, которые могут не только завершаться аварийно, но и активно препятствовать принятию общего решения.

В главе 6 мы рассмотрим *варианты консенсуса*, которые могут быть получены расширением или адаптацией абстракции консенсуса в соответствии с потребностями приложений. Сюда входят рассылка в общем порядке (total-order broadcast), надежная рассылка с ограничениями, атомарная (неблокирующая) фиксация, членство в группах и синхронные взаимодействия.

Распределенные алгоритмы, которые мы будем обсуждать, отличаются не только фактическими абстракциями, которые они реализуют, но также набором предположений и допущений о распределенной среде. Начальное множество абстракций, которое алгоритм считает самим собой разумеющимся, мы назвали мо-

делью *распределенной системы*. Существует множество аспектов, оказывающих фундаментальное влияние на проектирование алгоритмов, таких как надежность связей, степень синхронности системы, серьезность отказов и требования к детерминизму решения.

В разных местах в книге одни и те же базовые примитивы распределенного программирования будут реализованы в нескольких моделях распределенных систем. Такой подход имеет перед собой две цели: во-первых, дать представление о конкретной проблеме, свойственной данной модели, и, во-вторых, показать, как выбор модели влияет на реализацию примитива.

Материал всех глав и комплекс упражнений образуют обширное и исчерпывающее введение в тему распределенного программирования. Каждая глава, описывающая отдельный класс абстракций и алгоритмы их реализации в самом простом представлении, то есть в самой простой модели, предусматривающей только аварийные отказы, могла бы с успехом составить основу более короткого и простого курса. Такие курсы могли бы служить хорошим подспорьем для более практических курсов по распределенному программированию.

Изменения во втором издании

Это издание является полностью пересмотренной версией первого издания. Обновлениям подверглась большая часть книги. Но самое большое изменение произошло в области охватываемых тем, число которых увеличилось и дополнилось темой обеспечения *безопасности против злонамеренных действий*. Абстракции и алгоритмы в модели распределенных вычислений, где допускается вероятность злонамеренных действий, известны как Византийская модель сбоев (Byzantine fault-tolerance), или модель произвольных сбоев.

Первое издание книги называлось «Introduction to Reliable Distributed Programming». Добавление одного слова («secure» – «безопасное») в название и увеличение числа соавторов в жизни книги отражает развитие сферы распределенных систем в современном мире. Еще в 2006 году, когда вышло первое издание книги, было очевидно, что большинство действующих распределенных систем постоянно находится под угрозой вторжения и что штатных сотрудников, имеющих доступ к конфиденциальной информации, нельзя исключать из списка потенциальных источников злонамеренных атак. Для создания надежных распределенных систем в настоящее время требуются серьезные усилия, с привлечением знаний распределенных алгоритмов, методов обеспечения безопасности и из других областей.

С технической точки зрения, были изменены синтаксис модулей и имена некоторых событий, чтобы сделать алгоритмы более структурированными. Теперь модуль может существовать сразу в нескольких экземплярах внутри алгоритма, и в связи с этим каждый экземпляр получает свой уникальный идентификатор. Мы считаем, что это упростило представление некоторых важнейших алгоритмов.

В первое издание был включен вспомогательный набор действующих примеров реализации алгоритмов на языке Java, использующих фреймворк *Appia*. Эти примеры доступны для загрузки на веб-сайте книги.

Ресурсы в Интернете

Дополнительная информация о книге, включая реализацию многих протоколов из первого издания, учебные презентации, слайды и список опечаток, доступна на веб-сайте книги: <http://distributedprogramming.net>.

Ссылки

Мы занимаемся исследованием абстракций распределенного программирования уже почти два десятилетия. Основой для этой книги послужили работы многих исследователей в области распределенных вычислений. Но особенно нам хотелось бы отметить Лесли Лэмпорта (Leslie Lamport) и Нэнси Линч (Nancy Lynch), сформулировавших интереснейшие проблемы в распределенных вычислениях, и ученых из Корнельского университета, занимавшихся исследованиями надежности распределенных вычислений, включая Озалпа Бабайоглу (Özalp Babaoglu), Кена Бирмана (Ken Birman), Кейта Марцзулло (Keith Marzullo), Роберта ван Ренеше (Robbert van Renesse), Рика Шлихтинга (Rick Schlichting), Фреда Шнайде-ра (Fred Schneider) и Сэма Туера (Sam Toueg).

Прямое или косвенное отношение к этой книге имеют также многие другие исследователи. Мы постарались никого не забыть и добавить ссылки на их работы повсюду в этой книге. Все главы заканчиваются примечаниями, где приводятся информация по теме и исторические ссылки. Цель этих примечаний состоит в том, чтобы подсказать, какие еще работы можно прочесть в дальнейшем, проследить историю представленных концепций, а также отдать должное всем, кто причастен к их созданию и разработке. В конце книги приводится список книг для будущего чтения, посвященных другим аспектам распределенных вычислений.

Благодарности

Мы хотели бы выразить глубокую благодарность нашим студентам и аспирантам Федеральной политехнической школы Лозанны (École Polytechnique Fédérale de Lausanne, EPFL) и Лиссабонского университета (University of Lisboa, UL) за рецензирование рукописей этой книги. На самом деле они не имели выбора и в любом случае должны были сдавать экзамены нам! Но они были весьма снисходительны к ошибкам и опечаткам, которых было немало в первых версиях книги и сопроводительных слайдах, и они представили множество очень ценных замечаний.

Парта Дутта (Partha Dutta), Корин Хэри (Corine Hari), Михаль Капалка (Michal Kapalka), Петр Кузнецов (Petr Kouznetsov), Рон Леви (Ron Levy), Максим Монод (Maxime Monod), Бастиан Пошон (Bastian Pochon) и Джаспер Спринг (Jesper Spring), аспиранты школы информационных технологий и теории связи (School of Computer and Communication Sciences) университета EPFL, Филипе Араужо (Filipe Araujo) и Хьюго Миранда (Hugo Miranda), аспиранты из группы «Распределенные алгоритмы и сетевые протоколы» (Distributed Algorithms and Network Protocol, DIALNP) кафедры информатики (Departamento de Informatica)

факультета естественных наук Лиссабонского университета (Faculdade de Ciências da Universidade de Lisboa), Лейла Халил (Leila Khalil) и Роберт Басмаджан (Robert Basmadjian), аспиранты из Ливанского университета в Бейруте (Lebanese University), а также Али Годши (Ali Ghodsi), аспирант из Шведского института компьютерных наук (Swedish Institute of Computer Science, SICS) в Стокгольме предложили массу улучшений к алгоритмам, представленным в книге.

Несколько реализаций для «практической» части книги было разработано Александром Пинто (Alexandre Pinto) или с его помощью, ключевым членом проекта *Appia*, при участии некоторых преподавателей и студентов группы DIALNP, включая Нуно Крвальо (Nuno Carvalho), Марию Жоао Монтейро (Maria João Monteiro) и Луиса Сардинья (Luís Sardinha).

Наконец, мы хотели бы поблагодарить всех наших коллег, которые оказали нам любезность, передав свои комментарии к первым рукописям книги. В их числе: Феликс Гартнер (Felix Gaertner), Бенуит Гарбинато (Benoit Garbinato) и Мартен ван Стин (Maarten van Steen).

Благодарности ко второму изданию

Работа над вторым изданием книги началась в 2009 году, когда Кристиан Качин (Christian Cachin), работавший тогда в IBM Research (EPFL), находился в ежегодном отпуске. Мы благодарны EPFL и IBM Research за поддержку.

Мы снова хотим сказать спасибо студентам из EPFL и Лиссабонского университета, которым довелось работать с книгой, за предложенные улучшения к первому изданию. Мы также выражаем свою благодарность студентам из высшего технологического института (Instituto Superior Técnico, IST) Лиссабонского технического университета (Universidade Técnica de Lisboa), федерального технологического института в Цюрихе (ETH Zurich) и EPFL, на чей суд были представлены рукописи и дополнительные материалы, включенные во второе издание, за их ценные замечания.

Мы благодарны многим внимательным читателям первого издания и тем, кто прислал свои комментарии к рукописям второго издания, за подмеченные проблемы и предложения по улучшению. В частности, мы благодарны Зинаиде Бененсон (Zinaida Benenson), Элиссон Бессани (Alysson Bessani), Диего Бьюррун (Diego Buiurrun), Филиппе Кристову (Filipe Cristóvão), Дэну Добре (Dan Dobre), Феликсу Фрейлингу (Felix Freiling), Али Годши (Ali Ghodsi), Сейфу Хариди (Seif Haridi), Матусу Харвану (Matus Harvan), Рудигеру Капице (Rüdiger Kapitza), Николу Кнежевичу (Nikola Knežević), Андреасу Кнобелю (Andreas Knobel), Михаю Летиа (Mihai Letia), Томасу Локеру (Thomas Locher), Хейну Мелингу (Hein Meling), Хьюго Миранде (Hugo Miranda), Луису Пина (Luís Pina), Мартину Шобу (Martin Schaub) и Марко Вуколичу (Marko Vukolić).

*Кристиан Качин,
Рашид Гуерру,
Луис Родригес*

Введение

Я полностью использую свои возможности. Помоему, именно к этому стремится любое разумное существо.

(HAL 9000)

В этой главе сначала будет обоснована необходимость абстракций распределенного программирования. Особое внимание при этом мы уделим абстракциям, охватывающим задачи, которые лежат в основе устойчивых форм взаимодействий множественных процессов в распределенной системе и обычно называются *абстракциями соглашений* (agreement abstractions). Затем будет представлена модульная стратегия разработки распределенных программ с использованием упомянутых абстракций и специализированных прикладных программных интерфейсов (Application Programming Interfaces, API).

Также будет дан конкретный пример API для иллюстрации обозначений и схемы вызовов на основе событий, используемых на протяжении всей книги для описания алгоритмов, реализующих абстракции. Обозначения и схемы вызовов очень близки к тем, что можно найти в практических реализациях распределенных алгоритмов.

1.1. Обоснование

Под распределенными вычислениями понимаются алгоритмы для множества процессов, целью которых является организация некоторой формы сотрудничества. Помимо одновременного выполнения, некоторые процессы в распределенной системе могут прекратить работу, например в результате аварии или разрыва соединения, тогда как другие могут продолжать действовать. Само понятие такого *частичного отказа* является одной из характеристик распределенной системы. В действительности это понятие может пригодиться, когда возникнет необходимость обозначить отличие распределенной системы от параллельной системы. Здесь уместно процитировать Лесли Лэмпорта (Leslie Lamport):

«Распределенной называется система, в которой отказ компьютера, о существовании которого вы даже не подозревали, может сделать ваш собственный компьютер бесполезным».

Когда подмножество процессов терпит аварию или теряет связь, перед процессами, которые продолжают действовать или сохраняют связь с большинством, вста-

ет задача синхронизировать свою деятельность согласованным способом. Иными словами, сотрудничество процессов должно быть достаточно устойчивым, чтобы допускать возможность частичных отказов, а иногда и злонамеренных нападений. Это делает распределенные вычисления сложной, но чрезвычайно интересной задачей. Из-за асинхронности процессов, возможности отказов в инфраструктуре взаимодействий и, может быть, даже преднамеренных действий скомпрометированных процессов порой нельзя точно классифицировать отказы; в частности, невозможно отличить отказ процесса от отказа сети, о чем будет рассказываться далее в книге. Хуже того, процесс, находящийся под контролем злоумышленника, может целенаправленно нарушать связь между остальными процессами. Это делает задачу обеспечения согласованного сотрудничества еще более сложной. Задача распределенных вычислений как раз и заключается в создании алгоритма, который будет обеспечивать остающиеся процессы непротиворечивой информацией, благодаря чему они смогут взаимодействовать и решать общие задачи.

Фактически многие современные программы, которые мы используем, являются распределенными программами. Простые повседневные программы, такие как клиенты электронной почты или веб-браузеры, вовлекают распределенные вычисления в некоторой форме. Однако чаще в таких приложениях используется самая простая разновидность распределенных вычислений: вычисления в архитектуре *клиент-сервер*. В архитектуре клиент-сервер существует центральный процесс – *сервер*, – обслуживающий множество удаленных *клиентов*. Клиенты и сервер взаимодействуют, обмениваясь сообщениями и, как правило, следуя форме взаимодействий типа запрос-ответ. Например, чтобы отобразить веб-страницу, браузер посылает запрос веб-серверу и ожидает получить от него ответ с информацией для отображения. Основная проблема, связанная с достижением непротиворечивой формы сотрудничества при ненулевой вероятности частичных отказов, может проявиться даже в такой простой форме взаимодействий. Вернемся к примеру с веб-браузером. Было бы разумным ожидать, что пользователь сможет продолжить обзор Всемирной паутины, если запрошенный веб-сервер потерпит аварию (пользователь автоматически будет подключен к другому веб-серверу). Еще более разумным было бы ожидать, что серверный процесс продолжит передавать информацию другим клиентским процессам, даже если часть из них потерпит аварию или потеряет соединение с сетью.

Проблемы выше уже не тривиальны, когда распределенные вычисления ограничиваются всего двумя взаимодействующими сторонами, такими как клиент и сервер. Однако распределенные вычисления не ограничиваются схемой взаимодействий клиент-сервер. Очень часто во взаимодействиях участвует более двух процессов, деятельность которых требуется синхронизировать для решения общей задачи. Наличие множества процессов еще больше усложняет организацию распределенных вычислений. Иногда мы будем называть *многосторонние* взаимодействия общим случаем. Фактически оба шаблона могут вполне естественно сосуществовать друг с другом. Более того, многие распределенные приложения имеют части, взаимодействующие по схеме клиент-сервер, а другие – по схеме

многосторонних взаимодействий. Часто определение схемы (шаблона) зависит от точки зрения. Например, когда клиент подключается к серверу, чтобы получить услугу, он может не знать, что для оказания запрошенной услуги серверу придется самому послать запрос одному или нескольким вспомогательным серверам. Иногда, чтобы подчеркнуть наличие центрального сервера, используют термин *вычисления типа точка-точка* (peer-to-peer computing).

1.2. Абстракции распределенного программирования

Подобно улыбке, круг абстракций ограничен очень небольшим числом их естественных разновидностей. Определяя свойства, общие для большого или существенного диапазона систем, абстракции помогают отличать основное от второстепенного и избавляют проектировщиков систем и инженеров от необходимости повторно изобретать одни и те же решения для мало отличающихся вариантов одной и той же задачи.

От простого... Рассуждения о распределенных системах должны начинаться с абстрагирования физической системы, лежащей в основе: описания соответствующих элементов максимально абстрактным способом, описания их внутренних свойств и характеристики взаимодействий между ними, в результате чего должна получиться так называемая *модель системы* (system model). В этой книге в основном используются две абстракции, представляющие физическую систему: *процессы* (processes) и *связи* (links).

Процессы в абстракции распределенной программы – это активные сущности, выполняющие вычисления. Процессом может быть компьютер, процессор в компьютере или конкретное ядро в процессоре. В контексте сетевой безопасности процессом может быть также доверенный домен, или административный модуль. Для совместной работы над некоторой общей задачей процессам обычно требуется возможность обмениваться сообщениями посредством некоторой сети. Связи как раз являются абстракцией физических и логических сетей, обеспечивающих обмен информацией между процессами. Различные реализации распределенных систем можно представить этими двумя абстракциями, описав разные свойства процессов и связей, например как эти элементы могут действовать или отказывать при разных условиях.

В главе 2 будет представлено более глубокое обсуждение разных моделей распределенных систем, используемых в этой книге.

...к сложному. Получив модель системы, следующий шаг состоит в том, чтобы понять, как определить абстракции, описывающие повторяющиеся шаблоны взаимодействий в распределенных приложениях. В этой книге нас интересуют абстракции, решающие проблемы надежности взаимодействий в группах процессов, как наиболее важные. Взаимодействие процессов иногда можно смоделировать как задачу *распределенных соглашений* (distributed agreement). Например, процессам может потребоваться выяснить, произошло ли (или не произошло) неко-

торое событие, чтобы договориться о последовательности дальнейших действий (из числа нескольких альтернатив) или о порядке обработки входных данных. Желательно последовательно устанавливать более сложные формы соглашений, от решений до простых задач соглашений. Рассмотрим, к примеру, следующие ситуации.

Чтобы процессы имели возможность обмениваться информацией, они прежде должны договориться об идентификации себя (например, с использованием IP-адресов в Интернете) и общих форматах представления сообщений. Им может также потребоваться договориться о способах обмена сообщениями (например, использовать для взаимодействий надежный поток данных, такой как TCP в Интернете).

После обмена некоторыми сообщениями процессы могут столкнуться с несколькими альтернативными планами действий. Им может потребоваться добиться *консенсуса* в выборе одной из нескольких альтернатив, и каждый из договаривающихся процессов может первоначально иметь собственный план, отличный от планов других процессов.

В некоторых случаях сотрудничающие процессы могут пойти на некоторый шаг, только если все другие процессы согласятся выполнить его. Если это условие не соблюдается, все процессы должны согласиться не выполнять этого шага. Такая форма соглашений чрезвычайно важна для обработки распределенных транзакций, где данная задача известна под названием *атомарное подтверждение* (atomic commitment).

Процессам может потребоваться не только договориться (согласиться) о каких-то действиях, но и о порядке их выполнения. Эта форма соглашения является основой одного из фундаментальных приемов репликации вычислений для отказоустойчивости, который называется задачей *расылки в едином порядке* (total-order broadcast).

Данная книга рассказывает, как преодолеть сложности, кроющиеся в этих задачах, и определить *абстракции*, инкапсулирующие эти задачи. Сложность задач объясняется необходимостью координации действий множества процессов; учитывая, что процессы могут терпеть неудачу или даже проявлять злонамеренное поведение, создать такие абстракции не всегда просто. Ниже мы обоснуем уместность некоторых абстракций, рассматриваемых в книге. Мы будем различать случаи, когда абстракции вытекают из естественного распределения приложения, с одной стороны, и когда абстракции являются следствием технического решения о необходимости распределения – с другой.

1.2.1. Естественно распределенные системы

Приложения, требующие совместного использования или распределения информации между несколькими сотрудничающими процессами, являются благодатной почвой для проблем, решение которых требует применения абстракций распределенного программирования. Примерами таких приложений могут служить механизмы распределения информации, многопользовательские коллективные систе-

мы, распределенные разделяемые пространства, системы управления процессами (технологическими), коллективные редакторы, распределенные базы данных и распределенные системы хранения.

Распределение информации. В распределенных системах, опирающихся на распространение информации, процессы могут играть одну из следующих ролей: источники (производители) информации, которых также называют *издателями* (publishers), и потребители информации, которых также называют *подписчиками* (subscribers). Парадигма взаимодействий таких процессов часто называется *публикация–подписка* (publish–subscribe).

Издатели производят информацию в форме извещений. Подписчики регистрируют свою заинтересованность в получении определенных извещений. Существуют разные варианты парадигмы публикация–подписка, по-разному удовлетворяющие интересы подписчиков, в том числе на основе каналов, на основе тем, на основе содержимого или на основе типа подписки. Независимо от метода подписки весьма вероятно, что в получении одних и тех же извещений будет заинтересовано несколько подписчиков. В этом случае мы обычно заинтересованы в доставке всем подписчикам одной и той же информации – одного и того же множества сообщений. Иначе система будет несправедлива в оказании услуг, так как одни подписчики получат больше информации, чем другие.

Если только это не является свойством надежности, обеспечиваемым инфраструктурой по умолчанию (что встречается крайне редко), издатель и подписчики должны договориться о том, какие сообщения доставлять, а какие нет. Например, когда распространяется аудиопоток, процессы обычно заинтересованы в получении большей части информации и терпимо относятся к потере ограниченного числа сообщений, особенно если это позволит системе добиться большей пропускной способности. Соответствующую абстракцию обычно называют *негарантированной рассылкой*, или *рассылкой с максимально доступным качеством* (best-effort broadcast).

Для распространения информации, производимой фондовой биржей, может потребоваться более надежный способ доставки, который называют *надежной рассылкой* (reliable broadcast), так как желательно, чтобы все активные процессы приняли одну и ту же информацию. От инфраструктуры рассылки биржевой информации можно было бы также потребовать рассылки информации в определенном порядке. В некоторых приложениях типа издатель–подписчик источники и потребители информации взаимодействуют опосредованно, через промежуточных брокеров. В таких случаях абстракции соглашений (agreement abstractions) могут использоваться для организации взаимодействий между несколькими брокерами.

Управление процессами. Приложения управления (технологическими) процессами – это такие приложения, в которых несколько программных процессов управляют технологическими процессами. Программные процессы могут, например, управлять параметрами движения авиалайнера или поезда, температурой в ядерном реакторе или автоматизировать работу конвейера автомобильного завода.

Обычно каждый программный процесс подключен к некоторому датчику. Процессам может потребоваться, например, обмениваться значениями, прочитанными

ми с датчиков, и выводить некоторое общее значение, допустим координаты самолета, на экран дисплея пилота, даже при том что из-за погрешностей или отказов локальных датчиков они могут получать немного разные значения. Такое сотрудничество должно осуществляться всегда, несмотря на то что некоторые датчики (или соответствующие им программные процессы) вышли из строя или не получают никаких данных. Подобного рода сотрудничество можно упростить, если все процессы договорятся о едином множестве входных данных для управляющего алгоритма, – требование, зафиксированное абстракцией *консенсуса*.

Коллективные системы. Пользователи на разных узлах сети могут совместно работать над созданием общего программного обеспечения или документа, или просто участвовать в распределенном диалоге, например в чате или виртуальной конференции. Поддержку совместного труда в таких системах очень эффективно обеспечивает абстракция *общего рабочего пространства*. Такого рода абстракции общей распределенной памяти обычно предоставляют пользователям операции *чтения/записи* для эффективного обмена информацией и ее хранения. В простейшей форме общее рабочее пространство можно представить как единый виртуальный неструктурированный объект-хранилище. Более сложные реализации общих рабочих пространств могут добавлять структурные компоненты для создания пространств отдельных пользователей. Примеров таких реализаций можно привести множество: от разнообразных Вики (Wiki) до сложных многопользовательских распределенных файловых систем. Для поддержания целостности таких совместных рабочих пространств процессы должны договариваться об относительном порядке операций *чтения* и *записи* в пространстве.

Распределенные базы данных. Базы данных составляют еще один класс приложений, где могут пригодиться абстракции соглашений, чтобы гарантировать получение всеми диспетчерами транзакций непротиворечивого представления выполняющихся транзакций и помочь принять решение об упорядочении этих транзакций.

Кроме того, такие транзакции можно использовать для координации действий диспетчеров транзакций, когда решается вопрос о прямых результатах транзакций. То есть серверы баз данных, выполняющие данную распределенную транзакцию, должны координировать свои действия и решать, подтвердить транзакцию или прервать ее. Они могут решить прервать транзакцию, если какой-либо из серверов баз данных обнаружит угрозу целостности данных, несогласованное управление параллельным выполнением, дисковую ошибку или просто аварию какого-то другого сервера баз данных. Как уже отмечалось выше, такое распределенное взаимодействие обеспечивает абстракция *атомарного подтверждения* (atomic commit).

Распределенные системы хранения. Системы хранения огромной емкости распределяют данные между несколькими узлами хранения, каждый из которых является частью общего пространства хранилища. Доступ к хранимым данным обычно сопряжен с обращениями к нескольким узлам, потому что даже единственный элемент данных может оказаться разбросанным по нескольким узлам. Элемент данных может подвергаться сложным преобразованиям программными

компонентами определения и/или исправления ошибок, которые сами могут быть разбросаны по множеству узлов для обеспечения стойкости системы хранения против потери или повреждения некоторых узлов. Такие системы распределяют данные не только из-за ограниченной емкости отдельных узлов, но также для увеличения отказоустойчивости всей системы и снижения нагрузки на отдельные узлы.

Концептуально распределенная система хранения напоминает абстракцию общей памяти, доступ к которой осуществляется посредством операций *чтения* и *записи*, подобно общему рабочему пространству, как упоминалось выше. Но, так как распределенность используется также с целью увеличения общей надежности, распределенные системы хранения сочетают в себе черты естественно распределенных систем с чертами искусственно распределенных систем, о которых рассказывается в следующем разделе.

1.2.2. Искусственно распределенные системы

Часто для реализации приложений, не имеющих естественных предпосылок к распределению, также используются сложные абстракции распределенного программирования. Иногда это является следствием технического решения с целью удовлетворить такие требования, как *отказоустойчивость*, *балансировка нагрузки* или обеспечение *быстрого доступа* при большом числе клиентов.

Проиллюстрируем эту идею на примере *репликации автоматов* (state-machine replication) – мощного средства достижения отказоустойчивости в распределенных системах. Если говорить кратко, репликация автоматов используется для создания служб с высокой доступностью, реализованных в виде нескольких действующих копий службы на разных машинах, которые, как предполагается, не могут выйти из строя одновременно. За счет этого обеспечивается непрерывность обслуживания, даже в случае выхода из строя нескольких машин. Никакого специального аппаратного обеспечения для этого не требуется: отказоустойчивость репликации достигается исключительно программными средствами. В действительности прием репликации можно также использовать внутри информационной системы для ускорения доступа к данным, помещая действующие копии службы как можно ближе к процессам, которые наиболее вероятно будут выполнять запросы. Для служб, доступных для атак через Интернет, например, тот же подход может повысить устойчивость к злонамеренным вторжениям, компрометирующим ограниченное число узлов, участвующих в репликации.

Для большей эффективности репликации разные копии должны поддерживать состояние непротиворечивости. Если состояния копий могут отклоняться произвольно, бессмысленно говорить о репликации. Иллюзия *единой* службы с высокой степенью доступности быстро развалится, и взамен ее появится несколько распределенных служб, каждая из которых может терпеть неудачу независимо от других. Если перечень состояний копий четко определен, один из самых простых способов гарантировать целостность и непротиворечивость всех узлов заключается в передаче им всем одного и того же набора запросов в одном и том же порядке.