

$$2^{2^5} + 1 = 641 \cdot 6700417$$

$$\text{avg}(x, y) = (x \& y) + ((x \oplus y) \gg 1)$$



$$2^{2^6} + 1 = 274177 \cdot 67280421310721$$

$$x - y = x + \bar{y} + 1$$



$$\lfloor a \rfloor + \lfloor b \rfloor \leq \lfloor a + b \rfloor \leq \lfloor a \rfloor + \lfloor b \rfloor + 1 \quad \text{pop}(x) = \sum_{i=0}^{31} (x \ll i) \text{rot}$$

George Boole
1815 - 1864

$$\lfloor \sqrt{111111111} \rfloor = 1111$$

$$(x \neq 0) = (x | -x) \gg 31$$

$$\text{mux}(x, y, m) = ((x \oplus y) \& m) \oplus y$$

$$A(n, d) = A(n - 1, d - 1), d \text{ even}$$

$$-\bar{x} = x + 1$$

АЛГОРИТМИЧЕСКИЕ ТРЮКИ для программистов

$$\frac{1}{3} = 0.01010101\dots$$

Второе издание

$$1111^2 = 1110001$$

$$n = -2^{31} b_{31} + 2^{30} b_{30} + 2^{29} b_{29} + \dots + 2^0 b_0$$

$$\lfloor x \rfloor = - \lfloor -x \rfloor \quad f(x, y, z) = g(x, y) \oplus zh(x, y)$$

Num factors of 2 in $x =$
 $\log_2(x \& (-x)), x \neq 0$

$$\text{rjust}(x) = x \gg (x \& -x), x \neq 0$$

$$p_n = 1 + \sum_{m=1}^{2^n} \left[\frac{p_m}{\sqrt{m}} \left[\sum_{x=1}^m \left[\cos^2 \pi \frac{(x-1)! + 1}{x} \right] \right] \right]^{1/m}$$

$$x \oplus y = (x | y) - (x \& y)$$

$$x + y = (x | y) + (x \& y)$$

Генри Уоррен-мл.

Hacker's Delight

SECOND EDITION

HENRY S. WARREN, JR.

♣Addison-Wesley

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

АЛГОРИТМИЧЕСКИЕ ТРЮКИ для программистов

Второе издание

АЛГОРИТМИЧЕСКИЕ ТРЮКИ для программистов

Второе издание

ГЕНРИ УОРРЕН-МЛ.



Москва • СанктПетербург • Киев
2014

ББК (Ж/О)32.973.26-018.2.75

У64

УДК 681.3.06

Издательский дом "Вильямс"

Зав. редакцией *С.Н. Тригуб*

Перевод с английского и редакция канд. техн. наук *И.В. Красикова*

По общим вопросам обращайтесь в Издательский дом "Вильямс" по адресу:
info@williamspublishing.com, <http://www.williamspublishing.com>

Уоррен, Генри С.

У64 Алгоритмические трюки для программистов, 2-е изд. : Пер. с англ. — М. : ООО "И.Д. Вильямс", 2014. — 512 с. : ил. — Парал. тит. англ.

ISBN 978-5-8459-1838-3 (рус.)

ББК (Ж/О) 32.973.26-018.2.75

Все названия программных продуктов являются зарегистрированными торговыми марками соответствующих фирм.

Никакая часть настоящего издания ни в каких целях не может быть воспроизведена в какой бы то ни было форме и какими бы то ни было средствами, будь то электронные или механические, включая фотокопирование и запись на магнитный носитель, если на это нет письменного разрешения издательства Addison-Wesley Publishing Company, Inc.

Authorized translation from the English language edition published by Addison-Wesley Publishing Company, Inc © 2013.

All rights reserved. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise.

Russian language edition is published by Williams Publishing House according to the Agreement with R&I Enterprises International, Copyright © 2014.

Научно-популярное издание

Генри С. Уоррен

Алгоритмические трюки для программистов

2-е издание

Литературный редактор *Л.Н. Красножон*

Верстка *О.В. Мишутина*

Художественный редактор *Е.П. Дынный*

Корректор *Л.А. Гордиенко*

Подписано в печать 15.10.2013. Формат 70х100/16

Гарнитура Times. Печать офсетная

Усл. печ. л. 32,0. Уч.-изд. л. 41,2

Тираж 1500 экз. Заказ № 3835

Первая Академическая типография "Наука"

190034, Санкт-Петербург, 9-я линия, 12/28

ООО "И. Д. Вильямс", 127055, г. Москва, ул. Лесная, д. 43, стр. 1

ISBN 978-5-8459-1838-3 (рус.)

ISBN 978-0-321-84268-8 (англ.)

© Издательский дом "Вильямс", 2014

© Pearson Education, Inc., 2013

ОГЛАВЛЕНИЕ

ПРЕДИСЛОВИЕ	14
ВВЕДЕНИЕ	16
ГЛАВА 1. ВВЕДЕНИЕ	19
ГЛАВА 2. ОСНОВЫ	31
ГЛАВА 3. ОКРУГЛЕНИЕ К СТЕПЕНИ 2	81
ГЛАВА 4. АРИФМЕТИЧЕСКИЕ ГРАНИЦЫ	89
ГЛАВА 5. ПОДСЧЕТ БИТОВ	103
ГЛАВА 6. ПОИСК В СЛОВЕ	141
ГЛАВА 7. ПЕРЕСТАНОВКА БИТОВ И БАЙТОВ	155
ГЛАВА 8. УМНОЖЕНИЕ	197
ГЛАВА 9. ЦЕЛОЧИСЛЕННОЕ ДЕЛЕНИЕ	207
ГЛАВА 10. ЦЕЛОЕ ДЕЛЕНИЕ НА КОНСТАНТЫ	231
ГЛАВА 11. НЕКОТОРЫЕ ЭЛЕМЕНТАРНЫЕ ФУНКЦИИ	305
ГЛАВА 12. СИСТЕМЫ СЧИСЛЕНИЯ С НЕОБЫЧНЫМИ ОСНОВАНИЯМИ	325
ГЛАВА 13. КОД ГРЕЯ	337
ГЛАВА 14. ЦИКЛИЧЕСКИЙ ИЗБЫТОЧНЫЙ КОД	345
ГЛАВА 15. КОДЫ С КОРРЕКЦИЕЙ ОШИБОК	357
ГЛАВА 16. КРИВАЯ ГИЛЬБЕРТА	381
ГЛАВА 17. ЧИСЛА С ПЛАВАЮЩЕЙ ТОЧКОЙ	401
ГЛАВА 18. ФОРМУЛЫ ДЛЯ ПРОСТЫХ ЧИСЕЛ	419
ОТВЕТЫ К УПРАЖНЕНИЯМ	433
ПРИЛОЖЕНИЕ А. АРИФМЕТИЧЕСКИЕ ТАБЛИЦЫ ДЛЯ 4-БИТОВОЙ МАШИНЫ	483
ПРИЛОЖЕНИЕ Б. МЕТОД НЬЮТОНА	487
ПРИЛОЖЕНИЕ В. ГРАФИКИ ДИСКРЕТНЫХ ФУНКЦИЙ	489
СПИСОК ЛИТЕРАТУРЫ	501
ПРЕДМЕТНЫЙ УКАЗАТЕЛЬ	506

СОДЕРЖАНИЕ

Предисловие	14
Введение	16
Благодарности	17
ГЛАВА 1. ВВЕДЕНИЕ	19
1.1. Система обозначений	19
1.2. Система команд и модель оценки времени выполнения команд	23
Время выполнения	28
Упражнения	30
ГЛАВА 2. ОСНОВЫ	31
2.1. Манипуляции младшими битами	31
Расширенные законы де Моргана	33
Проверка выполнимости справа налево	33
Новое применение	35
2.2. Сложение и логические операции	36
2.3. Неравенства с логическими и арифметическими выражениями	38
2.4. Абсолютное значение	39
2.5. Среднее двух целых	39
2.6. Распространение знака	40
2.7. Знаковый сдвиг вправо на основе беззнакового сдвига	40
2.8. Функция sign	41
2.9. Трехзначная функция сравнения	42
2.10. Перенос знака	42
2.11. Декодирование поля "0 означает 2^{*n} "	43
2.12. Предикаты сравнения	43
Команды сравнения и бит переноса	47
Вычисление отношений	48
2.13. Обнаружение переполнения	49
Знаковое сложение и вычитание	49
Установка переполнения при знаковом сложении и вычитании	51
Беззнаковое сложение/вычитание	52
Умножение	53
Деление	55
2.14. Флаги условий после сложения, вычитания и умножения	58
2.15. Циклический сдвиг	59
2.16. Сложение/вычитание двойных слов	60
2.17. Сдвиг двойного слова	61
2.18. Сложение, вычитание и абсолютное значение многобайтовых величин	62
2.19. Функции doz, max, min	63
2.20. Обмен содержимого регистров	68
Обмен соответствующих полей регистров	68

Обмен двух полей одного регистра	69
Условный обмен	70
2.21. Выбор среди двух или более значений	70
2.22. Формула булева разложения	73
2.23. Реализация команд для всех 16 бинарных булевых операций	75
Исторические примечания	77
Упражнения	77
ГЛАВА 3. ОКРУГЛЕНИЕ К СТЕПЕНИ 2	81
3.1. Округление к кратному степени 2	81
3.2. Округление к ближайшей степени 2	82
Округление в меньшую сторону	83
Округление в большую сторону	84
3.3. Проверка пересечения границы степени 2	85
Упражнения	87
ГЛАВА 4. АРИФМЕТИЧЕСКИЕ ГРАНИЦЫ	89
4.1. Проверка границ целых чисел	89
4.2. Определение границ суммы и разности	92
Знаковые числа	94
4.3. Определение границ логических выражений	96
Знаковые границы	100
Упражнения	101
ГЛАВА 5. ПОДСЧЕТ БИТОВ	103
5.1. Подсчет единичных битов	103
Сумма и разность количества единичных битов в двух словах	110
Сравнение степени заполнения двух слов	111
Подсчет единичных битов в массиве	111
Применение	117
5.2. Четность	118
Вычисление четности слова	119
Добавление бита четности к 7-битовой величине	121
Применение	121
5.3. Подсчет ведущих нулевых битов	122
Методы с плавающей точкой	127
Сравнение количества ведущих нулевых битов двух слов	129
Связь с логарифмом	129
Применения	130
5.4. Подсчет завершающих нулевых битов	130
Применение	137
Упражнения	140
ГЛАВА 6. ПОИСК В СЛОВЕ	141
6.1. Поиск первого нулевого байта	141
Некоторые простые обобщения	145
Поиск значения из заданного диапазона	146
6.2. Поиск строки единичных битов заданной длины	147

6.3.	Поиск наидлиннейшей строки единичных битов	150
6.4.	Поиск кратчайшей строки единичных битов	152
	Упражнения	153
ГЛАВА 7. ПЕРЕСТАНОВКА БИТОВ И БАЙТОВ		155
7.1.	Реверс битов и байтов	155
	Обобщенный реверс битов	160
	Современные методы реверса битов	161
	Увеличение обращенного целого числа	163
7.2.	Перемешивание битов	165
7.3.	Транспонирование битовой матрицы	167
	Транспонирование битовой матрицы размером 32×32	171
7.4.	Сжатие, или Обобщенное извлечение	176
	Использование команд вставки и извлечения	181
	Сжатие влево	182
7.5.	Расширение, или Обобщенная вставка	182
7.6.	Аппаратные алгоритмы сжатия и расширения	183
	Сжатие	183
	Расширение	185
7.7.	Обобщенные перестановки	187
7.8.	Перегруппировки и преобразования индексов	191
7.9.	Алгоритм LRU	192
	Упражнения	195
ГЛАВА 8. УМНОЖЕНИЕ		197
8.1.	Умножение больших чисел	197
8.2.	Старшее слово 64-битового произведения	200
8.3.	Преобразование знакового и беззнакового	
	произведений одно в другое	200
	Беззнаковое произведение из знакового	201
8.4.	Умножения на константы	201
	Упражнения	205
ГЛАВА 9. ЦЕЛОЧИСЛЕННОЕ ДЕЛЕНИЕ		207
9.1.	Предварительные сведения	207
9.2.	Деление больших чисел	210
	Знаковое деление больших чисел	215
9.3.	Беззнаковое короткое деление на основе знакового	215
	Использование знакового длинного деления	215
	Использование знакового короткого деления	216
9.4.	Беззнаковое длинное деление	219
	Аппаратный алгоритм сдвига и вычитания	219
	Использование короткого деления	222
9.5.	Деление двойных слов из длинного деления	224
	Беззнаковое деление двойных слов	225
	Знаковое деление двойных слов	228
	Упражнения	229

ГЛАВА 10. ЦЕЛОЕ ДЕЛЕНИЕ НА КОНСТАНТЫ	231
10.1. Знаковое деление на известную степень 2	231
10.2. Знаковый остаток от деления на известную степень 2	232
10.3. Знаковое деление и вычисление остатка для других случаев	233
Деление на 3	233
Деление на 5	235
Деление на 7	235
10.4. Знаковое деление на делитель, не меньший 2	236
Алгоритм	238
Доказательство пригодности алгоритма	239
Доказательство корректности произведения	240
10.5. Знаковое деление на делитель, не превышающий -2	243
Для каких делителей $m(-d) \neq -m(d)$?	245
10.6. Встраивание в компилятор	246
10.7. Дополнительные вопросы	249
Единственность	249
Делители с лучшими программами	250
10.8. Беззнаковое деление	253
Беззнаковое деление на 3	253
Беззнаковое деление на 7	254
10.9. Беззнаковое деление на делитель, не меньший 1	256
Беззнаковый алгоритм	257
Доказательство пригодности беззнакового алгоритма	257
Доказательство корректности произведения в случае беззнакового деления	258
10.10. Встраивание в компилятор при беззнаковом делении	258
10.11. Дополнительные вопросы (беззнаковое деление)	260
Делители с лучшими программами (беззнаковое деление)	260
Использование знакового деления вместо беззнакового и наоборот	261
Более простой беззнаковый алгоритм	262
10.12. Применение к модульному делению и делению с округлением к меньшему значению	263
10.13. Другие похожие методы	263
10.14. Некоторые магические числа	265
10.15. Простой код на языке программирования Python	266
10.16. Точное деление на константу	266
Вычисление мультипликативного обратного по алгоритму Евклида	268
Вычисление мультипликативного обратного по методу Ньютона	271
Некоторые мультипликативные обратные	273
10.17. Проверка нулевого остатка при делении на константу	274
Беззнаковое деление	275
Знаковое деление, делитель ≥ 2	276
10.18. Методы, не использующие команды умножения со старшим словом	278
Беззнаковое деление	278
Знаковое деление	286
10.19. Получение остатка суммированием цифр	288
Беззнаковый остаток	289
Знаковый остаток	293

10.20. Получение остатка путем умножения и сдвига вправо	295
Беззнаковый остаток	295
Знаковый остаток	299
10.21. Преобразование в точное деление	301
10.22. Проверка времени выполнения	302
10.23. Аппаратная схема для деления на 3	303
Упражнения	304
ГЛАВА 11. НЕКОТОРЫЕ ЭЛЕМЕНТАРНЫЕ ФУНКЦИИ	305
11.1. Целочисленный квадратный корень	305
Метод Ньютона	305
Бинарный поиск	309
Аппаратный алгоритм	310
11.2. Целочисленный кубический корень	313
11.3. Целочисленное возведение в степень	314
Вычисление x^p бинарным разложением p	314
2^n в Fortran	315
11.4. Целочисленный логарифм	316
Целочисленный логарифм по основанию 2	317
Целочисленный логарифм по основанию 10	317
Упражнения	323
ГЛАВА 12. СИСТЕМЫ СЧИСЛЕНИЯ С НЕОБЫЧНЫМИ ОСНОВАНИЯМИ	325
12.1. Основание -2	325
12.2. Основание $-1+i$	332
12.3. Другие системы счисления	335
12.4. Какое основание наиболее эффективно	335
Упражнения	336
ГЛАВА 13. КОД ГРЕЯ	337
13.1. Код Грея	337
13.2. Увеличение чисел кода Грея	339
13.3. Отрицательно-двоичный код Грея	341
13.4. Краткая история и применение	341
Упражнения	343
ГЛАВА 14. ЦИКЛИЧЕСКИЙ ИЗБЫТОЧНЫЙ КОД	345
14.1. Введение	345
Основы	345
14.1. Теория	347
14.3. Практика	349
Аппаратное обеспечение	351
Программная реализация	353
Упражнения	356
ГЛАВА 15. КОДЫ С КОРРЕКЦИЕЙ ОШИБОК	357
15.1. Введение	357
15.2. Код Хэмминга	358
Код SEC-DED	360

Минимально необходимое количество проверочных битов	362
Заключительные замечания	362
15.3. Программная реализация SEC-DED для 32 информационных битов	364
15.4. Общее рассмотрение задачи коррекции ошибок	369
Расстояние Хэмминга	370
Основная задача теории кодирования	372
Сферы	374
Упражнения	379
Глава 16. КРИВАЯ ГИЛЬБЕРТА	381
16.1. Рекурсивный алгоритм построения кривой Гильберта	383
16.2. Преобразование расстояния вдоль кривой Гильберта в координаты	385
16.3. Преобразование координат в расстояние вдоль кривой Гильберта	393
16.4. Увеличение координат кривой Гильберта	395
16.5. Нерекурсивный алгоритм генерации кривой Гильберта	398
16.6. Другие кривые, заполняющие пространство	398
16.7. Применения	399
Упражнения	400
Глава 17. ЧИСЛА С ПЛАВАЮЩЕЙ ТОЧКОЙ	401
17.1. Формат IEEE	401
17.2. Преобразование чисел с плавающей точкой в целые и обратно	404
17.3. Сравнение чисел с плавающей точкой с использованием целых операций	408
17.4. Программа приближенного вычисления обратного к квадратному корню	410
17.5. Распределение ведущих цифр	413
17.6. Таблица различных значений	415
Упражнения	417
Глава 18. ФОРМУЛЫ ДЛЯ ПРОСТЫХ ЧИСЕЛ	419
18.1. Введение	419
18.2. Формулы Вилланса	421
Первая формула	421
Вторая формула	422
Третья формула	423
Четвертая формула	424
18.3. Формула Вормелла	424
18.4. Формулы для других сложных функций	425
Упражнения	431
ОТВЕТЫ К УПРАЖНЕНИЯМ	433
Глава 1. Введение	433
Глава 2. Основы	435
Глава 3. Округление к степени 2	444
Глава 4. Арифметические границы	444
Глава 5. Подсчет битов	445
Глава 6. Поиск в слове	446
Глава 7. Перестановка битов и байтов	451
Глава 8. Умножение	453

Глава 9. Целочисленное деление	456
Глава 10. Целое деление на константы	459
Глава 11. Некоторые элементарные функции	462
Глава 12. Системы счисления с необычными основаниями	464
Глава 13. Код Грея	468
Глава 14. Циклический избыточный код	470
Глава 15. Коды с коррекцией ошибок	470
Глава 16. Кривая Гильберта	475
Глава 17. Числа с плавающей точкой	475
Глава 18. Формулы для простых чисел	478
ПРИЛОЖЕНИЕ А. АРИФМЕТИЧЕСКИЕ ТАБЛИЦЫ ДЛЯ 4-БИТОВОЙ МАШИНЫ	483
ПРИЛОЖЕНИЕ Б. МЕТОД НЬЮТОНА	487
ПРИЛОЖЕНИЕ В. ГРАФИКИ ДИСКРЕТНЫХ ФУНКЦИЙ	489
В.1. Графики логических операций над целыми числами	489
В.2. Графики для сложения, вычитания и умножения	491
В.3. Графики функций, включающих деление	493
В.4. Графики функций сжатия, обобщенного упорядочения и циклического сдвига влево	494
В.5. Графики некоторых унарных функций	496
СПИСОК ЛИТЕРАТУРЫ	501
ПРЕДМЕТНЫЙ УКАЗАТЕЛЬ	506

*Джозефу У. Гауду (Joseph W. Gaud), моему школьному учителю алгебры,
который зажг в мне пламя восхищения математикой.*

ПРЕДИСЛОВИЕ

Около 30 лет назад, во время первой летней практики, мне посчастливилось участвовать в проекте MAC в Массачусетском технологическом институте. Тогда я впервые смог поработать на компьютере DEC PDP-10, который в то время предоставлял более широкие возможности при программировании на ассемблере, чем любой другой компьютер. Этот процессор имел большой набор команд для выполнения операций с отдельными битами, их тестирования, маскирования и для работы с полями и целыми числами. Несмотря на то что PDP-10 давно снят с производства, до сих пор есть энтузиасты, работающие на этих машинах либо на их программных эмуляторах. Некоторые из них даже пишут новые приложения; по крайней мере, в настоящее время есть как минимум один веб-узел, поддерживаемый эмулированным PDP-10. (Не смейтесь — поддержка такого узла ни в чем не уступает поддержке антикварного автомобиля “на ходу”.)

Тем же летом 1972 года я увлекся чтением серии отчетов по новым исследованиям, публиковавшихся в институте под общим названием НАКМЕМ, — довольно странный и эклектичный сборник разнообразных технических мелочей¹. Тематика статей могла быть любой: от электронных схем до теории чисел; но больше всего меня заинтриговал небольшой раздел, в котором публиковались программистские хитрости. Практически в каждой статье этого раздела содержалось описание некоторой (зачастую необычной) операции над целыми числами или битовыми строками (например, подсчет единичных битов в слове), которую можно было легко реализовать в виде длинной последовательности машинных команд либо цикла, а затем обсуждалось, как можно сделать то же самое, используя только четыре, три или две тщательно подобранные команды, взаимодействие между которыми становилось очевидным только после соответствующих объяснений или самостоятельных исследований. Я изучал эти маленькие программные самородки с таким же наслаждением, с каким другие пьют пиво или едят сладости. Я просто не мог остановиться. Каждая такая программа оказывалась для меня находкой, в каждой из них была интеллектуальная глубина, элегантность и даже своеобразная поэзия.

“Наверняка, — думал я тогда, — таких программ должно быть намного больше. Должна же где-то быть книга, посвященная таким штукам”.

Книга, которую вы держите в руках, меня просто потрясла. Автор систематически, на протяжении многих лет собирал такие программные перлы и в конце концов свел их воедино, тематически организовал и снабдил каждый четким и подробным описанием. Многие из них можно записать в машинных командах, но книга полезна не только для

¹ Почему “НАКМЕМ”? Сокращение от “hacks memo”; одно 36-битовое PDP-10-слово может содержать шесть 6-битовых символов, поэтому многие имена PDP-10, с которыми тогда работали программисты, были ограничены шестью символами. Аббревиатуры из шести символов широко использовались для сокращенных названий, иногда ими пользовались просто для удобства. Поэтому название сборника “НАКМЕМ” имело ясный смысл для специалистов, по крайней мере в то время.

тех, кто пишет программы на ассемблере. Главная тема книги — рассмотреть базовые структурные отношения среди целых чисел и битовых строк и эффективные приемы реализации операций над ними. Рассмотренные в книге методы так же полезны и практичны при программировании на языках высокого уровня, например С или Java, как и на языке ассемблера.

Во многих книгах по алгоритмам и структурам данных описаны сложные методы сортировки и поиска, поддержания хеш-таблиц и двоичных деревьев, а также работы с записями и указателями. Здесь же рассматривается, что можно сделать с очень крошечной частью данных — битами и массивами битов. Поразительно, сколько всего можно сделать, используя только операции двоичного сложения и вычитания вместе с некоторыми поразрядными операциями. Операции с переносом позволяют одному биту воздействовать на все биты, находящиеся слева от него, что делает сложение особенно мощной операцией при работе с данными способами, не получившими широкого распространения.

Сейчас у вас в руках книга именно о таких методах. Если вы работаете над оптимизирующим компилятором или создаете высокопроизводительный код, вам обязательно нужно прочесть эту книгу. Возможно, в своей повседневной работе вы нечасто будете использовать изложенные здесь приемы, но, если вам потребуется организовать цикл по всем битам слова или ускорить работу с отдельными битами во внутреннем цикле либо же вы окажетесь в ситуации, когда код получается слишком сложным, загляните в эту книгу — и, я уверен, вы найдете в ней помощь в решении стоящей перед вами проблемы. В любом случае чтение этой книги доставит вам огромное удовольствие.

Гай Л. Стил-мл. (Guy L. Steele, Jr.)
Берлингтон, Массачусетс
Апрель 2002

ВВЕДЕНИЕ

Caveat emptor²: стоимость сопровождения программного обеспечения пропорциональна квадрату творческих способностей программиста.

(Первый закон творческого программирования,
Роберт Д. Блис (Robert D. Bliss), 1992)

Перед вами сборник программных приемов, которые я собирал много лет. Большинство из них работают только на компьютерах, на которых целые числа представлены в дополнительном до 2 коде. Хотя в данной книге речь идет о 32-разрядных машинах с соответствующей длиной регистра, большую часть представленных здесь алгоритмов легко перенести на машины с регистрами других размеров.

В этой книге не рассматриваются сложные вопросы наподобие методов сортировки или оптимизации компилируемого кода. Основное внимание здесь уделяется приемам работы с отдельными машинными словами или командами, например подсчету количества единичных битов в заданном слове. В подобных приемах часто используется смесь арифметических и логических команд.

Предполагается, что прерывания, связанные с переполнением целых чисел, замаскированы и произойти не могут. Программы на C, Fortran и даже Java работают именно в таком окружении, но программистам на Pascal и ADA следует быть осторожными!

Представление материала в книге — неформальное. Доказательства приводятся только в том случае, если алгоритм неочевиден, а иногда не приводятся вообще. Все методы используют компьютерную арифметику, функции типа “пол”, комбинации арифметических и логических операций и тому подобные средства, а доказательства теорем в этой предметной области часто сложны и громоздки.

Чтобы свести к минимуму количество типографских ошибок и опечаток, многие алгоритмы реализованы на языке программирования высокого уровня, в качестве которого используется C. Это обусловлено его распространенностью и тем, что он позволяет непосредственно комбинировать операции с целыми числами и битовыми строками; кроме того, компиляторы языка C генерируют объектный код высокого качества.

Ряд алгоритмов написан на машинном языке. В книге применяется трехадресный формат команд, главным образом для повышения удобочитаемости. Использован язык ассемблера для некой абстрактной машины, которая является представителем современных RISC-компьютеров.

² Caveat emptor (лат.) — да будет осмотрителем покупатель.

Отсутствие ветвлений в программе всячески приветствуется. Это связано с тем, что на многих машинах наличие ветвлений замедляет выборку команд и блокирует их параллельное выполнение. Кроме того, наличие ветвлений может препятствовать выполнению оптимизаций компилятором — таких как расписание выполнения инструкций, распределение регистров и других. Оптимизирующий компилятор более эффективно работает с несколькими большими блоками кода, чем с множеством небольших.

Крайне желательно использовать малые величины при непосредственном задании операнда, сравнение с нулем (а не с некоторыми другими числами) и параллелизм на уровне команд. Хотя программу часто можно значительно сократить за счет использования поиска в таблице (расположенной в памяти), этот метод не слишком распространен. Дело в том, что по сравнению с выполнением арифметических команд загрузка данных из памяти занимает намного больше времени, а методы поиска в таблице зачастую не представляют большого интереса (хотя и весьма практичны).

Напоследок мне хотелось бы напомнить исходное значение слова “хакер”³. Хакер — это страстный любитель компьютеров: он создает что-то новое, переделывает или совершенствует то, что уже есть. Хакер очень хорошо разбирается в том, что делает, хотя часто не является профессиональным программистом или разработчиком. Обычно он пишет программы не ради выгоды, а ради собственного удовольствия. Такая программа может оказаться полезной, а может остаться всего лишь игрой интеллекта. Например, хакер может написать программу, которая во время выполнения выводит точную копию себя самой⁴. Таких людей называют хакерами, и эта книга написана именно для них, а не для тех, кто хочет получить совет о том, как взломать что-либо в компьютере.

Благодарности

Прежде всего я хочу поблагодарить Брюса Шрайвера (Bruce Shriver) и Денниса Эллисона (Dennis Allison), оказавших мне помощь в опубликовании книги. Я признателен коллегам из IBM, многие из которых упомянуты в списке литературы. Особой благодарности достоин Мартин Хопкинс (Martin E. Hopkins) из IBM, которого по праву можно назвать “Мистер Компьютер”. Этот человек неудержим в своем стремлении подсчитать в программе каждый такт, и многие его идеи нашли отражение в этой книге. Благодарю также обозревателей из Addison-Wesley, которые значительно улучшили мою книгу. Со многими из них я не знаком, но не могу не упомянуть выдающийся 50-страничный обзор одного из них, Гая Л. Стила-мл. (Guy L. Steele, Jr.), в котором, в частности, затронуты вопросы перемешивания битов, обоб-

³ В последнее время хакерами часто называют тех, кто получает несанкционированный доступ к банковским системам, взламывает веб-узлы и ведет прочую разрушительную деятельность либо ради получения денег, либо для демонстрации всем своей “крутости”. К настоящим хакерам таковы не имеют никакого отношения. — *Примеч. пер.*

⁴ Вот пример такой программы на языке программирования C.

```
main() {char *p="main() {char *p=%c%c; (void)printf(p, 34, p, 34, 10); } %c"; (void)printf(p, 34, p, 34, 10); }
```

щенного упорядочения и многие другие. Ряд предложенных им алгоритмов лучше тех, которые я использовал ранее. Помогла мне и его пунктуальность. Например, я ошибочно написал, что шестнадцатеричное число 0хAAAAAAAA можно разложить на множители $2 \cdot 3 \cdot 17 \cdot 257 \cdot 65537$; Гай указал, что 3 необходимо заменить на 5. Он не упустил ни одной из подобных мелочей и намного улучшил стиль изложения материала. Кроме того, весь материал, связанный с методом “параллельного префикса”, появился в книге исключительно благодаря Гаю.

Дополнительный материал по данной книге можно найти по адресу www.HackersDelight.org

Г.С. Уоррен-мл. (H.S. Warren, Jr.)

Йорктаун, Нью-Йорк

Июнь 2012

ГЛАВА I

ВВЕДЕНИЕ

1.1. Система обозначений

Повсюду в книге при описании машинных команд используются выражения, отличающиеся от обычных арифметических выражений. В "компьютерной арифметике" операнды представляют собой битовые строки (или векторы) некоторой фиксированной длины. Выражения компьютерной арифметики похожи на выражения обычной арифметики, но в этом случае переменные обозначают содержимое регистров компьютера. Значение выражения компьютерной арифметики представляет собой строку битов без какой-либо специальной интерпретации; операнды же могут интерпретироваться по-разному. Например, в команде сравнения операнды интерпретируются либо как двоичные целые числа со знаком, либо как беззнаковые целые числа. Поэтому, в зависимости от типа операндов, для обозначения оператора сравнения используются разные символы.

Основное различие между обычными арифметическими действиями и компьютерными сложением, вычитанием и умножением состоит в том, что результат действий компьютерной арифметики приводится по модулю 2^n , где n — размер машинного слова. Еще одно отличие состоит в том, что компьютерная арифметика содержит большее количество операций. Кроме основных четырех арифметических действий, машина способна выполнить множество других команд: логические *и*, *исключающее или*, *сравнение*, *сдвиг влево* и т.п.

Если не оговорено иное, везде в книге предполагается, что длина слова равна 32 бит, а целые числа со знаком представляются в дополнительном к 2 коде. Если длина слова иная, такие случаи оговариваются особо.

Все выражения компьютерной арифметики записываются аналогично обычным арифметическим выражениям с тем отличием, что переменные, обозначающие содержимое регистров компьютера, выделяются полужирным шрифтом (это обычное соглашение, принятое в векторной алгебре). Машинное слово интерпретируется как вектор, состоящий из отдельных битов. Константы также выделяются полужирным шрифтом, если обозначают содержимое регистра (в векторной алгебре аналогичного обозначения нет, так как там для записи констант используется только один способ — указание компонентов вектора). Если константа означает часть команды (например, непосредственно значение в команде сдвига), то она не выделяется.

Если оператор, например "+", складывает выделенные полужирным шрифтом операнды, то он подразумевает машинную операцию сложения ("векторное сложение"). Если операнды не выделены, то подразумевается обычное арифметическое сложение скалярных величин (скалярное сложение). Переменная x обозначает арифметическое значение выделенной полужирным шрифтом переменной x (интерпретируемое как знаковое или беззнаковое, что должно быть понятно из контекста). Таким образом, если $x = 0x80000000$,

$a = y = 0x80000000$, то при знаковой интерпретации $x = y = -2^{31}$, $x + y = -2^{32}$ и $x + y = 0$. (Здесь $0x80000000$ — шестнадцатеричная запись строки битов, состоящей из одного единичного бита и следующих за ним 31 нулевого бита.)

Биты нумеруются справа налево, причем крайний справа (младший) бит имеет номер 0. Длины бита, полубайта, байта, полуслова, слова и двойного слова равны соответственно 1, 4, 8, 16, 32 и 64 бит.

Небольшие и простые фрагменты кода представлены в виде алгебраических выражений с оператором присваивания (стрелка влево) и при необходимости с оператором *if*. Для таких задач использование математических выражений предпочтительнее, чем составление программы на языке ассемблера.

Более сложные и громоздкие алгоритмы представлены в виде программ на языке C, определяемом стандартом ISO 1999.

Полное описание языка C выходит за рамки данной книги. Однако для тех читателей, которые не знакомы с C, в табл. 1.1 приведено краткое описание основных элементов этого языка [57]. Это описание особенно полезно для читателей, которые не знакомы с C, но знают какой-либо иной процедурный язык программирования. В таблице перечислены также операторы, которые используются в алгебраических и арифметических выражениях этой книги. Все операторы упорядочены по приоритетам: в начале таблицы расположены операторы, имеющие наиболее высокий приоритет, в конце — операторы с самым низким приоритетом. Буква "L" в столбце "Приоритет" означает, что данный оператор левоассоциативен, а именно:

$$a \cdot b \cdot c = (a \cdot b) \cdot c.$$

Буква "R" означает, что оператор правоассоциативен. Приведенные в таблице и используемые в книге алгебраические операторы следуют правилам ассоциативности и приоритетам операторов языка C.

В дополнение к перечисленному в табл. 1.1 используется ряд обозначений из булевой алгебры и стандартной математики (соответствующие пояснения приводятся по мере необходимости).

ТАБЛИЦА 1.1. ЭЛЕМЕНТЫ ЯЗЫКА C И АЛГЕБРАИЧЕСКИЕ ОПЕРАТОРЫ

Приоритет	C	Алгебра	Описание
	$0x \dots$	$0x \dots, 0b \dots$	Шестнадцатеричные, двоичные константы
16	$a[k]$		Выбор k -го компонента
16		x_0, x_1, \dots	Различные переменные или выборка битов (зависит от контекста)
16	$f(x, \dots)$	$f(x, \dots)$	Вычисление функции
16		$abs(x)$	Абсолютное значение (однако $abs(-2^{31}) = -2^{31}$)

Продолжение табл. 1.1

Приоритет	C	Алгебра	Описание
16		$\text{abs}(x)$	Отрицательное абсолютное значение
15	$x++, x--$		Прибавление, вычитание единицы (постфиксная форма, постинкремент, постдекремент)
14	$++x, --x$		Прибавление, вычитание единицы (префиксная форма, преинкремент, предекремент)
14	$(\text{type name}) x$		Приведение типа
14 R		x^k	x в степени k
14	$\sim x$	$\sim x, \bar{x}$	Побитовое не (побитовое дополнение до единицы)
14	$!x$		Логическое отрицание (<i>if x = 0 then 1 else 0</i>)
14	$-x$	$-x$	Арифметическое отрицание
13 L	$x * y$	$x \cdot y$	Умножение по модулю размера слова
13 L	x / y	$x \div y$	Знаковое целочисленное деление
13 L	$x \ / \ y$	$x \overset{\cdot}{\div} y$	Беззнаковое целочисленное деление
13 L	$x \% y$	$\text{rem}(x, y)$	Остаток (может быть отрицательным) от деления ($x \div y$) знаковых аргументов
13 L	$x \ \% \ y$	$\text{remu}(x, y)$	Остаток от деления ($x \overset{\cdot}{\div} y$) беззнаковых аргументов
		$\text{mod}(x, y)$	Приведение x по модулю y к интервалу $[0, \text{abs}(y) - 1]$; знаковые аргументы
12 L	$x + y, x - y$	$x + y, x - y$	Сложение, вычитание
11 L	$x \ll y, x \gg y$	$x \ll y, x \overset{\circ}{\gg} y$	Сдвиг влево, вправо с заполнением нулями ("логический" сдвиг)
11 L	$x \gg y$	$x \overset{\circ}{\gg} y$	Сдвиг вправо со знаковым заполнением ("арифметический" или "алгебраический" сдвиг)
11 L		$\overset{\text{rot}}{x} \ll y, \overset{\text{rot}}{x} \gg y$	Циклический сдвиг влево, вправо

Окончание табл. 1.1

Приоритет	C	Алгебра	Описание
10 L	$x < y, x \leq y$ $x > y, x \geq y$	$x < y, x \leq y$ $x > y, x \geq y$	Сравнение знаковых аргументов
10 L	$x < y, x \leq y$ $x > y, x \geq y$	$\overset{\cdot}{x} < \overset{\cdot}{y}, \overset{\cdot}{x} \leq \overset{\cdot}{y}$ $\overset{\cdot}{x} > \overset{\cdot}{y}, \overset{\cdot}{x} \geq \overset{\cdot}{y}$	Сравнение беззнаковых аргументов
9 L	$x == y, x != y$	$x = y, x \neq y$	Равно, не равно
8 L	$x \& y$	$x \& y$	Побитовое и
7 L	$x \wedge y$	$x \oplus y$	Побитовое <i>исключающее или</i>
7 L		$x \equiv y$	Побитовая <i>эквивалентность</i> ($\neg(x \oplus y)$)
6 L	$x y$	$x y$	Побитовое <i>или</i>
5 L	$x \&\& y$	$x \bar{\&} y$	Логическое <i>и</i> (<i>if</i> $x = 0$ <i>then</i> 0 <i>else if</i> $y = 0$ <i>then</i> 0 <i>else</i> 1)
4 L	$x y$	$x \bar{ } y$	Логическое <i>или</i> (<i>if</i> $x \neq 0$ <i>then</i> 1 <i>else if</i> $y \neq 0$ <i>then</i> 1 <i>else</i> 0)
3 L		$x y$	Конкатенация
2 R	$x = y$	$x \leftarrow y$	Присваивание

Кроме функций "abs", "tem" и прочих, в книге используется множество других функций, которые будут определены позже.

При вычислении выражения $x < y < z$ в языке C сначала вычисляется выражение $x < y$ (результат равен 1, если выражение истинно, и 0, если выражение ложно), затем полученный результат сравнивается с z . Выражение компьютерной алгебры $x < y < z$ вычисляется как $(x < y) \& (y < z)$.

В языке C есть три оператора цикла: while, do и for. Цикл while имеет вид

while (*выражение*) *оператор*

Перед выполнением цикла вычисляется *выражение*. Если *выражение* истинно (не нуль), выполняется *оператор*. Затем снова вычисляется *выражение*. Цикл while завершается, когда *выражение* становится ложным (равным нулю).

Цикл do аналогичен циклу while, однако проверочное условие стоит после оператора цикла, а именно

do *оператор* while (*выражение*)

В этом цикле сначала выполняется *оператор*, затем вычисляется *выражение*. Если выражение истинно, операторы цикла выполняются еще раз, если выражение ложно, цикл завершается.

Оператор `for` имеет вид

`for (e1; e2; e3) оператор`

Сначала вычисляется выражение e_1 — как правило, это оператор присваивания (аналог выражения-инициализации). Затем вычисляется e_2 — оператор сравнения (или условное выражение). Если значение условного выражения равно нулю (условие ложно), цикл завершается, если не равно нулю (условие истинно), выполняется *оператор*. Затем вычисляется e_3 (также, как правило, оператор присваивания), и вновь вычисляется условное выражение e_2 . Таким образом, знакомый всем цикл “do $i = 1$ to n ” запишется как `for (i=1; i<=n; i++)` (это один из контекстов использования оператора постинкремента).

Стандарт ISO не определяет, распространяют ли правые сдвиги (оператор `>>`) знаковых величин бит 0 или бит знака. В коде C в данной книге предполагается, что если левый операнд знаковый, то в результате сдвига вправо будет распространяться бит знака (в случае беззнакового операнда биты слева будут заполняться нулями, как и требует ISO). Большинство современных компиляторов C работают именно таким образом.

В данной книге мы также полагаем, что левые сдвиги являются “логическими” (на ряде машин, в особенности на старых, выполняется “арифметический” левый сдвиг, при котором сохраняется знаковый бит).

Еще одна потенциальная проблема со сдвигами заключается в том, что стандарт ISO указывает, что если величина сдвига отрицательна или не менее ширины левого операнда, то результат операции не определен. Однако почти все 32-разрядные машины трактуют величины сдвигов по модулю 32 или 64. Представленный в книге код зависит от способа такой обработки величины сдвига; там, где это важно, в книге приводится соответствующее пояснение.

1.2. Система команд и модель оценки времени выполнения команд

Чтобы можно было хотя бы грубо сравнивать алгоритмы, представим, что они кодируются для работы на машине с набором команд, подобных современным RISC-компьютерам общего назначения (архитектуры типа IBM RS/6000, Oracle SPARC и ARM). Это трехадресная машина, имеющая достаточно большое количество регистров общего назначения — не менее 16. Если не оговорено иное, все регистры 32-разрядные. Регистр общего назначения с номером 0 всегда содержит нули, все другие регистры равноправны и могут использоваться для любых целей.

Для простоты будем считать, что в компьютере нет регистров “специального назначения”, в частности слова состояния процессора или регистра с битами состояний, например “переполнение”. Не рассматриваются также команды для работы с числами с плавающей точкой как выходящие за рамки тематики данной книги.

В книге описаны два типа RISC: “базовый RISC”, команды которого перечислены в табл. 1.2, и “RISC с полным набором команд”, в который, кроме основных RISC-команд, входят дополнительные команды, перечисленные в табл. 1.3.

ТАБЛИЦА 1.2. БАЗОВЫЙ НАБОР RISC-КОМАНД

Мнемокод команды	Операнды	Описание команды
add, sub, mul, div, divu, rem, remu	RT, RA, RB	$RT \leftarrow RA \text{ op } RB$ Здесь op — сложение (<i>add</i>), вычитание (<i>sub</i>), умножение (<i>mul</i>), знаковое деление (<i>div</i>), беззнаковое деление (<i>divu</i>), остаток от знакового деления (<i>rem</i>) или остаток от беззнакового деления (<i>remu</i>)
addi, muli	RT, RA, I	$RT \leftarrow RA \text{ op } I$ Здесь op — сложение (<i>addi</i>) или умножение (<i>muli</i>), I — непосредственно заданное 16-битовое знаковое значение
addis	RT, RA, I	$RT \leftarrow RA + (I \ll 16)$
and, or, xor	RT, RA, RB	$RT \leftarrow RA \text{ op } RB$ Здесь op — побитовое и (<i>and</i>), или (<i>or</i>) или исключающее или (<i>xor</i>)
andi, ori, xori	RT, RA, Iu	То же самое, но последний операнд является непосредственно заданным 16-битовым беззнаковым значением
beq, bne, blt, ble, bgt, bge	RT, target	Переход к метке target, если выполнено некоторое условие, т.е. если $RT=0$, $RT \neq 0$, $RT < 0$, $RT \leq 0$, $RT > 0$ или $RT \geq 0$ соответственно (RT — целое знаковое число)
bt, bf	RT, target	Переход к метке target, если выполнено некоторое условие (true/false); аналогичны командам bne/beq соответственно
cmpeq, cmrne, cmplt, cmple, cmrgt, cmrge, cmpltu, cmpleu, cmrgtu, cmrgeu	RT, RA, RB	RT содержит результат сравнения RA и RB; RT равен 0, если условие ложно, и 1, если условие истинно. Мнемокоды означают: сравнить на равенство, неравенство, меньше, не больше и так далее, как и в командах условного перехода. Суффикс u в названии обозначает сравнение беззнаковых величин
cmpeq, cmrne, cmplt, cmple, cmrgt, cmrge	RT, RA, I	Аналогичны командам серии cmreq, но второй операнд представляет собой непосредственно заданное 16-битовое знаковое значение

Окончание табл. 1.2

Мнемокод команды	Операнды	Описание команды
cmpiequ, cmpineu, cmpiltu, cmpileu, cmpigtu, cmpigeu	RT, RA, Iu	Аналогичны командам серии cmpiltu, но второй операнд представляет собой непосредственно заданное 16-битовое беззнаковое значение
ldbu, ldh, ldhu, ldw	RT, d(RA)	Загрузка беззнакового байта, знакового полуслова, беззнакового полуслова или слова в RT из ячейки памяти по адресу RA + d, где d — непосредственно заданное знаковое 16-битовое значение
mulhs, mulhu	RT, RA, RB	В RT помещаются старшие 32 бит результата умножения RA и RB (знакового и беззнакового)
not	RT, RA	В RT помещается побитовое дополнение RA до единицы
shl, shr, shrs	RT, RA, RB	В RT помещается значение RA, сдвинутое влево или вправо; величина сдвига задается шестью младшими битами второго операнда (RB). При выполнении команды shrs освобожденные биты заполняются содержимым знакового разряда, в остальных командах освобожденные биты заполняются нулями. (Значение величины сдвига вычисляется по модулю 64.)
shli, shri, shrsi	RT, RA, Iu	В RT помещается значение RA, сдвинутое влево или вправо на величину, задаваемую пятью младшими битами непосредственно заданного значения I
stb, sth, stw	RS, d(RA)	Сохранение байта, полуслова или слова из регистра RS в ячейку памяти по адресу RA + d, где d — непосредственно заданное 16-битовое знаковое значение

Везде в описаниях команд исходные операнды RA и RB представляют собой содержимое регистров.

В реальной машине обязательно должны быть команды ветвления и обращения к подпрограммам, команды перехода по адресу, содержащемуся в регистре (для возврата из подпрограмм и обработки оператора выбора альтернативы типа switch), а также, возможно, ряд команд для работы с регистрами специального назначения. Конечно же, должны быть привилегированные команды и команды вызова служб супервизора. Кроме того, вероятно наличие команд для работы с числами с плавающей точкой.

Краткое описание некоторых дополнительных команд, которые может иметь RISC-компьютер, приведено в табл. 1.3.

ТАБЛИЦА 1.3. ДОПОЛНИТЕЛЬНЫЕ КОМАНДЫ ПОЛНОГО НАБОРА RISC-КОМАНД

Мнемонакод команды	Операнды	Описание команды
abs, nabs	RT, RA	RT получает абсолютное значение (или отрицательное абсолютное значение) RA
andc, eqv, nand, nor, orc	RT, RA, RB	Побитовое <i>и с дополнением</i> , эквивалентность, отрицание <i>и</i> , отрицание <i>или</i> <i>и</i> или <i>с дополнением</i>
extr	RT, RA, I, L	Извлечение битов от I до I+L-1 из RA и размещение их в младших битах RT; остальные разряды заполняются нулями
extrs	RT, RA, I, L	Аналогична extr, но свободные биты заполняются содержимым бита знака
ins	RT, RA, I, L	Вставляет биты регистра RA от 0 до L-1 в биты от I до I+L-1 приемника RT
nlz	RT, RA	RT содержит количество старших нулевых битов RA (от 0 до 32)
pop	RT, RA	RT содержит количество единичных битов RA (от 0 до 32)
ldb	RT, d(RA)	Загрузка знакового байта из ячейки памяти по адресу RA + d в RT, где d — непосредственно задаваемое 16-битовое знаковое значение
moveq, movne, movlt, movle, movgt, movge	RT, RA, RB	В RT помещается значение RB, если RA=0 (RA≠0 и т.д.). Если условие не выполняется, содержимое RT не изменяется
shlr, shrr	RT, RA, RB	В RT помещается значение RA после циклического сдвига влево или вправо; величина сдвига задается пятью младшими битами RB
shlri, shrri	RT, RA, Iu	В RT помещается значение RA после циклического сдвига влево или вправо; величина сдвига задается пятью младшими битами непосредственно заданного значения
trpeq, trpne, trplt, trple, trpgt, trpge, trpltu, trpleu, trpgtu, trpgeu	RA, RB	Прерывание (<i>trap</i>), если RA=RB (RA≠RB и т.д.)

Окончание табл. 1.3

Мнемокод команды	Операнды	Описание команды
trpieq, trpine, trpilt, trpile, trpigt, trpige	RA, I	Аналогичны trreq и остальным командам серии с тем отличием, что второй операнд представляет собой непосредственно заданное 16-битовое знаковое значение
trpiequ, trpineu, trpiltu, trpileu, trpigtu, trpigeu	RA, Iu	Аналогичны trpltu и другим командам серии с тем отличием, что второй операнд представляет собой непосредственно заданное 16-битовое беззнаковое значение

На практике в языке ассемблера удобно иметь ряд “расширенных мнемоник”, похожих на макросы, которые обычно выполняются как одна команда. Некоторые из них представлены в табл. 1.4.

ТАБЛИЦА 1.4. РАСШИРЕННЫЕ МНЕМОНИКИ

Расширенная мнемоника	Расширение	Описание
b target	beq R0, target	Безусловный переход
li RT, I	См. пояснение в тексте	Загрузка непосредственно заданного числа, $-2^{31} \leq I < 2^{32}$
mov RT, RA	ori RT, RA, 0	Пересылка регистра RA в RT
neg RT, RA	sub RT, R0, RA	Отрицание (побитовое дополнение до двух)
subi RT, RA, I	addi RT, RA, -I	Вычитание непосредственно заданного значения ($I \neq -2^{15}$)

Команда загрузки непосредственно заданного значения расширяется до одной или двух команд, в зависимости от непосредственного значения I. Например, если $0 \leq I < 2^{16}$, можно применить команду ori с использованием регистра R0. Если $-2^{15} \leq I < 0$, можно обойтись командой addi с регистром R0. Если младшие 16 бит I — нулевые, используем команду addis. Для остальных значений I выполняются две команды; например, после команды ori выполняется addis. (Иначе говоря, в последнем случае выполняется загрузка из памяти, но при оценке времени и места, которое требуется на выполнение этой макрокоманды, будем считать, что выполняются две элементарные арифметические команды.)

Существуют разные мнения о том, к какому набору RISC-команд следует отнести ту или иную команду — к базовому или к расширенному. Команды беззнакового деления и получения остатка от деления вполне можно было бы включить в расширенный набор RISC-команд. Команда же загрузки знакового байта должна входить в базовый набор.

Команда входит в расширенный набор из-за редкого использования, а также из-за того, что в некоторых методах трудно распространить знаковый бит на много позиций, или из-за длительного времени выполнения.

Различия между основным и расширенным наборами RISC-команд приводят ко множеству подобных спорных моментов, но мы не будем подробно останавливаться на этой теме.

Действие команд ограничивается двумя входными регистрами и одним выходным, что упрощает работу компьютера. Упрощается также работа оптимизирующего компилятора — ему не приходится обрабатывать инструкции с несколькими целевыми регистрами. Однако за такое упрощение приходится платить: если в программе требуется вычислить и частное, и остаток от деления двух чисел, то приходится выполнять две команды (деление и получение остатка). Стандартный алгоритм деления вместе с частным одновременно вычисляет и остаток от деления, поэтому на многих машинах и частное, и остаток получаются в результате выполнения одной команды деления. Аналогичные замечания применимы и к случаю, когда при умножении двух слов получается двойное слово-произведение.

Создается впечатление, что команда *условной пересылки* (*moveq*) имеет только два входных операнда, хотя на самом деле их три. Поскольку результат выполнения этой команды зависит от содержимого регистров RT, RA и RB, при неупорядоченном выполнении команд RT следует интерпретировать и как *используемый*, и как *устанавливаемый* регистр одновременно. Представьте ситуацию, когда за командой, устанавливающей значение RT, следует команда *условной пересылки* и при этом нельзя потерять результат первой команды. Таким образом, разработчик машины может убрать команду *условной пересылки* из набора допустимых команд, тем самым исключив обработку команд с тремя (логическими) входными операндами. С другой стороны, команда *условной пересылки* снижает количество ветвлений в программе.

В данной книге формат команд не играет большой роли. Полное множество RISC-команд, описанное выше, вместе с командами для работы с числами с плавающей точкой и рядом команд супервизора может быть реализовано с помощью 32-битовых команд на компьютере с 32 регистрами общего назначения (5-битовые поля регистров). Если размер непосредственно заданных полей в командах сравнения, загрузки, сохранения и прерывания снизить до 14 бит, это же множество команд может быть реализовано на компьютере с 64 регистрами общего назначения (с использованием 6-битовых полей регистров).

Время выполнения

Предполагается, что все команды выполняются за один такт процессора, за исключением команд умножения, деления и получения остатка от деления, для которых невозможно точно определить время выполнения. Команды перехода занимают один такт независимо от того, был ли выполнен переход.

Команда *загрузки непосредственного значения* выполняется за один или два такта в зависимости от того, сколько элементарных арифметических команд требуется для формирования константы в регистре.

Хотя команды *загрузки* и *сохранения* не часто упоминаются в данной книге, будем считать, что они выполняются за один такт, пренебрегая при этом задержкой при загрузке (промежутком времени между моментом завершения команды в арифметическом модуле и моментом, когда данные становятся доступны следующей команде).

Однако зачастую одних только знаний о количестве тактов, используемых всеми арифметическими и логическими командами, недостаточно для правильной оценки времени выполнения программы. Выполнение программы часто замедляется вследствие задержек, возникающих при загрузке и выборке данных. Задержки такого рода не обсуждаются в книге, хотя и могут существенно влиять на скорость выполнения программ (причем это влияние постоянно возрастает). Другим источником сокращения времени выполнения является возможность распараллеливания вычислений на уровне команд, которая реализована практически на всех современных RISC-компьютерах, особенно на специализированных быстродействующих машинах.

В таких машинах имеется несколько исполнительных модулей и возможность диспетчеризации команд, что позволяет параллельное выполнение независимых команд (независимыми команды считаются тогда, когда ни одна из них не использует результаты других команд, команды не заносят значения в один и тот же регистр или бит состояния). Поскольку такие возможности компьютеров теперь уже не редкость, в книге много внимания уделяется независимым командам. Так, можно сказать, что некая формула может быть закодирована таким образом, что для вычисления потребуется выполнение восьми инструкций за пять тактов на машине с неограниченными возможностями распараллеливания вычислений на уровне команд. Это означает, что если команды расположены в правильном порядке и машина имеет достаточное количество регистров и арифметических модулей, то она в принципе способна выполнить такую программу за пять тактов.

Более подробно обсуждать этот вопрос не имеет смысла, так как возможности машин в плане распараллеливания вычислений на уровне команд существенно различаются. Например, процессор IBM RS/6000, созданный в 1992 году, имеет трехходовый сумматор и может параллельно выполнять две следующие друг за другом команды сложения, даже если одна из них использует результаты другой (например, когда команда сложения использует результат команды сравнения или основной регистр команды загрузки). В качестве обратного примера можно рассмотреть простейший дешевый встраиваемый в различные системы компьютер, у которого регистр ввода-вывода имеет только один порт для чтения. Естественно, что такой машине для выполнения команды с двумя операндами-источниками потребуется дополнительный такт для повторного чтения регистра ввода-вывода. Однако если команда выдает операнд для команды, непосредственно следующей за ней, то этот операнд оказывается доступным и без чтения регистра ввода-вывода. На машинах такого типа возможно ускорение работы, если каждая команда предоставляет данные для следующей команды, т.е. если параллелизмы в коде отсутствуют.

Упражнения

1. Выразите цикл
`for (e1; e2; e3) оператор`
через цикл `while`. Можно ли выразить его через цикл `do`?
2. Закодируйте на языке программирования C цикл, в котором управляющая переменная `i` принимает все значения от 0 до максимального беззнакового значения (на 32-разрядной машине) `0xFFFFFFFF`.
3. Задача для более опытных читателей: команды базового и расширенного наборов RISC, определенных в этой книге, могут быть выполнены не более чем с помощью двух чтений из регистров и одной записи в регистр. Можете ли вы указать RISC-команды, которым требуется либо более двух чтений, либо более одной записи в регистр?

ГЛАВА 2 ОСНОВЫ

2.1. Манипуляции младшими битами

Многие из приведенных в этом разделе формул используются в следующих главах.

Для того чтобы обнулить крайний справа единичный бит, а если такового нет — вернуть 0 (например, $01011000 \Rightarrow 01010000$), используется формула

$$x \& (x - 1)$$

Эту формулу можно применять для проверки того, является ли беззнаковое целое степенью 2 или нулем — путем проверки результата вычислений на равенство 0.

Чтобы установить крайний справа нулевой бит равным 1, а если такового нет — установить все биты равными 1 (например, $10100111 \Rightarrow 10101111$), можно воспользоваться формулой

$$x | (x + 1)$$

Приведенная ниже формула позволяет превратить все завершающие значение единичные биты в нулевые, а если таковых нет — просто вернуть исходное значение (например, $10100111 \Rightarrow 10100000$):

$$x \& (x + 1)$$

Эта формула может использоваться для проверки того, имеет ли беззнаковое целое число вид $2^n - 1$, равно нулю или все его биты равны 1: для этого следует результат вычислений проверить на равенство нулю.

Чтобы превратить все завершающие значение нулевые биты в единичные, а если таковых нет, просто вернуть исходное значение (например, $10101000 \Rightarrow 10101111$), можно выполнить следующую операцию:

$$x | (x - 1)$$

Приведенная далее формула позволяет создать слово с единственным битом 1 в позиции крайнего справа 0 в слове x , а если такового нет, то вернуть значение 0 (например, $10100111 \Rightarrow 00001000$):

$$\sim x \& (x + 1)$$

Небольшое изменение превращает эту формулу в другую, которая создает слово с единственным битом 0 в позиции крайнего справа единичного бита в слове x , а если такового нет, то возвращающую слово со всеми битами, установленными равными 1 (например, $10101000 \Rightarrow 11110111$):

$$\sim x \& (x - 1)$$

Для создания слова с единицами в позициях завершающих нулевых битов и с нулевыми битами во всех остальных местах можно воспользоваться одной из приведенных ниже формул (например, $01011000 \Rightarrow 00000111$):

$$\begin{aligned} & \neg x \& (x - 1) \text{ или} \\ & \neg(x | -x), \text{ или} \\ & (x \& -x) - 1 \end{aligned}$$

Первая из этих формул позволяет распараллелить вычисления на уровне команд.

Для создания слова с нулевыми битами в позициях завершающих единичных битов в x и единичными битами во всех остальных местах, или, если завершающих единичных битов нет, то слова, состоящего из единичных битов (например, $10100111 \Rightarrow 11111000$), используется формула

$$\neg x | (x + 1)$$

Формула

$$x \& (-x)$$

выделяет крайний справа единичный бит, возвращая 0, если такового нет (например, $01011000 \Rightarrow 00001000$).

Формула

$$x \oplus (x - 1)$$

создает слово, в котором в позициях крайнего справа единичного бита и следующих за ним нулевых битов x стоят единичные биты; если единичного бита нет — слово из единичных битов, и целое число, равное 1, если нет завершающих нулей (например, $01011000 \Rightarrow 00001111$).

Используйте приведенную далее формулу для создания слова с единичными битами в позициях крайнего справа нулевого бита и следующих за ним единичных битов x ; если нулевого бита нет — слово из единичных битов, и целое число, равное 1, если нет завершающих единиц (например, $01010111 \Rightarrow 00001111$):

$$x \oplus (x + 1)$$

Для того чтобы обнулить крайнюю справа непрерывную строку единиц (например, $01011100 \Rightarrow 01000000$) [112], можно воспользоваться одной из формул

$$\begin{aligned} & (((x | (x - 1)) + 1) \& x) \text{ и} \\ & ((x \& -x) + x) \& x \end{aligned}$$

Эти формулы можно применить для определения, имеет ли неотрицательное целое число вид $2^j - 2^k$ для некоторых $j \geq k \geq 0$. Для этого следует применить формулу и проверить результат на равенство нулю.

Расширенные законы де Моргана

Логические тождества, известные как законы де Моргана, можно рассматривать как распространение, или умножение, знака *не*. Эту идею можно применить к выражениям из этого раздела (и к некоторым другим), как показано ниже (первые два уравнения представляют собой исходные законы де Моргана).

$$\neg(x \& y) = \neg x \mid \neg y$$

$$\neg(x \mid y) = \neg x \& \neg y$$

$$\neg(x + 1) = \neg x - 1$$

$$\neg(x - 1) = \neg x + 1$$

$$\neg\neg x = x - 1$$

$$\neg(x \oplus y) = \neg x \oplus y = x \equiv y$$

$$\neg(x \equiv y) = \neg x \oplus y = x \oplus y$$

$$\neg(x + y) = \neg x - y$$

$$\neg(x - y) = \neg x + y$$

Вот пример применения этих формул:

$$\neg(x \mid \neg(x + 1)) = \neg x \& \neg\neg(x + 1) = \neg x \& ((x + 1) - 1) = \neg x \& x = 0$$

Проверка выполнимости справа налево

Имеется простая проверка того, можно ли реализовать заданную функцию в виде последовательности команд *сложения*, *вычитания*, *и*, *или* и *отрицания* [108]. К перечисленным выше командам, конечно, можно добавить другие команды из основного множества, которые, в свою очередь, являются комбинацией команд *сложения*, *вычитания*, *и*, *или* и *не*, например можно добавить команду *сдвиг влево* на определенную величину (которая эквивалентна последовательности команд *сложения*) или команду *умножения*. Команды, которые не могут быть выражены через команды *сложения*, *вычитания* и другие из основного списка, исключаются из рассмотрения. Критерий проверки следует из следующей теоремы.

ТЕОРЕМА. *Функция, отображающая слова в слова, может быть реализована посредством операций побитового сложения, вычитания, и, или, отрицания тогда и только тогда, когда каждый бит результата зависит только от битов исходных операндов в той же позиции и правее нее.*

Иными словами, представьте ситуацию, когда вы вычисляете младший бит результирующего значения, исходя только из младших битов операндов. Затем вычисляется второй справа бит результата, для чего используются только два младших бита операндов, и т.д. Если таким образом можно получить результирующее значение, то данная функция может быть вычислена с помощью последовательности команд *сложения*, *и* и т.п.

Если же таким образом (справа налево) вычислить функцию невозможно, то ее невозможно реализовать в виде последовательности указанных команд.

Наибольший интерес представляет утверждение, что любая из функций сложения, вычитания, *и*, *или* и *не* может быть вычислена справа налево, так что любая комбинация этих функций также обладает этим свойством, т.е. может быть вычислена справа налево.

Доказательство второй части теоремы достаточно громоздко, так что просто приведем конкретный пример. Предположим, что функция двух переменных, x и y , обладает свойством вычислимости справа налево, и пусть второй бит результата r вычисляется следующим образом.

$$r_2 = x_2 | (x_0 \& y_1) \quad (1)$$

Биты пронумерованы справа налево от 0 до 31. Так как бит 2 результата является функцией только крайних справа битов исходных операндов (бита номер 2 и битов младше его), он также вычисляется справа налево.

Запишем машинные слова x , x , сдвинутое на два бита влево, и y , сдвинутое на один бит влево, как показано ниже. Добавим к записи маску, которая выделяет второй бит.

x_{31}	x_{30}	...	x_3	x_2	x_1	x_0	
x_{29}	x_{28}	...	x_1	x_0	0	0	
y_{30}	y_{29}	...	y_2	y_1	y_0	0	
0	0	...	0	1	0	0	
0	0	...	0	r_2	0	0	

К строкам 2 и 3 применим операцию *побитового и*, затем, согласно формуле (1), к получившемуся слову применим операцию *побитового или* со строкой 1, после чего к полученному результату и маске (строка 4) применим операцию *побитового и*. В полученном таким образом слове все биты, кроме второго, будут нулевыми. Выполним подобные вычисления для получения остальных битов результата и объединим все 32 слова с помощью операции *или*. В итоге получим искомую функцию.

Такое построение не дает эффективной программы вычисления значений функции, а всего лишь показывает, что искомая функция может быть выражена с использованием команд из базового списка.

Используя эту теорему, можно сразу же увидеть, что не существует такой последовательности команд, которая обнуляла бы в слове крайний слева единичный бит, так как для этого необходимо просмотреть биты, находящиеся слева от него (чтобы точно знать, что это действительно крайний слева единичный бит). Аналогично не существует такой последовательности операций, которая бы давала сдвиг вправо, циклический сдвиг или сдвиг влево на переменную величину, а также могла подсчитать количество завершающих нулевых битов в слове (при подсчете завершающих нулевых битов младший бит результата должен быть равен 1, если это количество нечетно, и 0, если четно; так что необходимо просмотреть биты слева от младшего, чтобы выяснить его значение).

Новое применение

Одним из применений рассмотренной сортировки битов является задача поиска следующего числа, которое больше заданного, но имеет такое же количество единичных битов. Зачем это может понадобиться? Эта задача возникает при использовании битовых строк для представления подмножеств. Возможные члены множества перечислены в одномерном массиве, а подмножество представляет собой слово или последовательность слов, в которых бит i установлен равным 1, если i -й элемент принадлежит этому подмножеству. Тогда объединение множеств вычисляется с помощью *побитовой операции или* над битовыми строками, пересечение — с помощью операции *и* и т.д.

Иногда требуется выполнить определенные действия над всеми подмножествами заданной длины, что легко сделать с помощью функции, отображающей представленное число заданное подмножество в следующее большее число, содержащее такое же количество единичных битов, что и исходное (строка подмножества интерпретируется как целое число).

Компактный алгоритм для этой задачи составлен Р.У. Госпером (R.W. Gosper) [44, item 175]¹. Пусть имеется некоторое слово x , представляющее собой подмножество. Идея в том, чтобы сначала найти в x крайнюю справа группу смежных единичных битов, следующую за нулевыми битами, а затем “увеличить” представленную этой группой величину таким образом, чтобы количество единичных битов осталось неизменным. Например, строка $xxx011110000$, где xxx — произвольные биты, преобразуется в строку $xxx100000111$. Алгоритм работает следующим образом. Сначала в x определяется самый младший единичный бит, т.е. вычисляется $s = x \& -x$ (в результате получаем 000000010000). Затем эта строка складывается с x , давая $r = xxx100000000$, в которой получен первый единичный бит результата. Чтобы получить остальные биты, установим $n-1$ младших битов слова равными 1, где n — размер крайней справа группы единичных битов в исходном слове x . Для этого над r и x выполняется операция *исключающего или*, что в нашем примере дает строку 000111110000 .

В полученной строке содержится больше единичных битов, чем в исходной, поэтому выполним еще одно преобразование. Поделим искомое число на s , что эквивалентно его сдвигу вправо (s — степень 2), сдвинем его вправо еще на два разряда, что отбросит еще два нежелательных бита, и выполним операцию побитового *или* над полученным числом и r .

В записи компьютерной алгебры результат y получается следующим образом.

$$\begin{aligned} s &\leftarrow x \& -x \\ r &\leftarrow s + x \\ y &\leftarrow r \mid \left(\left((x \oplus r) \gg 2 \right) + s \right) \end{aligned} \quad (2)$$

¹ Вариант этого алгоритма имеется в [57] (раздел 7.6.7).

В листинге 2.1 представлена законченная процедура на языке C, выполняющая семь базовых RISC-команд, одной из которых является команда деления. (Данная процедура неприменима для $x = 0$, так как в этом случае получается деление на 0.)

**ЛИСТИНГ 2.1. Вычисление следующего большего числа
с тем же количеством единичных битов**

```
unsigned snob(unsigned x) {
    unsigned smallest, ripple, ones;

    smallest = x & -x;           // x = xxx0 1111 0000
    ripple   = x + smallest;     //      0000 0001 0000
    ones     = x ^ ripple;       //      xxx1 0000 0000
    ones     = x ^ ripple;       //      0001 1111 0000
    ones     = (ones >> 2)/smallest; // 0000 0000 0111
    return  ripple | ones;      //      xxx1 0000 0111
}
```

Если деление выполняется медленно, но у вас есть быстрый способ вычисления количества завершающих нулевых битов $\text{nlz}(x)$, количества ведущих нулевых битов $\text{nlz}(x)$ или функции степени заполнения $\text{pop}(x)$ (количество единичных битов в x), то последнюю строку в (2) можно заменить одной из приведенных ниже формул. (Первые два метода могут не работать на машине, которая выполняет сдвиги по модулю 32.)

$$y \leftarrow r | \left((x \oplus r) \gg (2 + \text{nlz}(x)) \right)$$

$$y \leftarrow r | \left((x \oplus r) \gg (33 - \text{nlz}(s)) \right)$$

$$y \leftarrow r | \left((1 \ll (\text{pop}(x \oplus r) - 2)) - 1 \right)$$

2.2. Сложение и логические операции

Предполагается, что читатель знаком с элементарными тождествами обычной и булевой алгебры. Ниже приводится набор аналогичных тождеств для операций сложения и вычитания в комбинации с логическими операциями.

$$\begin{aligned} \text{а)} \quad & -x = \neg x + 1 \\ \text{б)} \quad & = \neg(x - 1) \\ \text{в)} \quad & \neg x = -x - 1 \\ \text{г)} \quad & \neg \neg x = x + 1 \\ \text{д)} \quad & \neg \neg x = x - 1 \\ \text{е)} \quad & x + y = x - \neg y - 1 \\ \text{ж)} \quad & = (x \oplus y) + 2(x \& y) \\ \text{з)} \quad & = (x | y) + (x \& y) \\ \text{и)} \quad & = 2(x | y) - (x \oplus y) \\ \text{к)} \quad & x - y = x + \neg y + 1 \\ \text{л)} \quad & = (x \oplus y) - 2(\neg x \& y) \end{aligned}$$

$$\begin{aligned}
 \text{м)} & \quad = (x \& \neg y) - (\neg x \& y) \\
 \text{н)} & \quad = 2(x \& \neg y) - (x \oplus y) \\
 \text{о)} \quad x \oplus y & = (x | y) - (x \& y) \\
 \text{п)} \quad x \& \neg y & = (x | y) - y \\
 \text{р)} & \quad = x - (x \& y) \\
 \text{с)} \quad \neg(x - y) & = y - x - 1 \\
 \text{т)} & \quad = \neg x + y \\
 \text{у)} \quad x \equiv y & = (x \& y) - (x | y) - 1 \\
 \text{ф)} & \quad = (x \& y) + \neg(x | y) \\
 \text{х)} \quad x | y & = (x \& \neg y) + y \\
 \text{ц)} \quad x \& y & = (\neg x | y) - \neg x
 \end{aligned}$$

Равенство (г) можно повторно применять к самому себе, например $\neg\neg\neg x = x + 2$ и т.д. Аналогично после повторного использования равенства (д) получаем $\neg\neg\neg x = x - 2$. Таким образом, чтобы добавить к x или вычесть из него некоторую константу, можно воспользоваться только этими выражениями дополнения.

Равенство (е) дуально равенству (к), которое представляет собой известное отношение, позволяющее получить вычитающее устройство из сумматора.

Равенства (ж) и (з) взяты из [44, item 23]. В равенстве (ж) сумма x и y вычисляется следующим образом: сначала x и y суммируются без учета переносов ($x \oplus y$), а затем к полученному результату добавляются переносы. Равенство (з) изменяет операнды перед сложением так, что в любом разряде все комбинации типа $0 + 1$ заменяются на $1 + 0$.

Можно показать, что при обычном сложении двоичных чисел с независимыми битами вероятность возникновения переноса в любом разряде равна 0.5. Однако, если при создании сумматора используется предварительная подготовка входных данных, как при сложении по формуле (ж), вероятность появления переноса становится равной 0.25. При разработке сумматора вероятность переноса значения не имеет; самой важной характеристикой при этом является максимальное количество логических схем, через которые может потребоваться пройти переносу, а использование равенства (ж) позволяет свести это количество всего лишь к одной схеме.

Равенства (л) и (м) дуальны равенствам (ж) и (з) для вычитания. В (л) сначала выполняется вычитание без заемов ($x \oplus y$), а затем из полученного результата вычитаются заемы. Аналогично равенство (м) изменяет операнды перед вычитанием таким образом, что все комбинации типа $1 - 1$ заменяются на $0 - 0$.

Равенство (о) показывает, как реализовать *исключающее или*, используя всего три базовые RISC-команды. Использование логики *и-или-не* требует выполнения четырех команд $((x | y) \& \neg(x \& y))$. Аналогично равенства (х) и (ц) реализуют операции *и* и *или* с помощью трех элементарных команд (хотя по законам де Моргана (De Morgan) их требуется четыре).

2.3. Неравенства с логическими и арифметическими выражениями

Неравенства с двоичными логическими выражениями, значения которых интерпретируются как целые беззнаковые величины, выводятся практически тривиально. Вот пара примеров.

$$(x \oplus y) \leq (x | y)$$

$$(x \& y) \leq (x \equiv y)$$

Их легко получить из табл. 2.1, в которой представлены все логические операции.

ТАБЛИЦА 2.1. РЕЗУЛЬТАТЫ РАБОТЫ 16 ЛОГИЧЕСКИХ ОПЕРАЦИЙ

x	y	0	$x \& y$	$x \& \neg y$	x	$\neg x \& y$	y	$x \oplus y$	$x y$	$\neg(x y)$	$x \equiv y$	$\neg y$	$x \neg y$	$\neg x$	$\neg x y$	$\neg(x y)$	1
0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1
0	1	0	0	0	0	1	1	1	1	0	0	0	0	1	1	1	1
1	0	0	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1
1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1

Пусть функции $f(x, y)$ и $g(x, y)$ представлены двумя столбцами табл. 2.1. Если в каждой строке, где $f(x, y)$ равна 1, $g(x, y)$ также равна 1, то для любых значений (x, y) выполняется соотношение $f(x, y) \leq g(x, y)$. Очевидно, что это неравенство справедливо и для параллельных логических операций над словами. Легко вывести и более сложные соотношения (многие из которых тривиальны), как, например, $(x \& y) \leq x \leq (x | \neg y)$ и т.п. Если же значения функций $f(x, y)$ и $g(x, y)$ в одной строке таблицы равны 0 и 1, а в другой — 1 и 0 соответственно, то между этими логическими функциями не существует отношения неравенства. Таким образом, всегда можно выяснить, выполняется ли для каких-то функций f и g соотношение $f(x, y) \leq g(x, y)$.

При работе с неравенствами требуется внимание. Например, в обычной арифметике, если $x + y \leq a$ и $z \leq x$, то $z + y \leq a$. Но этот вывод становится неверным, если "+" заменить оператором *или*.

Неравенства, включающие как логические, так и арифметические выражения, более интересны. Ниже приведены некоторые из них.

а) $(x | y) \geq \max(x, y)$

б) $(x \& y) \leq \min(x, y)$

в) $(x | y) \leq x + y$ Если сложение не вызывает переполнения

г) $(x|y) \dot{>} x + y$ Если сложение вызывает переполнение

д) $|x - y| \dot{\leq} (x \oplus y)$

Доказать все эти неравенства достаточно просто, за исключением, возможно, только неравенства (д). Под $|x - y|$ подразумевается абсолютное значение $x - y$, которое может быть вычислено в области беззнаковых чисел как $\max(x, y) - \min(x, y)$. Это неравенство можно доказать по индукции по длинам x и y (доказать неравенство будет проще, если длины x и y расширять справа налево, а не наоборот).

2.4. Абсолютное значение

Если нет отдельной команды вычисления абсолютного значения, ее можно реализовать тремя или четырьмя командами (без ветвления). Сначала вычисляется значение $y \leftarrow x \dot{>} 31$, а затем одно из перечисленных ниже выражений (разумеется, $2x$ означает $x + x$, или $x \ll 1$).

abs	nabs
$(x \oplus y) - y$	$y - (x \oplus y)$
$(x + y) \oplus y$	$(y - x) \oplus y$
$x - (2x \& y)$	$(2x \& y) - x$

Если у вас есть возможность быстрого умножения на переменную, значение которой равно ± 1 , можно поступить следующим образом.

$$\left(\left((x \dot{>} 30) | 1 \right) * x \right)$$

2.5. Среднее двух целых

Приведенная далее формула может применяться для вычисления среднего значения двух беззнаковых целых чисел $\lfloor (x + y)/2 \rfloor$ (без переполнения) [25].

$$(x \& y) + \left((x \oplus y) \dot{>} 1 \right) \quad (3)$$

Приведенная далее формула вычисляет для беззнаковых целых чисел значение $\lceil (x + y)/2 \rceil$.

$$(x|y) - \left((x \oplus y) \dot{>} 1 \right)$$

Для вычисления тех же значений (среднее с применением пола и потолка), но для знаковых целых чисел, используются те же формулы, но беззнаковый сдвиг в них заменен знаковым.

В случае знаковых целых чисел может также потребоваться вычислить среднее делением на 2 с округлением к 0. Вычисление такого "усеченного среднего" (без переполнения) — несколько более сложная задача. Его можно вычислить, вычисляя среднее

с применением пола, а затем корректируя его. Коррекция заключается в добавлении 1, если арифметически значение $x + y$ отрицательно и нечетно. Но $x + y$ отрицательно тогда и только тогда, когда результат вычисления (3) с беззнаковым сдвигом, замененным знаковым, отрицателен. Все это приводит к следующему методу (семь базовых команд RISC с общим подвыражением $x \oplus y$).

$$r \leftarrow (x \& y) + \left((x \oplus y) \gg 1 \right) \\ r + \left(\left(r \gg 31 \right) \& (x \oplus y) \right)$$

В некоторых распространенных частных случаях вычисление можно выполнить более эффективно. Если x и y являются знаковыми целыми числами и если известно, что они не отрицательны, то среднее значение может быть вычислено просто как $(x + y) \gg 1$. Такая сумма может вызвать переполнение, но бит переполнения сохранится в регистре с суммой, так что беззнаковый сдвиг переместит этот бит в правильную позицию и добавит нулевой знаковый бит.

Если x и y представляют собой беззнаковые целые числа и $x \leq y$ или если x и y — знаковые целые числа и $x \leq y$ (сравнение знаковое), то их среднее можно получить как $x + \left((y - x) \gg 1 \right)$. Это — средние с применением пола, например среднее -1 и 0 равно -1 .

2.6. Распространение знака

Под распространением знака (*sign extension*) подразумевается наличие бита в определенной позиции слова, выступающего в роли бита знака, и распространение этого бита влево при игнорировании всех остальных битов слова. Стандартный способ решения этой задачи состоит в *логическом сдвиге влево*, за которым следует *знаковый сдвиг вправо*. Однако, если эти команды работают медленно или их вообще нет на вашем компьютере, можно поступить иначе. Ниже приведены примеры распространения влево бита в седьмой позиции.

$$\begin{aligned} & ((x + 0x00000080) \& 0x000000FF) - 0x00000080 \\ & ((x \& 0x000000FF) \oplus 0x00000080) - 0x00000080 \\ & (x \& 0x0000007F) - (x \& 0x00000080) \end{aligned}$$

Вместо “+” можно использовать “-” или “ \oplus ”. Если известно, что все ненужные старшие биты нулевые, вторая формула оказывается более практичной, так как в этом случае можно не выполнять команду *и*.

2.7. Знаковый сдвиг вправо на основе беззнакового сдвига

Если на машине нет команды *знакового сдвига вправо* (*shift right signed*), ее можно заменить другими командами. Первая формула взята из [39], вторая основана на той же

идее, что и первая. Эти формулы работают при $0 \leq n \leq 31$, а если машина выполняет сдвиги по модулю 64, то последняя формула работает при $0 \leq n \leq 63$. В принципе она применима для любых n , где под "применима" подразумевается "рассматривает величину сдвига по тому же модулю, что и логический сдвиг".

Для переменной n реализация каждой формулы требует пяти или шести команд из базового множества RISC-команд.

$$\begin{aligned} & \left((x + 0x80000000) \gg n \right) - \left(0x80000000 \gg n \right) \\ & t \leftarrow 0x80000000 \gg n; \left((x \gg n) \oplus t \right) - t \\ & t \leftarrow (x \& 0x80000000) \gg n; \left(x \gg n \right) - \left(t + t \right) \\ & \left(x \gg n \right) | \left(- \left(x \gg 31 \right) \ll 31 - n \right) \\ & t \leftarrow - \left(x \gg 31 \right); \left((x \oplus t) \gg n \right) \oplus t \end{aligned}$$

В первых двух формулах выражение $0x80000000 \gg n$ можно заменить выражением $1 \ll 31 - n$.

Если n — константа, то для реализации первых двух формул на большинстве машин требуется всего лишь три команды. Если $n = 31$, функция может быть вычислена с помощью двух команд: $- \left(x \gg 31 \right)$.

2.8. Функция sign

Функция *sign*, или *signum*, определяется как

$$\text{sign}(x) = \begin{cases} -1, & x < 0, \\ 0, & x = 0, \\ 1, & x > 0. \end{cases}$$

На большинстве машин ее можно вычислить с помощью четырех команд [52].

$$\left(x \gg 31 \right) | \left(-x \gg 31 \right)$$

Если на машине нет команды *знакового сдвига вправо*, вместо нее можно использовать последнее выражение из раздела 2.7. В результате получается элегантная симметричная формула (вычисляемая за пять команд).

$$- \left(x \gg 31 \right) | \left(-x \gg 31 \right)$$

Наличие команд сравнения допускает решения из трех инструкций.

$$\text{либо } (x > 0) - (x < 0), \text{ либо } (x \geq 0) - (x \leq 0) \quad (4)$$

И последнее замечание: почти всегда работает формула $(-x \gg 31) - (x \gg 31)$; ошибка возникает только при $x = -2^{31}$.

2.9. Трехзначная функция сравнения

Трехзначная функция сравнения представляет собой определенное обобщение функции *sign* и определяется следующим образом.

$$\text{str}(x, y) = \begin{cases} -1, & x < y, \\ 0, & x = y, \\ 1, & x > y. \end{cases}$$

Данное определение справедливо для любых переменных — и знаковых, и беззнаковых. Все, что говорится в этом разделе, если явно не оговорено иное, применимо к обоим случаям.

Использование команд сравнения позволяет реализовать эту функцию за три команды как очевидное обобщение выражений (4):

$$\text{либо } (x > y) - (x < y), \text{ либо } (x \geq y) - (x \leq y)$$

Решение для беззнаковых целых чисел на компьютере PowerPC приведено ниже [20]. На этой машине “перенос” означает “не заем”.

```
subf  R5, Ry, Rx    # R5 <-- Rx-Ry
subfc R6, Rx, Ry    # R6 <-- Ry- Rx, установлен перенос
subfe R7, Ry, Rx    # R7 <-- Rx-Ry+перенос, установить перенос
subfe R8, R7, R5    # R8 <-- R5-R7+перенос, (установить перенос)
```

Если ограничиться командами из базового набора RISC-команд, то эффективного метода вычисления этой функции нет. Для вычисления $x < y$, $x \leq y$ и других команд сравнения потребуется выполнить пять команд (см. раздел 2.12), что приводит к решению, содержащему 12 команд (при использовании определенной общности вычислений $x < y$ и $y < x$). Вероятно, на RISC-машинах с базовым набором команд лучшим путем решения будет использование команд сравнения и ветвления (в худшем случае потребуется выполнить шесть команд).

2.10. Перенос знака

Функция *переноса знака* (известная в Fortran как ISIGN) определяется следующим образом.

$$\text{ISIGN}(x, y) = \begin{cases} \text{abs}(x), & y \geq 0, \\ -\text{abs}(x), & y < 0. \end{cases}$$

На большинстве машин эта функция может быть вычислена (по модулю 2^{32}) за четыре команды.

$$t \leftarrow y \gg 31;$$

$$\begin{aligned} \text{ISIGN}(x, y) &= (\text{abs}(x) \oplus t) - t \\ &= (\text{abs}(x) + t) \oplus t \end{aligned}$$

$$t \leftarrow (x \oplus y) \gg 31;$$

$$\begin{aligned} \text{ISIGN}(x, y) &= (x \oplus t) - t \\ &= (x + t) \oplus t \end{aligned}$$

2.11. Декодирование поля "0 означает 2^{**n} "

Иногда 0 или отрицательное значение для некоторой величины не имеет смысла. Поэтому она кодируется нулевым значением в n -м поле, которое интерпретируется как 2^n , а ненулевое значение имеет обычный двоичный смысл. Например, у PowerPC длина поля в команде загрузки непосредственно заданной строки слов (lswi) занимает 5 бит. Нелогично использовать эту команду для загрузки нуля байтов, поскольку длина представляет собой непосредственно заданную величину, но было бы очень полезно иметь возможность загрузить 32 байта. Длина поля кодируется значениями от 0 до 31, которые могут, например, означать длины от 1 до 32, но соглашение "нуль означает 32" приводит к более простой логике, в особенности когда процессор должен также поддерживать соответствующую команду с переменной (задаваемой в регистре) длиной, которая выполняет простое бинарное кодирование (например, команду lswx на PowerPC).

Кодирование "0 означает 2^n " целых чисел от 1 до 2^n тривиально: нужно просто маскировать биты целого числа с помощью маски $2^n - 1$. Выполнить же декодирование без проверок и переходов не так-то просто. Ниже приведено несколько возможных вариантов для работы с третьим битом. Все эти варианты требуют выполнения трех команд, не считая возможных загрузок констант.

$$\begin{array}{llll} ((x-1) \& 7) + 1 & ((x-1) | -8) + 9 & ((x+7) | 8) - 7 & 8 - (-x \& 7) \\ ((x+7) \& 7) + 1 & ((x+7) | -8) + 9 & ((x-1) \& 8) + x & -(-x | -8) \end{array}$$

2.12. Предикаты сравнения

Предикаты сравнения представляют собой функции, которые сравнивают две величины и возвращают однобитовый результат, равный 1, если проверяемое отношение истинно, и 0, если ложно. В этом разделе приводятся несколько методов вычисления результата сравнения с размещением его в бите знака (эти методы не используют команд ветвления). Чтобы получить значение 1 или 0, используемое в некоторых языках (например, в C), после вычисления следует использовать команду *сдвига вправо* на 31 бит. Для получения результата $-1/0$, используемого в некоторых других языках (например, Basic), выполняется команда *знакового сдвига вправо* на 31 бит.

На таких машинах, как MIPS и наша модель RISC, есть команды сравнения, которые непосредственно вычисляют многие отношения и помещают однобитовый результат 0/1 в регистр общего назначения, а потому для таких машин приводимые ниже формулы не представляют особого интереса.

$$\begin{aligned}
 x = y : & \quad \text{abs}(x - y) - 1 \\
 & \quad \text{abs}(x - y + 0x80000000) \\
 & \quad \text{nlz}(x - y) \ll 26 \\
 & \quad - \left(\text{nlz}(x - y) \gg 5 \right) \\
 & \quad \neg(x - y | y - x) \\
 x \neq y : & \quad \text{nabs}(x - y) \\
 & \quad \text{nlz}(x - y) - 32 \\
 & \quad x - y | y - x \\
 x < y : & \quad (x - y) \oplus [(x \oplus y) \& ((x - y) \oplus x)] \\
 & \quad (x \& \neg y) | ((x = y) \& (x - y)) \\
 & \quad \text{nabs}(\text{doz}(y, x)) \qquad [43] \\
 x \leq y : & \quad (x | \neg y) \& ((x \oplus y) | \neg(y - x)) \\
 & \quad \left((x = y) \gg 1 \right) + (x \& \neg y) \qquad [43] \\
 x <^{\circ} y : & \quad (\neg x \& y) | ((x = y) \& (x - y)) \\
 & \quad (\neg x \& y) | ((\neg x | y) \& (x - y)) \\
 x \leq^{\circ} y : & \quad (\neg x | y) \& ((x \oplus y) | \neg(y - x))
 \end{aligned}$$

При вычислении функций сравнения двух операндов весьма полезной оказывается машинная команда “nabs”, вычисляющая отрицательное абсолютное значение. В отличие от функции абсолютного значения, она никогда не приводит к переполнению. Если среди команд машины нет функции “nabs”, но есть более привычная команда “abs”, то вместо $\text{nabs}(x)$ можно использовать $-\text{abs}(x)$. Если x — максимальное отрицательное число, дважды возникает переполнение, но конечный результат получается правильный (предполагается, что абсолютное значение и отрицание максимального отрицательного числа дают одинаковые результаты). Поскольку на многих машинах нет ни команды “abs”, ни команды “nabs”, приведены и альтернативные формулы, в которых эти функции не применяются.

Функция “nlz” возвращает количество ведущих нулевых битов аргумента. Описание функции doz (difference or zero — *разность или ноль*) приводится в разделе 2.19. Для $x > y$, $x \geq y$ и прочих отношений необходимо в соответствующих формулах поменять

местами x и y . Команду сложения с $0x80000000$ можно заменить любой другой командой, которая инвертирует значение старшего бита.

Другой класс формул основан на том, что предикат $x < y$ вычисляется по знаку величины $x/2 - y/2$, а вычисление данной разности не может привести к переполнению. В тех случаях, когда сдвиги приводят к потере информации, результат можно уточнить путем вычитания единицы.

$$x < y: \quad \left(x \overset{\cdot}{\gg} 1 \right) - \left(y \overset{\cdot}{\gg} 1 \right) - \left(-x \& y \& 1 \right)$$

$$x \overset{\cdot}{<} y: \quad \left(x \overset{\cdot}{\gg} 1 \right) - \left(y \overset{\cdot}{\gg} 1 \right) - \left(-x \& y \& 1 \right)$$

Большинство машин вычисляют эти отношения с помощью семи команд (или шести при наличии команды *и-не*), что ничем не лучше результата, достигнутого в предыдущих формулах (от пяти до семи команд в зависимости от полноты используемого множества команд).

Формулы, в которых используется функция "nlz", взяты из [102]. Его формула для $x = y$ оказывается особенно полезной, поскольку небольшое ее изменение позволяет вычислить результат в виде 1 или 0 всего за три команды.

$$\text{nlz}(x - y) \overset{\cdot}{\gg} 5$$

Команды знакового сравнения с 0 встречаются достаточно часто, чтобы заслужить отдельного рассмотрения. Ниже приводится несколько формул, которые в основном непосредственно выведены из предыдущих выражений. Так же, как и ранее, результат сравнения сохраняется в позиции бита знака.

$$\begin{aligned}
 x = 0: & \quad \text{abs}(x) - 1 \\
 & \quad \text{abs}(x + 0x80000000) \\
 & \quad \text{nlz}(x) \ll 26 \\
 & \quad - \left(\text{nlz}(x) \overset{\cdot}{\gg} 5 \right) \\
 & \quad \neg(x | -x) \\
 & \quad -x \& (x - 1) \\
 x \neq 0: & \quad \text{nabs}(x) \\
 & \quad \text{nlz}(x) - 32 \\
 & \quad x | -x \\
 & \quad \left(x \overset{\cdot}{\gg} 1 \right) - x \quad [20] \\
 x < 0: & \quad x \\
 x \leq 0: & \quad x | (x - 1) \\
 & \quad x | \neg -x
 \end{aligned}$$

$$\begin{aligned}
 x > 0: & \quad x \oplus \text{nabs}(x) \\
 & \quad \left(x \gg 1 \right) - x \\
 & \quad -x \& \neg x \\
 x \geq 0: & \quad \neg x
 \end{aligned}$$

Знаковое сравнение можно получить из беззнакового, увеличивая сравниваемые величины на 2^{31} и рассматривая получающиеся числа как беззнаковые. Обратное преобразование также справедливо². Таким образом, получаем

$$\begin{aligned}
 x < y &= x + 2^{31} < y + 2^{31} \\
 x <^n y &= x - 2^{31} < y - 2^{31}
 \end{aligned}$$

Аналогичные соотношения выполняются для \leq , \leq^n и т.д. В данном случае прибавление и вычитание числа 2^{31} (а также *исключающее или*) — эквивалентные действия, так как все эти операции инвертируют значение знакового разряда.

Еще один путь получения знакового сравнения из беззнакового основан на том, что если x и y имеют одинаковые знаки, то $x < y = x <^n y$, в то время как если их знаки различны, то $x < y = x >^n y$ [76]. В этом случае также справедливы обратные преобразования; таким образом, получаем следующее.

$$\begin{aligned}
 x < y &= \left(x <^n y \right) \oplus x_{31} \oplus y_{31} \\
 x <^n y &= \left(x < y \right) \oplus x_{31} \oplus y_{31}
 \end{aligned}$$

Здесь x_{31} и y_{31} представляют собой знаковые биты x и y соответственно. Аналогичные выражения получаются для отношений \leq , \leq^n и др.

Применение описанных способов позволяет вычислить все стандартные команды сравнения (кроме $=$ и \neq) через другие команды сравнения с использованием не более трех дополнительных команд на большинстве машин. Например, возьмем в качестве исходной команды сравнения отношение $x \leq^n y$, поскольку реализовать его проще других команд (как бит переноса $y - x$). Тогда остальные команды сравнения можно реализовать следующим образом.

$$\begin{aligned}
 x < y &= \neg \left(y + 2^{31} \leq^n x + 2^{31} \right) \\
 x \leq y &= x + 2^{31} \leq^n y + 2^{31} \\
 x > y &= \neg \left(x + 2^{31} \leq^n y + 2^{31} \right) \\
 x \geq y &= y + 2^{31} \leq^n x + 2^{31}
 \end{aligned}$$

² Это очень полезно в случае сравнений в Java, где нет беззнаковых целых чисел.

$$x \overset{\cdot}{<} y = \neg \left(y \overset{\cdot}{\leq} x \right)$$

$$x \overset{\cdot}{>} y = \neg \left(x \overset{\cdot}{\leq} y \right)$$

$$x \overset{\cdot}{\geq} y = y \overset{\cdot}{\leq} x$$

Команды сравнения и бит переноса

Если процессор может легко поместить бит переноса в регистр общего назначения, то для некоторых операций отношения это позволяет получить очень лаконичный код. Ниже приводятся выражения такого рода для некоторых отношений. Запись $\text{carry}(\text{выражение})$ означает, что разряд переноса генерируется при выполнении команды *выражение*. Предполагается, что бит переноса для разности $x - y$ тот же, что и на выходе сумматора для выражения $x + \bar{y} + 1$, и является дополнением “заема”.

$$x = y: \text{carry}(0 - (x - y)) \text{ или}$$

$$\text{carry}((x + \bar{y}) + 1), \text{ или}$$

$$\text{carry}((x - y - 1) + 1)$$

$$x \neq y: \text{carry}((x - y) - 1) = \text{carry}((x - y) + (-1))$$

$$x < y: \neg \text{carry}(((x + 2^{31}) - (y + 2^{31})))$$

$$\neg \text{carry}(x - y) \oplus x_{31} \oplus y_{31}$$

$$x \leq y: \text{carry}(((y + 2^{31}) - (x + 2^{31})))$$

$$\text{carry}(y - x) \oplus x_{31} \oplus y_{31}$$

$$x \overset{\cdot}{<} y: \neg \text{carry}(x - y)$$

$$x \overset{\cdot}{\leq} y: \text{carry}(y - x)$$

$$x = 0: \text{carry}(0 - x), \text{ или } \text{carry}(\bar{x} + 1)$$

$$x \neq 0: \text{carry}(x - 1) = \text{carry}(x + (-1))$$

$$x < 0: \text{carry}(x + x)$$

$$x \leq 0: \text{carry}(2^{31} - (x + 2^{31}))$$

Для $x > y$ используется дополнение выражения для $x \leq y$; то же справедливо и для других типов отношения “больше”.

При вычислении условных выражений на IBM RS/6000 и близком к нему PowerPC была использована программа GNU Superoptimizer [35]. У RS/6000 имеются команды $\text{abs}(x)$, $\text{pabs}(x)$, $\text{doz}(x, y)$ и ряд различных команд сложения и вычитания с использованием переноса. Как выяснилось, RS/6000 может вычислить все целочисленные условные выражения с помощью не более чем трех элементарных (выполняющихся за один такт)

команд — результат, удививший даже самих создателей машины. В состав “всех условных выражений” входят шесть команд знакового сравнения и четыре команды беззнакового сравнения, а также их версии, в которых второй операнд равен нулю (при этом каждая из команд имеет два варианта — возвращающий результат 1/0 и возвращающий результат -1/0). Power PC, где нет функций $\text{abs}(x)$, $\text{pabs}(x)$ и $\text{doz}(x, y)$, вычисляет все перечисленные условные выражения не более чем за четыре элементарные команды.

Вычисление отношений

Большинство компьютеров возвращают однобитовый результат операции сравнения, который обычно размещается в “регистре условий”, а на некоторых машинах (как, например, в нашей RISC-модели) — в регистре общего назначения. В любом случае вычисление отношений зачастую реализуется вычитанием операндов и некоторыми действиями над битами полученного результата, чтобы получить окончательный однобитовый результат работы команды сравнения.

Рассмотрим логику описанных выше действий. Пусть вместо разности $x - y$ компьютер вычисляет сумму $x + \bar{y} + 1$ и в результате этих вычислений получаются следующие величины.

- C_0 Перенос из старшего разряда в разряд переноса
- C_1 Перенос в старший разряд
- N Бит знака результата
- Z Величина, равная 1, если полученный результат (за исключением C_0) нулевой, и 0 — в противном случае

Тогда в системе обозначений булевой алгебры для различных отношений получим приведенные ниже выражения (расположение рядом двух величин подразумевает операцию и , а оператор “+” означает операцию или).

$$\begin{aligned}
 V: & C_1 \oplus C_0 \text{ (знаковое переполнение)} \\
 x = y: & Z \\
 x \neq y: & \bar{Z} \\
 x < y: & N \oplus V \\
 x \leq y: & (N \oplus V) + Z \\
 x > y: & (N \text{ и } V) \bar{Z} \\
 x \geq y: & N \text{ и } V \\
 x \overset{\circ}{<} y: & \bar{C}_0 \\
 x \overset{\circ}{\leq} y: & \bar{C}_0 + Z \\
 x \overset{\circ}{>} y: & C_0 \bar{Z} \\
 x \overset{\circ}{\geq} y: & C_0
 \end{aligned}$$

2.13. Обнаружение переполнения

“Переполнение” означает, что результат арифметической операции либо слишком велик, либо слишком мал, чтобы его можно было корректно представить в целевом регистре. В этом разделе обсуждаются методы, которые позволяют определить, когда возникнет переполнение, не используя при этом биты состояния, предназначенные для этой цели. Эти методы имеют особое значение на тех машинах (например, MIPS), где биты состояния отсутствуют, и даже при наличии таких битов — так как обращение к ним из языков высокого уровня, как правило, затруднено или попросту невозможно.

Знаковое сложение и вычитание

Если в результате сложения и вычитания целых чисел произошло переполнение, то, как правило, старший бит результата отбрасывается и сумматор выдает только младшие биты. Переполнение при сложении целых знаковых чисел возникает тогда и только тогда, когда операнды имеют одинаковые знаки, а знак суммы противоположен знаку операндов. Удивительно, но это правило справедливо даже тогда, когда в сумматоре был выполнен перенос, т.е. при вычислении суммы $x + y + 1$. Данное правило играет важную роль при сложении знаковых целых чисел, состоящих из нескольких слов, так как в этом случае последнее сложение представляет собой знаковое сложение двух полных слов и бита переноса, значение которого может быть равно 0 или +1.

Для доказательства этого правила предположим, что слагаемые x и y представляют собой целые знаковые значения, состоящие из одного слова, а бит переноса c равен 0 или 1. Предположим также для простоты, что сложение выполняется на 4-разрядной машине. Если x и y имеют разные знаки, то $-8 \leq x \leq -1$ и $0 \leq y \leq 7$ (аналогичные границы получаются и при отрицательном y и неотрицательном x). В любом случае после сложения этих неравенств и необязательного добавления c получим $-8 \leq x + y + c \leq 7$.

Как видите, сумма может быть представлена как 4-разрядное целое число со знаком; следовательно, при сложении операндов с разными знаками переполнения не будет.

Теперь предположим, что x и y имеют одинаковые знаки. Здесь возможны два случая.

$$\begin{array}{ll} \text{(a)} & \text{(б)} \\ -8 \leq x \leq -1 & 0 \leq x \leq 7 \\ -8 \leq y \leq -1 & 0 \leq y \leq 7 \end{array}$$

Таким образом:

$$\begin{array}{ll} \text{(a)} & \text{(б)} \\ -16 \leq x + y + c \leq -1 & 0 \leq x + y + c \leq 15 \end{array}$$

Переполнение возникает тогда, когда сумма не может быть представлена в виде 4-битового целого числа со знаком, т.е. когда

$$\begin{array}{ll} \text{(a)} & \text{(б)} \\ -16 \leq x + y + c \leq -9 & 8 \leq x + y + c \leq 15 \end{array}$$

В случае (а) старший разряд 4-битовой суммы равен 0, т.е. противоположен знакам x и y . В случае (б) старший бит 4-битовой суммы равен 1, т.е. также противоположен знакам x и y .

При вычитании целых чисел, состоящих из нескольких слов, нас интересует разность $x - y - c$, где c равно 0 или 1, причем 1 означает, что был выполнен заем. Проводя анализ, аналогичный только что рассмотренному, можно увидеть, что при вычислении выражения $x - y - c$ переполнение будет тогда и только тогда, когда x и y имеют разные знаки и знак разности противоположен знаку x (или, что то же самое, имеет знак, одинаковый со знаком y).

Это приводит нас к следующему выражению предиката наличия переполнения (помещающему результат в бит знака). Беззнаковый или знаковый сдвиг вправо на 31 бит дает нам результат в виде чисел 1/0 или -1/0.

$$\begin{array}{ll} x + y + c & x - y - c \\ (x \equiv y) \& ((x + y + c) \oplus x) & (x \oplus y) \& ((x - y - c) \oplus x) \\ ((x + y + c) \oplus x) \& ((x + y + c) \oplus y) & ((x - y - c) \oplus x) \& ((x - y - c) \equiv y) \end{array}$$

Если взять вторую формулу для суммы и первую формулу для разности (не содержащую операции эквивалентности), то при использовании базового множества RISC-команд для проверки на переполнение потребуется выполнить три команды в дополнение к тем, которые вычисляют саму сумму (или разность). Четвертой может стать команда перехода к коду обработки переполнения.

Если при возникновении переполнения генерируется прерывание, то может потребоваться проверка, не вызовет ли сложение или вычитание переполнения (причем сама проверка ни в коем случае не должна вызвать переполнения). Для этого можно воспользоваться приведенными ниже формулами (в которых нет команд ветвления и которые заведомо не генерируют переполнения).

$$\begin{array}{ll} x + y + c & x - y - c \\ z \leftarrow (x \equiv y) \& 0x80000000 & z \leftarrow (x \oplus y) \& 0x80000000 \\ z \& (((x \oplus z) + y + c) \equiv y) & z \& (((x \oplus z) - y - c) \oplus y) \end{array}$$

Присваивание в левом столбце устанавливает значение z равным $0x80000000$, если знаки x и y одинаковы, и 0 в противном случае. Таким образом, во втором выражении слагаемые имеют противоположные знаки и переполнение возникнуть не может. Если x и y неотрицательны, значение знакового бита в результате во втором выражении будет равным 1 тогда и только тогда, когда $(x - 2^{31}) + y + c \geq 0$, т.е. если $x + y + c \geq 2^{31}$, что и представляет собой условие возникновения переполнения при вычислении $x + y + c$. Если x и y отрицательны, то значение бита знака во втором выражении будет равным 1 тогда и только тогда, когда $(x + 2^{31}) + y + c < 0$, т.е. если $x + y + c < -2^{31}$, что также представляет собой условие переполнения при вычислении $x + y + c$. Выполнение z гарантирует корректность результата (равный 0 знаковый бит), если x и y имеют разные знаки. Условия переполнения для разности (приведенные в правом столбце) получаются

аналогично. Для реализации этих формул понадобится выполнить девять команд из базового множества RISC.

Может показаться, что в случае легкодоступности значения разряда переноса это может помочь в разработке эффективного алгоритма проверки на переполнение в случае знакового сложения, но на самом деле это не так. Однако есть еще один метод.

Если x — целое знаковое число, то $x + 2^{31}$ представляет собой корректное беззнаковое число, которое получается в результате инвертирования старшего бита x . При сложении положительных знаковых чисел переполнение возникает тогда, когда $x + y \geq 2^{31}$, т.е. если $(x + 2^{31}) + (y + 2^{31}) \geq 3 \cdot 2^{31}$. Последнее выражение означает, что при сложении беззнаковых операндов был перенос (так как сумма не меньше 2^{32}) и старший бит суммы равен 1. Аналогично при сложении отрицательных чисел переполнение возникает тогда, когда значение разряда переноса и значение старшего разряда суммы равны 0.

Это дает нам следующий алгоритм обнаружения переполнения при знаковом сложении.

Вычисляем $(x \oplus 2^{31}) + (y \oplus 2^{31})$ и получаем сумму s и значение бита переноса c .

Переполнение возникает тогда и только тогда, когда значение разряда переноса c равно значению старшего разряда суммы s .

Полученная сумма представляет собой корректное значение знакового сложения, поскольку инвертирование старшего бита обоих операндов не изменяет их суммы.

При вычитании используется такой же алгоритм, однако на первом шаге вместо суммы вычисляется разность. Мы считаем, что перенос в данном случае генерируется при вычислении $x - y$ как $x + \bar{y} + 1$. Полученная разность представляет собой корректный результат знакового вычитания.

Несмотря на то что рассмотренные методы крайне интересны, на многих компьютерах они работают не так эффективно, как методы без использования переноса (например, вычисление переполнения как $(x \equiv y) \& (s \oplus x)$ для сложения и $(x \oplus y) \& (d \oplus x)$ для вычитания, где s и d — соответственно сумма и разность x и y).

Установка переполнения при знаковом сложении и вычитании

Зачастую переполнение при сложении знаковых операндов устанавливается посредством правила: "перенос в знаковый разряд отличается от переноса из знакового разряда". Как ни удивительно, но эта логика позволяет корректно определять переполнение как при сложении, так и при вычитании, считая, что разность $x - y$ вычисляется как $x + \bar{y} + 1$. Более того, это утверждение корректно независимо от того, был ли перенос или заем или не был. Конечно, то, что перенос в знаковый разряд легко вычислить, вовсе не означает, что существуют эффективные методы определения знакового переполнения. В случае сложения и вычитания перенос/заем формируется в знаковом разряде после вычисления приведенных ниже выражений (здесь c равно 0 или 1).

$$\begin{array}{l} \text{Перенос} \\ (x + y + c) \oplus x \oplus y \end{array}$$

$$\begin{array}{l} \text{Заем} \\ (x - y - c) \oplus x \oplus y \end{array}$$

В действительности для каждого разряда i эти выражения дают перенос/заем в разряд i .

Беззнаковое сложение/вычитание

При беззнаковом сложении/вычитании для вычисления предикатов переполнения можно использовать следующий метод без ветвления, результат работы которого помещается в знаковый разряд. Методы, использующие команды сдвига вправо, эффективно работают только тогда, когда известно, что $c = 0$. Выражения в скобках вычисляют перенос или заем, который появляется из младшего разряда.

Беззнаковое вычисление $x + y + c$

$$(x \& y) | ((x | y) \& -(x + y + c)) \\ (x \gg 1) + (y \gg 1) + \left[\left((x \& y) | ((x | y) \& c) \right) \& 1 \right]$$

Беззнаковое вычисление $x - y - c$

$$(-x \& y) | ((x \equiv y) \& (x - y - c)) \\ (-x \& y) | ((-x | y) \& (x - y - c)) \\ (x \gg 1) - (y \gg 1) - \left[\left((-x \& y) | ((-x | y) \& c) \right) \& 1 \right]$$

В [84] для беззнаковых сложения и вычитания есть существенно более простые формулы, использующие команды сравнения. При беззнаковом сложении переполнение (перенос) возникает, если полученная сумма меньше (при беззнаковом сравнении) любого из операндов. Эта и подобные формулы приведены ниже. К сожалению, нет возможности обобщить эти формулы для использования в них переменной c , которая представляет собой перенос или заем. Вместо этого приходится сначала проверять значение c , а затем, в зависимости от того, равно оно 0 или 1, использовать различные типы сравнений (напомним, что в приведенной ниже таблице вычисления сумм и разностей беззнаковые).

$x + y$	$x + y + 1$	$x - y$	$x - y - 1$
$\overset{\cdot}{-}x < \overset{\cdot}{y}$	$\overset{\cdot}{-}x \leq \overset{\cdot}{y}$	$\overset{\cdot}{x} < \overset{\cdot}{y}$	$\overset{\cdot}{x} \leq \overset{\cdot}{y}$
$x + y < x$	$x + y + 1 \leq x$	$x - y > x$	$x - y - 1 \geq x$

В каждом из приведенных случаев перед выполнением сложения/вычитания вычисляется первая формула, которая проверяет, возможно ли переполнение, не вызывая при этом самого переполнения. Вторая формула вычисляется после выполнения сложения/вычитания, при которых возможно переполнение.

Для случая знакового сложения/вычитания аналогичного (с использованием сравнений) простого метода проверки, похоже, не существует.

Умножение

При умножении переполнение означает, что результат не помещается в 32 разрядах (результат как знакового, так и беззнакового умножения 32-битовых чисел всегда может быть представлен в виде 64-битового числа). Если доступны старшие 32 разряда произведения, то определить, имеет ли место переполнение, очень просто. Пусть $hi(x \times y)$ и $lo(x \times y)$ — старшая и младшая половины 64-битового произведения. Тогда условия переполнения можно записать следующим образом [84].

Беззнаковое произведение $x \times y$

$$hi(x \times y) \neq 0$$

Знаковое произведение $x \times y$

$$hi(x \times y) \neq \left(lo(x \times y) \gg 31 \right)$$

Один из способов проверки наличия переполнения состоит в контрольном делении. Однако при этом необходимо следить, чтобы не произошло деления на 0, а кроме того, при знаковом умножении возникают дополнительные сложности.

Переполнение при умножении имеет место, если приведенные ниже выражения истинны.

Беззнаковое умножение

$$z \leftarrow x * y$$

$$y \neq 0 \bar{\&} z + y \neq x$$

Знаковое умножение

$$z \leftarrow x * y$$

$$(y < 0 \& x = -2^{31}) | (y \neq 0 \bar{\&} z + y \neq x)$$

Затруднения возникают при $x = -2^{31}$ и при $y = -1$. В этом случае при умножении будет переполнение, но результат может получиться равным -2^{31} . Такой результат приводит к переполнению при делении, и, таким образом, оказывается возможным любой результат (на некоторых машинах). Поэтому данный частный случай должен проверяться отдельно, что и делается добавлением члена ($y < 0 \& x = -2^{31}$). Чтобы в приведенных выше выражениях предотвратить деление на 0, используется оператор “условное и” (в языке C — оператор $\&\&$).

Можно также использовать проверку делением без выполнения самого умножения (т.е. без генерации переполнения). При беззнаковом умножении переполнение возникает тогда и только тогда, когда $xu > 2^{32} - 1$, или $x > ((2^{32} - 1)/y)$, или, поскольку x — целое число, $x > \lfloor (2^{32} - 1)/y \rfloor$. В обозначениях компьютерной арифметики это выглядит следующим образом.

$$y \neq 0 \bar{\&} x > \left(0xFFFFFFFF + y \right)$$

Для целых знаковых чисел переполнение в результате умножения x и y определить сложнее. Если операнды имеют одинаковые знаки, то переполнение будет тогда и только тогда, когда $xu > 2^{31} - 1$. Если x и y имеют противоположные знаки, то переполнение будет тогда и только тогда, когда $xu < -2^{31}$. Эти условия могут быть проверены с помощью знакового деления (табл. 2.2). Как видно из таблицы, при проверке необходимо учесть четыре возможных случая. Вследствие проблем переполнения и представления числа $+2^{31}$ с унификацией этих выражений могут возникнуть трудности.

ТАБЛИЦА 2.2. ПРОВЕРКА НА ПЕРЕПОЛНЕНИЕ ПРОИЗВЕДЕНИЯ ЗНАКОВЫХ ОПЕРАНДОВ

	$y > 0$	$y \leq 0$
$x > 0$	$x > 0x7FFFFFFF + y$	$y < 0x80000000 + x$
$x \leq 0$	$x < 0x80000000 + y$	$x \neq 0 \ \& \ y < 0x7FFFFFFF + x$

В случае доступности беззнакового деления тест можно упростить. Если использовать абсолютные значения x и y , которые корректно представимы в виде целых чисел без знака, то весь тест можно записать так, как показано ниже. Переменная c имеет значение $2^{31} - 1$, если x и y имеют одинаковые знаки, и 2^{31} в противном случае.

$$c \leftarrow \left((x = y) \gg 31 \right) + 2^{31}$$

$$x \leftarrow \text{abs}(x)$$

$$y \leftarrow \text{abs}(y)$$

$$y \neq 0 \ \& \ x > \left(c + y \right)$$

Чтобы выяснить, будет ли переполнение при вычислении произведения $x * y$, можно использовать команду, вычисляющую количество ведущих нулей у операндов. Способ, основанный на использовании этой функции, позволяет более точно определить условие переполнения. Сначала рассмотрим умножение беззнаковых чисел. Легко показать, что если x и y — 32-разрядные величины, содержащие m и n ведущих нулей соответственно, то в 64-разрядном произведении будет $m + n$ или $m + n + 1$ ведущих нулей (или 64, если $x = 0$ или $y = 0$). Переполнение возникает тогда, когда в 64-разрядном произведении количество ведущих нулей оказывается меньше 32. Из этого следует:

если $\text{nlz}(x) + \text{nlz}(y) \geq 32$, переполнения точно не будет;
 если $\text{nlz}(x) + \text{nlz}(y) \leq 30$, переполнение точно будет.

Если же $\text{nlz}(x) + \text{nlz}(y) = 31$, то переполнение может быть, а может и не быть. Чтобы узнать, будет ли переполнение в этом частном случае, вычисляется значение $t = x \lfloor y/2 \rfloor$ (переполнения при этом вычислении не возникает). Поскольку xu равно $2t$ (или $2t + x$, если y нечетно), переполнение при вычислении xu будет в том случае, если $t \geq 2^{31}$. Все это приводит к коду, показанному в листинге 2.2, который вычисляет значение произведения xu с проверкой переполнения (в этом случае выполняется переход к метке `overflow`).

ЛИСТИНГ 2.2. Беззнаковое умножение с проверкой переполнения

```
unsigned x, y, z, m, n, t;

m = nlz(x);
n = nlz(y);
if (m + n <= 30) goto overflow;
t = x*(y >> 1);
if ((int)t < 0) goto overflow;
```

```

z = t * 2;
if (y & 1) {
    z = z + x;
    if (z < x) goto overflow;
}
// В z содержится корректное произведение x и y.

```

При знаковом умножении условие переполнения можно определить через количество ведущих нулей положительных аргументов и количество ведущих единиц отрицательных аргументов. Пусть

$$m = \text{nlz}(x) + \text{nlz}(\bar{x}) \text{ и } n = \text{nlz}(y) + \text{nlz}(\bar{y}).$$

Тогда,

если $m + n \geq 34$, переполнения точно не будет;
 если $m + n \leq 31$, переполнение точно будет.

Остались два неоднозначных случая — 32 и 33. В случае $m + n = 33$ переполнение возникает только тогда, когда оба аргумента отрицательны и произведение равно в точности 2^{31} (машинный результат равен -2^{31}), так что распознать эту ситуацию можно с помощью проверки корректности знака произведения (т.е. переполнение будет, если $m \oplus n \oplus (m * n) < 0$). Для случая $m + n = 32$ определить наличие переполнения не так просто.

Не будем больше задерживаться на этой теме; заметим только, что определение наличия переполнения может основываться на значении $\text{nlz}(\text{abs}(x)) + \text{nlz}(\text{abs}(y))$, но и при этом возможны неоднозначные случаи, когда значение данного выражения равно 31 или 32.

Деление

При знаковом делении $x + y$ переполнение возникает, если истинно выражение

$$y = 0 | (x = 0x80000000 \& y = -1)$$

Большинство компьютеров сообщают о переполнении (или генерируют прерывание) при возникновении неопределенности типа $0 + 0$.

Простой код для вычисления этого выражения, включая завершающую команду перехода к коду обработки переполнения, состоит из семи команд (три из которых являются командами ветвления). Создается впечатление, что этот код улучшить нельзя, однако рассмотрим некоторые возможности для этого.

$$[\text{abs}(y \oplus 0x80000000) | (\text{abs}(x) \& \text{abs}(y \neq 0x80000000))] < 0$$

Здесь сначала вычисляется выражение в скобках, а затем выполняется ветвление, если полученный результат меньше 0. Это выражение на машине, имеющей все необходимые команды (включая сравнение с 0), вычисляется примерно за девять команд с учетом загрузки константы и завершающего перехода.

Еще один метод предполагает первоначальное вычисление значения

$$z \leftarrow (x \oplus 0x80000000) | (y + 1)$$

(что на большинстве машин выполняется за три команды), а затем выполнение проверки и ветвления по условию $y = 0 \mid z = 0$ одним из следующих способов.

$$((y | -y) \& (z | -z)) \geq 0$$

$$(\text{nabs}(y) \& \text{nabs}(z)) \geq 0$$

$$\left((nlz(y) | nlz(z)) \gg 5 \right) \neq 0$$

Вычисление этих выражений может быть реализовано за девять, семь и восемь команд соответственно (при наличии всех необходимых для вычисления команд). Для PowerPC наиболее эффективным оказывается последний из перечисленных методов.

При беззнаковом делении $x \div y$ переполнение возникает тогда и только тогда, когда $y = 0$.

У некоторых машин имеется команда “длинного деления” (с. 219), и с помощью элементарных действий можно предсказать, возникнет ли переполнение. Мы обсудим этот вопрос в терминах команды, которая делит двойное слово на полное слово, приводя к полному слову частного и, возможно, полному слову остатка.

Такая команда приводит к переполнению тогда, когда либо делитель равен 0, либо частное не может быть представлено 32 битами. Обычно при таком переполнении оказываются неверными и частное, и остаток. Остаток не может вызвать переполнения в том смысле, что он слишком велик для размещения в 32 битах (по величине он меньше делителя), так что проверка корректности остатка — то же самое, что и проверка корректности частного.

Мы считаем, что либо машина имеет 64-битовые регистры общего назначения, либо она имеет 32-битовые регистры и без проблем в состоянии выполнять элементарные операции (сдвиг, сложение и т.д.) с 64-битовыми величинами. Например, компилятор может реализовывать целочисленный тип, размещающийся в двойном слове.

В беззнаковом случае проверка тривиальна: для $x + y$, где x — двойное слово, а y — полное слово, деление не приводит к переполнению тогда (и только тогда), когда любое из следующих эквивалентных выражений истинно.

$$y \neq 0 \& x < (y \ll 32)$$

$$y \neq 0 \& \left(x \gg 32 \right) < y$$

На 32-разрядной машине выполнять сдвиги не требуется; достаточно просто сравнить y с регистром, содержащим старшую половину x . Чтобы гарантировать корректность результата на 64-разрядной машине, необходимо также проверить, является ли делитель y 32-разрядной величиной (например, проверить $\left(y \gg 32 \right) = 0$).

Знаковый случай более интересен. Сначала нужно проверить, что $y \neq 0$ и, на 64-разрядной машине, что значение y корректно представимо 32 битами (проверить, что $((y \ll 32) \gg 32) = y$). В предположении, что все указанные тесты выполнены, в приведенной далее таблице показано, как можно точно определить, представимо ли частное 32 битами, рассматривая четыре отдельных случая, связанных со знаками делимого и делителя. Выражения в таблице используют обычную, а не компьютерную арифметику.

$x \geq 0, y \geq 0$	$x \geq 0, y < 0$	$x < 0, y > 0$	$x < 0, y < 0$
$\lfloor x/y \rfloor < 2^{31}$	$\lfloor x/y \rfloor \geq -2^{31}$	$\lfloor x/y \rfloor \geq -2^{31}$	$\lfloor x/y \rfloor < 2^{31}$
$x/y < 2^{31}$	$\lfloor x/y \rfloor > -2^{31} - 1$	$\lfloor x/y \rfloor > -2^{31} - 1$	$x/y < 2^{31}$
$x < 2^{31} y$	$x/y > -2^{31} - 1$	$x/y > -2^{31} - 1$	$x > 2^{31} y$
	$x < -2^{31} y - y$	$x > -2^{31} y - y$	$-x < 2^{31} (-y)$
	$x < 2^{31} (-y) + (-y)$	$-x < 2^{31} y + y$	

В каждом столбце каждое отношение следует из приведенного выше по правилу “тогда и только тогда, когда”. Для устранения функций “пол” и “потолок” были использованы некоторые из соотношений теоремы D1 (с. 209).

В качестве примера интерпретации данной таблицы рассмотрим крайний слева столбец. Он применяется, когда $x \geq 0$ и $y > 0$. В этом случае частное равно $\lfloor x/y \rfloor$, и для того, чтобы оно могло быть представлено 32 битами, оно должно быть строго меньше 2^{31} . Отсюда следует, что действительное число x/y должно быть меньше 2^{31} или x должно быть меньше $2^{31} y$. Эта проверка может быть реализована путем сдвига y влево на 31 разряд и сравнения полученной величины с x .

Когда знаки x и y различны, частное при обычном делении равно $\lfloor x/y \rfloor$. Поскольку частное отрицательно, оно может достигать значения -2^{31} .

В нижней строке все сравнения имеют один и тот же вид — “меньше, чем”. Поскольку имеется возможность, что x является наибольшим отрицательным числом, в третьем и четвертом столбцах требуется применять беззнаковое сравнение. В первых двух столбцах сравниваются значения с нулевым старшим битом, так что здесь можно использовать как знаковое, так и беззнаковое сравнения.

Эти проверки, конечно, могут быть реализованы с применением условного ветвления для разделения четырех случаев, выполнения указанных арифметических действий и еще одного ветвления к коду для случая переполнения и для отсутствия такового. Однако количество ветвлений можно снизить за счет использования $-y$ при отрицательном значении y (аналогично и для x). Следовательно, проверки можно сделать более унифицированными, если воспользоваться абсолютными значениями x и y . Применение, кроме того, стандартного подхода для необязательного сложения во втором и третьем столбцах приводит к следующей схеме.

$$\begin{aligned}
 x' &= |x| \\
 y' &= |y| \\
 \delta &= \left((x \oplus y) \gg 63 \right) \& y' \\
 \text{if } \left(x' < (y' \ll 31) + \delta \right) &\text{ then \{переполнения нет\}}
 \end{aligned}$$

Используя трехкомандный метод вычисления абсолютного значения (с. 39) на 64-разрядной машине с базовым набором команд RISC, все вычисления можно выполнить с помощью 12 команд, не считая условного перехода.

2.14. Флаги условий после сложения, вычитания и умножения

Во многих процессорах в результате выполнения арифметических операций над целочисленными данными устанавливаются так называемые флаги условий, или код условия (condition code). Зачастую имеется только одна команда *сложения*, и указанная характеристика отражает результат как для беззнаковой, так и для знаковой интерпретации операндов (но не для смешанных типов). Как правило, на флаги условий оказывают влияние следующие события:

- был ли перенос (беззнаковое переполнение);
- было ли знаковое переполнение;
- является ли полученный 32-битовый результат, интерпретируемый как знаковое целое число в дополнительном коде, без учета переноса и переполнения, отрицательным, положительным или нулевым.

На некоторых старых моделях машин в случае переполнения обычно выдается сообщение о том, какой именно результат (в случае *сложения* и *вычитания* — 33-битовый) получен: положительный, отрицательный или нулевой. Однако использовать такое сообщение в компиляторах языков высокого уровня — непростая задача, так что данная возможность нами не рассматривается.

При сложении возможны только 9 из 12 комбинаций этих трех событий. Невозможны следующие три комбинации: “нет переноса, переполнение, результат > 0 ”, “нет переноса, переполнение, результат $= 0$ ” и “перенос, переполнение, результат < 0 ”. Таким образом, для кода условия требуется четыре бита. Две из комбинаций уникальны в том смысле, что могут произойти только при единственном значении исходных аргументов: комбинация “нет переноса, нет переполнения, результат $= 0$ ” возникает только при сложении 0 с 0; комбинация “перенос, переполнение, результат $= 0$ ” возникает только при сложении двух максимальных отрицательных чисел.

В случае вычитания полагаем, что для вычисления разности $x - y$ компьютер на самом деле вычисляет сумму $x + \bar{y} + 1$ с переносом, который образуется так же, как и при *сложении* (в этой схеме понятие “перенос” трактуется наоборот: перенос, равный 1, означает, что результат помещается в одно слово, а перенос, равный 0, означает, что ре-

зультат в одно слово не помещается). Тогда при вычитании могут произойти только семь из всех возможных комбинаций событий. Невозможны те же три комбинации, что и при сложении, а кроме того, еще две комбинации: “нет переноса, нет переполнения, результат = 0” и “перенос, переполнение, результат = 0”.

Если умножение на компьютере может давать произведение в виде двойного слова, то желательно иметь две команды умножения: одну для знаковых операндов, а другую — для беззнаковых. (На 4-разрядной машине в результате умножения знаковых операндов произведение $F \times F = 01$, а в случае перемножения беззнаковых операндов — $F \times F = E1$ (в шестнадцатеричной записи)). Ни переноса, ни переполнения при выполнении этих команд умножения не возникает, так как результат умножения всегда может разместиться в двойном слове.

Если же произведение дает в результате одно слово (младшее слово результата из двух слов), то в случае, когда операнды и результат интерпретируются как целые беззнаковые числа, перенос означает, что результат не помещается в одно слово. Если же операнды и результат интерпретируются как знаковые целые числа в дополнительном коде, то, если результат не помещается в слово, устанавливается флаг переполнения. Таким образом, при умножении возможны только девять комбинаций событий. Невозможны комбинации “нет переноса, переполнение, результат > 0 ”, “нет переноса, переполнение, результат = 0” и “перенос, нет переполнения, результат = 0”. Таким образом, при рассмотрении сложения, вычитания и умножения получаем, что реально могут встретиться только 10 комбинаций флагов условий из 12 возможных.

2.15. Циклический сдвиг

Здесь все тривиально. Возможно, это покажется удивительным, но приведенные ниже выражения справедливы для любых n от 0 до 32 включительно, даже если сдвиги выполняются по модулю 32.

$$\text{Циклический сдвиг влево на } n \text{ разрядов: } y \leftarrow (x \ll n) \mid (x \gg (32 - n))$$

$$\text{Циклический сдвиг вправо на } n \text{ разрядов: } y \leftarrow (x \gg n) \mid (x \ll (32 - n))$$

Если ваш компьютер оснащен сдвигами двойной длины, ими можно воспользоваться для выполнения циклических сдвигов. Эти команды могут быть записаны как

```
shldi  RT, RA, RB, I
shrdi  RT, RA, RB, I
```

Они рассматривают конкатенацию RA и RB как единую величину двойной длины и выполняют ее сдвиг влево или вправо на непосредственно указанное количество битов I. (Если величина сдвига находится в регистре, такую команду очень сложно реализовать на большинстве RISC-машин, поскольку при этом требуется чтение трех регистров.) Результатом левого сдвига является старшее слово сдвинутого “двухсловного” значения, а результатом правого сдвига — младшее слово.

При использовании команды `shldi` левый циклический сдвиг регистра R_x можно выполнить как

```
shldi RT, Rx, Rx, I
```

Правый циклический сдвиг выполняется аналогично с помощью команды `shrdi`.

Левый циклический сдвиг на одну позицию можно выполнить путем сложения содержимого регистра с самим собой с "циклическим переносом" (прибавление переноса, получающегося в результате сложения в младшем разряде). У большинства машин такой команды нет, но на многих машинах это действие можно выполнить с помощью двух команд: (1) сложить содержимое регистра с самим собой, генерируя перенос (в регистре состояния), а затем (2) добавить перенос к получившейся сумме.

2.16. Сложение/вычитание двойных слов

Взяв за основу любое из приведенных на с. 52 выражений для переполнения при беззнаковом сложении и вычитании, можно легко реализовать сложение и вычитание двойных слов, не используя при этом значение бита переноса. Для иллюстрации применения сложения с двойной длиной слова примем, что (z_1, z_0) — результат сложения операндов (x_1, x_0) и (y_1, y_0) . Индекс 1 обозначает старшее полуслово, а индекс 0 — младшее. Считаем также, что при сложении используются все 32 бита регистров. Младшие слова интерпретируются как беззнаковые величины. Тогда

$$\begin{aligned} z_0 &\leftarrow x_0 + y_0 \\ c &\leftarrow [(x_0 \& y_0) | ((x_0 | y_0) \& \neg z_0)] \overset{\#}{\gg} 31 \\ z_1 &\leftarrow x_1 + y_1 + c. \end{aligned}$$

Всего для вычисления потребуется девять команд. Во второй строке можно использовать выражение $c \leftarrow (z_0 \overset{\#}{<} x_0)$, допускающее получение конечного результата за четыре команды на машинах, в наборе команд которых имеется оператор сравнения в форме, помещающей результат (0 или 1) в регистр, как, например, команда `SLTU` (*Set on Less Than Unsigned* — беззнаковая проверка "меньше, чем") в MIPS [84].

Аналогичный код используется для вычисления разности двойных слов.

$$\begin{aligned} z_0 &\leftarrow x_0 - y_0 \\ b &\leftarrow [(-x_0 \& y_0) | ((x_0 = y_0) \& z_0)] \overset{\#}{\gg} 31 \\ z_1 &\leftarrow x_1 - y_1 - b \end{aligned}$$

На машинах с полным набором логических команд для вычисления потребуется выполнить восемь команд. Если вторую строку записать как $b \leftarrow (x_0 \overset{\#}{<} y_0)$, то на машинах, имеющих команду `SLTU`, для вычисления разности понадобится всего четыре команды.

На большинстве машин сложение и вычитание двойных слов реализуется за пять команд посредством представления данных в младшем слове с помощью 31 бита (старший бит всегда равен 0, за исключением временного хранения в нем бита переноса или заема).

2.17. Сдвиг двойного слова

Пусть (x_1, x_0) — два 32-разрядных слова, которые требуется сдвинуть вправо или влево как единое 64-разрядное слово, в котором x_1 представляет старшее полуслово, и пусть (y_1, y_0) представляет собой аналогично интерпретируемый результат сдвига. Будем считать, что переменная величина сдвига n может принимать любые значения от 0 до 63. Предположим также, что машина может выполнять команду сдвига по модулю 64 или большему. Это означает, что если величина сдвига лежит в диапазоне от 32 до 63 или от -32 до -1 , то в результате сдвига получится слово, состоящее из одних нулевых битов. Исключением является команда знакового сдвига вправо, так как в этом случае 32 разряда результата будут заполнены значением знакового разряда слова. (Рассматриваемый код не работает на процессорах Intel x86, сдвиг в которых выполняется по модулю 32.)

При этих предположениях команду *сдвига двойного слова влево* можно реализовать, как показано ниже (для вычисления требуется восемь команд).

$$\begin{aligned} y_1 &\leftarrow x_1 \ll n \mid x_0 \gg (32 - n) \mid x_0 \ll (n - 32) \\ y_0 &\leftarrow x_0 \ll n \end{aligned}$$

Чтобы при $n = 32$ получился корректный результат, в первом присваивании должна обязательно использоваться команда *или*, а не *плюс*. Если известно, что $0 \leq n \leq 32$, то последний член в первом присваивании можно опустить, тем самым реализуя сдвиг двойного слова за шесть команд.

Аналогично *беззнаковый сдвиг двойного слова вправо* можно вычислить так.

$$\begin{aligned} y_0 &\leftarrow x_0 \gg n \mid x_1 \ll (32 - n) \mid x_1 \gg (n - 32) \\ y_1 &\leftarrow x_1 \gg n \end{aligned}$$

Знаковый сдвиг двойного слова вправо реализовать сложнее из-за нежелательного распространения знакового разряда в одном из членов. Вот как выглядит простейший код для выполнения этой операции.

$$\begin{aligned} &\text{if } n < 32 \\ &\quad \text{then } y_0 \leftarrow x_0 \gg n \mid x_1 \ll (32 - n) \\ &\quad \text{else } y_0 \leftarrow x_1 \gg (n - 32) \\ & y_1 \leftarrow x_1 \gg n \end{aligned}$$

Если в компьютере есть команда *условной пересылки* (*conditional move*), то код реализуется за восемь команд без переходов. Если же команды условной пересылки нет, то

понадобится выполнить десять команд с использованием знакомой конструкции создания маски с помощью *знакового сдвига вправо на 31 разряд* для устранения нежелательного распространения знака в одном из членов.

$$y_0 \leftarrow x_0 \gg n \mid x_1 \ll (32 - n) \mid \left[\left(x_1 \gg (n - 32) \right) \& \left((32 - n) \gg 31 \right) \right]$$

$$y_1 \leftarrow x_1 \gg n$$

2.18. Сложение, вычитание и абсолютное значение многобайтовых величин

Ряд приложений работает с массивами коротких целых чисел (обычно байтов или полуслов), и часто работу программы можно ускорить, если начать выполнять действия не над короткими числами, а над целыми словами. Пусть для определенности слово содержит четыре однобайтовых целых числа (хотя все методы, которые рассматриваются в этом разделе, легко приспособить к другим типам упаковки, например к словам, содержащим одно 12-битовое и два 10-битовых целых числа и т.д.). Описываемые методы играют большую роль при работе на 64-разрядных машинах, поскольку в этом случае повышается степень возможной параллельности вычислений.

Сложение должно быть организовано таким образом, чтобы не возникало переносов из одного байта в другой; для этого можно воспользоваться приведенным ниже двухшаговым алгоритмом.

1. Замаскировать старший бит в каждом байте операнда и выполнить сложение (избежав таким образом переносов через границы байтов).
2. Учесть старшие биты каждого байта, т.е. выполнить сложение старших битов операндов (и при необходимости перенос в этот бит).

Перенос в старший бит каждого байта определяется старшим битом каждого байта суммы, вычисленной на первом шаге алгоритма. Аналогичный метод используется и при вычитании.

Сложение

$$s \leftarrow (x \& 0x7F7F7F7F) + (y \& 0x7F7F7F7F)$$

$$s \leftarrow ((x \oplus y) \& 0x80808080) \oplus s$$

Вычитание

$$d \leftarrow (x \mid 0x80808080) - (y \& 0x7F7F7F7F)$$

$$d \leftarrow ((x \oplus y) \mid 0x7F7F7F7F) \oplus d$$

На компьютере, оснащенном полным набором логических команд, реализация описанных выше вычислений потребует выполнения восьми команд (с учетом загрузки чис-

ла $0x7F7F7F7F$). (В приведенных формулах команды *и* и *или* с числом $0x80808080$ можно заменить соответственно командами *и-не* и *или-не* с числом $0x7F7F7F7F$.)

Существует и другая методика для случая, когда слово состоит только из двух полей. Сложение при этом можно реализовать посредством сложения 32-битовых чисел с последующим вычитанием нежелательного переноса из полученной суммы. Ранее отмечалось, что выражение $(x + y) \oplus x \oplus y$ даст переносы в каждом разряде. Использование этой и аналогичной формул для вычитания дает нам приведенный ниже код для сложения/вычитания двух полуслов по модулю 2^{16} (реализуется семью командами).

Сложение	Вычитание
$s \leftarrow x + y$	$d \leftarrow x - y$
$c \leftarrow (s \oplus x \oplus y) \& 0x00010000$	$b \leftarrow (d \oplus x \oplus y) \& 0x00010000$
$s \leftarrow s - c$	$d \leftarrow d + b$

Абсолютное значение многобайтовой величины можно вычислить путем дополнения этого числа и прибавления 1 к каждому байту, содержащему отрицательное целое число (т.е. число, старший бит которого равен 1). Приведенный код делает каждый байт y равным абсолютному значению соответствующего байта x (реализуется восемью командами).

$a \leftarrow x \& 0x80808080$	Выделяем биты знаков
$b \leftarrow a \gg 7$	Целое число, равное 1, если x отрицательно
$m \leftarrow (a - b) a$	$0xFF$ для отрицательного x
$y \leftarrow (x \oplus m) + b$	Дополнение и прибавление 1 для отрицательных чисел

Вместо третьей формулы можно использовать $m \leftarrow a + a - b$. Прибавление b в четвертой строке не может вызвать переноса через границу байта, так как в каждом байте значение старшего бита $x \oplus m$ равно 0.

2.19. Функции `doz`, `max`, `min`

Функция `doz` для знаковых аргументов определяется следующим образом.

Знаковая функция	Беззнаковая функция
$\text{doz}(x, y) = \begin{cases} x - y, & x \geq y \\ 0, & x < y \end{cases}$	$\text{dozu}(x, y) = \begin{cases} x - y, & x \geq y \\ 0, & x < y \end{cases}$

Функцию doz называют также вычитанием в первом классе, так как, если вычитаемое больше уменьшаемого, результат равен 0³. Эту функцию удобно использовать при вычислении двух других функций — $\max(x, y)$ и $\min(x, y)$ (как в знаковом, так и в беззнаковом случае). Аппаратная реализация $\max(x, y)$ и $\min(x, y)$ трудна, поскольку машине может потребоваться путь передачи из выходного порта блока регистров во входной порт, минуя сумматор, а обычно такие пути отсутствуют. Ситуация проиллюстрирована на рис. 2.1. Сумматор используется (командой) для вычисления $x - y$. Старший бит результата вычитания (бит знака и переносы, как описано на с. 48) определяет справедливость $x \geq y$ или $x < y$. Результат сравнения передается мультиплексору (MUX), который выбирает либо x , либо y для записи в целевой регистр. Такие пути из выхода блока регистров в мультиплексор обычно отсутствуют. Команды “разность или 0” могут быть реализованы и без таких путей, поскольку в блок регистров передается либо выход сумматора, либо нулевое значение.

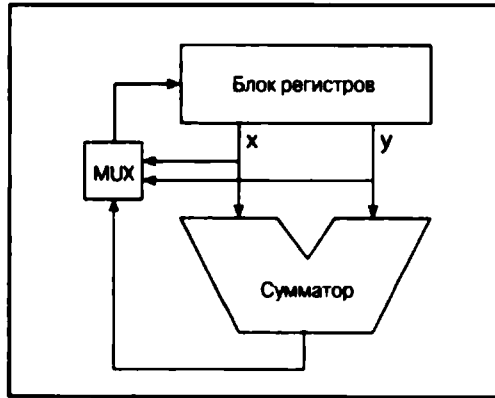


Рис. 2.1. Реализация $\max(x, y)$ и $\min(x, y)$

Используя функцию “разность или 0”, вычисление $\max(x, y)$ и $\min(x, y)$ можно реализовать с помощью двух команд следующим образом.

Знаковое вычисление	Беззнаковое вычисление
$\max(x, y) = y + \text{doz}(x, y)$	$\max_u(x, y) = y + \text{dozu}(x, y)$
$\min(x, y) = x - \text{doz}(x, y)$	$\min_u(x, y) = x - \text{dozu}(x, y)$

В знаковом случае результат команды “разность или 0” может быть отрицательным. Это происходит в случае переполнения при вычитании. Переполнение следует игнорировать; сложение с y или вычитание из x также приведет к переполнению, и результат будет правильным. Если значение $\text{doz}(x, y)$ отрицательное, в действительности она дает верное значение, если интерпретировать его как беззнаковое целое.

³ Математическое название операции — *минус*, и она обозначается как $-$.

Предположим, что ваш компьютер не оснащен командой “разность или \bar{f} ”, но вы хотите закодировать вычисление $\text{doz}(x, y)$, $\text{max}(x, y)$ и так далее эффективным способом без ветвлений. В нескольких следующих абзацах мы покажем, как выполнить такое кодирование на вашей машине с применением команд *условного перемещения*, предикатов сравнения, эффективной работы с битом переноса или даже без всего перечисленного.

Если ваша машина оснащена командой *условного перемещения*, она может вычислить $\text{doz}(x, y)$ за три команды, а деструктивные⁴ $\text{max}(x, y)$ и $\text{min}(x, y)$ — за две команды. Например, на машине с полным набором команд RISC $z \leftarrow \text{doz}(x, y)$ можно вычислить следующим образом ($r0$ — постоянный нулевой регистр).

```
sub    z, x, y      Установить z = x - y
cmplt  t, x, y      Установить t = 1, если x < y, иначе 0
movne  z, t, r0     Установить z = 0, если x < y
```

На машине с полным набором команд RISC можно также вычислить $\text{max}(x, y)$.

```
cmplt  t, x, y      Установить t = 1, если x < y, иначе 0
movne  x, t, y      Установить x = y, если x < y
```

Функция min и ее беззнаковый двойник получаются путем изменения условий сравнения.

С использованием предикатов сравнения эти функции можно вычислить четырьмя или пятью командами (тремя или четырьмя, если предикаты сравнения дают -1 как значение “истинно”).

$$\begin{aligned}\text{doz}(x, y) &= (x - y) \& -(x \geq y) \\ \text{max}(x, y) &= y + \text{doz}(x, y) \\ &= ((x \oplus y) \& -(x \geq y)) \oplus y \\ \text{min}(x, y) &= x - \text{doz}(x, y) \\ &= ((x \oplus y) \& -(x \leq y)) \oplus y\end{aligned}$$

На некоторых машинах бит переноса может оказаться хорошим подспорьем для вычисления беззнаковых версий этих функций. Обозначим как $\text{carry}(x - y)$ бит на выходе сумматора при операции $x + \bar{y} + 1$, помещенный в регистр общего назначения. Таким образом, $\text{carry}(x - y) = 1$ тогда и только тогда, когда $x \geq y$. Тогда мы имеем следующее.

$$\begin{aligned}\text{dozu}(x, y) &= ((x - y) \& -(\text{carry}(x - y) - 1)) \\ \text{maxu}(x, y) &= x - ((x - y) \& (\text{carry}(x - y) - 1)) \\ \text{minu}(x, y) &= y + ((x - y) \& (\text{carry}(x - y) - 1))\end{aligned}$$

⁴ Деструктивной называется операция, которая перезаписывает один или несколько своих аргументов.

На большинстве машин, оснащенных командой *вычитания*, генерирующей перенос или заем, и другой формой вычитания, которая использует перенос или заем в качестве входных данных, выражение $\text{sgtu}(x - y) - 1$ можно вычислить за еще одну дополнительную команду после вычитания y из x . Например, на машинах Intel x86 $\text{minu}(x, y)$ можно вычислить с помощью четырех команд.

```
sub eax,ecx ; x и y находятся в eax и ecx
sbb edx,edx ; edx = 0, если x >= y, иначе -1
and eax,edx ; 0, если x >= y, иначе x - y
add eax,ecx ; Прибавить y, что дает y, если x >= y, иначе x
```

Таким способом все три функции можно вычислить за четыре команды (три для $\text{dozu}(x, y)$, если машина имеет команду *и с дополнением*).

Метод, применимый почти на всех машинах RISC, заключается в использовании одного из приведенных выше выражений, использующих предикат сравнения, и подстановке в качестве предиката одного из выражений на с. 43. Например

$$d \leftarrow x - y$$

$$\text{doz}(x, y) = d \& \left[\left(d \equiv ((x \oplus y) \& (d \oplus x)) \right) \gg 31 \right]$$

$$\text{dozu}(x, y) = d \& \neg \left[\left((-x \& y) | ((x = y) \& d) \right) \gg 31 \right]$$

Такое вычисление требует от семи до десяти команд, в зависимости от имеющегося на машине набора команд; для вычисления функций max и min необходимо выполнить еще одну команду.

Эти операции можно выполнить за четыре команды из базового набора команд RISC без ветвления, если известно, что $-2^{31} \leq x - y \leq 2^{31} - 1$ (обратите внимание на то, что это выражение обычной, а не компьютерной арифметики). Один и тот же код работает и для знаковых, и для беззнаковых целых чисел с одним и тем же ограничением на x и y . Чтобы приведенные формулы были корректны, достаточным условием для знаковых чисел является $-2^{30} \leq x, y \leq 2^{30} - 1$, а для беззнаковых — $0 \leq x, y \leq 2^{31} - 1$.

$$\text{doz}(x, y) = \text{dozu}(x, y) = (x - y) \& \neg \left((x - y) \gg 31 \right)$$

$$\text{max}(x, y) = \text{maxu}(x, y) = x - \left((x - y) \& \left((x - y) \gg 31 \right) \right)$$

$$\text{min}(x, y) = \text{minu}(x, y) = y + \left((x - y) \& \left((x - y) \gg 31 \right) \right)$$

Далее приводятся некоторые возможные применения команды *"разность или 0"*. Здесь результат вычисления $\text{doz}(x, y)$ должен рассматриваться как беззнаковое целое число.

1. Она непосредственно реализует функцию IDIM в Fortran.
2. Она используется для вычисления абсолютного значения разности [72].

$$\begin{aligned} |x - y| &= \text{doz}(x, y) + \text{doz}(y, x), && \text{знаковые аргументы,} \\ &= \text{dozu}(x, y) + \text{dozu}(y, x), && \text{беззнаковые аргументы} \end{aligned}$$

Следствие: $|x| = \text{doz}(x, 0) + \text{doz}(0, x)$ (еще одно трехкомандное решение приведено в разделе 2.4 на с. 39).

3. Для фиксации верхней границы истинной суммы беззнаковых целых x и y максимальным положительным значением $(2^{32} - 1)$ [72] можно использовать

$$\neg \text{dozu}(-x, y).$$

4. Она используется для вычисления некоторых предикатов сравнения (четыре команды для каждого):

$$\begin{aligned} x > y &= (\text{doz}(x, y) | \neg \text{doz}(x, y)) \gg 31, \\ x \dot{>} y &= (\text{dozu}(x, y) | \neg \text{dozu}(x, y)) \gg 31. \end{aligned}$$

5. С ее помощью можно вычислить бит переноса при сложении $x + y$ (пять команд):

$$\text{carry}(x + y) = x \dot{>} \neg y = (\text{dozu}(x, \neg y) | \neg \text{dozu}(x, \neg y)) \gg 31.$$

Выражение $\text{doz}(x, -y)$, результат вычисления которого рассматривается как беззнаковое целое число, в большинстве случаев представляет собой истинную сумму $x + y$ с фиксированной нулевой нижней границей. Однако это не так, когда y является максимальным отрицательным числом.

В компьютере IBM RISC/6000 (и его предшественнике 801) есть отдельная команда для вычисления знаковой *разности* или *0*. Компьютер MMIX Кнута [72] имеет беззнаковую версию этой команды (включая некоторые разновидности, работающие с частями слова параллельно). Это приводит к вопросу о том, как получить знаковую версию функции из беззнаковой и наоборот. Данная задача может быть решена следующим образом (здесь сложение и вычитание просто дополняют бит знака).

$$\begin{aligned} \text{doz}(x, y) &= \text{dozu}(x + 2^{31}, y + 2^{31}) \\ \text{dozu}(x, y) &= \text{doz}(x - 2^{31}, y - 2^{31}) \end{aligned}$$

Могут оказаться полезными и некоторые другие тождества.

$$\begin{aligned} \text{doz}(\neg x, \neg y) &= \text{doz}(y, x) \\ \text{dozu}(\neg x, \neg y) &= \text{dozu}(y, x) \end{aligned}$$

Соотношение $\text{doz}(-x, -y) = \text{doz}(y, x)$ не выполняется, если либо x , либо y (но не оба одновременно) представляет собой максимальное отрицательное число.

2.20. Обмен содержимого регистров

Существует старый хорошо известный способ обмена содержимым двух регистров без использования третьего [59].

$$x \leftarrow x \oplus y$$

$$y \leftarrow y \oplus x$$

$$x \leftarrow x \oplus y$$

Этот метод хорошо работает на двухадресной машине. Он остается работоспособным и при замене оператора \oplus дополнительным к нему оператором эквивалентности \equiv . Возможно также использование различных наборов операций сложения и вычитания.

$$x \leftarrow x + y \quad x \leftarrow x - y \quad x \leftarrow y - x$$

$$y \leftarrow x - y \quad y \leftarrow y + x \quad y \leftarrow y - x$$

$$x \leftarrow x - y \quad x \leftarrow y - x \quad x \leftarrow x + y$$

К сожалению, в каждом из вариантов обмена есть команда, неприемлемая для двухадресной машины (если, конечно, среди команд машины нет “обратного вычитания”).

Вот еще одна маленькая хитрость, которая может оказаться полезной в приложениях с двойной буферизацией, в которых меняются местами два указателя. Первая команда выносится за пределы цикла, в котором выполняется обмен (хотя при этом и теряется преимущество сохранения регистра).

$$\text{Вне цикла:} \quad t \leftarrow x \oplus y$$

$$\text{Внутри цикла:} \quad x \leftarrow x \oplus t$$

$$y \leftarrow y \oplus t$$

Обмен соответствующих полей регистров

Зачастую требуется обменять биты регистров x и y , если i -й бит маски $m_i = 1$, и оставить их неизменными, если $m_i = 0$. Под обменом “соответствующими” полями подразумевается обмен без сдвигов. В маске m единичные биты не обязательно должны быть смежными. Приведем простейший метод решения поставленной задачи.

$$x' \leftarrow (x \& \bar{m}) | (y \& m)$$

$$y \leftarrow (y \& \bar{m}) | (x \& m)$$

$$x \leftarrow x'$$

Использование промежуточных временных данных в четырех выражениях с командой $|$ требует выполнения семи команд при условии, что загрузка \bar{m} или m занимает одну команду и имеется команда $|$ -не. Если машина в состоянии вычислить четыре независимые команды $|$ параллельно, то время выполнения равно всего трем тактам.

Более эффективный метод (пять команд, требующих для выполнения четыре такта на машине с неограниченными возможностями распараллеливания вычислений на уровне команд) показан ниже, в столбце (а). Этот алгоритм использует три команды *исключающего или*.

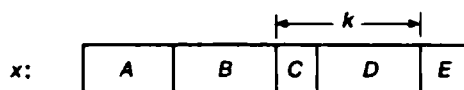
(а)	(б)	(в)
$x \leftarrow x \oplus y$	$x \leftarrow x \equiv y$	$t \leftarrow (x \oplus y) \& m$
$y \leftarrow y \oplus (x \& m)$	$y \leftarrow y \equiv (x \bar{m})$	$x \leftarrow x \oplus t$
$x \leftarrow x \oplus y$	$x \leftarrow x \equiv y$	$y \leftarrow y \oplus t$

Метод, предложенный в столбце (б), лучше использовать в тех случаях, когда m не помещается в поле непосредственно задаваемого значения, но в него можно поместить \bar{m} , а в наборе команд имеется команда *эквивалентности*.

Метод из столбца (в) предложен в [37]. Он также требует выполнения пяти команд (вновь в предположении, что для загрузки m в регистр используется одна команда), но на машинах с достаточными возможностями распараллеливания вычислений на уровне команд выполнение этого кода занимает всего три такта.

Обмен двух полей одного регистра

Предположим, что в регистре x необходимо поменять два поля (одинаковой длины), не изменяя при этом прочие биты регистра, т.е. нам требуется поменять местами поля B и D в пределах одного машинного слова, не изменяя при этом содержимое полей A , C и E . Поля расположены друг от друга на расстоянии k бит.



Следующий простейший код сдвигает поля D и B на новые позиции и затем комбинирует полученные слова с помощью команд *и* и *или*.

$$\begin{aligned}
 t_1 &= (x \& m) \ll k \\
 t_2 &= \left(x \overset{n}{\gg} k \right) \& m \\
 x' &= (x \& m') | t_1 | t_2
 \end{aligned}$$

Здесь маска m содержит единичные биты в поле D (и нулевые биты в других полях), а маска m' содержит единичные биты в полях A , C и E . На машине с неограниченными возможностями распараллеливания вычислений на уровне команд этот код требует выполнения одиннадцати команд за шесть тактов, отводя четыре команды для генерации двух масок.

Ниже описан метод обмена полями [37], который в тех же предположениях требует выполнения восьми команд и выполняется за пять тактов. Он аналогичен алгоритму обмена соответствующими полями двух регистров, описанному ранее в подразделе

“Обмен соответствующих полей регистров” в столбце (в). Здесь, как и ранее, m — маска, выделяющая поле D .

$$t_1 = \left[x \oplus \left(x \gg k \right) \right] \& m$$

$$t_2 = t_1 \ll k$$

$$x' = x \oplus t_1 \oplus t_2$$

Идея состоит в том, что t_1 в позициях поля D содержит биты $B \oplus D$ (и нули в остальных полях), а t_2 содержит биты $B \oplus D$ в позициях поля B . Этот код, как и приведенный ранее простейший, работает корректно, даже если B и D представляют собой поля, разбитые на части, т.е. если в маске m единичные биты не являются смежными.

Условный обмен

Методы обмена из предыдущих двух разделов, использующие команду *исключающего или*, вырождаются в отсутствие каких-либо действий, если все биты маски m нулевые. Следовательно, эти методы могут выполнять условный обмен содержимым регистров, соответствующих полей регистров или полей в одном регистре, если некоторое условие c истинно, и не выполнять его при ложности условия. Для этого достаточно, чтобы все биты маски m были нулевыми, если условие ложно, и имелись корректно установленные единичные биты при выполнении условия c . Если формирование маски m выполняется без команд ветвления, код условного обмена также не будет содержать этих команд.

2.21. Выбор среди двух или более значений

Предположим, что переменная x может принимать только два возможных значения — a и b . Пусть требуется присвоить переменной значение, отличающееся от текущего, причем алгоритм не должен зависеть от конкретных значений a и b . Например, в компиляторе x может быть кодом одной из операций — *переход при выполнении условия* или *переход при невыполнении условия*, и вам требуется изменить текущий код на противоположный. Значения кодов произвольны и могут быть определены, например, с помощью макроопределения `#define` или перечисления `enum` в заголовочном файле.

Вот простейший код переключения.

```
if (x == a) x = b;
else x = a;
```

А это его эквивалент в языке C.

```
x = (x == a) ? b : a;
```

Однако существенно лучше (по крайней мере эффективнее) использовать другой способ.

$$x \leftarrow a + b - x, \text{ или}$$

$$x \leftarrow a \oplus b \oplus x$$

Если a и b — константы, то вычисления требуют всего одной или двух команд из базового множества RISC-команд. Естественно, переполнение при вычислении $a + b$ можно игнорировать.

Возникает закономерный вопрос: существует ли эффективный метод циклического перебора трех и более значений? Иными словами, пусть заданы три произвольные, не равные друг другу константы a , b и c . Требуется найти легко вычисляемую функцию f , которая удовлетворяет условиям

$$f(a) = b; \quad f(b) = c; \quad f(c) = a.$$

Любопытно отметить, что всегда существует полиномиальная функция, удовлетворяющая указанным условиям. В случае трех констант это функция

$$f(x) = \frac{(x-a)(x-b)}{(c-a)(c-b)}a + \frac{(x-b)(x-c)}{(a-b)(a-c)}b + \frac{(x-c)(x-a)}{(b-c)(b-a)}c. \quad (5)$$

Идея заключается в том, что при $x = a$ первый и последний члены функции обращаются в нуль и остается только средний член при коэффициенте b и т.д. Для вычисления этой функции необходимо выполнить 14 арифметических операций, при этом в общем случае промежуточные результаты превышают размер машинного слова, хотя мы имеем дело всего лишь с квадратичной функцией. Однако если для вычисления полинома применить схему Горнера⁵, то понадобится выполнить только пять арифметических операций (четыре для вычисления квадратов с целыми коэффициентами и заключительное деление). После преобразований (5) получим следующий вид функции

$$f(x) = \frac{1}{(a-b)(a-c)(b-c)} \{ [(a-b)a + (b-c)b + (c-a)c] x^2 + [(a-b)b^2 + (b-c)c^2 + (c-a)a^2] x + [(a-b)a^2b + (b-c)b^2c + (c-a)ac^2] \}$$

Полученное выражение слишком громоздко и непригодно для практического использования.

Еще один метод аналогичен формуле (5) в том смысле, что при вычислении также остается только один из трех членов функции.

$$f(x) = ((-(x=c)) \& a) + ((-(x=a)) \& b) + ((-(x=b)) \& c)$$

Для вычисления этого выражения, при условии наличия в наборе команд предиката равенства, требуется выполнить 11 команд, не считая команды загрузки констант. Так как две команды сложения суммируют два нулевых значения с одним ненулевым, их можно заменить командами *и* или *исключающего или*.

⁵ Правило Горнера состоит в вынесении x за скобки. Например, полином четвертой степени $ax^4 + bx^3 + cx^2 + dx + e$ вычисляется как $x(x(x(ax+b)+c)+d)+e$. Для вычисления полинома степени n требуются n умножений и n сложений; при этом наличие команды *умножения со сложением* существенно повышает эффективность вычислений.

Данную формулу можно упростить, если предварительно вычислить значения $a - c$ и $b - c$ [37].

$$f(x) = (((-x = c)) \& (a - c)) + (((-x = a)) \& (b - c)) + c \text{ или}$$

$$f(x) = (((-x = c)) \& (a \oplus c)) \oplus (((-x = a)) \& (b \oplus c)) \oplus c$$

Каждое из этих выражений вычисляется восемью командами. Но на большинстве компьютеров это, вероятно, ничуть не эффективнее простого кода на языке С, для выполнения которого требуется от четырех до шести команд для малых значений a , b и c .

```
if (x == a) x = b;
else if (x == b) x = c;
else x = a;
```

В [37] предлагается еще один оригинальный метод без переходов с циклом выборки из трех значений, работающий даже там, где нет команд сравнения. На большинстве машин его можно выполнить за восемь команд.

Пусть a , b и c не равны друг другу. Следовательно, имеются два разряда, n_1 и n_2 , такие, что биты чисел a , b и c в этих разрядах различны, и два числа, оба бита которых отличаются от битов другого числа (при этом в каждой позиции имеется одно число, бит которого отличается от битов двух других чисел). Проиллюстрируем это на примере чисел 21, 31 и 20, записанных в двоичном представлении.

$$\begin{array}{cccccc} 1 & 0 & 1 & 0 & 1 & c \\ 1 & 1 & 1 & 1 & 1 & a \\ 1 & 0 & 1 & 0 & 0 & b \\ & n_1 & & n_2 & & \end{array}$$

Без потери общности переименуем a , b и c так, чтобы бит a отличался от битов двух других чисел в разряде n_1 , а бит b — в разряде n_2 , как показано выше. Тогда существует два возможных набора битов в позиции n_1 , а именно $(a_{n_1}, b_{n_1}, c_{n_1}) = (0, 0, 1)$ или $(1, 0, 0)$. Аналогично для битов в позиции n_2 также возможны два варианта: $(a_{n_2}, b_{n_2}, c_{n_2}) = (0, 1, 0)$ или $(1, 0, 1)$. Это приводит нас к рассмотрению четырех возможных случаев, формулы для которых приведены ниже.

Случай 1. $(a_{n_1}, b_{n_1}, c_{n_1}) = (0, 1, 1)$, $(a_{n_2}, b_{n_2}, c_{n_2}) = (0, 1, 0)$:
 $f(x) = x_{n_1} * (a - b) + x_{n_2} * (c - a) + b$

Случай 2. $(a_{n_1}, b_{n_1}, c_{n_1}) = (0, 1, 1)$, $(a_{n_2}, b_{n_2}, c_{n_2}) = (1, 0, 1)$:
 $f(x) = x_{n_1} * (a - b) + x_{n_2} * (a - c) + (b + c - a)$

Случай 3. $(a_{n_1}, b_{n_1}, c_{n_1}) = (1, 0, 0)$, $(a_{n_2}, b_{n_2}, c_{n_2}) = (0, 1, 0)$:

$$f(x) = x_{n_1} * (b - a) + x_{n_2} * (c - a) + a$$

Случай 4. $(a_{n_1}, b_{n_1}, c_{n_1}) = (1, 0, 0)$, $(a_{n_2}, b_{n_2}, c_{n_2}) = (1, 0, 1)$:

$$f(x) = x_{n_1} * (b - a) + x_{n_2} * (a - c) + c$$

В этих формулах левый операнд каждого умножения представляет собой отдельный бит. Умножение на 0 или 1 можно заменить операцией u с 0 или с числом, все биты которого равны 1. Таким образом, эти формулы можно переписать иначе, например для первого случая получим следующее.

$$f(x) = \left((x \ll (31 - n_1)) \gg 31 \right) \& (a - b) + \left((x \ll (31 - n_2)) \gg 31 \right) \& (c - a) + b$$

Поскольку все переменные, кроме x , являются константами, функцию можно вычислить восемью базовыми RISC-командами. Как и ранее, команды сложения и вычитания можно заменить командами *исключающего или*.

Данная идея может быть распространена на случай циклической выборки из четырех и более констант. Главное — найти такие номера битов n_1, n_2, \dots , которые бы идентифицировали константы единственным образом. Для четырех констант всегда достаточно трех позиций битов. Затем (в случае четырех констант) решаются уравнения для s, t, u и v (т.е. решается система из четырех линейных уравнений, в которой $f(x)$ принимает значения a, b, c и d , а коэффициенты x_{n_i} равны 0 или 1).

$$f(x) = x_{n_1} s + x_{n_2} t + x_{n_3} u + v$$

Если все четыре константы определяются двумя битами единственным образом, то уравнение принимает следующий вид.

$$f(x) = x_{n_1} s + x_{n_2} t + x_{n_1} x_{n_2} u + v$$

2.22. Формула булева разложения

В этом разделе мы рассмотрим вопрос о минимальном количестве бинарных булевых операций, или команд, которых достаточно для реализации любой булевой функции от трех, четырех или пяти переменных. Под булевой функцией мы понимаем функцию от булевых аргументов, возвращающую булево значение.

В наших обозначениях булевой алгебры символ “+” используется для *или*, соединение двух переменных — для \oplus — для *исключающего или* и надчеркивание, или символ “¬”, — для *не*. Эти операторы можно применять как к однобитовым операндам, так и “побитово” к машинным словам. Основной результат приведен в следующей теореме.

ТЕОРЕМА. Если $f(x, y, z)$ — булева функция от трех переменных, то она может быть переписана в виде $g(x, y) \oplus zh(x, y)$, где g и h являются булевыми функциями от двух переменных⁶.

Доказательство [26]. Функция $f(x, y, z)$ может быть выражена в виде суммы минитермов, после чего из них могут быть выделены \bar{z} и z , так что

$$f(x, y, z) = \bar{z}f_0(x, y) + zf_1(x, y).$$

Поскольку операнды оператора “+” одновременно быть равны 1 не могут, или можно заменить *исключающим или*, и получить следующее.

$$\begin{aligned} f(x, y, z) &= \bar{z}f_0(x, y) \oplus zf_1(x, y) \\ &= (1 \oplus z)f_0(x, y) \oplus zf_1(x, y) \\ &= f_0(x, y) \oplus zf_0(x, y) \oplus zf_1(x, y) \\ &= f_0(x, y) \oplus z(f_0(x, y) \oplus f_1(x, y)) \end{aligned}$$

Здесь дважды использовано тождество $(a \oplus b)c = ac \oplus bc$.

Это и есть искомый вид записи с $g(x, y) = f_0(x, y)$ и $h(x, y) = f_0(x, y) \oplus f_1(x, y)$. Кстати, $f_0(x, y)$ представляет собой $f(x, y, z)$ при $z = 0$, а $f_1(x, y)$ представляет собой $f(x, y, z)$ при $z = 1$.

СЛЕДСТВИЕ. Если набор команд компьютера включает команды для каждой из 16 булевых функций от двух переменных, то любая булева функция от трех переменных может быть реализована с помощью не более четырех команд.

Одна команда реализует $g(x, y)$, другая — $h(x, y)$, и они комбинируются с помощью *или* *или* *исключающего или*.

В качестве примера рассмотрим булеву функцию, которая равна 1, когда единице равны ровно две из переменных x, y, z .

$$f(x, y, z) = xy\bar{z} + x\bar{y}z + \bar{x}yz$$

Перед тем как читать дальше, заинтересованный читатель может попытаться реализовать функцию f с помощью четырех команд, не прибегая к доказанной выше теореме.

Из доказательства теоремы следует

$$\begin{aligned} f(x, y, z) &= f_0(x, y) \oplus z(f_0(x, y) \oplus f_1(x, y)) \\ &= xy \oplus z(xy \oplus (\bar{x}\bar{y} + \bar{x}y)) \\ &= xy \oplus z(x + y), \end{aligned}$$

что можно реализовать четырьмя командами.

⁶ Это разложение булевой функции известно также как разложение Рида-Маллера (Reed-Muller) или положительное разложение Давио (Davio). Согласно Кнуту [69, раздел 7.1.1] это разложение было открыто И.И. Жегалкиным [Математический сборник, 35 (1928), 311–369] и иногда упоминается как “русское разложение”.

Понятно, что теорема может быть распространена на функции четырех и более переменных, т.е. любая булева функция $f(x_1, x_2, \dots, x_n)$ может быть приведена к виду $g(x_1, x_2, \dots, x_{n-1}) \oplus x_n h(x_1, x_2, \dots, x_{n-1})$. Таким образом, можно выполнить разложение функции четырех переменных следующим образом.

$$\begin{aligned} f(w, x, y, z) &= g(w, x, y) \oplus zh(w, x, y), \text{ где} \\ g(w, x, y) &= g_1(w, x) \oplus yh_1(w, x) \text{ и} \\ h(w, x, y) &= g_2(w, x) \oplus yh_2(w, x) \end{aligned}$$

Тем самым показано, что компьютер, который оснащен командами для каждой из 16 бинарных булевых функций, может реализовать любую функцию от четырех переменных с помощью десяти команд. Аналогично для реализации функции от пяти переменных достаточно 22 команд.

Однако можно достичь гораздо лучшего. Для функций от четырех или более переменных, вероятно, нет простых уравнений, таких, как дает нам рассмотренная теорема, но можно провести исчерпывающий поиск. В результате выяснилось, что любая булева функция от четырех переменных может быть реализована с помощью семи бинарных булевых команд, а любая такая функция от пяти переменных может быть реализована с помощью 12 таких команд [69, раздел 7.1.2].

В случае пяти переменных только 1920 из $2^5 = 4\,294\,967\,296$ функций реализуются с помощью 12 команд, и эти 1920 функций, по сути, являются одной и той же функцией. Все они получаются одна из другой путем перестановки аргументов, замены некоторых из аргументов их дополнениями или дополнением значения функции.

2.23. Реализация команд для всех 16 бинарных булевых операций

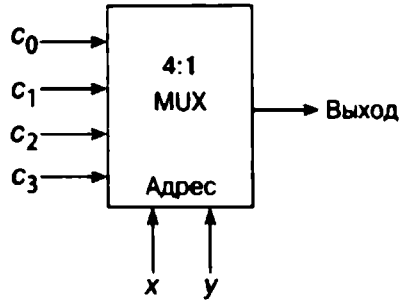
Набор команд некоторых компьютеров включает все 16 бинарных булевых операций. Многие из них бесполезны, в том смысле что они могут быть выполнены с помощью других команд. Например, функция $f(x, y) = 0$ просто очищает регистр, и на большинстве машин это можно сделать разными способами. Тем не менее одна из причин, по которым разработчики компьютеров могут захотеть реализовать все 16 команд, заключается в том, что имеется простая и обычная схема для решения поставленной задачи.

Обратимся к табл. 2.1 на с. 38, в которой показаны все 16 бинарных булевых функций. Чтобы реализовать эти функции как команды, выберем четыре бита кода операции такими же, как и значения соответствующей функции, приведенные в таблице. Обозначим эти биты кода операции как c_0 , c_1 , c_2 и c_3 , считывая их в таблице снизу вверх, а входные регистры — как x и y . Тогда схема для реализации всех 16 бинарных булевых операций описывается логическим выражением

$$c_0xy + c_1x\bar{y} + c_2\bar{x}y + c_3\bar{x}\bar{y}.$$

Например, в случае $c_0 = c_1 = c_2 = c_3 = 0$ указанная команда вычисляет нулевую функцию $f(x, y) = 0$. Если в коде операции $c_0 = 1$, а все остальные биты нулевые, мы получаем функцию x . При $c_0 = c_3 = 0$ и $c_1 = c_2 = 1$ получается *исключающее или*, и т.д.

Все описанное можно реализовать с помощью n мультиплексоров 4:1, где n — размер машинного слова. Биты данных x и y представляют собой адресные входы, а четыре бита кода операции — входные данные для каждого мультиплексора. В современной технологии мультиплексор представляет собой стандартный строительный блок и обычно является очень быстрой схемой. Она проиллюстрирована ниже.



Функция схемы состоит в выборе c_0 , c_1 , c_2 или c_3 в качестве выхода в зависимости от того, равны ли биты x и y 00, 01, 10 или 11 соответственно. Это чем-то похоже на четырехпозиционный поворотный переключатель.

Это элегантное решение несколько дорого с точки зрения кодов операций — в том смысле, что используются целых 16 кодов. Имеется ряд способов реализации всех 16 булевых операций с применением только восьми кодов ценой применения менее распространенной логики. Одна такая схема проиллюстрирована в табл. 2.3.

Таблица 2.3. НАБОР ИЗ ВОСЬМИ БУЛЕВЫХ КОМАНД

Значения функции	Формула	Название команды
0001	xy	<i>and</i> (и)
0010	$x\bar{y}$	<i>andc</i> (и с дополнением)
0110	$x \oplus y$	<i>xor</i> (исключающее или)
0111	$x + y$	<i>or</i> (или)
1110	\overline{xy}	<i>nand</i> (отрицательное и)
1101	\overline{xy} , или $\bar{x} + y$	<i>cor</i> (дополнение и или)
1001	$\overline{x \oplus y}$, или $x \equiv y$	<i>eqv</i> (эквивалентность)
1000	$\overline{x + y}$	<i>nor</i> (отрицательное или)

Восемь операций, не показанных в таблице, могут быть выполнены с помощью указанных восьми операций путем обмена входов или использования одного и того же регистра в качестве обоих входов (см. упр. 13).

Эта схема используется архитектурой IBM POWER с минимальным отличием, заключающимся в том, что вместо команды “дополнение и или” в ней имеется команда “или с дополнением”. Схема, показанная в табл. 2.3, позволяет реализовать последние четыре команды через дополнение соответствующей команды из первой четверки.

Исторические примечания

Алгебра логики, изложенная в работе Джорджа Буля (George Boole) *Исследование законов мышления (An Investigation of the Laws of Thought)*⁷ в 1854 году, несколько отличается от того, что мы называем булевой алгеброй сегодня. Буль для представления истины и лжи использовал целые числа 1 и 0 соответственно и показал, как можно работать с ними методами обычной числовой алгебры, чтобы формализовать утверждения обычного языка, включающие “и”, “или” и “за исключением”. Он также воспользовался обычной алгеброй для формализации утверждений теории множеств, включая пересечения, объединения непересекающихся множеств и дополнения. Кроме того, он формализовал утверждения теории вероятности, в которой переменные принимают действительные значения от 0 до 1. В работе часто затрагиваются вопросы философии, религии и юриспруденции.

Буль является великим мыслителем в области логики, потому что сумел формализовать ее, позволяя безупречно работать со сложными утверждениями чисто механически, используя знакомые методы обычной алгебры.

Забегая вперед, заметим, что имеется несколько языков программирования, оснащенных всеми 16 булевыми операциями. IBM PL/I (ок. 1966 года) включает встроенную функцию BOOL. В вызове BOOL(x, y, z) аргумент z представляет собой четырехбитовую строку (или при необходимости преобразуется в нее), а x и y — битовые строки равной длины (или преобразуемые в них при необходимости). Аргумент z определяет булеву операцию, выполняемую над x и y . Бинарное значение 0000 определяет нулевую функцию, 0001 — xy , 0010 — $x\bar{y}$ и т.д.

Другим таким языком программирования является Basic для компьютера Wang System 2200B (ок. 1974 года), в котором имеется версия BOOL, работающая вместо битовых строк или целых чисел с символьными строками [89].

Еще одним таким языком программирования является MIT PDP-6 Lisp, позже названный MacLisp [37].

Упражнения

1. Дэвид де Клот (David de Kloet) предложил следующий код для функции snoob для $x \neq 0$, где результат получается при последнем присваивании y .

⁷ Полностью эта 335-страничная работа доступна по адресу www.gutenberg.org/etext/15114.

```

y ← x + (x & -x)
x ← x & -y
while((x & 1) = 0) x ← x >> 1
x ← x >> 1
y ← y | x

```

По сути это то же, что и код Госпера (с. 35), с тем отличием, что вместо команды деления выполняется правый сдвиг в цикле `while`. Поскольку обычно деление является дорогостоящей в смысле времени выполнения операцией, такой код может составить конкуренцию коду Госпера при не слишком большом количестве итераций цикла `while`. Пусть n — длина битовых строк x и y , k — число единичных битов в строках, и предположим, что код выполняется для всех значений x , в которых ровно по k единичных битов. Сколько итераций цикла `while` выполняется в среднем при таких условиях?

2. В тексте упоминается, что левый сдвиг на переменную величину позиций не является вычислимым справа налево. Рассмотрим функцию $x \ll (x \& 1)$ [73]. Это левый сдвиг на переменную величину, но он может быть вычислен следующим образом:

$$x + (x \& 1) * x, \quad \text{или}$$

$$x + (x \& -(x \& 1)),$$

т.е. с помощью операций, вычисляемых справа налево. Что здесь происходит? Можете ли вы придумать другую такую функцию?

3. Выведите формулу Дитца (Dietz) для среднего двух беззнаковых целых

$$(x \& y) + \left((x \oplus y) \gg 1 \right).$$

4. Разработайте метод вычисления среднего значения четырех беззнаковых целых чисел $\lfloor (a + b + c + d) / 4 \rfloor$, при котором не происходит переполнения.
5. Многие предикаты сравнения, показанные на с. 43, могут быть существенно упрощены, если известен 31-й бит x или y . Покажите, как можно упростить семикомандное выражение для $x \stackrel{2}{\leq} y$ до трех базовых команд RISC, среди которых нет сравнений, если известно, что $y_{31} = 0$.
6. Покажите, что если два (возможно, различных) числа суммируются с "циклическим переносом", то прибавление бита переноса не может привести к генерации другого переноса из старшего разряда.
7. Покажите, как использовать циклический перенос для сложения в случае представления отрицательных чисел в записи с дополнением до единицы. Чему равно

максимальное число разрядов, на которое может распространиться перенос (из произвольного разряда)?

8. Покажите, что мультиплексорная операция $(x \& m) | (y \& \sim m)$ может быть выполнена тремя командами из базового набора RISC (в котором отсутствует команда *и с дополнением*).
9. Покажите, как реализовать $x \oplus y$ за четыре команды логики *и-или-не*.
10. Покажите, как для заданного 32-битового слова x и двух целых переменных i и j закодировать копирование бита x из позиции i в позицию j . Значения i и j не связаны одно с другим, но считаем, что $0 \leq i, j \leq 31$.
11. Какого количества бинарных булевых команд достаточно для вычисления любой булевой функции от n переменных, если воспользоваться рекурсивным разложением методом, описанным в теореме в тексте раздела.
12. Покажите, что альтернативными разложениями булевых функций от трех переменных являются
 - (а) $f(x, y, z) = g(x, y) \oplus \bar{z}h(x, y)$ (отрицательное разложение Давио) и
 - (б) $f(x, y, z) = g(x, y) \oplus (z + h(x, y))$.
13. В тексте упоминалось, что все 16 булевых операций можно осуществить с помощью 8 команд из табл. 2.3 путем обмена входных аргументов или использования одного и того же регистра в качестве обоих входных регистров. Покажите, как это сделать.
14. Предположим, что нас не интересуют шесть булевых функций, которые в действительности являются константами или унарными функциями, а именно — $f(x, y) = 0, 1, x, y, \bar{x}$ и \bar{y} , но мы хотим, чтобы наш набор команд мог вычислить остальные десять функций одной командой. Можно ли добиться этого с использованием менее чем восьми типов бинарных булевых команд?
15. В упр. 13 показано, что восьми типов команд достаточно для вычисления любой из 16 булевых операций с двумя операндами одной командой типа R-R (регистр-регистр). Покажите, что в случае команд R-I (регистр-непосредственное значение) достаточно шести типов команд. При использовании команд R-I входные операнды нельзя обменивать или приравнять, но второй входной операнд (поле непосредственного значения) может быть дополнен (или, в действительности, получить любое значение) без затрат времени. Предположим для простоты, что поля непосредственных значений имеют ту же длину, что и регистры общего назначения.
16. Покажите, что не все булевы функции от трех переменных могут быть реализованы с помощью трех бинарных логических команд.

ГЛАВА 3 ОКРУГЛЕНИЕ К СТЕПЕНИ 2

3.1. Округление к кратному степени 2

Округление целого беззнакового числа x в меньшую сторону, например к ближайшему числу, кратному 8, тривиально: для этого можно использовать выражение $x \& -8$ или $\left(x \overset{\cdot}{\gg} 3\right) \ll 3$. Этот способ работает и со знаковыми целыми числами, просто здесь “округление в меньшую сторону” означает округление к меньшему отрицательному числу, например $(-37) \& (-8) = -40$.

Округление в большую сторону выполняется практически так же просто. Например, округление беззнакового целого x в большую сторону к следующему кратному 8 реализуется одной из следующих формул.

$$\begin{aligned}(x + 7) \& -8 \\ x + (-x \& 7)\end{aligned}$$

Эти формулы верны и для знаковых целых чисел, причем “округление в большую сторону” означает округление в положительном направлении. Второе слагаемое во втором выражении может быть полезным, если вы хотите знать, какую величину нужно прибавить к x , чтобы сделать его кратным 8 [40].

Для округления знаковых целых чисел к ближайшему кратному 8 по направлению к 0 можно просто скомбинировать приведенные ранее выражения.

$$\begin{aligned}t \leftarrow \left(x \overset{\cdot}{\gg} 31\right) \& 7 \\ (x + t) \& -8\end{aligned}$$

Первую формулу можно переписать несколько иначе: $t \leftarrow \left(x \overset{\cdot}{\gg} 2\right) \overset{\cdot}{\gg} 29$, что может пригодиться при округлении на компьютере с отсутствующей командой *и* с непосредственно заданным значением или если константа оказывается слишком велика для размещения в отводимом для нее поле непосредственного значения.

Иногда при округлении указывается множитель округления, представляющий собой \log_2 от величины округления (так, значение 3 означает округление к числу, кратному 8, так как $\log_2 8 = 3$). В этом случае можно использовать приведенные ниже формулы, где k — множитель округления.

Округление к меньшему числу: $x \& ((-1) \ll k)$

$$\left(x \gg k \right) \ll k$$

Округление к большему числу: $t \leftarrow (1 \ll k) - 1; (x + t) \& -t$

$$t \leftarrow (-1) \ll k; (x - t - 1) \& t$$

3.2. Округление к ближайшей степени 2

Определим функции $\text{flr2}(x)$ и $\text{clp2}(x)$, подобные функциям “пол” и “потолок” и округляющие число x к ближайшей целой степени 2, а не к ближайшему целому, следующим образом.

$$\text{flr2}(x) = \begin{cases} \text{не определена,} & x < 0, \\ 0, & x = 0, \\ 2^{\lfloor \log_2 x \rfloor}, & x > 0; \end{cases} \quad \text{clp2}(x) = \begin{cases} \text{не определена,} & x < 0, \\ 0, & x = 0, \\ 2^{\lceil \log_2 x \rceil}, & x > 0. \end{cases}$$

Первые буквы в названиях функций отражают их подобие с функциями “пол” (floor) и “потолок” (ceil). Значение функции $\text{flr2}(x)$ представляет собой наибольшую целую степень 2, не превосходящую значение x , а функции $\text{clp2}(x)$ — наименьшую целую степень 2, не меньшую, чем x . Эти определения имеют смысл и для нецелых x (например, $\text{flr2}(0.1) = 0.0625$). Эти функции удовлетворяют ряду отношений, аналогичных отношениям между функциями “пол” и “потолок”, в частности приведенным ниже, где n — целое число.

$$\begin{aligned} \lfloor x \rfloor &= \lceil x \rceil, \text{ только если } x \text{ — целое} & \text{flr2}(x) &= \text{clp2}(x), \text{ только если } x \text{ — степень 2 или 0} \\ \lfloor x + n \rfloor &= \lfloor x \rfloor + n & \text{flr2}(2^n x) &= 2^n \text{flr2}(x) \\ \lceil x \rceil &= -\lfloor -x \rfloor & \text{clp2}(x) &= 1/\text{flr2}(1/x), \quad x \neq 0 \end{aligned}$$

При вычислениях используются только целые беззнаковые значения x , так что обе функции определены для любых возможных значений x . Кроме того, используется требование, чтобы вычисленное значение представляло собой значение по модулю 2^{32} (таким образом, $\text{clp2}(x) = 0$ при $x > 2^{31}$). Далее представлены значения функций для некоторых значений x .

x	$\text{flr2}(x)$	$\text{clp2}(x)$
0	0	0
1	1	1
2	2	2
3	2	4
4	4	4

x	$\text{flp2}(x)$	$\text{clp2}(x)$
5	4	8
...
2^{31}	2^{30}	2^{31}
-1		
2^{31}	2^{31}	2^{31}
2^{31}	2^{31}	0
+1		
...
2^{32}	2^{31}	0
-1		

Функции $\text{flp2}(x)$ и $\text{clp2}(x)$ связаны между собой определенными соотношениями, которые позволяют вычислить значения одной функции из значений другой (с учетом указанных ограничений).

$$\begin{aligned} \text{clp2}(x) &= 2 \text{flp2}(x-1), & x \neq 1 \\ &= \text{flp2}(2x-1), & 1 \leq x \leq 2^{31} \\ \text{flp2}(x) &= \text{clp2}\left(x \overset{\cdot}{+} 2 + 1\right), & x \neq 0 \\ &= \text{clp2}(x+1) \overset{\cdot}{+} 2, & x < 2^{31} \end{aligned}$$

Функции округления в большую и меньшую сторону можно достаточно легко вычислить с помощью функции, вычисляющей количество ведущих нулевых битов, как показано ниже. Однако, чтобы эти формулы можно было использовать для $x = 0$ и $x > 2^{31}$, команды сдвига на вашем компьютере должны давать 0 при сдвиге на величины -1, 32 и 63. На многих машинах (например, PowerPC) есть команда "сдвига по модулю 64", которая работает именно так. Отрицательный сдвиг эквивалентен сдвигу в противоположном направлении (т.е. сдвиг влево на -1 — это просто сдвиг вправо на 1).

$$\begin{aligned} \text{flp2}(x) &= 1 \ll (31 - \text{nlz}(x)) \\ &= 1 \ll (\text{nlz}(x) \oplus 31) \\ &= 0x80000000 \overset{\cdot}{\gg} \text{nlz}(x) \\ \text{clp2}(x) &= 1 \ll (32 - \text{nlz}(x-1)) \\ &= 0x80000000 \overset{\cdot}{\gg} (\text{nlz}(x-1) - 1) \end{aligned}$$

Округление в меньшую сторону

В листинге 3.1 приведен алгоритм без ветвления, который полезен в случае отсутствия функции, вычисляющей *количество ведущих нулевых битов*. Алгоритм основан на распространении вправо крайнего слева единичного бита и требует для реализации выполнения 12 команд.

ЛИСТИНГ 3.1. Наибольшая степень 2, не превосходящая значение x , вычисление без ветвления

```
unsigned flp2(unsigned x)
{
    x = x | (x >> 1);
    x = x | (x >> 2);
    x = x | (x >> 4);
    x = x | (x >> 8);
    x = x | (x >> 16);
    return x - (x >> 1);
}
```

В листинге 3.2. показаны два простых цикла, вычисляющих одну и ту же функцию. Все переменные в приведенных фрагментах представляют собой беззнаковые целые числа. Цикл справа обнуляет крайний справа единичный бит переменной x до тех пор, пока она не становится равной 0, после чего возвращается предыдущее значение x .

ЛИСТИНГ 3.2. Округление x в меньшую сторону, простые циклы

```
y = 0x80000000;
while (y > x)
    y = y >> 1;
return y;

do {
    y = x;
    x = x & (x - 1);
} while(x != 0);
return y;
```

При вычислении левого фрагмента требуется $4\text{nlz}(x)+3$ команды; цикл справа при $x \neq 0$ требует выполнения $4\text{pop}(x)$ команд¹, если сравнение с 0 имеет нулевую стоимость.

Округление в большую сторону

Использование распространения старшего бита вправо дает хороший алгоритм для округления к ближайшей большей степени 2. Представленный в листинге 3.3 алгоритм не имеет команд ветвления и выполняется за 12 команд.

ЛИСТИНГ 3.3. Наименьшая степень 2, не уступающая значению x

```
unsigned clp2(unsigned x)
{
    x = x - 1;
    x = x | (x >> 1);
    x = x | (x >> 2);
    x = x | (x >> 4);
}
```

¹ Функция $\text{pop}(x)$ возвращает количество единичных битов в x .

```

x = x | (x >> 8);
x = x | (x >> 16);
return x + 1;
)

```

Попытка провести вычисление с использованием цикла работает не так хорошо, как хотелось бы.

```

y = 1;
while (y < x)    // Беззнаковое сравнение
    y = 2*y;
return y;

```

Для $x=0$ данный алгоритм дает 1, что не совсем то, что требуется; а при $x \geq 2^{31}$ происходит заикливание. Вычисления требуют выполнения $4n+3$ команд, где n — степень 2 возвращаемого кодом числа. Таким образом, этот алгоритм будет работать медленнее (в смысле количества выполняемых команд), чем алгоритм без ветвления, уже при $n \geq 3$ ($x \geq 8$).

3.3. Проверка пересечения границы степени 2

Предположим, что вся память разделена на блоки, причем размер каждого блока представляет собой степень 2. Нумерация адресов начинается с 0. Блоки могут быть словами, двойными словами, страницами и т.д. Пусть даны некоторый начальный адрес a и длина l . Требуется определить, пересекает ли диапазон адресов от a до $a+l-1$, $l \geq 2$, границы блока (величины a и l — целые беззнаковые, полностью размещаемые в регистрах).

Если $l=0$ или 1, то независимо от значения a пересечения границ блока не будет. Если l превышает размер блока, то, каким бы ни было значение a , произойдет выход за границы блока. Для очень больших значений l (с возможным циклическим возвратом) выход за границы блока может произойти, даже если первый и последний байты диапазона адресов расположены в одном блоке.

Для IBM System/370 определить, произойдет ли выход за границы блока, на удивление просто [16]. Вот как этот метод работает для размера блока, равного 4096 байт (распространенный размер страницы).

```

O    RA, -A(-4096)
ALR  RA, RL
BO   CROSSES

```

В первой команде над регистром RA (в котором содержится начальный адрес a) и числом 0xFFFF000 выполняется команда *побитового или*. Вторая команда добавляет к регистру длину, в результате чего устанавливается двухбитовый код условия. Первый флаг устанавливается равным 1, если при сложении произошел перенос, а второй оказывается установленным, если 32-битовый результат сложения ненулевой. Переход в последней команде выполняется только тогда, когда оба флага условия равны 1. После

выполнения перехода регистр RA будет содержать длину диапазона адресов, выходящую за пределы первой страницы (дополнительная возможность, которая не запрашивалась).

Если, например, $a = 0$ и $l = 4096$, перенос будет иметь место, однако результат в регистре RA окажется нулевым, так что перехода к метке CROSSES не будет.

Теперь посмотрим, нельзя ли адаптировать этот метод для RISC-компьютеров, на которых нет команды *перехода при наличии переноса и ненулевого результата команды*. Примем для простоты размер блока равным 8. Согласно [16] переход к метке CROSSES произойдет только в том случае, если был перенос (т.е. $(a|-8)+l \geq 2^{32}$) и если результат не равен нулю (т.е. $(a|-8)+l \neq 2^{32}$), что эквивалентно условию

$$(a|-8)+l > 2^{32}.$$

Это условие, в свою очередь, эквивалентно появлению переноса в последнем сложении при вычислении $((a|-8)-1)+l$. Если в компьютере имеется команда *переход при наличии переноса*, то данный алгоритм можно использовать непосредственно; при этом потребуется пять команд с учетом загрузки константы -8 .

Если же такой команды нет, то можно воспользоваться тем, что перенос при суммировании $x+y$ происходит тогда и только тогда, когда $\neg x \overset{n}{<} y$ (см. раздел "Беззнаковое сложение/вычитание" на с. 52), чтобы получить выражение

$$\neg((a|-8)-1) \overset{n}{<} l.$$

Использование различных тождеств типа $\neg(x-1) = -x$ дает ряд эквивалентных выражений для предиката "выход за границы блока".

$$\begin{aligned} -(a|-8) \overset{n}{<} l \\ \neg(a|-8)+1 \overset{n}{<} l \\ (-a \& 7)+1 \overset{n}{<} l \end{aligned}$$

На большинстве RISC-компьютеров эти выражения вычисляются за пять или шесть команд.

С другой стороны, очевидно, что выход за границы 8-байтового блока произойдет тогда и только тогда, когда

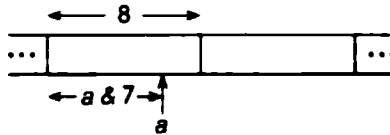
$$(a \& 7)+l-1 \geq 8.$$

Непосредственному вычислению это выражение не поддается в силу возможного переполнения (которое возникает при очень больших l). Однако если переписать выражение в виде $8-(a \& 7) < l$, то его можно вычислить, не вызывая переполнения. Это дает нам выражение

$$8-(a \& 7) \overset{n}{<} l,$$

которое на большинстве RISC-компьютеров вычисляется с помощью пяти команд (или четырех, если компьютер имеет команду *вычитания из непосредственно заданного значения*). При наличии выхода за границы блока одна дополнительная команда вычитания позволяет вычислить значение $1 - (8 - (a \& 7))$, которое дает нам длину поддиапазона адресов, выходящего за пределы первого блока.

Эту формулу легко понять из приведенного ниже рисунка [75], который иллюстрирует, что $a \& 7$ представляет собой смещение a в пределах своего блока, так что $8 - (a \& 7)$ представляет собой размер оставшейся части блока.



Упражнения

1. Покажите, как округлить беззнаковое целое число до ближайшего кратного 8, когда среднее значение округляется (а) в сторону увеличения, (б) в сторону уменьшения, (в) в большую или меньшую сторону — так, чтобы следующий слева бит стал нулевым (“несмещенное” округление).
2. Покажите, как округлить беззнаковое целое число до ближайшего кратного 10, когда среднее значение округляется (а) в сторону увеличения, (б) в сторону уменьшения, (в) в большую или меньшую сторону — так, чтобы результат был четным кратным 10. Можете использовать деление, получение остатка при делении и умножение и не беспокоиться о том, что получающиеся в процессе вычисления значения слишком близки к наибольшему беззнаковому целому числу.
3. Напишите на языке программирования С функцию, выполняющую “невыровненную загрузку”. Функция получает адрес a и загружает четыре байта, располагающиеся от адреса a до $a + 3$, в 32-битовый регистр общего назначения, как если бы эти четыре байта содержали целое число. Параметр a адресует младший байт. Функция не должна содержать ветвлений, должна выполнять не более двух команд загрузки и, если a выровнено на границу слова, не должна пытаться выполнять загрузку из адреса $a + 4$, так как по этому адресу может содержаться защищенный от чтения блок памяти.

ГЛАВА 4

АРИФМЕТИЧЕСКИЕ ГРАНИЦЫ

4.1. Проверка границ целых чисел

Под *проверкой границ* (bounds checking) подразумевается проверка того, находится ли целое значение x в пределах диапазона от a до b , т.е. выполняется ли соотношение

$$a \leq x \leq b.$$

Считаем, что все упомянутые здесь величины представляют собой целые знаковые числа.

Важным применением данного действия является проверка индексов элементов массива. Например, пусть объявлен одномерный массив A , нижняя и верхняя границы которого соответственно равны 1 и 10. При обращении к элементу массива $A(i)$ компилятор может сгенерировать код для проверки корректности индекса $1 \leq i \leq 10$. Если i выходит за пределы объявленного диапазона индексов, выполняется команда *перехода* или *прерывания* для обработки ошибки. В этом разделе показано, что такая проверка может быть выполнена посредством одного сравнения [93].

$$i - 1 \leq 9$$

Очевидно, что этот метод лучше предыдущего, поскольку содержит только одну команду *условного ветвления*; кроме того, значение $i - 1$ в дальнейшем можно использовать при адресации элемента массива.

Возникает вопрос: всегда ли справедлива реализация

$$a \leq x \leq b \Rightarrow x - a \leq b - a$$

(включая случай переполнения при вычитании)? Да, всегда при выполнении условия $a \leq b$. В случае проверки границ массива правила языка программирования могут потребовать, чтобы массив не мог иметь нулевого или отрицательного числа элементов, причем проверка выполнения этого правила может выполняться во время компиляции программы или при выделении памяти для динамических массивов. В этом случае рассмотренное выше преобразование вполне корректно.

Для начала докажем следующую лемму.

ЛЕММА. Если a и b — знаковые целые числа, причем $a \leq b$, то вычисленное значение $b - a$ корректно представляет арифметическое значение $b - a$, если рассматривать вычисленное значение как беззнаковое.

Доказательство. (Предполагается, что используется 32-разрядный компьютер.) Так как $a \leq b$, значения арифметической разности $b - a$ находятся в диапазоне от 0 до $(2^{31} - 1) - (-2^{31}) = 2^{32} - 1$. Если истинная разность находится в диапазоне от 0 до $2^{31} - 1$,

то вычисленная компьютером разность корректна (так как результат представим в знаковой интерпретации), а бит знака равен 0. Следовательно, вычисленный результат будет корректен независимо от того, как он интерпретируется — как знаковое или как беззнаковое целое число.

Если значение истинной разности находится между 2^{31} и $2^{32} - 1$, вычисленный компьютером результат будет отличаться от истинного на некоторое кратное 2^{32} (поскольку это число не может быть представлено при знаковой интерпретации). Таким образом, при знаковой интерпретации будет получен результат, лежащий в диапазоне от -2^{31} до -1 . Результат машинного вычисления оказывается уменьшенным на 2^{32} , при этом бит знака становится равным 1. Однако если рассматривать вычисленный результат как беззнаковое целое число, то его значение за счет использования бита знака (который при этом имеет вес $+2^{31}$, а не -2^{31}) увеличивается на 2^{32} . Следовательно, при рассмотрении результата как беззнакового целого числа он корректно представляет истинное значение.

Теперь мы можем перейти к "теореме границ".

ТЕОРЕМА. Если a и b — целые знаковые числа и $a \leq b$, то

$$a \leq x \leq b = x - a \overset{n}{\leq} b - a \quad (1)$$

Доказательство. В зависимости от значения x возможны три случая. Согласно рассмотренной ранее лемме, поскольку в любом из случаев $a \leq b$, значение $b - a$, рассматриваемое как беззнаковое целое число (что и происходит в уравнении (1)), равно арифметическому значению разности b и a .

Случай 1: $x < a$. При этом разность $x - a$, рассматриваемая как беззнаковое число, равна $x - a + 2^{32}$. Какими бы ни были значения x и b (в диапазоне 32-битовых чисел), выполняется неравенство

$$x + 2^{32} > b.$$

Таким образом,

$$x - a + 2^{32} > b - a,$$

а следовательно,

$$x - a \overset{n}{>} b - a,$$

и обе части уравнения (1) ложны.

Случай 2: $a \leq x \leq b$. В этом случае арифметически $x - a \leq b - a$. Так как $a \leq x$, согласно лемме истинное значение разности $x - a$ равно вычисленному значению $x - a$, если рассматривать последнее как беззнаковое. Следовательно,

$$x - a \overset{n}{\leq} b - a,$$

и обе части уравнения (1) истинны.

Случай 3: $x > b$. В таком случае $x - a > b - a$. Так как $b \geq a$, в этом случае $x > a$, и согласно лемме истинное значение разности $x - a$ равно вычисленному значению $x - a$, рассматриваемому как беззнаковое число. Следовательно,

$$x - a \overset{\circ}{>} b - a,$$

т.е. обе части уравнения (1) ложны.

Рассмотренная теорема справедлива и для *беззнаковых* целых a и b . Для беззнаковых чисел рассмотренная лемма выполняется тривиально, а значит, приведенное выше доказательство теоремы справедливо и для беззнаковых чисел.

Ниже приведено несколько аналогичных преобразований, полученных из теоремы границ и применимых как для знаковых, так и для беззнаковых целых a , b и x .

$$\text{Если } a \leq b, \text{ то } a \leq x \leq b = x - a \overset{\circ}{\leq} b - a = b - x \overset{\circ}{\leq} b - a.$$

$$\text{Если } a \leq b, \text{ то } a \leq x < b = x - a \overset{\circ}{<} b - a.$$

$$\text{Если } a \leq b, \text{ то } a < x \leq b = b - x \overset{\circ}{<} b - a.$$

$$\text{Если } a < b, \text{ то } a < x < b = x - a - 1 \overset{\circ}{<} b - a - 1 = b - x - 1 \overset{\circ}{<} b - a - 1.$$

(2)

В последнем правиле $b - a - 1$ можно заменить на $b + \neg a$.

При проверке граничных условий вида $-2^{n-1} \leq x \leq 2^{n-1} - 1$ лучше использовать другие преобразования. Проверка таких условий определяет, может ли знаковое число x быть корректно представлено как n -разрядное целое число в дополнительном до 2 коде. Ниже для иллюстрации приводится несколько равноценных методов проверки для $n = 8$.

$$\text{а) } -128 \leq x \leq 127$$

$$\text{б) } x + 128 \overset{\circ}{\leq} 255$$

$$\text{в) } \left(x \overset{\circ}{\gg} 7 \right) + 1 \overset{\circ}{\leq} 1$$

$$\text{г) } x \overset{\circ}{\gg} 7 = x \overset{\circ}{\gg} 31$$

$$\text{д) } \left(x \overset{\circ}{\gg} 7 \right) + \left(x \overset{\circ}{\gg} 31 \right) = 0$$

$$\text{е) } \left(x \ll 24 \right) \overset{\circ}{\gg} 24 = x$$

$$\text{ж) } x \oplus \left(x \overset{\circ}{\gg} 31 \right) \leq 127$$

Соотношения (б) и (в) получены непосредственно из предыдущего материала главы (в формуле (в) преобразование применено к значению x , сдвинутому вправо на 7 бит). Выражения (в)–(е) (а возможно, и (ж)), вероятно, следует использовать, только если константы из (а) и (б) превышают размеры полей непосредственно задаваемых значений в командах *сравнения* и *сложения*.

Для еще одного частного случая проверки границ, включающего степень 2, применимо преобразование

$$0 \leq x \leq 2^n - 1 \Leftrightarrow (x \dot{\gg} n) = 0$$

или, в более общем виде,

$$a \leq x \leq a + 2^n - 1 \Leftrightarrow ((x - a) \dot{\gg} n) = 0.$$

4.2. Определение границ суммы и разности

Ряд оптимизирующих компиляторов в процессе работы выполняет анализ области значений выражений, который представляет собой процесс определения верхней и нижней границ значений арифметического выражения в программе. Хотя такая оптимизация реально не приводит к большому выигрышу, она позволяет вносить в программу определенные улучшения, например избегать проверки диапазона в операторе выбора `switch` языка C или проверки значений индексов в режиме отладки.

Предположим, что диапазон значений двух переменных, x и y , задается следующим образом (все величины беззнаковые):

$$\begin{aligned} a \leq x \leq b, \\ c \leq y \leq d. \end{aligned} \quad (3)$$

Каковы в этом случае компактные границы значений выражений $x + y$, $x - y$ и $-x$? Очевидно, что арифметически $a + c \leq x + y \leq b + d$. Но проблема в том, что при вычислении вполне может возникнуть переполнение.

Один из способов определения диапазона значений основан на следующей теореме.

ТЕОРЕМА. Если a, b, c, d, x и y — беззнаковые целые числа, причем

$$\begin{aligned} a \dot{\leq} x \dot{\leq} b \quad \text{и} \\ c \dot{\leq} y \dot{\leq} d, \end{aligned}$$

то

$$\begin{aligned} 0 \dot{\leq} x + y \dot{\leq} 2^{32} - 1, \quad \text{если } a + c \dot{\leq} 2^{32} - 1 \text{ и } b + d \geq 2^{32}, \\ a + c \dot{\leq} x + y \dot{\leq} b + d \quad \text{в противном случае;} \end{aligned} \quad (4)$$

$$\begin{aligned} 0 \dot{\leq} x - y \dot{\leq} 2^{32} - 1, \quad \text{если } a - d < 0 \text{ и } b - c \geq 0, \\ a - d \dot{\leq} x - y \dot{\leq} b - c \quad \text{в противном случае;} \end{aligned} \quad (5)$$

$$\begin{aligned} 0 \overset{\cdot}{\leq} x \overset{\cdot}{\leq} 2^{32} - 1, & \text{ если } a = 0 \text{ и } b \neq 0, \\ -b \overset{\cdot}{\leq} x \overset{\cdot}{\leq} a & \text{ в противном случае.} \end{aligned} \quad (6)$$

Неравенство (4) гласит, что значения суммы $x + y$ *обычно* находятся в диапазоне от $a + c$ до $b + d$. Однако если при вычислении $b + d$ было переполнение, а при вычислении $a + c$ переполнения *не* было, то нижней и верхней границами суммы $x + y$ будут соответственно 0 и максимальное беззнаковое целое число. Выражения (5) интерпретируются аналогично, но значение истинной разности, меньшее 0, означает, что имело место переполнение (в отрицательном направлении).

Доказательство. Пусть значения x и y находятся в указанных выше диапазонах. Если при вычислении обеих сумм — $a + c$ и $b + d$ — переполнения нет, то и при вычислении суммы $x + y$ переполнения не будет, и вычисленный результат будет правильным (так что выполняется второе неравенство из выражения (4)). Если же переполнение возникло при вычислении как $a + c$, так и $b + d$, то при вычислении $x + y$ также произойдет переполнение. С арифметической точки зрения ясно, что

$$a + c - 2^{32} \leq x + y - 2^{32} \leq b + d - 2^{32}.$$

Но именно так и выполняются вычисления в случае, когда переполнение происходит во всех трех суммах. Следовательно, и в этом случае

$$a + c \overset{\cdot}{\leq} x + y \overset{\cdot}{\leq} b + d.$$

Если при вычислении $a + c$ переполнения нет, а при вычислении $b + d$ — есть, то

$$a + c \leq 2^{32} - 1 \text{ и } b + d \geq 2^{32}.$$

Поскольку $x + y$ может принимать все значения из диапазона от $a + c$ до $b + d$, значение этой суммы находится между $2^{32} - 1$ и 2^{32} , т.е. вычисленное значение суммы $x + y$ будет находиться между $2^{32} - 1$ и 0 (хотя и не принимает *все* значения из указанного диапазона).

Случай, когда при вычислении $a + c$ возникает переполнение, а при вычислении $b + d$ не возникает, невозможен, поскольку $a \leq b$ и $c \leq d$.

Доказательство неравенств (4) завершено. Неравенства (5) доказываются аналогично, с тем отличием, что здесь "переполнение" означает, что истинное значение разности меньше 0.

При доказательстве неравенств (6) можно воспользоваться неравенствами (5), в которых $a = b = 0$, а затем переименовать переменные. (Выражение $-x$, где x — беззнаковое число, означает вычисление значения $2^{32} - x$ или $-x + 1$, как вам больше нравится.)

Поскольку беззнаковое переполнение легко распознается (см. раздел "Беззнаковое сложение/вычитание" на с. 52), полученные результаты несложно перевести в код, показанный в листинге 4.1, в котором вычисляются верхняя (переменная s) и нижняя (переменная t) границы суммы и разности.

Листинг 4.1. Верхняя и нижняя границы беззнаковых суммы и разности

```

s = a + c;           s = a - d;
t = b + d;           t = b - c;
if (s >= a && t < b) if (s > a && t <= b)
{                     {
    s = 0;             s = 0;
    t = 0xFFFFFFFF;  t = 0xFFFFFFFF;
}                     }

```

Знаковые числа

При сложении и вычитании целых знаковых чисел не все так просто. Как и прежде, предположим, что диапазоны значений x и y задаются следующим образом (все величины здесь — знаковые).

$$a \leq x \leq b$$

$$c \leq y \leq d$$

Мы хотим определить точные границы значений $x + y$, $x - y$ и $-x$. Доказательство формул для знаковых величин очень похоже на доказательство в беззнаковом случае, так что приведем только конечный результат для сложения.

$$\begin{aligned}
 a+c < -2^{31}, b+d < -2^{31} &: a+c \leq x+y \leq b+d \\
 a+c < -2^{31}, b+d \geq -2^{31} &: -2^{31} \leq x+y \leq 2^{31}-1 \\
 -2^{31} \leq a+c < 2^{31}, b+d < 2^{31} &: a+c \leq x+y \leq b+d \\
 -2^{31} \leq a+c < 2^{31}, b+d \geq 2^{31} &: -2^{31} \leq x+y \leq 2^{31}-1 \\
 a+c \geq 2^{31}, b+d \geq 2^{31} &: a+c \leq x+y \leq b+d
 \end{aligned} \tag{8}$$

Первая строка означает, что если при вычислении обеих сумм — $a+c$ и $b+d$ — было переполнение в отрицательном направлении, то вычисленная сумма $x+y$ будет находиться между вычисленными суммами $a+c$ и $b+d$, так как все три вычисленные суммы оказываются слишком велики на одну и ту же величину (2^{32}). Вторая строка означает, что если при вычислении суммы $a+c$ происходит переполнение в отрицательном направлении, а при вычислении $b+d$ переполнения либо не было, либо оно было в положительном направлении, то вычисленное значение суммы $x+y$ находится между крайним отрицательным и крайним положительным числами (хотя и может принимать не все значения из указанного диапазона). Остальные строки интерпретируются аналогично.

Если представить выражение для границ y в виде

$$-d \leq -y \leq -c$$

и воспользоваться уже имеющимися выражениями для суммы, то получим формулы, по которым определяется диапазон значений в результате вычитания знаковых чисел.

$$\begin{aligned}
 a-d < -2^{31}, b-c < -2^{31} &: a-d \leq x-y \leq b-c \\
 a-d < -2^{31}, b-c \geq -2^{31} &: -2^{31} \leq x-y \leq 2^{31}-1 \\
 -2^{31} \leq a-d < 2^{31}, b-c < 2^{31} &: a-d \leq x-y \leq b-c \\
 -2^{31} \leq a-d < 2^{31}, b-c \geq 2^{31} &: -2^{31} \leq x-y \leq 2^{31}-1 \\
 a-d \geq 2^{31}, b-c \geq 2^{31} &: a-d \leq x-y \leq b-c
 \end{aligned}$$

Формулы для отрицания (изменения знака) можно получить из выражений для разности знаковых чисел в предположении $a = b = 0$, отбросив ряд невозможных комбинаций, переименовав переменные и проведя некоторые упрощения.

$$\begin{aligned}
 a = -2^{31}, b = -2^{31} &: -x = -2^{31} \\
 a = -2^{31}, b \neq -2^{31} &: -2^{31} \leq -x \leq 2^{31}-1 \\
 a \neq -2^{31} &: -b \leq -x \leq -a
 \end{aligned}$$

Программа на языке C, вычисляющая границы арифметических выражений, несколько сложновата, поэтому рассмотрим только код определения диапазона значений суммы двух знаковых операндов. По-видимому, проще всего выполнить проверку границ в двух случаях из (7) — когда нижняя и верхняя границы равны крайним отрицательному и положительному числам соответственно. Переполнение в отрицательном направлении имеет место, если оба операнда отрицательны, а их сумма — неотрицательна (см. раздел “Знаковое сложение и вычитание” на с. 49). Таким образом, для проверки условия $a+c < -2^{31}$ сначала необходимо вычислить сумму $s = a + c$; а затем использовать код наподобие `if (a < 0 && c < 0 && s >= 0) ...`. Однако программа будет более эффективной¹, если все логические операции будут выполняться непосредственно над арифметическими переменными, при этом результат будет находиться в бите знака. Тогда рассмотренное условие можно представить иначе: `if ((a & c & ~s) < 0) ...`. Таким образом, получаем фрагмент кода, приведенный в листинге 4.2.

Листинг 4.2. Определение границ знакового сложения

```

s = a + c;
t = b + d;
u = a & c & ~s & ~(b & d & ~t);
v = ((a ^ c) | ~(a ^ s)) & (~b & ~d & t);
if ((u | v) < 0)
{
    s = 0x80000000;
    t = 0x7FFFFFFF;
}

```

Переменная `u` имеет значение `true` (т.е. бит знака равен 1), если при вычислении $a + c$ происходит переполнение в отрицательном направлении, а при вычислении $b + d$

¹ Другими словами, более компактной, с минимальным количеством переходов. В предположении, что границы обычно не слишком велики, проверку следует начинать для случая отсутствия переполнения, так как велика вероятность получить результат сразу же, в ходе этой проверки.

переполнения в отрицательном направлении *нет*. Переменная v имеет значение *true*, если при вычислении $a + c$ переполнения *нет*, а при вычислении $b + d$ имеет место переполнение в положительном направлении. Первое условие можно записать как “ a и c имеют разные знаки или a и c имеют один и тот же знак”. Проверка $\text{if}((u|v) < 0)$ эквивалентна проверке $\text{if}(u < 0 | v < 0)$, т.е. если истинно или u , или v (или оба одновременно).

4.3. Определение границ логических выражений

Как и в предыдущем разделе, предположим, что границы двух переменных, x и y , задаются следующим образом (все величины беззнаковые).

$$\begin{aligned} a \leq x \leq b \\ c \leq y \leq d \end{aligned} \quad (8)$$

Необходимо определить точные границы выражений $x|y$, $x \& y$, $x \oplus y$ и $\neg x$.

Комбинируя неравенства (8) с неравенствами из раздела 2.3 на с. 38 и учитывая, что $\neg x = 2^{31} - 1 - x$, получим следующее.

$$\begin{aligned} \max(a, c) \leq (x|y) \leq b + d \\ 0 \leq (x \& y) \leq \min(b, d) \\ 0 \leq (x \oplus y) \leq b + d \\ \neg b \leq \neg x \leq \neg a \end{aligned}$$

Здесь предполагается, что при сложении $b + d$ переполнения не возникает. Эти формулы легко вычисляются и достаточно эффективны для использования в компиляторах, о чем упоминалось в предыдущем разделе. Однако границы в первых двух неравенствах не являются точными. Пусть, например, диапазон значений переменных x и y в двоичной записи задан следующим образом.

$$\begin{aligned} 00010 \leq x \leq 00100 \\ 01001 \leq y \leq 10100 \end{aligned} \quad (9)$$

После проверки (перебора всех 36 возможных комбинаций значений x и y) получим, что $01010 \leq (x|y) \leq 10111$, т.е. нижняя граница не равна ни $\max(a, c)$, ни $a|c$; верхняя граница также не равна ни $b + d$, ни $b|d$.

Можно ли определить точные границы логических выражений, если известны значения a , b , c и d из (8)? Давайте сначала найдем минимальное значение выражения $x|y$. Разумно предположить, что значение этого выражения будет минимальным при минимальных значениях x и y , равным $a|c$. Однако, как видно из примера (9), минимальное значение на самом деле может оказаться меньше предполагаемого.

Для поиска минимума начнем со значений $x = a$ и $y = c$ и будем искать значения x и y , при которых значение выражения $x|y$ будет минимальным. Этот результат и будет искомым минимальным значением. Вместо того чтобы присваивать значения a и c пере-

менным x и y , будем работать непосредственно со значениями a и c , увеличивая одно из них так, чтобы снизить значение $a | c$.

Процедура состоит в сканировании битов a и c слева направо. Если оба бита равны 0, то в соответствующем разряде результата будет 0. Если оба бита равны 1, то и соответствующий бит результата будет равен 1 (очевидно, никакие значения x и y не могут уменьшить полученный результат). В обоих этих случаях переходим к сравнению следующих разрядов. Если в следующем разряде одной из границ содержится единичный бит, а у другой границы — нулевой, то, возможно, замена 0 на 1 в этом разряде и установка всех следующих за ним битов равными 0 снизит значение $a | c$. Такое изменение границ не может увеличить значение $a | c$, так как в любом случае в этом разряде результат будет равным 1 за счет другой границы. Итак, мы формируем число с нулевым битом, замененным единичным, а все последующие биты этого числа заменяем нулями. Если полученное таким образом число меньше или равно соответствующей верхней границе, такая замена возможна. Выполним эту замену, а затем применим операцию *или* к полученному значению и другой нижней границе. Если же после описанной замены битов новое значение нижней границы превосходит соответствующую верхнюю границу, то такая замена невозможна, и мы переходим к следующим битам.

Вот и все. Казалось бы, после замены следует продолжить сканирование, чтобы искать дальнейшие возможности улучшения полученного результата — еще меньшее значение $a | c$. Однако, даже если в одном из следующих разрядов можно заменить 0 на 1, установка следующих за ним битов равными нулю не даст меньшего значения $a | c$, так как в этих разрядах уже содержатся нули.

Программа на языке C для данного алгоритма представлена в листинге 4.3. Предполагается, что оптимизирующий компилятор вынесет вычисление подвыражений $\sim a \& c$ и $a \& \sim c$ за пределы цикла. Если у нас есть функция вычисления *количества ведущих нулевых битов*, код можно ускорить, инициализируя m следующим образом.

```
m = 0x80000000 >> nlz(a^c);
```

При этом в $a \wedge c$ опускаются те разряды, в которых изначально у обеих констант оба бита равны 0 или оба равны 1. Это ускорение работы программы будет особенно эффективным при $a \wedge c$ равном 0 (т.е. если $a = c$); при этом сдвиг вправо должен выполняться по модулю 64. Если функции `nlz` в вашем распоряжении нет, можно воспользоваться одной из версий функции `flr2` (см. с. 82) с аргументом $a \wedge c$.

ЛИСТИНГ 4.3. Вычисление нижней границы $x | y$

```
unsigned minOR(unsigned a, unsigned b,
               unsigned c, unsigned d)
```

```
{
    unsigned m, temp;

    m = 0x80000000;
    while (m != 0)
    {
        if (~a & c & m)
        {
```

```

        temp = (a | m) & ~m;
        if (temp <= b) {a = temp; break;}
    }
    else if (a & ~c & m)
    {
        temp = (c | m) & ~m;
        if (temp <= d) {c = temp; break;}
    }
    m = m >> 1;
}
return a | c;
}

```

Рассмотрим теперь алгоритм поиска *максимального* значения выражения $x|y$, в котором, как и ранее, значения переменных ограничены неравенствами (8). Этот алгоритм аналогичен алгоритму поиска минимального значения. Верхние границы поразрядно сканируются слева направо до тех пор, пока в некоторой позиции биты в b и d не окажутся равными единицам. Если такая позиция обнаружена, алгоритм пытается увеличить значение $c|d$, уменьшая одну из границ заменой в соответствующем разряде 1 на 0 и заменяя все последующие биты единичными. Если такая замена корректна (т.е. полученное таким образом новое значение верхней границы не ниже соответствующей нижней границы), то изменения принимаются, и результатом является $c|d$ с использованием модифицированной границы. Если же замена оказывается некорректной, пытаемся изменить другую границу. Если и эта замена неприемлема, сканируются следующие биты. Программа на языке C, реализующая данный алгоритм, приведена в листинге 4.4. Здесь за пределы цикла можно вынести вычисление подвыражения $b \& d$; кроме того, работу программы можно ускорить следующей инициализацией.

```
m = 0x80000000 >> nlz(b & d);
```

Листинг 4.4. Вычисление верхней границы $x|y$

```

unsigned maxOR(unsigned a, unsigned b,
               unsigned c, unsigned d)
{
    unsigned m, temp;

    m = 0x80000000;
    while (m != 0)
    {
        if (b & d & m)
        {
            temp = (b - m) | (m - 1);
            if (temp >= a) {b = temp; break;}
            temp = (d - m) | (m - 1);
            if (temp >= c) {d = temp; break;}
        }
        m = m >> 1;
    }
    return b | d;
}

```

Определить диапазон значений выражения $x \& y$, где границы x и y определены неравенствами (8), можно с помощью двух методов: алгебраического и непосредственного вычисления. В основе алгебраического метода лежит закон де Моргана (DeMorgan).

$$x \& y = \neg(\neg x \mid \neg y)$$

Поскольку мы знаем, как определить точные границы выражения $x \mid y$, границы $x \& y$ вычисляются тривиально с помощью команды *ne* ($a \leq x \leq b \Leftrightarrow \neg b \leq \neg x \leq \neg a$).

$$\text{minAND}(a, b, c, d) = \neg \text{maxOR}(\neg b, \neg a, \neg d, \neg c)$$

$$\text{maxAND}(a, b, c, d) = \neg \text{minOR}(\neg b, \neg a, \neg d, \neg c)$$

Код вычисления этих функций, приведенный в листингах 4.5 и 4.6, очень похож на код вычисления границ выражения *c и* *и*.

ЛИСТИНГ 4.5. Вычисление нижней границы $x \& y$

```
unsigned minAND(unsigned a, unsigned b,
                unsigned c, unsigned d)
{
    unsigned m, temp;

    m = 0x80000000;
    while (m != 0)
    {
        if (~a & ~c & m)
        {
            temp = (a | m) & ~m;
            if (temp <= b) {a = temp; break;}
            temp = (c | m) & ~m;
            if (temp <= d) {c = temp; break;}
        }
        m = m >> 1;
    }
    return a & c;
}
```

ЛИСТИНГ 4.6. Вычисление верхней границы $x \& y$

```
unsigned maxAND(unsigned a, unsigned b,
                unsigned c, unsigned d)
{
    unsigned m, temp;

    m = 0x80000000;
    while (m != 0)
    {
        if (b & ~d & m)
        {
            temp = (b & ~m) | (m - 1);
            if (temp >= a) {b = temp; break;}
        }
        else if (~b & d & m)
        {

```

```

temp = (d & ~m) | (m - 1);
if (temp >= c) (d = temp; break;)
}
m = m >> 1;
}
return b & d;
}

```

Алгебраический метод поиска границ можно использовать для любых побитовых логических выражений, кроме *исключающего или* и *эквивалентности*. Причина в том, что при выражении этих операций через команды *и*, *или* и *не* получаются два члена, содержащие *x* и *y*. Например, для *исключающего или* имеем

$$\min_{csysd} (x \oplus y) = \min_{csysd} ((x \& \neg y) | (\neg x \& y)).$$

Указанные операнды команды *или* нельзя минимизировать по отдельности (не доказав предварительно, что так можно делать), поскольку мы ищем одно значение *x* и одно значение *y*, которые минимизируют выражение *или* в целом.

Для вычисления границ выражения с *исключающим или* можно использовать следующие выражения.

$$\begin{aligned} \min\text{XOR}(a, b, c, d) &= \min\text{AND}(a, b, \neg d, \neg c) | \min\text{AND}(\neg b, \neg a, c, d) \\ \max\text{XOR}(a, b, c, d) &= \max\text{OR}(0, \max\text{AND}(a, b, \neg d, \neg c), 0, \max\text{AND}(\neg b, \neg a, c, d)) \end{aligned}$$

Значения *minXOR* и *maxOR* несложно вычислить непосредственно. Код для вычисления *minXOR* такой же, как и для вычисления *minOR* (см. листинг 4.3), нужно лишь удалить два оператора *break*, а в последней строке заменить возвращаемое значение на *a ^ c*. Код для вычисления *maxXOR* тот же, что и для вычисления *maxOR* (см. листинг 4.4), за исключением того, что четыре строки после оператора *if* необходимо заменить таким фрагментом.

```

temp = (b - m) | (m - 1);
if (temp >= a) b = temp;
else
{
temp = (d - m) | (m - 1);
if (temp >= c) d = temp;
}

```

Кроме того, возвращаемое значение нужно заменить выражением *b ^ d*.

Знаковые границы

Вычисление целых *знаковых* границ для логических выражений значительно усложняется. Вычисления границ зависят от расположения числа 0 по отношению к заданным диапазонам (от *a* до *b* и от *c* до *d*), и для каждого частного случая используются свои формулы. В табл. 4.1 приведен один из методов вычисления верхней и нижней границ выражения *x | y*. Если значение переменной больше или равно 0, в соответствующей

графе таблицы стоит знак “плюс”, если значение переменной меньше 0 — знак “минус”. Под заголовком `minOR` приведены выражения, вычисляющие нижнюю границу $x|y$ для данных диапазонов, а под заголовком `maxOR` — формулы для вычисления верхней границы. Вычислить верхнюю и нижнюю границы выражения $x|y$ можно следующим образом: из битов знаков чисел a , b , c и d составляется 4-битовое число, значение которого находится в пределах от 0 до 15, и используется оператор выбора `switch`. Обратите внимание, что будут задействованы не все значения от 0 до 15, так как невозможны ситуации, когда $a > b$ или $c > d$.

Для знаковых чисел справедливо отношение

$$a \leq x \leq b \Leftrightarrow -b \leq -x \leq -a,$$

так что для распространения результатов из табл. 4.1 на другие логические выражения (кроме *исключающего или* и *эквивалентности*) может использоваться алгебраический метод. Выражения, определяющие границы других логических команд, читателю предлагается получить самостоятельно.

Таблица 4.1. Вычисление знаковых функций `minOR` и `maxOR` с помощью беззнаковых

a	b	c	d	<code>minOR</code> (знаковый)	<code>maxOR</code> (знаковый)
-	-	-	-	<code>minOR(a, b, c, d)</code>	<code>maxOR(a, b, c, d)</code>
-	-	-	+	a	-1
-	-	+	+	<code>minOR(a, b, c, d)</code>	<code>maxOR(a, b, c, d)</code>
-	+	-	-	c	-1
-	+	-	+	<code>min(a, c)</code>	<code>maxOR(0, b, 0, d)</code>
-	+	+	+	<code>minOR(a, 0xFFFFFFFF, c, d)</code>	<code>maxOR(0, b, c, d)</code>
+	+	-	-	<code>minOR(a, b, c, d)</code>	<code>maxOR(a, b, c, d)</code>
+	+	-	+	<code>minOR(a, b, c, 0xFFFFFFFF)</code>	<code>maxOR(a, b, 0, d)</code>
+	+	+	+	<code>minOR(a, b, c, d)</code>	<code>maxOR(a, b, c, d)</code>

Упражнения

1. Найдите границы $x - y$ для беззнаковых чисел, если

$$0 \leq x \leq b \quad \text{и} \\ 0 \leq y \leq d.$$

2. Покажите, как можно упростить функцию `maxOR` (листинг 4.4) на машине, оснащенной командой вычисления количества ведущих нулевых битов, в случае, когда либо $a = 0$, либо $c = 0$.

ГЛАВА 5

ПОДСЧЕТ БИТОВ

5.1. Подсчет единичных битов

На компьютере IBM Stretch (выпущенном в 1960 году) имелась возможность подсчета количества единичных битов в слове и количества ведущих нулевых битов слова, причем эти величины являлись побочным результатом выполнения любой логической операции! Функция подсчета количества единичных битов иногда называлась *степенью заполнения* (population count), например на машинах Stretch или SPARCv9.

Если на машине нет отдельной команды для вычисления этой функции, то для подсчета единичных битов в слове можно воспользоваться другими методами. Обычно исходное слово делится на 2-разрядные поля и в каждое поле помещается количество имевшихся в нем единичных битов. Затем значения, содержащиеся в соседних 2-разрядных полях, складываются и результат помещается в 4-разрядное поле и т.д. Более подробно этот метод обсуждается в [97] и иллюстрируется на рис. 5.1. Слово, в котором подсчитывается количество единичных битов, находится в первой строке, в последней строке содержится результат (равный 23 в десятичной системе счисления).

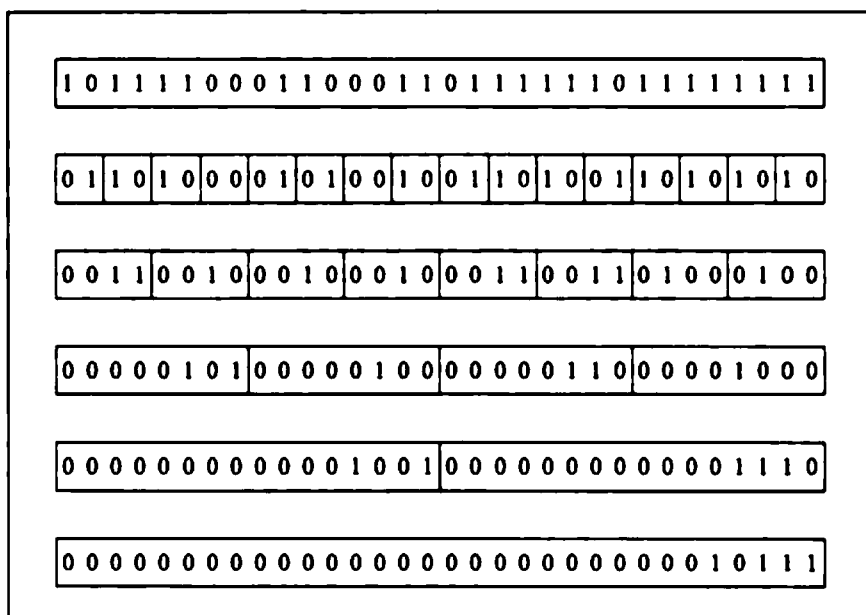


Рис. 5.1. Подсчет количества единичных битов; стратегия "разделяй и властвуй"

Рассмотренный метод — пример стратегии “разделяй и властвуй”: исходная задача (суммирование в 32-битовом слове) разбивается на две подзадачи (суммирование в 16-битовом слове), каждая из которых решается независимо от другой. Затем результаты объединяются (в нашем случае суммируются). Данная стратегия является рекурсивной: 16-битовые слова обрабатываются, как два 8-битовых, и т.д.

Выполнение подзадач на самом нижнем уровне (суммирование значений соседних битов) можно осуществлять параллельно, как и суммирование значений в соседних полях, которое можно выполнить параллельно за фиксированное число шагов на любом этапе. Весь алгоритм может быть выполнен за $\log_2(32) = 5$ шагов.

Другими примерами стратегии “разделяй и властвуй” являются хорошо всем известные алгоритмы бинарного поиска, быстрой сортировки, а также метод реверса битов слова, рассматриваемый в разделе 7.1 на с. 155.

Метод, приведенный на рис. 5.1, может быть выражен на языке C следующим образом.

```
x = (x & 0x55555555) + ((x >> 1) & 0x55555555);
x = (x & 0x33333333) + ((x >> 2) & 0x33333333);
x = (x & 0x0F0F0F0F) + ((x >> 4) & 0x0F0F0F0F);
x = (x & 0x00FF00FF) + ((x >> 8) & 0x00FF00FF);
x = (x & 0x0000FFFF) + ((x >> 16) & 0x0000FFFF);
```

В первой строке можно было бы написать, пожалуй, более естественное выражение $(x \& 0xAAAAAAAA) \gg 1$, но в данном случае лучше использовать приведенное выражение $(x \gg 1) \& 0x55555555$, так как при этом удастся избежать загрузки в регистры двух больших констант, что на машине с отсутствующей командой *и-не* будет стоить дополнительной команды. Аналогичные замечания справедливы и для остальных строк кода.

Очевидно, что последняя команда *и* не является необходимой; кроме того, если точно известно, что при суммировании полей не произойдет переноса в соседнее поле, остальные команды *и* также не нужны. Этот код можно еще немного усовершенствовать, а именно — преобразовать первую строку так, чтобы в ней выполнялось на одну команду меньше. Такой упрощенный код показан в листинге 5.1; он не содержит условных переходов и выполняется за 21 команду.

Листинг 5.1. Подсчет количества единичных битов в слове

```
int pop(unsigned x)
{
    x = x - ((x >> 1) & 0x55555555);
    x = (x & 0x33333333) + ((x >> 2) & 0x33333333);
    x = (x + (x >> 4)) & 0x0F0F0F0F;
    x = x + (x >> 8);
    x = x + (x >> 16);
    return x & 0x0000003F;
}
```

Первое присвоение *x* основано на первых двух членах одной удивительной формулы.

$$\text{pop}(x) = x - \left\lfloor \frac{x}{2} \right\rfloor - \left\lfloor \frac{x}{4} \right\rfloor - \dots - \left\lfloor \frac{x}{2^{31}} \right\rfloor \quad (1)$$

В формуле (1) значение x должно быть не меньше 0. Если x — целое беззнаковое число, формула (1) может быть реализована в виде последовательности из 31 команды *сдвига вправо на единицу* и 31 команды *вычитания*. Процедура в листинге 5.1 использует параллельно первые два члена этой формулы для каждого двухбитового поля.

Существует простое доказательство уравнения (1), приведенное ниже для случая 4-битового слова. Пусть имеется слово $b_3b_2b_1b_0$, в котором каждое b_i равно 0 или 1. Тогда можно записать следующее.

$$\begin{aligned} x - \left\lfloor \frac{x}{2} \right\rfloor - \left\lfloor \frac{x}{4} \right\rfloor - \left\lfloor \frac{x}{8} \right\rfloor &= b_3 \cdot 2^3 + b_2 \cdot 2^2 + b_1 \cdot 2^1 + b_0 \cdot 2^0 \\ &\quad - (b_3 \cdot 2^2 + b_2 \cdot 2^1 + b_1 \cdot 2^0) \\ &\quad - (b_3 \cdot 2^1 + b_2 \cdot 2^0) \\ &\quad - (b_3 \cdot 2^0) \\ &= b_3 (2^3 - 2^2 - 2^1 - 2^0) + b_2 (2^2 - 2^1 - 2^0) + b_1 (2^1 - 2^0) + b_0 (2^0) \\ &= b_3 + b_2 + b_1 + b_0 \end{aligned}$$

Формулу (1) можно получить и иначе, например заметив, что i -й бит в двоичном представлении неотрицательного числа x вычисляется по формуле

$$b_i = \left\lfloor \frac{x}{2^i} \right\rfloor - 2 \left\lfloor \frac{x}{2^{i+1}} \right\rfloor,$$

и вычислив сумму всех b_i , где i принимает значения от 0 до 31. Последний член будет равен нулю, так как $x < 2^{32}$.

Формулу (1) можно обобщить для других систем счисления. Например, для системы счисления с основанием 10 получим следующее выражение.

$$\text{sum_digits}(x) = x - 9 \left\lfloor \frac{x}{10} \right\rfloor - 9 \left\lfloor \frac{x}{100} \right\rfloor - \dots$$

Здесь слагаемые вычисляются до тех пор, пока их значение не становится равным нулю. Полученная формула доказывается аналогично предыдущей.

Если применить описанный алгоритм для системы счисления по основанию 4, можно получить другой код для вычисления функции $\text{pop}(x)$, который аналогичен приведенному в листинге 5.1, но вторая строка которого заменяется строкой

```
x = x - 3 * ((x >> 2) & 0x33333333);
```

Этот код выполняется за то же количество команд, что и замененный, но требует наличия команды *быстрого умножения на 3*.

При подсчете количества единичных битов в слове по алгоритму, предложенному в НАКМЕМ [44, item 169], используются только три первых члена из формулы (1) для формирования трехбитовых полей, в каждом из которых помещается сумма его единичных битов. Затем значения, содержащиеся в соседних трехбитовых полях, складываются

ся, а результат помещается в шестибитовое поле. После этого суммируются значения, содержащиеся в шестибитовых полях, вычисляя значение слова по модулю 63. На языке C этот алгоритм записывается следующим образом (обратите внимание, что длинные константы здесь записаны в восьмеричной системе счисления).

```
int pop(unsigned x)
{
    unsigned n;

    n = (x >> 1) & 033333333333;           // Подсчет битов
    x = x - n;                               // в 3-битовых
    n = (n >> 1) & 033333333333;           // полей
    x = x - n;                               //
    x = (x + (x >> 3)) & 030707070707;     // 6-битовые суммы
    return x%63;                             // Сложение 6-битовых сумм
}
```

В последней строке используется беззнаковая функция получения значения x по модулю 63 (это может быть как знаковое, так и беззнаковое число, если длина слова кратна 3). Данная функция складывает значения, содержащиеся в 6-битовых полях слова x , что становится понятно, если рассматривать слово x как целое число, записанное в системе счисления с основанием 64. Остаток от деления целого числа в системе счисления с основанием b на $b-1$ для $b \geq 3$ сравним по модулю $b-1$ с суммой цифр этого числа и конечно, меньше $b-1$. Так как сумма цифр в нашем случае не превышает 32, значение $\text{mod}(x, 63)$ должно быть равно сумме цифр x , т.е. количеству единичных битов в исходном слове.

На компьютерах DEC PDP-10 есть команда, вычисляющая остаток от деления, второй операнд которой содержит адрес слова в памяти, а потому на этих машинах подсчет единичных битов в слове выполняется за 10 команд. На компьютере RISC с базовым набором команд требуется около 13 команд, если в их числе имеется отдельная команда *беззнакового остатка от деления по модулю*. Однако, вероятно, это не очень быстрый метод, так как деление — почти всегда операция довольно медленная. Кроме того, простое расширение констант не позволяет работать с 64-битовыми словами, хотя вполне применимо для слов длиной до 62 бит.

Инструкция `return` в приведенном выше коде может быть заменена следующей, которая на большинстве машин работает быстрее, хотя и выглядит менее элегантно (здесь также используются восьмеричные константы).

```
return ((x*0404040404) >> 26) + (x >> 30);
```

В [47] предложена вариация алгоритма НАКМЕМ. Сначала с использованием формулы (1) параллельно подсчитывается количество единичных битов в каждом из восьми 4-битовых полей. Полученные 4-битовые суммы преобразуются в 8-битовые суммы, а затем четыре байта суммируются с помощью умножения на число `0x01010101`. Вот как это выглядит на языке программирования C.

```
int pop(unsigned x)
{
    unsigned n;
```

```

n = (x >> 1) & 0x77777777;           // Подсчет битов
x = x - n;                             // в 4-битовых
n = (n >> 1) & 0x77777777;           // полях
x = x - n;
n = (n >> 1) & 0x77777777;
x = x - n;
x = (x + (x >> 4)) & 0x0F0F0F0F;      // Вычисление сумм
x = x*0x01010101;                     // Сложение байтов
return x >> 24;
}

```

Здесь для вычислений требуется 19 базовых RISC-команд. Этот код эффективно работает на двухадресной машине, так как в первых шести строках выполняется только одна команда *пересылки регистра*. Многократно используемую маску 0x77777777 можно однократно загрузить в регистр и обращаться к ней с помощью команд, работающих только с регистрами. Кроме того, большинство сдвигов в этом коде — на один разряд.

Совсем другой метод подсчета битов предложен в [109, 97] (листинг 5.2). Здесь циклически сбрасывается крайний справа единичный бит исследуемого слова до тех пор, пока это слово не станет равным 0. Этот метод, для работы которого требуется выполнение $2 + 5 \text{pop}(x)$ команд, эффективно работает со словами, в которых содержится только небольшое количество единичных битов.

ЛИСТИНГ 5.2. Подсчет единичных битов в малозаполненных словах

```

int pop(unsigned x)
{
    int n;

    n = 0;
    while (x != 0)
    {
        n = n + 1;
        x = x & (x - 1);
    }
    return n;
}

```

Для слов с большим количеством единичных битов применяется дуальный алгоритм. В слове циклически устанавливается крайний справа нулевой бит ($x = x | (x+1)$) до тех пор, пока во всех разрядах слова не оказываются единицы (т.е. пока слово не станет равным -1). После этого возвращается значение $32 - n$. (В других вариантах этого алгоритма биты исходного значения x могут быть инвертированы, начальное значение n может быть задано равным 32 и при вычислении уменьшаться, а не увеличиваться.)

Еще один интересный алгоритм подсчета единичных битов предложен в [85]. В этом алгоритме вычисляется сумма всех 32 слов, полученных в результате циклического сдвига слова влево на один разряд. Итоговая сумма равна значению $\text{pop}(x)$ со знаком минус! Т.е.

$$\text{pop}(x) = -\sum_{i=0}^{31} \left(x \ll i \right), \quad (2)$$

где все сложения выполняются по модулю размера слова, а итоговая сумма рассматривается как целое, дополненное до 2. Этот способ — не более чем красивая, но не эффективная для практического применения формула: цикл выполняется 31 раз, так что для вычисления требуется выполнение 63 команд, не считая команд управления циклом.

Чтобы понять, как работает уравнение (2), посмотрим, что происходит с каждым отдельным битом слова x . Этот бит при циклическом сдвиге оказывается в каждом разряде, и при сложении всех 32 чисел, полученных в результате 31 циклического сдвига, в итоговом слове в каждом разряде будет единица. Но это число представляет собой -1 . Для иллюстрации рассмотрим 6-битовое слово $x = 001001$.

$$\begin{array}{r}
 001001 \quad x \\
 010010 \quad x \ll 1 \\
 100100 \quad x \ll 2 \\
 001001 \quad x \ll 3 \\
 010010 \quad x \ll 4 \\
 100100 \quad x \ll 5
 \end{array}$$

Разумеется, циклический сдвиг вправо работает не менее корректно.

Метод подсчета единичных битов по формуле (1) очень похож на метод циклического сдвига и сложения, что становится понятно, если уравнение (1) переписать следующим образом.

$$\text{pop}(x) = x - \sum_{i=1}^{31} (x \gg i)$$

Вычисление функции $\text{pop}(x)$ по этой формуле будет происходить немного быстрее по сравнению с подсчетом единичных битов по формуле (2) по двум причинам. Во-первых, в нем используется команда сдвига вправо, которая, в отличие от команды циклического сдвига, есть практически на всех машинах; во-вторых, если при сдвиге слова получается нулевое значение, цикл завершается. Таким образом, это позволяет упростить цикл и сэкономить несколько итераций. Сравнение алгоритмов показано в листинге 5.3.

ЛИСТИНГ 5.3. Циклические алгоритмы подсчета количества единичных битов

```

int pop(unsigned x)
{
    int i, sum;

    // Циклический сдвиг и сложение
    sum = x;
    for(i = 1; i <= 31; i++)
    {
        x = rotatel(x,1);
        sum = sum + x;
    }

    // Сдвиг вправо и вычитание
    // sum = x;
    // while(x != 0)
    // {
    //     x = x >> 1;
    //     sum = sum - x;
    // }

```

```

)                                     // }
    return -sum;                       // return sum;
)

```

Алгоритм, вычисляющий функцию $\text{pop}(x)$ с использованием массива значений, менее оригинален по сравнению с предыдущими алгоритмами, но вполне может с ними конкурировать. Например, пусть для значений от 0 до 255 имеется массив `table` значений функции $\text{pop}(x)$. При подсчете единичных битов в слове выполняются четыре обращения к массиву `table` и затем складываются четыре полученных числа. Вот один из вариантов этого алгоритма (без условных переходов).

```

int pop(unsigned x)           // Поиск в таблице
{
    static char table[256] = {
        0, 1, 1, 2, 1, 2, 2, 3, 1, 2, 2, 3, 2, 3, 3, 4,
        ...
        4, 5, 5, 6, 5, 6, 6, 7, 5, 6, 6, 7, 6, 7, 7, 8};

    return table [x           & 0xFF] +
           table [(x >> 8) & 0xFF] +
           table [(x >> 16) & 0xFF] +
           table [(x >> 24)];
}

```

В [44, item 167] приведен краткий алгоритм подсчета количества единичных битов в 9-битовом значении, выровненном вправо и размещенном в регистре. Этот алгоритм работает только на машинах с длиной регистра 36 и больше бит. Ниже приведен вариант этого алгоритма для работы на 32-разрядной машине, но только для 8-битовых величин.

```

x = x * 0x08040201;           // Создаются 4 копии
x = x >> 3;                   // Удаление соответствующих битов
x = x & 0x11111111;          // Каждый 4-й бит
x = x * 0x11111111;          // Сумма цифр (0 или 1)
x = x >> 28;                  // Позиция результата

```

Вот версия алгоритма для 7-битовых величин.

```

x = x * 0x02040810;          // Создаются 4 копии
x = x & 0x11111111;          // Каждый 4-й бит
x = x * 0x11111111;          // Сумма цифр (0 или 1)
x = x >> 28;                  // Позиция результата

```

Здесь последние два шага могут быть заменены вычислением остатка от x по модулю 15.

Эти методы работают не так уж эффективно. Большинство программистов предпочитают при подсчете единичных битов пользоваться таблицами значений этой функции. Последний из рассмотренных алгоритмов, однако, имеет версию с использованием 64-битовой арифметики и может оказаться полезным на 64-разрядных машинах, у которых есть команда быстрого умножения. Аргументом данного алгоритма является 15-битовая величина (я не думаю, что есть похожий алгоритм для 16-битовых величин, кроме случая, когда заранее известно, что не все 16 бит являются единичными). Тип данных `long long` представляет собой расширение языка C во множестве старых

и новых компиляторов и означает целое число длиной 64 бита. Этот тип является официально признанным в стандарте C99. Константы типа `unsigned long long` указываются с помощью суффикса `ULL`.

```
int pop(unsigned x)
{
    unsigned long long y;
    y = x * 0x0002000400080010ULL;
    y = y & 0x11111111111111ULL;
    y = y * 0x11111111111111ULL;
    y = y >> 60;
    return y;
}
```

Сумма и разность количества единичных битов в двух словах

При вычислении $\text{pop}(x) + \text{pop}(y)$ (если ваш компьютер не оснащен командой вычисления *степени заполнения*) можно сэкономить немного времени, используя первые две строки листинга 5.1 для x и y отдельно, после чего вычисляя следующие три стадии алгоритма для суммы. После выполнения первых двух строк листинга 5.1 x и y состоят из восьми 4-битовых полей, каждое из которых содержит значение, не превышающее 4. Таким образом, x и y могут быть безопасно суммированы, так как максимальное значение в каждом 4-битовом поле не превысит 8, так что никакое переполнение при этом произойти не может. (Фактически таким способом можно комбинировать целых три слова.)

Эта идея применима и для вычитания. Чтобы вычислить $\text{pop}(x) - \text{pop}(y)$, используйте

$$\begin{aligned} \text{pop}(x) - \text{pop}(y) &= \text{pop}(x) - (32 - \text{pop}(\bar{y})) \\ &= \text{pop}(x) + \text{pop}(\bar{y}) - 32 \end{aligned}$$

После этого только что описанная методика применяется для вычисления $\text{pop}(x) + \text{pop}(\bar{y})$. Соответствующий код показан в листинге 5.4. В нем используются 32 команды, в то время как выполнение двух кодов из листинга 5.1 с последующим вычитанием требует 43 команд.

ЛИСТИНГ 5.4. Вычисление $\text{pop}(x) - \text{pop}(y)$

```
int popDiff(unsigned x, unsigned y)
{
    x = x - ((x >> 1) & 0x55555555);
    x = (x & 0x33333333) + ((x >> 2) & 0x33333333);
    y = ~y;
    y = y - ((y >> 1) & 0x55555555);
    y = (y & 0x33333333) + ((y >> 2) & 0x33333333);
    x = x + y;
    x = (x & 0x0F0F0F0F) + ((x >> 4) & 0x0F0F0F0F);
    x = x + (x >> 8);
    x = x + (x >> 16);
    return (x & 0x0000007F) - 32;
}
```

Сравнение степени заполнения двух слов

Иногда требуется выяснить, в каком из двух слов больше единичных битов, безотносительно к их фактическим значениям. Можно ли выяснить это, не прибегая к вычислению степени заполнения каждого из слов? Можно вычислить разность степеней заполнения двух слов методом, представленным в листинге 5.4, и сравнить полученный результат с 0, однако есть и другой способ, более предпочтительный в случае, если степени заполнения малы или если имеется сильная корреляция между установленными битами в двух словах.

Идея заключается в сбросе по одному биту в каждом слове, пока одно из слов не станет нулевым. Этот процесс работает быстрее в наихудшем и среднем случаях, если биты, равные 1, находятся в одинаковых позициях в каждом слове. Соответствующий код приведен в листинге 5.5. Процедура возвращает отрицательное целое число, если $\text{pop}(x) < \text{pop}(y)$, нуль, если $\text{pop}(x) = \text{pop}(y)$, и положительное целое число (1), если $\text{pop}(x) > \text{pop}(y)$.

Листинг 5.5. Сравнение $\text{pop}(x)$ и $\text{pop}(y)$

```
int popCmpr(unsigned xp, unsigned yp)
{
    unsigned x, y;
    x = xp & ~yp;          // Сброс в обоих словах битов,
    y = yp & ~xp;          // где они оба равны 1
    while (1)
    {
        if (x == 0) return y | -y;
        if (y == 0) return 1;
        x = x & (x - 1);    // Сброс одного бита
        y = y & (y - 1);    // в каждом слове
    }
}
```

После сброса общих единичных битов в каждом 32-битовом слове максимально возможное количество единичных битов в обоих словах равно 32. Таким образом, слово с наименьшим количеством единичных битов содержит их не более 16, т.е. цикл в листинге 5.5 выполняется не более 16 раз, что в наихудшем случае дает нам 119 команд базового набора RISC ($16 \cdot 7 + 7$). Вычислительный эксперимент с использованием равномерно распределенных случайных 32-битовых чисел показал, что средняя степень заполнения слова с меньшим количеством единичных битов после сброса общих единичных битов составляет около 6.186. Это дает нам среднее время выполнения, равное около 50 команд для случайных 32-битовых входных данных, что не столь хорошо, как в случае применения кода из листинга 5.4. Чтобы описанная процедура была эффективнее приведенной в листинге 5.4, количество единичных битов в x или y после сброса общих единичных битов должно быть не более трех.

Подсчет единичных битов в массиве

Простейший способ подсчета единичных битов в массиве (векторе) полных слов при отсутствии команды $\text{pop}(x)$ состоит в подсчете количества единичных битов в каждом слове массива, например, с использованием процедуры из листинга 5.1 на с. 104, после

чего полученные значения просто складываются. Назовем такой метод “наивным”. Без учета управляющих команд цикла и загрузки массива этот метод требует выполнения 16 команд на слово: 15 — для выполнения кода из листинга 5.1 и одной — для сложения. Предполагается, что код процедуры встраивается в код программы (без осуществления вызова), маски загружаются вне циклов, а машина имеет достаточное количество регистров для хранения всех используемых при вычислениях величин.

Другой путь, который может оказаться более быстрым, состоит в использовании первых двух строк этой процедуры для групп из трех слов с последующим сложением трех промежуточных значений. Поскольку максимальное значение каждого промежуточного результата в каждом 4-разрядном поле равно 4, после сложения этих трех промежуточных результатов в каждом 4-разрядном поле будет максимум число 12, так что переполнения при сложении возникнуть не может. Эта же идея может быть применена к 8- и 16-битовым полям. Кодирование и компиляция этого метода указывает, что он дает около 20% выигрыша по сравнению с наивным методом в смысле количества команд из базового набора RISC. В основном экономия “съедается” дополнительными накладными расходами. Но мы не будем рассматривать этот метод подробно, поскольку имеется *гораздо* лучший метод подсчета единичных битов в массиве.

Этот лучший способ изобретен около 1996 года Робертом Харли (Robert Harley) и Дэвидом Силом (David Seal) [100]. Он основан на схеме, именуемой *сумматором с сохранением переноса* (carry-save adder — CSA). CSA представляет собой последовательность независимых полных сумматоров¹ [56] и часто применяется в схемах бинарных умножителей.

Используя булеву алгебру, логику каждого полного сумматора можно записать так.

$$h \leftarrow ab + ac + bc = ab + (a + b)c = ab + (a \oplus b)c,$$

$$l \leftarrow (a \oplus b) \oplus c$$

Здесь a , b и c представляют собой однобитовые входы, l является младшим выходным битом (суммой), а h — старшим битом (переносом). Замена $a + b$ в первой строке на $a \oplus b$ оправдана, поскольку, когда a и b равны 1, член ab делает значение всего выражения равным 1. Присваивая значение $a \oplus b$ временной переменной, можно вычислить логику полного сумматора за пять логических команд, работающих параллельно с 32 битами (на 32-разрядной машине). Будем обозначать эти пять команд как $CSA(h, l, a, b, c)$. Это “макрос”, в котором h и l являются выходными данными.

Один из способов использования операции CSA заключается в обработке элементов массива A группами по три, сводя каждую группу из трех слов к двум, а затем применяя операцию подсчета единичных битов к этим двум словам. Эти две степени заполнения суммируются в цикле. По окончании работы цикла общая степень заполнения массива равна удвоенной степени заполнения старших выходных битов CSA плюс накопленная степень заполнения младших выходных битов.

¹ Полный сумматор представляет собой схему с тремя однобитовыми входами и двумя однобитовыми выходами (сумма и перенос).

Пусть n_c — количество команд, требующихся для выполнения шага CSA, а n_p — количество команд, необходимое для подсчета единичных битов в одном слове. Для типичной RISC-машины $n_c = 5$ и $n_p = 15$. Игнорируя загрузку из массива и команды управления циклом (этот код может несколько изменяться при переходе от одной машины к другой), находим, что описанный выше цикл выполняет $(n_c + 2n_p + 2)/3 \approx 12.33$ команды на слово массива (“+2” требуется для двух сложений в цикле). Эта величина существенно отличается от 16 команд на слово в наивном методе.

Имеется другой способ использования операции CSA, который приводит к более эффективной и несколько более компактной программе, показанной в листинге 5.6. В ней на одно слово цикла требуется $(n_c + n_p + 1)/2 = 10.5$ команды (игнорируя управление циклом и загрузки). Здесь операция CSA раскрывается в следующий код.

```
u = ones ^ A[i];
v = A[i+1];
twos = (ones & A[i]) | (u & v);
ones = u ^ v;
```

ЛИСТИНГ 5.6. Количество единичных битов в массиве при обработке элементов группами по два

```
#define CSA(h,l, a,b,c) \
(unsigned u = a ^ b; unsigned v = c; \
 h = (a & b) | (u & v); l = u ^ v;)

int popArray(unsigned A[], int n)
{
    int tot, i;
    unsigned ones, twos;

    tot = 0; // Инициализация
    ones = 0;
    for (i = 0; i <= n - 2; i = i + 2)
    {
        CSA(twos, ones, ones, A[i], A[i+1])
        tot = tot + pop(twos);
    }
    tot = 2*tot + pop(ones);

    // При наличии необработанного последнего
    // элемента обрабатываем его отдельно
    if (n & 1)
        tot = tot + pop(A[i]);
    return tot;
}
```

Имеются возможности применения операции CSA для дальнейшего уменьшения количества команд, необходимых для вычисления количества единичных битов в массиве. Легче всего понять эти способы, если воспользоваться принципиальными схемами. Так, на рис. 5.2 показан способ кодирования цикла, который берет по восемь элементов массива одновременно и сжимает их в четыре величины, помеченные как *eights*, *fours*, *twos* и *ones* (восьмерки, четверки, двойки и единицы). Величины *fours*, *twos* и *ones* возвращаются

в сумматоры CSA на следующей итерации цикла, а единичные биты в *eights* подсчитываются функцией вычисления степени заполнения на уровне слова. После того как будет обработан весь массив, общая степень заполнения вычисляется следующим образом.

$$8 \text{ pop}(\text{eights}) + 4 \text{ pop}(\text{fours}) + 2 \text{ pop}(\text{twos}) + \text{pop}(\text{ones})$$

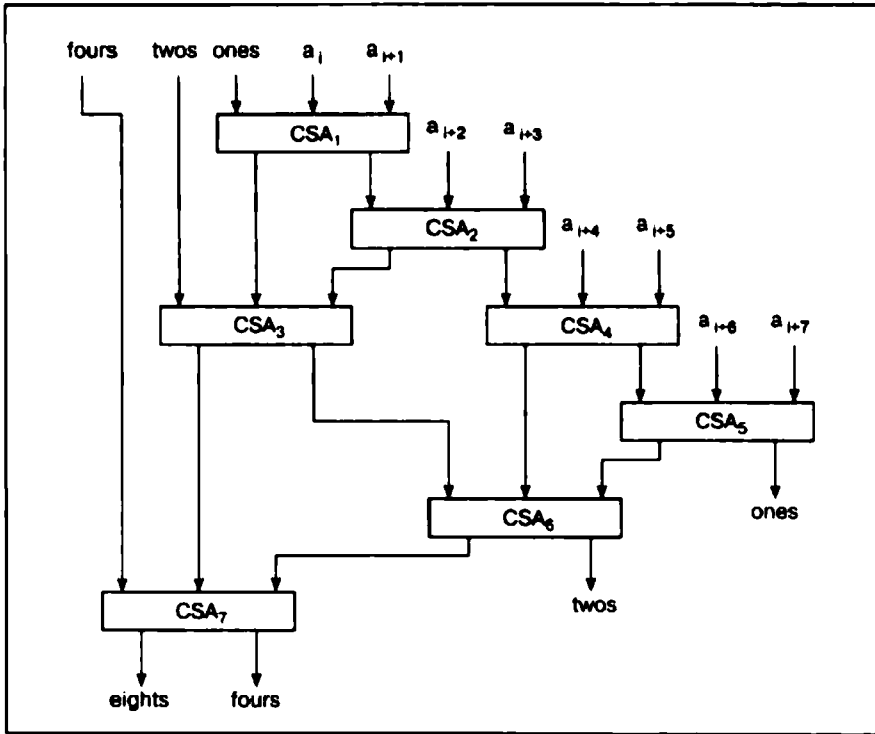


Рис. 5.2. Схема для подсчета количества единичных битов

Соответствующий код, показанный в листинге 5.7, использует макрос CSA, определенный в листинге 5.6. Нумерация блоков на рис. 5.2 соответствует порядку вызовов макросов CSA в листинге 5.7. Время выполнения цикла, без учета загрузки массива и управления циклом, составляет $(7n_c + n_p + 1)/8 = 6.375$ команды на слово массива.

Листинг 5.7. Количество единичных битов в массиве при обработке элементов группами по восемь

```
int popArray(unsigned A[], int n)
{
    int tot, i;
    unsigned ones, twos, twosA, twosB,
            fours, foursA, foursB, eights;

    tot = 0; // Инициализация
    fours = twos = ones = 0;

    for (i = 0; i <= n - 8; i = i + 8)
```

```

(
    CSA(twosA, ones, ones, A[i], A[i+1])
    CSA(twosB, ones, ones, A[i+2], A[i+3])
    CSA(foursA, twos, twos, twosA, twosB)
    CSA(twosA, ones, ones, A[i+4], A[i+5])
    CSA(twosB, ones, ones, A[i+6], A[i+7])
    CSA(foursB, twos, twos, twosA, twosB)
    CSA(eights, fours, fours, foursA, foursB)
    tot = tot + pop(eights);
)
tot = 8*tot + 4*pop(fours) + 2*pop(twos) + pop(ones);

// Просто добавляем последние 0-7 элементов
for (i = 1; i < n; i++)
    tot = tot + pop(A[i]);

return tot;
)

```

Сумматоры CSA могут быть соединены не только так, как показано на рис. 5.2, но и иными способами. Например, можно достичь более высокой степени параллельности при передаче первых трех элементов массива в один CSA, а следующих трех — во второй, что позволит этим двум сумматорам работать параллельно. Для повышения степени параллельности можно также переставить три входных операнда макроса CSA. План, показанный на рис. 5.2, позволяет легко увидеть, каким образом использовать для построения программы, обрабатывающей элементы массива группами по четыре, только три первых сумматора и как расширить его для построения программы, работающей с группами по 16 и более элементов массива. Показанный план несколько распределяет загрузки, что оказывается преимуществом при работе на машинах с относительно малым количеством одновременно выполняемых загрузок.

Схему на рис. 5.2 можно обобщить таким образом, чтобы выполнялось малое количество подсчетов единичных битов в отдельных словах. Набросаем примерный план построения такой программы. Для этого нам потребуется массив размером $m \times 2$ слов для хранения двух экземпляров каждой из переменных *ones*, *twos*, *fours* и т.д. В случае массива размером n достаточно выбрать $m \geq \lceil \log_2(n+1) \rceil + 1$ (для массива любого размера, который может храниться в адресуемой памяти 32-разрядной машины, достаточно выбора $m = 31$). Необходим также байтовый массив размером m для отслеживания того, сколько (0, 1 или 2) значений хранится в каждой строке массива размером $m \times 2$. Программа обрабатывает элементы массива группами по два. Для каждой группы вызывается сумматор CSA, который сжимает эти два элемента массива с сохраненным значением *ones*, которое наиболее удобно хранить в позиции $[0, 0]$ массива размером $m \times 2$. Во внутреннем цикле в массиве сохраняется получающееся в результате значение *twos* путем сканирования (обычно не очень длинного) для поиска строки с менее чем двумя элементами. Если строка *twos* заполнена, ее два значения комбинируются со значением *twos* с помощью сумматора CSA. Выходное значение *twos* помещается в массив, при этом устанавливая количество его строк равным 1. Сканирование продолжается, теперь для поиска места для размещения выходного значения *fours* и т.д.

По завершении прохода по входному массиву программа выполняет проход по (гораздо более короткому) массиву размером $m \times 2$, сжимая все заполненные строки, так что теперь в каждой строке содержится одно значащее значение. Наконец программа вызывает операцию подсчета количества единичных битов на уровне слова для первого элемента каждой строки до тех пор, пока не встретится строка с нулевым количеством. Общее количество единичных битов в массиве подсчитывается так.

$$\text{pop}(\text{строка } 0) + 2 \text{pop}(\text{строка } 1) + 4 \text{pop}(\text{строка } 2) + \dots$$

Предложенное выше значение m гарантирует, что последняя строка будет иметь нулевое количество единичных битов, так что это условие можно использовать для прекращения сканирования.

Получающаяся в результате программа выполняет ровно $\lceil \log_2(n+3) \rceil$ операций подсчета единичных битов в слове. К сожалению, этот метод нельзя назвать практичным из-за того, что накладные расходы на загрузку и сохранение промежуточных результатов превышают полученную в этом методе экономию вычислительных команд. Экспериментальная программа (при разработке которой не прилагались особые усилия по ее оптимизации) выполняла около 29 команд на слово массива (с учетом всех команд цикла). Это значительно хуже, чем при использовании наивного метода.

В табл. 5.1 подытожено число команд при выполнении описанной схемы для разных размеров групп. Значения в двух средних столбцах игнорируют команды загрузки и управления циклом. Четвертый столбец дает общее количество команд цикла на слово входного массива, сгенерированных компилятором для машины RISC с базовым набором команд (не имеющей команды индексированной загрузки).

ТАБЛИЦА 5.1. Количество команд на одно слово при вычислении количества единичных битов массива

Программа	Количество команд (без учета загрузок и управления циклом)		Все команды цикла (выход компилятора)
	Формула	Для $n_c = 5$, $n_p = 15$	
Наивный метод	$n_p + 1$	16	21
Группы по 2	$(n_c + n_p + 1)/2$	10.5	14
Группы по 4		7.75	10
Группы по 8	$(7n_c + n_p + 1)/8$	6.38	8
Группы по 16	$(15n_c + n_p + 1)/16$	5.69	7
Группы по 32	$(31n_c + n_p + 1)/32$	5.34	6.5
Группы по 2^n	$n_c + \frac{n_p - n_c + 1}{2^n}$	$5 + \frac{11}{2^n}$	—

Для малых массивов имеются схемы, которые лучше изображенной на рис. 5.2. Например, для массива из семи слов достаточно эффективная схема изображена на рис. 5.3 [100]. В ней выполняется $4n_c + 3n_p + 4 = 69$ команд, или 9.86 команды на слово. Имеется аналогичная схема, применимая к массивам размером $2^k - 1$ слов для любого положительного целого k . Схема для 15 слов выполняет $11n_c + 4n_p + 6 = 121$ команду, или 8.07 команды на слово.

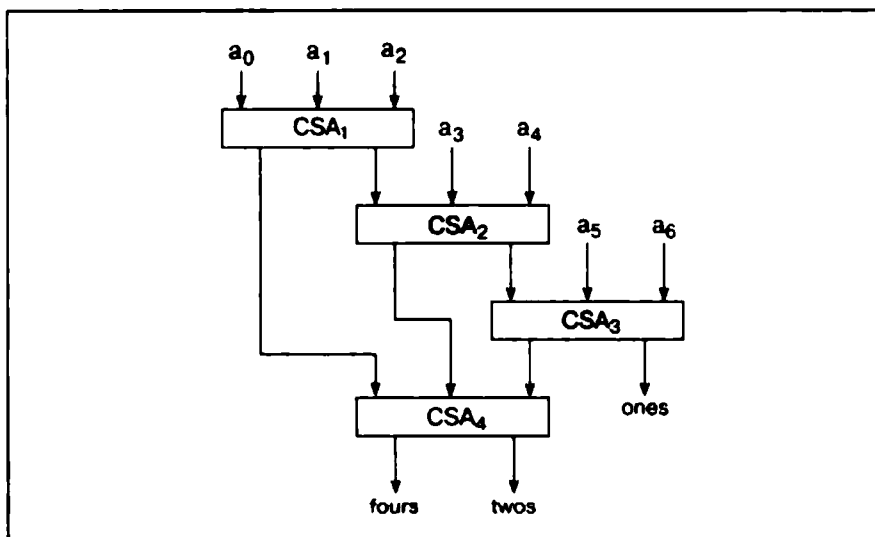


Рис. 5.3. Схема для подсчета единичных битов в семи словах

Применение

Функция $\text{pop}(x)$ используется, например, при вычислении “расстояния Хэмминга” (Hamming distance) между двумя битовыми векторами. Расстояние Хэмминга (концепция из теории кодов с исправлением ошибок) представляет собой количество разрядов, в которых эти векторы отличаются, т.е.

$$\text{dist}(x, y) = \text{pop}(x \oplus y)$$

(Например, обратитесь к главе о кодах с исправлением ошибок в [24].)

Другое применение — обеспечение быстрого индексного доступа к элементам разреженного массива A , представленного некоторым компактным образом. В компактном представлении сохраняются только определенные, ненулевые элементы массива. Пусть есть вспомогательный массив битовых строк $bits$ (составленный из 32-битовых слов), имеющий единичный бит для каждого индекса i , для которого определен элемент $A[i]$. Для ускорения доступа имеется также массив слов $bitsum$, такой, что $bitsum[j]$ представляет собой общее количество во всех словах массива $bits$ единичных битов, предшествующих j -му. Проиллюстрируем сказанное на массиве, в котором определены 0-, 2-, 32-, 47-, 48- и 95-й элементы.

<i>bits</i>	<i>bitsum</i>	Данные
0x00000005	0	A[0]
0x00018001	2	A[2]
0x80000000	5	A[32]
		A[47]
		A[48]
		A[95]

Для заданного индекса i , $0 \leq i \leq 95$, соответствующий индекс *sparse_i* в массиве данных задается количеством единичных битов элементов массива *bits*, предшествующих биту, соответствующему i . Его можно вычислить следующим образом.

```

j = i >> 5;           // j = i/32
k = i & 31;          // k = rem(i, 32)
mask = 1 << k;       // "1" в позиции k
if ((bits[j] & mask) == 0)
    goto no_such_element;
mask = mask - 1;     // единицы справа от k
sparse_i =
    bitsum[j] + pop(bits[j] & mask);

```

Затраты памяти при таком представлении разреженного массива — два бита на элемент полного массива.

Функция $\text{pop}(x)$ может также использоваться для генерации биномиально распределенных случайных целых чисел. Для генерации целого числа из совокупности $\text{BINOMIAL}(t, p)$, где t — число попыток, а $p = 1/2$, генерируем t случайных битов и подсчитываем количество единиц среди них. Этот подход можно обобщить для вероятности p , отличной от $1/2$ (см., например, [67, разд. 3.4.1, упр. 27]).

Еще одно применение рассматриваемой функции — при вычислении количества завершающих слово нулей (см. раздел “Подсчет завершающих нулевых битов” на с. 130).

Как утверждает программистский фольклор, эта функция имеет важное значение для секретных служб. Никто (исключая сотрудников этих самых служб) не знает, для чего ее там применяют, но, возможно, это применение как-то связано с криптографией или выполнением поиска в материале огромного объема.

5.2. Четность

Под четностью (*parity*) битовой строки понимается, какое количество единичных битов — четное или нечетное — содержит эта строка. Строка обладает “нечетной четностью”, если она содержит нечетное количество единичных битов; в противном случае строка имеет “четную четность”².

² В данном разделе значения терминов “четное” и “нечетное” отличаются от обычных. Везде в разделе, где не оговорено иное, “четность” означает наличие четного числа единичных битов в рассматриваемом значении, а “нечетность” — нечетное количество битов. — *Примеч. ред.*

Вычисление четности слова

Результат проверки на четность равен 1, если некоторое слово x имеет нечетную четность, и равен 0, если слово имеет четную четность. Этот результат представляет собой сумму всех битов слова по модулю 2, т.е. над всеми битами x выполняется операция *исключающего или*.

В качестве одним из способов проверки слова на четность можно использовать функцию $\text{pop}(x)$ — в таком случае четность представляет собой младший бит значения, возвращаемого этой функцией. Этот метод вполне применим при наличии команды для вычисления $\text{pop}(x)$. Если же такой команды нет, то вместо использования кода для вычисления $\text{pop}(x)$ лучше воспользоваться более эффективными специализированными методами.

Один из таких методов состоит в вычислении

$$y \leftarrow \bigoplus_{i=0}^{n-1} (x \gg i)$$

Здесь n — длина слова, а результат проверки на четность слова x помещается в младший бит y . (Оператор \oplus означает *исключающее или*, но в данном случае его можно заменить обычным сложением.)

Ниже приводится более быстрый способ проверки на четность для не слишком больших значений n (в примере $n = 32$; сдвиги могут быть как знаковыми, так и беззнаковыми).

$$\begin{aligned} y &= x \wedge (x \gg 1); \\ y &= y \wedge (y \gg 2); \\ y &= y \wedge (y \gg 4); \\ y &= y \wedge (y \gg 8); \\ y &= y \wedge (y \gg 16); \end{aligned} \tag{3}$$

Всего выполняется 10 команд, что значительно меньше, чем в первом методе (62 команды), даже при полном развертывании цикла. Здесь бит четности также оказывается в младшем разряде y . В действительности всегда при выполнении беззнаковых сдвигов i -й бит результата y дает четность битовой строки, состоящей из всех битов x от i -го до старшего. Кроме того, так как *исключающее или* представляет собой свою собственную инверсию, $y, \oplus y$, представляет собой четность подслова, состоящего из битов x с $i-1$ по j ($i \geq j$).

Приведенный выше метод является примером метода "параллельного префикса", или "сканирования", который применяется в параллельных вычислениях [53, 74]. При наличии достаточного количества процессоров он позволяет преобразовать некоторые кажущиеся последовательными процессы в параллельные, сокращая время вычислений от $O(n)$ до $O(\log n)$. Например, если необходимо выполнить операцию поразрядного *исключающего или* над целым массивом, вы можете сначала выполнить действия (3) над всем массивом, а затем продолжить со сдвигом на 32, 64 бита и так далее, выполняя операцию *исключающего или* над словами массива. При таком подходе выполняется большее количество элементарных (с отдельными словами) команд *исключающего или*, чем

в простом алгоритме, работающем слева направо, а следовательно, для однопроцессорного компьютера применять данный метод не стоит. Но на компьютере с достаточным количеством параллельных процессоров проверка на четность по описанному выше алгоритму занимает время $O(\log n)$, а не $O(n)$ (где n — количество слов в массиве).

Непосредственным приложением метода (3) является преобразование целых чисел в код Грея (см. с. 339).

Если изменить код (3), заменив все сдвиги вправо сдвигами влево, то результат проверки на четность целого слова x отображается в старшем бите y , а i -й бит слова y содержит четность битов слова x , расположенных в i -й позиции и *правее* нее. Это операция “параллельного суффикса”, так как каждый бит является функцией от самого себя и следующих за ним.

При использовании *циклических сдвигов* в результирующем слове во всех разрядах будут либо единицы, если x имеет нечетную четность, либо нули, если его четность четная.

Пять присваиваний в (3) можно выполнить в произвольном порядке (исходная переменная x в любом случае используется в первой строке). Если выполнить их в обратном порядке и если вас интересует только получение четности в младшем бите y , то две последние строки

```
y = y ^ (y >> 2);
y = y ^ (y >> 1);
```

можно заменить строкой [55]

```
y = 0x6996 >> (y & 0xF)
```

Это операция “табличного поиска в регистре”. При использовании базового набора команд RISC она позволяет сэкономить одну команду (или две, если не учитывать загрузку константы). Младший бит y содержит четность исходного слова, но в остальных битах y нет никакой полезной информации.

В следующем методе четность слова x вычисляется за девять команд, при этом результат оказывается равным 0 или 1 (все сдвиги в коде беззнаковые).

```
x = x ^ (x >> 1);
x = (x ^ (x >> 2)) & 0x11111111;
x = x*0x11111111;
p = (x >>28) & 1;
```

После выполнения второй строки каждая шестнадцатеричная цифра в x примет значение 0 или 1, в зависимости от четности исходной шестнадцатеричной цифры. Команда *умножения* суммирует все цифры, помещая сумму в старший шестнадцатеричный разряд. Переносов при этом возникнуть не может, поскольку максимальная сумма любого шестнадцатеричного разряда равна 8.

Команды *умножения* и *сдвига* можно заменить командой, которая вычисляет остаток от деления x на 15, что дает (медленную) реализацию алгоритма за восемь команд, если, конечно, в компьютере есть команда получения *остатка от деления на непосредственно заданное значение*.

На 64-разрядной машине код с применением умножения дает корректный результат после внесения очевидных изменений (расширение шестнадцатеричных констант до 16 полубайтов, каждый со значением 1, и замена последнего сдвига с 28 до 60). В этом случае максимальная сумма в любом шестнадцатеричном разряде частичных произведений, кроме старшего, равна 15, так что вновь не возникает никаких переполнений, влияющих на старший шестнадцатеричный разряд. С другой стороны, вариант, вычисляющий остаток при делении на 15, на 64-разрядной машине работать *не* будет, поскольку остаток является суммой полубайтов по модулю 15, а эта сумма может достигать значения 16.

Добавление бита четности к 7-битовой величине

В [44, item 167] есть интересный алгоритм, добавляющий бит четности к 7-битовой величине (выровненной вправо и помещенной в регистр). Под этим подразумевается установка бита четности слева от семи битов таким образом, чтобы получить восьмибитовую величину с четной четностью. В указанной работе приведен пример кода для 36-разрядной машины, но он корректно работает и на 32-разрядных машинах.

$$\text{modu}((x * 0x10204081) \& 0x888888FF, 1920)$$

Здесь $\text{modu}(a, b)$ означает остаток от беззнакового деления a на b , при этом аргументы и результат интерпретируются как беззнаковые целые числа, “*” означает умножение по модулю 2^{32} , а константа 1920 равна $15 \cdot 2^7$. Фактически здесь вычисляется сумма битов в слове x с размещением этой суммы слева от семи битов x . Например, это выражение преобразует число $0x0000007F$ в $0x000003FF$, а $0x00000055$ в $0x00000255$.

В [44] есть еще одна остроумная формула для добавления битов четности к 7-битовому числу (при этом образуется 8-битовое число с нечетным количеством единичных битов).

$$\text{modu}((x * 0x00204081) | 0x3DB6DB00, 1152)$$

Здесь $1152 = 9 \cdot 2^7$. Понять, как работает эта формула, поможет тот факт, что степень 8 по модулю 9 равна ± 1 . Если заменить значение $0x3DB6DB00$ значением $0xBDB6DB00$, эта формула будет давать восьмибитовое число с четным количеством единичных битов.

На современных машинах все эти методы недостаточно практичны, так как при том, что память дешевет, деление все еще выполняется довольно медленно. Поэтому большинство программистов предпочитают использовать вместо вычислений простой поиск в таблице.

Применение

Операция четности широко используется для добавления проверочных битов к данным. Она также полезна при перемножении битовых матриц в поле $GF(2)$ (где операция сложения представляет собой *исключающее или*).

5.3. Подсчет ведущих нулевых битов

Существует ряд простых способов подсчета ведущих нулевых битов слова, реализуемых с помощью метода бинарного поиска. Ниже приведена модель, имеющая ряд вариаций. На процессоре с базовым набором RISC-команд вычисления выполняются за 20–29 команд. Все сравнения — “логические” (беззнаковые).

```
if (x == 0) return(32);
n = 0;
if (x <= 0x0000FFFF) {n = n + 16; x = x << 16;}
if (x <= 0x00FFFFFF) {n = n + 8; x = x << 8;}
if (x <= 0x0FFFFFFF) {n = n + 4; x = x << 4;}
if (x <= 0x3FFFFFFF) {n = n + 2; x = x << 2;}
if (x <= 0x7FFFFFFF) {n = n + 1;}
return n;
```

В другом варианте все сравнения заменяются командами *и*.

```
if ((x & 0xFFFF0000) == 0) {n = n + 16; x = x << 16;}
if ((x & 0xFF000000) == 0) {n = n + 8; x = x << 8;}
...
```

Еще один вариант, позволяющий избежать больших непосредственно задаваемых значений, использует команды *сдвига вправо*.

Последний из операторов *if* просто прибавляет единицу к *n*, если старший бит *x* равен 0, так что в варианте без использования условного перехода можно прибегнуть к инструкции

```
n = n + 1 - (x >> 31);
```

Здесь “+1” можно опустить, если инициализировать переменную *n* значением не 0, а 1. Это наблюдение приводит к алгоритму, для работы которого потребуется выполнение 12–20 базовых RISC-команд (листинг 5.8). Этот код можно улучшить еще немного, если *x* начинается с единичного бита: замените первую строку следующей.

```
if ((int)x <= 0) return (~x >> 26) & 32;
```

Листинг 5.8. Подсчет ведущих нулевых битов бинарным поиском

```
int nlz(unsigned x)
{
    int n;

    if (x == 0) return(32);
    n = 1;
    if ((x >> 16) == 0) {n = n + 16; x = x << 16;}
    if ((x >> 24) == 0) {n = n + 8; x = x << 8;}
    if ((x >> 28) == 0) {n = n + 4; x = x << 4;}
    if ((x >> 30) == 0) {n = n + 2; x = x << 2;}
    n = n - (x >> 31);
    return n;
}
```

В листинге 5.9 показан вариант подсчета одним из методов, работающих в направлении, противоположном приведенному выше. Он требует выполнения тем меньшего количества операций, чем больше ведущих нулей имеется в слове, и позволяет избежать использования больших непосредственно задаваемых чисел и больших сдвигов. Всего выполняется 12–20 базовых RISC-команд.

**Листинг 5.9. Подсчет ведущих нулевых битов бинарным поиском
в обратном направлении**

```
int nlz(unsigned x)
{
    unsigned y;
    int n;

    n = 32;
    y = x >> 16; if (y != 0) {n = n - 16; x = y;}
    y = x >> 8;  if (y != 0) {n = n - 8;  x = y;}
    y = x >> 4;  if (y != 0) {n = n - 4;  x = y;}
    y = x >> 2;  if (y != 0) {n = n - 2;  x = y;}
    y = x >> 1;  if (y != 0) return n - 2;
    return n - x;
}
```

Этот алгоритм позволяет использовать метод поиска в таблице, при этом последние четыре выполняемые строки кода заменяются следующими строками.

```
static char table[256] = {0,1,2,2,3,3,3,3,4,4,..., 8};
return n - table[x];
```

Использовать вспомогательный поиск в таблице позволяют многие алгоритмы, хотя упоминается об этом нечасто.

Для компактности два последних алгоритма можно кодировать с использованием циклов. Например, алгоритм из листинга 5.10 является вариантом алгоритма из листинга 5.9 с использованием цикла. Выполняется этот вариант алгоритма за 23–33 базовые RISC-команды, 10 из которых представляют собой команды условного ветвления.

**Листинг 5.10. Подсчет ведущих нулевых битов бинарным поиском
с использованием цикла**

```
int nlz(unsigned x)
{
    unsigned y;
    int n, c;

    n = 32;
    c = 16;
    do {
        y = x >> c;
        if (y != 0) {n = n - c; x = y;}
        c = c >> 1;
    } while (c != 0);
    return n - x;
}
```



```

y = x - 0x100;      // Если разряды 8-15 нулевые,
m = (y >> 16) & 8; // к n добавляется 8 и выполняется
n = n + m;         // сдвиг x влево на 8
x = x << m;

y = x - 0x1000;    // Если разряды 12-15 нулевые,
m = (y >> 16) & 4; // увеличиваем n на 4 и
n = n + m;         // сдвигаем x влево на 4
x = x << m;

y = x - 0x4000;    // Если разряды 14-15 нулевые,
m = (y >> 16) & 2; // увеличиваем n на 2 и
n = n + m;         // сдвигаем x влево на 2
x = x << m;

y = x >> 14;       // Устанавливаем y = 0, 1, 2 или 3
m = y & ~(y >> 1); // m = 0, 1, 2 или 2 соответственно
return n + 2 - m;
}

```

Если ваш компьютер оснащен командой для вычисления функции $\text{pop}(x)$, то существует эффективный способ подсчета ведущих нулевых битов, показанный в листинге 5.13. Пять присваиваний значений переменной x в действительности можно выполнять в любом порядке. Этот алгоритм выполняется за 11 команд и не использует команд ветвления. Он полезен даже в том случае, когда специальной команды для вычисления функции $\text{pop}(x)$ в компьютере нет. В этой ситуации можно использовать программу из листинга 5.1 на с. 104, которая выполняется за 21 команду, так что поиск количества ведущих нулевых битов выполняется с помощью 32 базовых RISC-команд, причем среди них нет ни одной команды условного перехода.

ЛИСТИНГ 5.13. Подсчет ведущих нулевых битов с использованием функции $\text{pop}(x)$

```

int nlz(unsigned x)
{
    int pop(unsigned x);

    x = x | (x >> 1);
    x = x | (x >> 2);
    x = x | (x >> 4);
    x = x | (x >> 8);
    x = x | (x >> 16);
    return pop(~x);
}

```

Роберт Харли (Robert Harley) [46] разработал алгоритм вычисления $\text{nlz}(x)$, очень похожий на алгоритм Сила (Seal) для вычисления $\text{ntz}(x)$ (см. листинг 5.22 на с. 134). Метод Харли распространяет старший единичный бит вправо с применением сдвигов и команд *или* и выполняет умножение по модулю 2^{32} на специальную константу, в результате чего получается произведение, старшие шесть битов которого однозначно идентифицируют количество ведущих нулей в x . Затем выполняются сдвиг вправо и поиск в таблице (индексированная загрузка) для преобразования шестибитового значения в фактическое число ведущих нулей. Как показано в листинге 5.14, метод состоит из 14 команд, вклю-

чая *умножение*, и индексированной загрузки. Значения таблицы, показанные как *u*, не используются.

Листинг 5.14. Алгоритм Харли вычисления количества ведущих нулевых битов

```
int nlz(unsigned x)
{
    static char table[64] =
    {
        32, 31, u, 16, u, 30, 3, u,
        15, u, u, u, 29, 10, 2, u,
        u, u, 12, 14, 21, u, 19, u,
        u, 28, u, 25, u, 9, 1, u,
        17, u, 4, u, u, u, 11, u,
        13, 22, 20, u, 26, u, u, 18,
        5, u, u, 23, u, 27, u, 6,
        u, 24, 7, u, 8, u, 0, u);

    x = x | (x >> 1); // Распространение старшего
    x = x | (x >> 2); // единичного бита вправо
    x = x | (x >> 4);
    x = x | (x >> 8);
    x = x | (x >> 16);
    x = x*0x06EB14F9; // Множитель равен 7*255**3
    return table[x >> 26];
}
```

Множитель равен $7 \cdot 255^3$, так что умножение можно выполнить так, как показано ниже. В этом случае функция выполняется 19 элементарными командами, плюс индексированная загрузка.

```
x = (x << 3) - x; // Умножение на 7
x = (x << 8) - x; // Умножение на 255
x = (x << 8) - x; // ...и еще раз...
x = (x << 8) - x; // ...и еще
```

Имеется много множителей, обладающих требуемым свойством однозначности, сомножители которых имеют вид $2^k \pm 1$. Наименьшим таким множителем является $0x045BCED1 = 17 \cdot 65 \cdot 129 \cdot 513$. Множителя, состоящего из трех сомножителей, для таблиц размером 64 или 128 элементов не имеется, но они есть для таблицы из 256 элементов, и таких множителей много. Наименьший из них равен $0x01033CBF = 65 \cdot 255 \cdot 1025$ (его применение экономит две команды ценой большего размера таблицы).

Юлиус Горявски (Julius Goryavsky) [42] нашел ряд вариаций алгоритма Харли, в которых размер таблицы снижен ценой нескольких команд либо повышена степень параллельности вычислений, либо эти вариации обладают еще какими-то желательными свойствами. В листинге 5.15 приведен очевидный победитель, если выполнять умножение с помощью сдвигов и суммирования. В этом коде изменены только таблица и строки, содержащие *сдвиг вправо* на 16 и последующее *умножение* из листинга 5.14. Если машина оснащена командой *и-не*, то это позволяет сэкономить две команды, поскольку множитель можно разложить как $511 \cdot 2047 \cdot 16383 \pmod{2^{32}}$, так что вычисления можно

выполнить за шесть элементарных команд вместо восьми. Если же машина не оснащена командой *и-не*, то экономия составляет одну команду.

ЛИСТИНГ 5.15. Вариация Горявски алгоритма Харли вычисления количества ведущих нулевых битов

```
...
static char table[64] =
{
    32,20,19, u, u,18, u, 7, 10,17, u, u,14, u, 6, u,
    u, 9, u,16, u, u, 1,26, u,13, u, u,24, 5, u, u,
    u,21, u, 8,11, u,15, u, u, u, u, 2,27, 0,25, u,
    22, u,12, u, u, 3,28, u, 23, u, 4,29, u, u,30,31
};
...
x = x & ~(x >> 16);
x = x*0xFD7049FF;
...
```

Методы с плавающей точкой

При подсчете ведущих нулевых битов слова можно использовать постнормализацию чисел с плавающей точкой. Эта методика достаточно хорошо работает с числами с плавающей точкой в IEEE-формате. Идея заключается в преобразовании данного беззнакового целого числа в число с плавающей точкой двойной точности, выделении из него степенной части и вычитании ее из константы. Соответствующий код приведен в листинге 5.16.

ЛИСТИНГ 5.16. Подсчет ведущих нулевых битов с использованием чисел с плавающей точкой в IEEE-формате

```
int nlz(unsigned k) {
    union {
        unsigned asInt[2];
        double asDouble;
    };
    int n;

    asDouble = (double)k + 0.5;
    n = 1054 - (asInt[LE] >> 20);
    return n;
}
```

Этот код использует анонимное (безымянное) объединение C++ для перекрытия целого числа и величины с плавающей точкой двойной точности. Переменная LE должна быть равна 1, если процедура выполняется на машинах, на которых адрес младшего байта меньше адреса старшего, и 0 в противном случае. Слагаемое 0.5 (или другое малое число) необходимо для корректной работы программы при $k = 0$.

Оценить время выполнения этой программы практически невозможно, так как на разных машинах числа с плавающей точкой обрабатываются по-разному. Например, на многих машинах, кроме регистров для хранения целых чисел, имеются специальные регистры для хранения чисел с плавающей точкой. На таких машинах, в частности, может

потребуется пересылка данных в память для преобразования целого числа в число с плавающей точкой и наоборот.

Код в листинге 5.16 не соответствует ANSI-стандарту языков C и C++, поскольку обращается к одной и той же области памяти как к двум различным типам данных. Следовательно, нет никаких гарантий, что этот код будет правильно работать на всех машинах или со всеми компиляторами. Тем не менее он работает с компилятором XLC IBM на платформе AIX и с компилятором GCC на платформах AIX и Windows 2000 и XP при всех уровнях оптимизации (по крайней мере, на момент написания книги). Если изменить код так, как показано ниже, то в режиме оптимизации в указанных системах код работать не станет.

```
xx = (double)k + 0.5;
n = 1054 - (*(unsigned *)&xx + LE) >> 20;
```

Кстати, этот код нарушает еще один стандарт ANSI, а именно то, что арифметические действия могут выполняться только с указателями на элементы массива [17]. Проблема, тем не менее, связана с первым нарушением стандарта.

Несмотря на то что данный код не совсем корректен³, далее приведено несколько его вариантов.

```
asDouble = (double)k;
n = 1054 - (asInt[LE] >> 20);
n = (n & 31) + (n >> 9);
```

```
k = k & ~(k >> 1);
asFloat = (float)k + 0.5f;
n = 158 - (asInt >> 23);
```

```
k = k & ~(k >> 1);
asFloat = (float)k;
n = 158 - (asInt >> 23);
n = (n & 31) + (n >> 6);
```

В первом варианте проблема с $k = 0$ решена не путем добавления слагаемого 0.5, а с помощью дополнительных арифметических действий над результатом n (который без применения коррекции будет равен 1054 (0x41E)).

В следующих двух вариантах этого кода используются числа с плавающей точкой одинарной точности, которые преобразуются обычным методом с использованием анонимного объединения. Правда, здесь возникает новая проблема: при округлении результата возможны неточности, как при округлении к ближайшему числу (наиболее распространенный метод округления), так и при округлении в большую сторону. В режиме округления к ближайшему числу проблемы с округлением k возникают в диапазонах (в шестнадцатеричной записи) от FFFFFFF80 до FFFFFFFF, от 7FFFFFFC0 до 7FFFFFFF, от 3FFFFFFE0 до 3FFFFFFF и т.д. При округлении этих чисел добавление единицы при-

³ Некорректен в силу способа использования языка программирования C. Эти методы совершенно корректно будут работать при кодировании на машинном языке либо будучи сгенерированными компилятором для конкретного типа компьютера.

водит к появлению переноса в левых разрядах, что влечет за собой изменение позиции старшего единичного бита. Показанные выше способы коррекции сбрасывают бит справа от старшего единичного бита, позволяя тем самым избежать переноса. Если k является 64-битовой величиной, такая коррекция необходима и для кода в листинге 5.16, и в первом из трех приведенных выше вариантов.

Компилятор GNU C/C++ дает уникальную возможность кодировать любой из этих методов в виде макроса, позволяя вместо вызова функции использовать встроенный код [103]. Компилятор разрешает использовать в макросе различные языковые конструкции, включая объявления, а последнее выражение среди операторов макроса дает возвращаемое макросом значение. Ниже приводится пример макроопределения для алгоритма с применением чисел с одинарной точностью. (В языке C в названиях макросов обычно используются прописные буквы.)

```
#define NLZ(kp) \
    ({union {unsigned _asInt; float _asFloat;}; \
     unsigned _k = (kp), _kk = _k & ~(_k >> 1); \
     _asFloat = (float)_kk + 0.5f; \
     158 - (_asInt >> 23);})
```

Символы подчеркивания в именах нужны для того, чтобы избежать конфликтов имен с параметром kp (обычно определенные пользователем имена переменных не начинаются с символа подчеркивания).

Сравнение количества ведущих нулевых битов двух слов

Имеется простой способ определить, какое из слов x и y имеет большее количество ведущих нулевых битов [70] без реального вычисления $nlz(x)$ и $nlz(y)$. Эти методы приведены ниже. Три не приведенные здесь отношения, разумеется, получаются путем дополнения смысла сравнения в правой части.

$$\begin{array}{lll}
 nlz(x) = nlz(y) & \text{тогда и только тогда, когда} & (x \oplus y) \leq (x \& y) \\
 nlz(x) < nlz(y) & \text{тогда и только тогда, когда} & (x \& \neg y) > y \\
 nlz(x) \leq nlz(y) & \text{тогда и только тогда, когда} & (y \& \neg x) \leq x
 \end{array}$$

Связь с логарифмом

По существу, nlz представляет собой функцию “целого логарифма по основанию 2”. Для беззнакового $x \neq 0$ справедливы следующие соотношения (см. также раздел 11.4, “Целочисленный логарифм” на с. 316).

$$\begin{array}{l}
 \lfloor \log_2(x) \rfloor = 31 - nlz(x) \\
 \lceil \log_2(x) \rceil = 32 - nlz(x - 1)
 \end{array}$$

Другой тесно связанной с ней функцией является функция *bitsize* — количество битов, которое требуется для представления ее аргумента как знаковой величины в дополнительном к 2 коде. Эта функция определяется следующим образом.

$$\text{bitsize}(x) = \begin{cases} 1, & x = -1 \text{ или } 0 \\ 2, & x = -2 \text{ или } 1 \\ 3, & -4 \leq x \leq -3 \text{ или } 2 \leq x \leq 3 \\ 4, & -8 \leq x \leq -5 \text{ или } 4 \leq x \leq 7 \\ \dots & \dots \\ 32, & -2^{31} \leq x \leq -2^{30} + 1 \text{ или } 2^{30} \leq x \leq 2^{31} - 1 \end{cases}$$

Из определения ясно, что $\text{bitsize}(x) = \text{bitsize}(-x-1)$. Но, поскольку $-x-1 = -x$, можно записать следующий код для вычисления функции *bitsize* (сдвиг в коде — знаковый).

```
x = x ^ (x >> 31); // Если (x < 0), x = -x - 1;
return 32 - nlz(x);
```

Может использоваться и альтернативная функция — такая же, как и *bitsize(x)*, но отличающаяся тем, что при $x = 0$ она дает значение 0.

```
32 - nlz(x ^ (x << 1));
```

Применения

Функции вычисления ведущих нулевых битов слова обычно используются при моделировании арифметических действий над числами с плавающей точкой и в различных алгоритмах деления (см. листинги 9.1 на с. 211 и 9.3 на с. 222). Кроме того, команда *nlz(x)* имеет множество других применений.

Она может использоваться при вычислении отношения “ $x = y$ ” за три команды (см. раздел “Предикаты сравнения” на с. 43) и ряда элементарных функций (см. разделы 11.1–11.4).

Еще одно интересное применение этой функции — для генерации экспоненциально распределенных случайных целых чисел. Для этого генерируются случайные равномерно распределенные целые числа и к результату применяется функция *nlz* [37]. Вероятность возврата функцией значения 0 равна 1/2, вероятность значения 1 равна 1/4, 2 — 1/8 и т.д. Кроме того, функция применяется в качестве вспомогательной при поиске в слове подстроки определенной длины, состоящей из единичных (или нулевых) битов (процесс, используемый в некоторых алгоритмах распределения дискового пространства). В последних двух задачах используется также функция, вычисляющая количество завершающих нулевых битов.

5.4. Подсчет завершающих нулевых битов

Если доступна команда подсчета ведущих нулевых битов, то подсчет количества завершающих нулевых битов лучше всего свести к использованию этой команды.

$$32 - \text{nlz}(-x \& (x - 1))$$

При наличии команды вычисления количества единичных битов слова существует более эффективный метод, состоящий в формировании маски для выделения всех завершающих нулевых битов и подсчете единиц в этой маске [48], например так.

$$\text{pop}(-x \& (x - 1)) \text{ или}$$

$$32 - \text{pop}(x | -x)$$

Сформировать маску для выделения завершающих нулевых битов можно по формулам, которые приводятся в разделе 2.1, "Манипуляции младшими битами", на с. 31. Эти методы применимы даже в случае, когда среди команд компьютера команда подсчета количества единичных битов слова отсутствует. Например, если для подсчета единичных битов воспользоваться алгоритмом из листинга 5.1 на с. 104, то первое из приведенных выше выражений вычисляется за $3 + 21 = 24$ команды, среди которых нет команд ветвления.

В листинге 5.17 показан алгоритм, который выполняет описанные действия непосредственно, требуя от 12 до 20 базовых RISC-команд (для $x \neq 0$).

Листинг 5.17. Подсчет завершающих нулевых битов бинарным поиском

```
int ntz(unsigned x)
{
    int n;

    if (x == 0) return(32);
    n = 1;
    if ((x & 0x0000FFFF) == 0) {n = n + 16; x = x >> 16;}
    if ((x & 0x000000FF) == 0) {n = n + 8; x = x >> 8;}
    if ((x & 0x0000000F) == 0) {n = n + 4; x = x >> 4;}
    if ((x & 0x00000003) == 0) {n = n + 2; x = x >> 2;}
    return n - (x & 1);
}
```

Выражение $n + 16$ можно упростить до 17, если компилятор недостаточно интеллектуален, чтобы сделать это самостоятельно (это не влияет на подсчитанное нами количество команд).

Еще один вариант представлен в листинге 5.18. Здесь используются малые непосредственно задаваемые значения и выполняются более простые действия. При вычислениях требуется выполнить от 12 до 21 базовой RISC-команды. В отличие от первой процедуры, при меньшем количестве завершающих нулей выполняется большее количество команд (но, впрочем, и большее количество условных переходов, пропускающих вычисления).

Листинг 5.18. Подсчет завершающих нулевых битов,
малые непосредственные значения

```
int ntz(unsigned x)
{
    unsigned y;
```

```

int n;

if (x == 0) return 32;
n = 31;
y = x << 16; if (y != 0) {n = n - 16; x = y;}
y = x << 8; if (y != 0) {n = n - 8; x = y;}
y = x << 4; if (y != 0) {n = n - 4; x = y;}
y = x << 2; if (y != 0) {n = n - 2; x = y;}
y = x << 1; if (y != 0) {n = n - 1;}
return n;
}

```

Строку перед оператором `return` можно заменить строкой, которая экономит один условный переход (но не команду).

```
n = n - ((x << 1) >> 31);
```

В смысле количества выполняемых команд трудно превзойти поиск по дереву [7]. В листинге 5.19 представлена соответствующая процедура для 8-битового аргумента. В каждом пути этой процедуры выполняется семь команд (за исключением последних двух путей — `return 7` и `return 8`, которые требуют девяти команд). Для 32-битового аргумента требуется выполнить от 11 до 13 команд. К сожалению, для больших размеров слов программа быстро становится все большей и большей. Так, для 8-битового аргумента исходный код занимает 12 выполняемых строк (и компилируется в 41 команду); 32-битовый аргумент требует 48 выполняемых строк кода и около 164 команд, а 64-битовый — еще вдвое больше.

Листинг 5.19. Подсчет завершающих нулевых битов, поиск по бинарному дереву

```

int ntz(char x)
{
    if (x & 15) {
        if (x & 13) {
            if (x & 1) return 0;
            else return 1;
        }
        else if (x & 4) return 2;
        else return 3;
    }
    else if (x & 0x30) {
        if (x & 0x10) return 4;
        else return 5;
    }
    else if (x & 0x40) return 6;
    else if (x) return 7;
    else return 8;
}

```

Если ожидается малое или, напротив, большое число завершающих нулевых битов, то лучше использовать простые циклические алгоритмы (листинг 5.20). Код в левой части листинга выполняется за $5 + 3\text{ntz}(x)$ команд, а в правой — за $3 + 3(32 - \text{ntz}(x))$ базовых RISC-команд.

Листинг 5.20. Подсчет завершающих нулевых битов с использованием циклов

```

int ntz(unsigned x)
{
    int n;

    x = ~x & (x - 1);
    n = 0;           // n = 32;
    while(x != 0) { // while (x != 0) {
        n = n + 1;  //     n = n - 1;
        x = x >> 1; //     x = x + x;
    }              // }
    return n;      // return n;
}

```

Дин Годо (Dean Gaudet) [33] разработал алгоритм, интересный тем, что при использовании правильных команд он не содержит ветвлений, загрузок (не использует поиска в таблице) и обеспечивает высокую степень параллельности. Он показан в листинге 5.21.

Листинг 5.21. Алгоритм Годо подсчета количества завершающих нулевых битов

```

int ntz(unsigned x)
{
    unsigned y, bz, b4, b3, b2, b1, b0;
    y = x & -x; // Выделение крайнего справа единичного бита
    bz = y ? 0 : 1; // 1, если y = 0.
    b4 = (y & 0x0000FFFF) ? 0 : 16;
    b3 = (y & 0x00FF00FF) ? 0 : 8;
    b2 = (y & 0x0F0F0F0F) ? 0 : 4;
    b1 = (y & 0x33333333) ? 0 : 2;
    b0 = (y & 0x55555555) ? 0 : 1;
    return bz + b4 + b3 + b2 + b1 + b0;
}

```

Как видно, код использует "условное выражение" языка программирования C в шести местах. Эта конструкция имеет вид $a ? b : c$, а ее значение равно b , если a истинно (имеет ненулевое значение), и c , если a ложно (имеет значение 0). Хотя в общем случае условное выражение компилируется в сравнение с ветвлением, для простых случаев в листинге 5.21 ветвления можно избежать, если машина оснащена командой *сравнения на равенство нулю*, которая устанавливает целевой регистр равным 1, если операнд равен 0, и 0, если операнд ненулевой. Ветвления также можно избежать с помощью команд *условного перемещения*. С помощью *сравнения* присвоение $b3$ может быть скомпилировано в пять базовых команд RISC: две для генерации шестнадцатеричной константы, и, *сравнение* и *сдвиг влево* на 3. (Первое, второе и последнее условные выражения требуют соответственно одной, трех и четырех команд.)

Код может быть скомпилирован в 30 команд. Все шесть строк с условными выражениями могут выполняться параллельно. На машине с достаточной степенью параллелизма этот код выполняется за десять тактов. Современные машины такой степенью параллелизма не располагают, так что в практическом плане может помочь замена первых двух использований y на x , что позволит параллельно выполнить три первых оператора.

Дэвид Сил [101] разработал алгоритм для вычисления $\text{ntz}(x)$, который основан на идее сжатия 2^{32} возможных значений x в малое плотное множество целых чисел и поиска в таблице. Он использует выражение $x \& -x$, чтобы привести количество возможных значений к небольшому числу. Значением этого выражения является слово, которое содержит единственный единичный бит в позиции младшего единичного бита в x или равно 0, если $x = 0$. Таким образом, имеется только 33 возможных значения выражения $x \& -x$. Однако они не плотные и находятся в диапазоне от 0 до 2^{31} .

Чтобы получить плотное множество из 33 значений $x \& -x$, генерируется идентифицирующее значение в шести старших битах младшей половины произведения константы на $x \& -x$. Поскольку $x \& -x$ является целочисленной степенью 2 или 0, умножение представляет собой либо сдвиг константы влево, либо ее умножение на 0. Пяти старших битов произведения недостаточно, так как нам надо получить 33 различных значения.

Соответствующий код показан в листинге 5.22, в котором неиспользуемые элементы таблицы обозначены как `u`.

Листинг 5.22. Алгоритм Сила подсчета количества завершающих нулевых битов

```
int ntz(unsigned x)
{
    static char table[64] =
        {32, 0, 1,12, 2, 6, u,13, 3, u, 7, u, u, u, u,14,
         10, 4, u, u, 8, u, u,25, u, u, u, u, u,21,27,15,
         31,11, 5, u, u, u, u, u, 9, u, u,24, u, u,20,26,
         30, u, u, u, u,23, u,19,29, u,22,18,28,17,16, u};

    x = (x & -x)*0x0450FBAF;
    return table[x >> 26];
}
```

Рассмотрим в качестве примера нечетное кратное 16. Тогда $x \& -x = 16$, так что умножение представляет собой простой сдвиг влево на 4 разряда. Старшие шесть битов младшей половины произведения представляют собой 010001, или 17 в десятичной системе счисления. Поиск в таблице превращает 17 в 4, которое является корректным значением количества завершающих нулевых битов нечетного кратного 16.

Имеются тысячи констант, обладающих необходимым свойством единственности. Наименьшая из них равна `0x0431472F`, а наибольшая — `0xFDE75C6D`. Сил выбрал константу, для которой умножение можно выполнить с помощью небольшого количества сдвигов и сложений. Поскольку $0x0450FBAF = 17 \cdot 65 \cdot 65535$, умножение можно выполнить следующим образом.

```
x = (x << 4) + x; // x = x*17
x = (x << 6) + x; // x = x*65
x = (x << 16) - x; // x = x*65535
```

При такой подстановке код в листинге 5.22 состоит из девяти элементарных команд плюс индексированная загрузка. Сила интересовал набор команд ARM, в котором *сдвиг* и *сложение* выполняются одной командой. При использовании этой архитектуры код состоит из шести команд, включая индексированную загрузку.

Чтобы сделать умножение с помощью сдвигов и сложений еще проще, хотелось бы найти обладающую требуемыми свойствами константу вида $(2^k \pm 1)(2^k \pm 1)$. Для таблицы размером 64 элемента такой константы нет и есть только еще одна константа, являющаяся произведением трех сомножителей такого вида: $0x08A1FBAF = 17 \cdot 65 \cdot 131071$. Применение таблиц размером 128 или 256 также не помогает. Только для 512-элементной таблицы имеется четыре подходящие константы вида $(2^k \pm 1)(2^k \pm 1)$; наименьшая из них — $0x0080FF7F = 129 \cdot 65536$. Соответствующую 512-элементную таблицу читатель может построить самостоятельно.

Имеется вариант алгоритма Сила, основанный на циклах де Брейна (de Bruijn) [80]. Существуют такие циклические последовательности на заданном алфавите, которые содержат в качестве подпоследовательностей каждую из последовательностей букв алфавита заданной длины ровно один раз. Например, цикл, содержащий все подпоследовательности $\{a, b, c\}$ длиной 2, — *aabacbbcc*.

Обратите внимание, что последовательность *ca* получается при "защикливании" последовательности от конца к началу. Если размер алфавита — k , а длина подпоследовательности — n , имеется k^n таких подпоследовательностей. Чтобы цикл содержал все такие подпоследовательности, он должен иметь длину не менее k^n , которая получается, если каждая подпоследовательность начинается в очередной позиции последовательности. Можно показать, что всегда существует цикл такой минимально возможной длины, содержащий все k^n последовательностей.

Для наших целей воспользуемся алфавитом $\{0,1\}$, а для работы с 32-битовыми словами нас интересует цикл, который содержит все 32 последовательности 00000, 00001, 00010, ..., 11111. Если имеется такой цикл, который начинается по крайней мере с четырех нулей, можно вычислить $ntz(x)$, сначала преобразуя x в слово, содержащее единственный единичный бит в позиции младшего единичного бита x , как в алгоритме Сила. Затем с использованием умножения мы выбираем 5-битовое поле цикла де Брейна, которое для каждого множителя будет иметь свое единственное значение. Затем количество завершающих нулей может быть найдено с помощью поиска в таблице. Вот как выглядит данный алгоритм. Используемый цикл де Брейна имеет следующий вид.

```
0000 0100 1101 0111 0110 0101 0001 1111
```

Это действительно цикл, так как при работе с ним после 32 показанных битов идут завершающие нулевые биты, как если бы эта последовательность была "свернута" в кольцо и за последними битами шли первые.

Имеется 33 возможных значения $ntz(x)$, а в цикле де Брейна — только 32 пятибитовые последовательности. Следовательно, два слова с разными значениями $ntz(x)$ должны отображаться на одно и то же число при поиске в таблице. Конфликтующими словами являются нулевое слово и слова, заканчивающиеся единичным битом. Для разрешения этой неоднозначности код проверяет слово на равенство нулю и возвращает 32, если слово нулевое. Способ разрешения без ветвлений, полезный, если ваш компьютер имеет команды предикатов сравнения, заключается в замене последнего оператора следующим оператором.

```
return table[x >> 26] + 32*(x == 0);
```


Сравнивая эти два алгоритма, мы видим, что алгоритм Сила не требует проверки на равенство нулю и допускает альтернативный вариант, в котором умножение выполняется с помощью шести элементарных команд. Алгоритм деБрейна использует меньшую по размеру таблицу. Цикл деБрейна, использованный в листинге 5.23 и открытый Дэни Дьюбом (Danny Dubé) [27], хорош тем, что в нем соответствующее умножение выполняется с помощью восьми элементарных команд. Соответствующая константа равна $0x04D7651F = (2047 \cdot 5 \cdot 256 + 1) \cdot 31$, так что умножение можно выполнить посредством сдвигов, сложений и вычитаний.

Листинг 5.23. Подсчет количества завершающих нулевых битов с использованием цикла деБрейна

```
int ntz(unsigned x)
{
    static char table[32] =
        { 0, 1, 2, 24, 3, 19, 6, 25, 22, 4, 20, 10, 16, 7, 12, 26,
          31, 23, 18, 5, 21, 9, 15, 11, 30, 17, 8, 14, 29, 13, 28, 27 };

    if (x == 0) return 32;
    x = (x & -x) * 0x04D7651F;
    return table[x >> 27];
}
```

Джон Райзер (John Reiser) [95] заметил, что имеется другой способ отображения 33 значений множителя x & $-x$ в алгоритме Сила на плотное множество уникальных целых чисел: деление и использование остатка. Наименьший делитель, обладающий требуемым свойством уникальности, равен 37. В результате получается код, показанный в листинге 5.24; здесь неиспользуемые элементы таблицы указаны как u .

Листинг 5.24. Алгоритм Райзера для подсчета количества завершающих нулевых битов

```
int ntz(unsigned x)
{
    static char table[37] = { 32, 0, 1, 26, 2, 23, 27,
                             u, 3, 16, 24, 30, 28, 11, u, 13, 4,
                             7, 17, u, 25, 22, 31, 15, 29, 10, 12,
                             6, u, 21, 14, 9, 5, 20, 8, 19, 18 };

    x = (x & -x) % 37;
    return table[x];
}
```

Интересно отметить, что при равномерном распределении чисел x среднее количество завершающих нулевых битов очень близко к 1.0. Чтобы увидеть это, сложим произведения $p_i n_i$, где p_i — вероятность того, что в слове содержится n_i завершающих нулевых битов. Таким образом,

$$S \cong \frac{1}{2} \cdot 0 + \frac{1}{4} \cdot 1 + \frac{1}{8} \cdot 2 + \frac{1}{16} \cdot 3 + \frac{1}{32} \cdot 4 + \frac{1}{64} \cdot 5 + \dots \cong \sum_{n=0}^{\infty} \frac{n}{2^{n+1}}$$

Чтобы вычислить эту сумму, рассмотрим следующий массив.

$$\begin{array}{cccccc}
 1/4 & 1/8 & 1/16 & 1/32 & 1/64 & \dots \\
 & 1/8 & 1/16 & 1/32 & 1/64 & \dots \\
 & & 1/16 & 1/32 & 1/64 & \dots \\
 & & & 1/32 & 1/64 & \dots \\
 & & & & 1/64 & \dots \\
 & & & & & \dots
 \end{array}$$

Сумма каждого столбца является членом ряда S . Следовательно, S — это сумма всех элементов массива. Вычислим сумму каждой строки.

$$\begin{aligned}
 \frac{1}{4} + \frac{1}{8} + \frac{1}{16} + \frac{1}{32} + \dots &= \frac{1}{2} \\
 \frac{1}{8} + \frac{1}{16} + \frac{1}{32} + \frac{1}{64} + \dots &= \frac{1}{4} \\
 \frac{1}{16} + \frac{1}{32} + \frac{1}{64} + \frac{1}{128} + \dots &= \frac{1}{8} \\
 &\dots
 \end{aligned}$$

Полная сумма, таким образом, равна $\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots = 1$. Абсолютная сходимость исходного ряда оправдывает использование перестановки.

Иногда требуется функция, аналогичная $\text{ntz}(x)$, но рассматривающая аргумент, равный 0, как специальный (возможно, ошибочный), так что значение функции в этой точке должно отличаться от “обычных” значений. Например, определим функцию “количество множителей 2 в x ” как

$$\text{nfact2}(x) = \begin{cases} \text{ntz}(x), & x \neq 0, \\ -1, & x = 0. \end{cases}$$

Эту функцию можно вычислить по формуле $31 - \text{nlz}(x \& -x)$.

Применение

В [37] рассматривается несколько интересных применений функции, вычисляющей количество завершающих нулевых битов. Эрик Дженсен (Eric Jensen) назвал ее *линейчной функцией* (ruler function), так как она задает отметки на линейке, которая разделена на половины, четверти, восьмые части и т.д.

Функция ntz используется в алгоритме Р.В. Госпера (R.W. Gosper), обнаруживающем циклы и заслуживающем детального рассмотрения в силу своей элегантности и функциональности, которую, на первый взгляд, от него трудно ожидать.

Пусть имеется последовательность X_0, X_1, X_2, \dots , определяемая правилом $X_{n+1} = f(X_n)$. Если область значений функции f конечна, то эта последовательность

является периодической. Она состоит из некоторой ведущей последовательности $X_0, X_1, X_2, \dots, X_{\mu-1}$, за которой идет периодически повторяющаяся последовательность $X_\mu, X_{\mu+1}, \dots, X_{\mu+\lambda-1}$ ($X_\mu = X_{\mu-\lambda}$, $X_{\mu+1} = X_{\mu-\lambda+1}$ и так далее, где λ — ее период). Нам требуется найти индекс первого элемента периодической подпоследовательности μ и период λ . Обнаружение цикла часто используется при тестировании генераторов случайных чисел и для определения, имеется ли цикл в связанном списке.

Конечно, можно сохранять все значения последовательности по мере их вычисления и сравнивать новый элемент со всеми предыдущими. Таким образом можно непосредственно вычислить начало второго цикла. Однако имеются и более эффективные алгоритмы — как в смысле времени выполнения, так и в смысле необходимой памяти.

Возможно, самый простой метод предложен Р.В. Флойдом (R.W. Floyd) [67, раздел. 3.1, задача 6]. В этом алгоритме итеративно выполняется следующий процесс.

$$\begin{aligned}x &= f(x) \\ y &= f(f(y))\end{aligned}$$

(Здесь начальные значения x и y равны X_0 .) После n -го шага $x = X_n$ и $y = X_{2n}$. Затем эти значения сравниваются. Если они равны, значит, элементы последовательности X_n и X_{2n} отделены друг от друга количеством элементов, кратным периоду λ , т.е. $2n - n = n$ кратно λ . Чтобы найти индекс μ , последовательность генерируется заново и X_0 сравнивается с X_n , X_1 — с X_{n+1} и так до тех пор, пока не будет найден элемент X_μ , равный $X_{n+\mu}$. И наконец λ определяется путем дальнейшей генерации элементов и сравнением X_μ с элементами $X_{\mu+1}, X_{\mu+2}, \dots$. Этот алгоритм требует небольшого количества памяти, но функция f при этом вычисляется многократно.

В алгоритме Госпера [44, item 132; 65, ответы к упражнениям из раздела 3.1, задача 7] определяется только период λ , но не начальная точка μ первого цикла. Главное достоинство данного алгоритма в том, что в нем не требуются повторное вычисление функции f , а также скорость работы и экономное использование памяти. Алгоритм не ограничен каким-либо количеством памяти; для работы ему требуется таблица размером $\log_2(\Lambda) + 1$, где Λ — максимально возможный период. Это не так много, например если заранее известно, что $\Lambda \leq 2^{32}$, то достаточно массива из 33 слов.

В листинге 5.25 приведен код на языке С, реализующий алгоритм Госпера. В качестве параметров функции передаются указатель на рассматриваемую функцию f и начальное значение X_0 . Функция возвращает верхнюю и нижнюю границы μ и период λ . (Хотя алгоритм Госпера и не позволяет определить точное значение μ , он позволяет вычислить его нижнюю и верхнюю границы μ_l и μ_u , такие, что $\mu_u - \mu_l + 1 \leq \max(\lambda - 1, 1)$.) Алгоритм работает путем сравнения X_n (для $n = 1, 2, \dots$) с подмножеством (размером $\lfloor \log_2 n \rfloor + 1$) предшествующих значений. Элементами данного подмножества являются ближайшие предшествующие X_i , такие, что $i+1$ заканчивается единичным битом (т.е. i — четное число, предшествующее n), $i+1$ заканчивается ровно одним нулевым битом, $i+1$ заканчивается ровно двумя нулевыми битами и т.д.

ЛИСТИНГ 5.25. Алгоритм Госпера

```

void ld_Gosper(int (*f) (int), int X0, int *mu_l,
              int *mu_u, int *lambda)
{
    int Xn, k, m, kmax, n, lgl;
    int T[33];

    T[0] = X0;
    Xn = X0;
    for (n = 1; ; n++)
    {
        Xn = f(Xn);
        kmax = 31 - nlz(n);           // Floor(log2 n)
        for (k = 0; k <= kmax; k++)
        {
            if (Xn == T[k]) goto L;
        }
        T[ntz(n+1)] = Xn;           // Нет совпадений
    }
L:
    // Вычисляем m = max{i | i < n и ntz(i+1) = k}

    m = (((n >> k) - 1) | 1) << k) - 1;
    *lambda = n - m;
    lgl = 31 - nlz(*lambda - 1);    // Ceil(log2 lambda)-1
    *mu_u = m;                      // Верхняя граница mu
    *mu_l = m - max(1, 1 << lgl) + 1; // Нижняя граница mu
}

```

Таким образом, сравнения выполняются так.

$X_1 : X_0$	$X_7 : X_6, X_5, X_3$	$X_{13} : X_{12}, X_9, X_{11}, X_7$
$X_2 : X_0, X_1$	$X_8 : X_6, X_5, X_3, X_7$	$X_{14} : X_{12}, X_{13}, X_{11}, X_7$
$X_3 : X_2, X_1$	$X_9 : X_8, X_5, X_3, X_7$	$X_{15} : X_{14}, X_{13}, X_{11}, X_7$
$X_4 : X_2, X_1, X_3$	$X_{10} : X_8, X_9, X_3, X_7$	$X_{16} : X_{14}, X_{13}, X_{11}, X_7, X_{15}$
$X_5 : X_4, X_1, X_3$	$X_{11} : X_{10}, X_9, X_3, X_7$	$X_{17} : X_{16}, X_{13}, X_{11}, X_7, X_{15}$
$X_6 : X_4, X_3, X_3$	$X_{12} : X_{10}, X_9, X_{11}, X_7$	$X_{18} : X_{16}, X_{17}, X_{11}, X_7, X_{15}$

Можно показать, что вычисления всегда завершаются во втором цикле, т.е. при $n < \mu + 2\lambda$. Более подробно этот алгоритм описан в [67].

Линейная функция используется в задаче о ханойских башнях. Перенумеруем n дисков от 0 до $n-1$. При каждом перемещении k , где k принимает значения от 1 до $2^n - 1$, циклически перемещаем диск $\text{ntz}(k)$ на минимальное разрешенное расстояние вправо, циклическим образом.

Линейная функция используется также при генерации отраженного бинарного кода Грея (см. раздел 13.1 на с. 337). Начиная с произвольного n -битового слова на каждом шаге k (где k принимает значения от 1 до $2^n - 1$), изменяем бит $\text{ntz}(k)$.

Упражнения

1. Разработайте код алгоритма Дьюба с преобразованным в элементарные операции умножением.
2. Закодируйте функцию "правого выравнивания" $x \gg \text{ptz}(x)$, $x \neq 0$, тремя базовыми командами RISC.
3. Являются ли параллельные суффикс и префикс (с использованием исключающего или) обратимыми операциями? Если да, как можно вычислить соответствующие обратные функции?

ГЛАВА 6

ПОИСК В СЛОВЕ

6.1. Поиск первого нулевого байта

Необходимость в отдельной функции поиска нулевого байта возникает главным образом из-за способа представления символьных строк в языке C. Строки не содержат явно указанной длины; вместо этого в конце строки помещается нулевой байт. В C для вычисления длины строки используется функция `strlen`. Эта функция сканирует строку слева направо до тех пор, пока не находит нулевой байт, и возвращает количество просканированных байтов без учета нулевого.

Быстрая реализация функции `strlen` может загружать в регистр целое слово за раз и искать в нем нулевой байт. На компьютерах, где первым хранится старший байт слова, функция должна возвращать смещение первого нулевого байта слева. Удобно использовать значения от 0 до 3, которые соответствуют номеру нулевого байта в слове: если нулевого байта в слове нет, возвращается значение 4. Это значение по ходу поиска добавляется к длине строки. На компьютерах, где первым хранится младший байт, соответствующая проверка должна возвращать номер первого нулевого байта справа, поскольку порядок байтов при загрузке слова в регистр на таких компьютерах изменяется на обратный. В частности, нас интересуют функции `zbyte1(x)` и `zbyte2(x)`, определенные ниже (здесь "00" обозначает нулевой байт, "nn" — ненулевой байт, байт "xx" может быть как нулевым, так и ненулевым).

$$\text{zbyte1}(x) = \begin{cases} 0, & x = 00 \text{ xxxxxx}, \\ 1, & x = nn 00 \text{ xxxx}, \\ 2, & x = nnnn 00 \text{ xx}, \\ 3, & x = nnnnnn 00, \\ 4, & x = nnnnnnnn. \end{cases} \quad \text{zbyte2}(x) = \begin{cases} 0, & x = \text{xxxxxx} 00, \\ 1, & x = \text{xxxx} 00 \text{ nn}, \\ 2, & x = \text{xx} 00 \text{ nnnn}, \\ 3, & x = 00 \text{ nnnnnn}, \\ 4, & x = \text{nnnnnnnn}. \end{cases}$$

В листинге 6.1 приведен код функции `zbyte1(x)`, которая ищет *первый нулевой байта слева*. Проверка выполняется последовательно, слева направо, после чего возвращается номер первого слева нулевого байта (если таковой обнаружен).

Листинг 6.1. Поиск первого слева нулевого байта с использованием ветвления

```
int zbyte1 (unsigned x)
{
    if ((x >> 24) == 0) return 0;
    else if ((x & 0x00FF0000) == 0) return 1;
    else if ((x & 0x0000FF00) == 0) return 2;
    else if ((x & 0x000000FF) == 0) return 3;
    else return 4;
}
```

Для вычисления функции $\text{zbyte1}(x)$ требуется выполнить от 2 до 11 базовых RISC-команд (11 команд в случае, если в слове нет нулевых байтов, что является наиболее распространенной ситуацией для функции strlen). Код для функции $\text{zbyter}(x)$ практически аналогичен коду $\text{zbyte1}(x)$.

В листинге 6.2 функция $\text{zbyte1}(x)$ вычисляется без использования команд ветвления. Идея заключается в том, чтобы преобразовать каждый нулевой байт в 0x80, а каждый ненулевой — в нулевой, после чего достаточно воспользоваться функцией подсчета ведущих нулевых битов. Если в компьютере есть команда для вычисления количества ведущих нулевых битов и команда *или-не*, то для поиска нулевого байта потребуется восемь команд. Похожий алгоритм описан в [76].

Листинг 6.2. Поиск первого слева нулевого байта без использования ветвления

```
int zbyte1 (unsigned x)
{
    unsigned y;
    int n;

    // Исходный байт: 00 80 другие
    y = (x & 0x7F7F7F7F) + 0x7F7F7F7F; // 7F 7F 1xxxxxxx
    y = ~(y | x | 0x7F7F7F7F); // 80 00 00000000
    n = nlz(y) >> 3; // n = 0 ... 4; 4, если в x
    return n; // нет нулевого байта
}
```

Позицию первого нулевого байта справа можно получить, разделив нацело количество завершающих нулевых битов в y на 8 (с отбрасыванием дробной части). Если воспользоваться выражением для вычисления завершающих нулевых битов с использованием команды nlz (см. раздел 5.4, “Подсчет завершающих нулевых битов” на с. 130), то в приведенном выше коде строку, в которой вычисляется n , можно заменить следующей строкой.

```
n = (32 - nlz(~y & (y - 1))) >> 3;
```

Таким образом, если в компьютере имеются команды *или-не* и *и-не*, решение поставленной задачи будет получено за 12 команд.

Следует отметить, что на PowerPC в большинстве случаев процедура поиска первого справа нулевого байта не нужна; вместо нее можно воспользоваться командой *загрузки слова с обратным порядком байтов* lwbgrx .

Процедура из листинга 6.2 более ценна для 64-разрядных компьютеров, чем для 32-разрядных. Дело в том, что на 64-разрядном компьютере эта процедура (с необходимыми вполне очевидными изменениями) выполняется примерно за то же количество команд (семь или десять — в зависимости от способа генерации констант), что и на 32-разрядном. Процедура с использованием команд ветвления (листинг 6.1) в худшем случае потребует для выполнения 23 команд.

Если требуется только проверить, есть ли в данном слове нулевой байт, то после второго присваивания y можно добавить команду *ветвления при нулевом значении* (или, напротив, *ненулевом*).

Метод, аналогичный приведенному в листинге 6.2, но находящий крайний *справа* нулевой байт в слове x ($zbyte1(x)$) имеет следующий вид [88].

```
y = (x - 0x01010101) & ~x & 0x80808080;
n = ntz(y) >> 3;
```

Этот код выполняется за пять команд, не считая команды загрузки константы, если машина имеет команды *и-не* и *количество завершающих нулевых битов*. Этот метод нельзя использовать для вычисления $zbyte1(x)$ из-за проблем, связанных с заемами. Более полезным он может оказаться для поиска первого нулевого байта строки на машине с прямым порядком байтов либо для проверки наличия нулевого байта (при этом используется только присваивание y) на машине с любым порядком байтов.

Если команда `ntz` недоступна, поиск первого нулевого байта усложняется. В листинге 6.3 показан один из вариантов поиска нулевого байта в слове без использования функции `ntz` (приведена только выполняемая часть кода).

ЛИСТИНГ 6.3. Поиск первого слева нулевого байта без использования функции `ntz`

```
        // Исходный байт: 00 80 другие
y = (x & 0x7F7F7F7F) + 0x7F7F7F7F; //7F 7F 1xxxxxxx
y = ~(y | x | 0x7F7F7F7F);        //80 00 00000000
//Эти шаги отображают:
if (y == 0) return 4;             //00000000 ==> 4,
else if (y > 0x0000FFFF)         //80xxxxxx ==> 0,
    return (y >> 31) ^ 1;        //0080xxxx ==> 1,
else                               //000080xx ==> 2,
    return (y >> 15) ^ 3;        //00000080 ==> 3,
```

Этот код требует выполнения от 10 до 13 базовых RISC-команд (10 команд в случае, когда в слове нет нулевого байта). Таким образом, этот алгоритм, пожалуй, уступает коду из листинга 6.1, хотя и имеет меньшее количество команд ветвления. К сожалению, масштабирование данного алгоритма для 64-разрядных компьютеров неэффективно.

Рассмотрим еще одну возможность избежать использования функции `ntz`. Значение y , вычисленное в листинге 6.3, состоит из четырех байтов, каждый из которых равен либо $0x00$, либо $0x80$. Остаток от деления такого числа на $0x7F$ представляет собой исходное число, единичные биты которого (их не более четырех) перемещены в четыре крайние справа позиции. Таким образом, остаток от беззнакового деления в диапазоне от 0 до 15 однозначно идентифицирует исходное число, например так.

```
getu(0x80808080, 127) = 15
getu(0x80000000, 127) = 8
getu(0x00008080, 127) = 3 и т.д.
```

Это значение можно использовать при поиске нулевого байта в качестве индекса в таблице размером 16 байт. Таким образом, часть кода, которая начинается со строки `if (y == 0)`, можно заменить следующим (здесь y — беззнаковое).


```
static char table[16] = {4, 3, 2, 2, 1, 1, 1, 1,
                        0, 0, 0, 0, 0, 0, 0, 0};
return table[y%127];
```

Вместо 127 можно воспользоваться числом 31, но, разумеется, с другим массивом `table`.

На практике методы, использующие остаток от деления на 127 или 31, не применяются, так как функция *остатка от деления* требует 20 и более тактов даже при непосредственной аппаратной реализации этой команды. Тем не менее код из листинга 6.3 можно усовершенствовать, заменив часть кода, начинающуюся со строки `if (y == 0)`, одной из двух следующих строк.

```
return table[hopu(y, 0x02040810) & 15];
return table[y*0x00204081 >> 28];
```

Здесь `hopu(a, b)` означает старшие 32 бита беззнакового произведения a и b . Во второй строке в соответствии с соглашением, принятым в языках высокого уровня, предполагается, что значение произведения представляет собой младшие 32 бита результата. Этот метод вполне практичен, особенно на компьютерах с быстрым умножением и на компьютерах с командой *сдвига-и-сложения*, так как в последнем случае умножение на `0x00204081` реализуется следующим образом.

$$y(1 + 2^7 + 2^{14} + 2^{21}) = y(1 + 2^7)(1 + 2^{14})$$

При использовании такого 4-тактowego умножения общее время вычислений составляет 13 тактов (7 тактов для вычисления y , 4 — для выполнения команд *сдвига-и-сложения*, 2 — для *сдвига вправо* на 28 и индексирования массива `table`); кроме того, этот код не содержит команд условного перехода.

Этот алгоритм достаточно просто масштабируется для применения на 64-разрядных компьютерах. В алгоритме с остатком от деления используем следующую инструкцию.

```
return table[y%511];
```

В этом случае `table` содержит 256 элементов со значениями 8, 0, 1, 0, 2, 0, 1, 0, 3, 0, 1, 0, 2, 0, 1, 0, 4, ... (т.е. `table[i]` равно количеству завершающих нулевых битов в i).

В мультипликативных методах используется одна из двух инструкций.

```
return table[hopu(y, 0x0204081020408100) & 255];
return table[(y*0x0002040810204081)>>56];
```

Здесь `table` имеет размер 256 и значения 8, 7, 6, 6, 5, 5, 5, 5, 4, 4, 4, 4, 4, 4, 4, 4, 3,

Умножение на `0x2040810204081` можно выполнить за 13 тактов следующим образом.

$$\begin{aligned}t_1 &\leftarrow y(1 + 2^7) \\t_2 &\leftarrow t_1(1 + 2^{14}) \\t_3 &\leftarrow t_2(1 + 2^{28})\end{aligned}$$

Очевидно, все варианты поиска с использованием таблиц легко реализуют поиск первого справа нулевого байта путем простого изменения данных в таблице.

Если на 32-разрядном компьютере нет отдельной команды *или-не*, то вместо команды *не* во втором присваивании *y* в листинге 6.3 можно использовать один из трех описанных выше операторов *return* с массивом $table[i] = \{0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 2, 2, 3, 4\}$. Эта схема не работает для 64-разрядного компьютера.

Рассмотрим еще один интересный вариант процедуры из листинга 6.2, ориентированный на машины, в которых нет функции *nlz*. Пусть *a*, *b*, *c* и *d* — однобитовые переменные, означающие предикаты “первый байт *x* ненулевой”, “второй байт *x* ненулевой” и т.д. Тогда получаем следующее.

$$zbyte1(x) = a + ab + abc + abcd$$

Здесь умножение реализуется с помощью команды *и*, что дает процедуру, показанную в листинге 6.4 (приведена только выполняемая часть кода).

Листинг 6.4. Поиск первого слева нулевого байта вычислением полинома

```

y = (x & 0x7F7F7F7F) + 0x7F7F7F7F;
y = y | x;           // 1 в ненулевых байтах

t1 = y >> 31;       // t1 = a
t2 = (y >> 23) & t1; // t2 = ab
t3 = (y >> 15) & t2; // t3 = abc
t4 = (y >> 7) & t3;  // t4 = abcd
return t1 + t2 + t3 + t4;

```

Всего требуется выполнить 15 базовых RISC-команд. Этот алгоритм работает не очень быстро, зато часть вычислений можно распараллелить. На суперскалярном компьютере, который способен выполнять параллельно до трех независимых арифметических команд, выполнение алгоритма займет 10 тактов.

Внесение небольших изменений позволяет легко найти первый справа нулевой байт.

$$zbyte1(x) = abcd + bcd + cd + d$$

(Здесь для вычислений потребуется на одну команду *и* больше, чем в коде из листинга 6.4.)

Некоторые простые обобщения

Функции *zbyte1* и *zbyte1* могут использоваться при поиске байта, содержащего некоторое заданное число. Для этого над аргументом *x* и словом, в каждом байте которого содержится нужное значение, производится операция *исключающего или* и в получившемся слове осуществляется поиск нужного байта. Например, чтобы найти ASCII-пробел (0x20) в слове *x*, необходимо найти нулевой байт в слове $x \oplus 0x20202020$.

Аналогично для поиска одинаковых байтов в словах *x* и *y* выполняется поиск нулевого байта в слове $x \oplus y$.

Код из листинга 6.2 и его модификации с небольшими изменениями можно использовать и для поиска, не связанного с границами байта. Например, пусть требуется найти нулевое значение (если таковое есть) в первых четырех битах, следующих за ними 12 битах или в последних 16 битах. Для этого можно использовать код из листинга 6.2 с маской 0x77FF7FFF [92] (если длина поля равна 1, то в соответствующем разряде маски устанавливается 0).

Поиск значения из заданного диапазона

Код из листинга 6.2 можно легко модифицировать для поиска байта, значение которого находится в диапазоне от 0 до некоторого заданного значения, меньшего 128. В приведенном ниже фрагменте вычисляется индекс первого слева байта, значение которого находится в диапазоне от 0 до 9.

```
y = (x & 0x7F7F7F7F) + 0x76767676;
y = y | x;
y = y | 0x7F7F7F7F;          // Байты > 9 = 0xFF
y = ~y;                      // Байты > 9 = 0x00
                              // Байты <= 9 = 0x80
n = nlz(y) >> 3;
```

Рассмотрим более общий случай. Предположим, что в слове требуется найти крайний слева байт, значение которого находится в диапазоне от a до b , причем разность верхней и нижней границ этого диапазона меньше 128. Например, в ASCII прописные буквы имеют значения от 0x41 до 0x5A. Чтобы найти первую прописную букву в слове, вычтем из него число 0x41414141 так, чтобы при этом не возникло распространения заемов за пределы байта. После этого поиск байта, содержащего значение из диапазона от 0 до 0x19 (разность 0x5A – 0x41), выполняется так же, как и в приведенном выше коде. Используя формулы для вычитания из раздела 2.18, “Сложение, вычитание и абсолютное значение многобайтовых величин”, на с. 62 с очевидными упрощениями для $y = 0x41414141$, получим следующее.

```
d = (x | 0x80808080) - 0x41414141;
d = ~(x | 0x7F7F7F7F) ^ d;
y = (d & 0x7F7F7F7F) + 0x66666666;
y = y | d;
y = y | 0x7F7F7F7F; // Байты, не равные 41-5A, становятся FF
y = ~y;             // Байты, не равные 41-5A, становятся 00
                    // Байты, равные 41-5A, становятся 80
n = nlz(y) >> 3;
```

Для некоторых диапазонов значений возможен более простой код. Например, чтобы найти первый байт, содержащий значение в диапазоне от 0x30 до 0x39 (десятичная цифра в кодировке ASCII), необходимо выполнить операцию *исключающего или* с исходным словом и числом 0x30303030, а затем воспользоваться приведенным выше кодом для поиска байта со значением из диапазона от 0 до 9 (это упрощение применимо в тех случаях, когда n старших битов верхней и нижней границ диапазона совпадают и нижняя граница заканчивается 8 – n нулевыми битами).

Эти методы могут быть адаптированы для работы с большими диапазонами без дополнительных команд. Например, для поиска индекса первого слева байта, значение которого находится в диапазоне от 0 до 137 (0x89), строку программы поиска байта со значением от 0 до 9 $y = y | x$ следует заменить строкой $y = y \& x$.

Аналогично, заменив строку $y = y | d$ в программе поиска первого слева байта со значением от 0x41 до 0x5A строкой $y = y \& d$, получим код для поиска первого слева байта, содержащего значение в диапазоне от 0x41 до 0xDA.

6.2. Поиск строки единичных битов заданной длины

Здесь задача заключается в поиске в слове, находящемся в регистре, первой строки из единичных битов, длина которой не менее n . В результате поиска возвращается позиция первого бита строки; если такой строки в слове не обнаружено, выдается некоторое специальное значение. Вариантами этой задачи являются задача ответа на вопрос о существовании данной строки (возвращается значение *да* или *нет*) и задача поиска строки, состоящей ровно из n единичных битов. Задача находит применение в программах распределения дискового пространства, в частности при дефрагментации диска (перемещения данных на диске таким образом, чтобы в результате все блоки, в которых хранится некоторый файл, располагались непрерывно, друг за другом). Эта задача была предложена мне Альбертом Чангом (Albert Chang), который указал, что при ее решении можно использовать команду определения количества ведущих нулевых битов.

Далее предполагается, что компьютер имеет команду (либо подпрограмму) вычисления функции `niz`, определяющей количество ведущих нулевых битов.

Первым в голову приходит простейший алгоритм: вычислить количество ведущих нулевых битов и выполнить сдвиг влево на полученное число. Затем вычисляется количество ведущих единичных битов (инвертируя слово и подсчитывая количество ведущих нулевых битов). Если количество нулевых битов в этой группе не меньше заданной длины строки, значит, искомая группа найдена. Если количество нулевых битов меньше требуемого, выполняется сдвиг слова влево на количество нулевых битов и все действия повторяются сначала. Соответствующий код приведен ниже. Если найдены n идущих подряд единичных битов, возвращается число от 0 до 31, указывающее позицию первого бита первой такой строки. В противном случае возвращается значение 32, указывающее на то, что соответствующая строка не найдена.

```
int ffstr1(unsigned x, int n)
{
    int k, p;

    p = 0; // Инициализация возвращаемой позиции
    while (x != 0)
    {
        k = niz(x); // Удаление ведущих нулевых битов
        x = x << k; // (если они есть)
        p = p + k;
        k = niz(~x); // Подсчет группы единиц
        if (k >= n) // Если единиц достаточно,
```

```

        return p; // возвращается найденная позиция
    x = x << k;   // Если единиц мало,
    p = p + k;   // пропускаем их
}
return 32;
}

```

Алгоритм вполне применим, если цикл выполняется всего несколько раз, например когда в слове x есть относительно большие группы нулевых и единичных битов. Это условие вполне обычно для приложений, выполняющих распределение дискового пространства. Однако в худшем случае выполнение этого алгоритма занимает довольно много времени. Например, при $x = 0x55555555$ и $n \geq 2$ выполняется около 178 команд из полного набора RISC-команд.

Алгоритм, который лучше справляется с наихудшим случаем, основан на последовательности команд *сдвига влево* и *и*. Чтобы понять, как он работает, рассмотрим поиск строки из восьми или более последовательных единичных битов в 32-битовом слове x . Ниже показано, как можно выполнить этот поиск.

$$x \leftarrow x \& (x \ll 1)$$

$$x \leftarrow x \& (x \ll 2)$$

$$x \leftarrow x \& (x \ll 4)$$

После первого присваивания единичные биты в x указывают начальные положения строк единичных битов длиной 2. После второго присваивания единичные биты в x указывают начальные положения строк длиной 4 (строк длиной 2, следующих за другими строками длиной 2). После третьего присваивания единичные биты в x указывают начальные позиции строк длиной 8. Применение функции `nlz` к получившемуся слову возвращает положение первого бита строки, длина которой не меньше 8. Если такой строки нет, возвращается 32.

Для разработки алгоритма, который бы работал с любыми длинами n от 1 до 32, следует взглянуть на задачу несколько иначе. Прежде всего заметим, что три приведенные выше присваивания можно выполнять в любом порядке. Обратный порядок может быть даже более удобным. Для иллюстрации общего метода рассмотрим случай $n = 10$.

$$x_1 \leftarrow x \& (x \ll 5)$$

$$x_2 \leftarrow x_1 \& (x_1 \ll 2)$$

$$x_3 \leftarrow x_2 \& (x_2 \ll 1)$$

$$x_4 \leftarrow x_3 \& (x_3 \ll 1)$$

Первая инструкция выполняет сдвиг на $n/2$. После этого задача сводится к поиску строки из пяти последовательных единичных битов в x_1 . Это может быть сделано путем сдвига влево на $\lfloor 5/2 \rfloor = 2$ бита и выполнения команды *и*, после чего выполняется поиск строки из трех единичных битов ($3 = 5 - 2$). Последние две инструкции определяют начальные позиции строк длиной 3 в слове x_2 . Сумма величин всех сдвигов всегда равна $n - 1$. Соответствующий код показан в листинге 6.5. Время работы кода при n от 1 до 32 составляет от 3 до 36 команд из полного набора RISC-команд.

**Листинг 6.5. Поиск первой строки из n единичных битов
с помощью команд сдвига и и**

```
int ffstr1(unsigned x, int n)
{
    int s;

    while (n > 1) {
        s = n >> 1;
        x = x & (x << s);
        n = n - s;
    }
    return nlz(x);
}
```

Если n не очень велико, имеет смысл развернуть цикл (на 32-разрядном компьютере достаточно повторить тело цикла пять раз) и опустить проверку $n > 1$. Код с развернутым циклом не содержит условных переходов и выполняется за 20 команд (последнее присваивание n можно опустить). Тем не менее для малых значений n лучше использовать код с циклом. В случае малых значений n поиск по алгоритму с развернутым циклом выполняется медленнее, потому что, как только переменная n становится равной 1, ее значение больше не изменяется, и, таким образом, все последующие действия не изменят значений ни x , ни n . В смысле количества выполняемых инструкций развернутый цикл работает быстрее обычного цикла при $n \geq 5$.

Поиск строки, содержащей точно n единичных битов, требует на шесть команд больше (на четыре, если есть команда *и-не*). После выполнения всех действий листинга 6.5 единичные биты в x находятся в позициях, где в исходном слове начинается строка из n или более единичных битов. Следовательно, после подстановки вычисленного значения x выражение

$$x \& \neg(x \gg 1) \& \neg(x \ll 1)$$

содержит единичный бит в той позиции, в которой в x находится изолированный единичный бит, т.е. с этой позиции в исходном слове x начинается строка ровно из n единичных битов.

Данный алгоритм можно легко адаптировать и для поиска строк длиной n , которые начинаются в определенных позициях. Например, чтобы найти строку, первый бит которой находится на границе байта, необходимо применить команду *и* к конечному значению x и к числу 0x80808080.

Этот алгоритм можно использовать и для поиска строк из нулевых битов. При этом либо вначале инвертируются все биты x и затем выполняются все обычные вычисления, либо команды *и* заменяются командами *или*, а x инвертируется непосредственно перед вызовом функции *nlz*. Ниже показан алгоритм поиска первого (крайнего слева) нулевого байта (точное определение этой задачи приведено в разделе 6.1, "Поиск первого нулевого байта" на с. 141).

```

x ← x|(x << 4)
x ← x|(x << 2)
x ← x|(x << 1)
x ← 0x7F7F7F7F|x
p ← nlz(¬x) >> 3

```

Всего при этом требуется выполнение 12 команд из полного набора RISC-команд (т.е. этот алгоритм оказывается менее эффективным, чем алгоритм из листинга 6.2 на с. 142, в котором выполняется восемь команд).

6.3. Поиск наидлиннейшей строки единичных битов

Красивая и краткая функция, приведенная в листинге 6.6, возвращает длину самой длинной строки из единичных битов в слове x [54].

Листинг 6.6. Поиск наидлиннейшей строки единичных битов

```

int maxstr1(unsigned x)
{
    int k;
    for (k = 0; x != 0; k++) x = x & 2*x;
    return k;
}

```

Данная функция выполняет $4l + 3$ команды базового набора RISC, где l — длина наидлиннейшей строки, состоящей из единичных битов, или 131 команду в наихудшем случае.

Чтобы сократить время работы в наихудшем случае, можно прибегнуть к “логарифмической” версии, которая работает путем распространения нулей на одну, две, четыре, восемь и шестнадцать позиций влево, останавливаясь по достижении последнего ненулевого слова, а затем выполняя откат для поиска длины наидлиннейшей непрерывной строки единичных битов.

Пусть, например,

```
x = 0011 1111 1111 0011 1111 0011 1111 1000
```

Тогда

```

x2 = 0011 1111 1110 0011 1110 0011 1111 0000
x4 = 0011 1111 1000 0011 1000 0011 1100 0000
x8 = 0011 1000 0000 0000 0000 0000 0000 0000
x16 = 0000 0000 0000 0000 0000 0000 0000 0000

```

В этом случае последним ненулевым словом является $x8$. Заметим, что каждый единичный бит в $x8$ указывает позицию строки из восьми единиц. Таким образом, наидлиннейшая строка единичных битов начинается в позиции крайнего слева единичного бита $x8$, в данном примере — в 29-й позиции. Для проверки строк длиной 12 можно протестировать бит в 21-й позиции ($21 = 29 - 8$) в $x4$. Поскольку он равен нулю, строк

единичных битов длиной 12 в слове нет. Для проверки строк длиной 10 тестируем бит в 21-й позиции в $x2$. Так как этот бит единичный, позиция 29 является началом строки из 10 (или более) единичных битов. Наконец выполняем проверку наличия строки длиной 11, проверяя бит в позиции 19 ($19 = 21 - 2$) в x . Так как этот бит нулевой, наидлиннейшая строка единичных битов имеет длину 10 и начинается в позиции 29.

Эта схема закодирована в листинге 6.7, но приведенный в этом листинге код использует только две переменные — x и y — вместо пяти (x , $x2$, $x4$, $x8$ и $x16$). Этот код находит как длину, так и начальную позицию наидлиннейшей строки единичных битов, причем отсчет начальной позиции ведется с левого конца строки. Эта схема не работает, если x равно нулю или состоит только из единичных битов. Это частные случаи, обрабатываемые отдельно, причем последний из них обрабатывается в редко выполняемом месте, так как это достаточно редкая ситуация на практике.

Листинг 6.7. Поиск длины и начальной позиции наидлиннейшей строки единичных битов

```
int fmaxstr1(unsigned x, int *apos)
{
    unsigned y;
    int s;

    if (x == 0) (*apos = 32; return 0;)
    y = x & (x << 1);
    if (y == 0) {s = 1; goto L1;}
    x = y & (y << 2);
    if (x == 0) {s = 2; x = y; goto L2;}
    y = x & (x << 4);
    if (y == 0) {s = 4; goto L4;}
    x = y & (y << 8);
    if (x == 0) {s = 8; x = y; goto L8;}
    if (x == 0xFFFF8000) (*apos = 0; return 32;)
    s = 16;
L16: y = x & (x << 8);
    if (y != 0) {s = s + 8; x = y;}
L8: y = x & (x << 4);
    if (y != 0) {s = s + 4; x = y;}
L4: y = x & (x << 2);
    if (y != 0) {s = s + 2; x = y;}
L2: y = x & (x << 1);
    if (y != 0) {s = s + 1; x = y;}
L1: *apos = nlz(x);
    return s;
}
```

В наихудшем случае код выполняет 39 команд базового набора RISC плюс команды, необходимые для вычисления функции `nlz`. Если нас интересует только длина наидлиннейшей строки единичных битов, существенной экономии за счет отказа от сохранения позиции достичь не удастся (однако она возможна за счет отказа от вычисления функции `nlz`).

6.4. Поиск кратчайшей строки единичных битов

Гораздо труднее найти в слове *кратчайшую строку единичных битов*. Один из способов решения этой задачи заключается в том, чтобы пометить начала всех строк единичных битов в слове *b* и концы всех таких строк в слове *e*. Тогда, если *b&e* не равно нулю, длина кратчайшей строки равна 1. В противном случае сдвигаем *e* влево на одну позицию и повторяем проверку. Например, если

```
x = 0011 1111 1111 0011 1111 0011 1111 1000
```

то

```
b = 0010 0000 0000 0010 0000 0010 0000 0000
```

```
e = 0000 0000 0001 0000 0001 0000 0000 1000
```

После сдвига *e* влево на пять позиций получаем ненулевое значение *b&e*, что означает, что кратчайшая строка единичных битов имеет длину 6.

Эта идея реализована в листинге 6.8. Как и ранее, позиция строки отсчитывается от левого конца слова, и если две или больше строк имеют одинаковую длину, то функция находит ту, которая расположена левее всех прочих. Например, для $x = 0x00FF0FF0$ функция возвращает длину кратчайшей строки 8 и позицию, также равную 8.

ЛИСТИНГ 6.8. Поиск длины и начальной позиции кратчайшей строки единичных битов

```
int fminstr1(unsigned x, int *apos)
{
    int k;
    unsigned b, e;           // Начала и окончания строк

    if (x == 0) (*apos = 32; return 0;);
    b = ~(x >> 1) & x;      // Переходы 0-1
    e = x & ~(x << 1);     // Переходы 1-0

    for (k = 1; (b & e) == 0; k++)
        e = e << 1;
    *apos = nlz(b & e);
    return k;
}
```

При $n \geq 2$, где n — длина кратчайшей непрерывной строки единичных битов в x , эта функция выполняется за $8 + 4n$ команд базового набора RISC плюс время вычисления функции `nlz`.

Пожалуй, последней задачей в этом классе является поиск длины и позиции кратчайшей строки единичных битов в слове x , которая имеет длину не менее заданного целого числа $n > 0$. В терминах алгоритмов распределения памяти это алгоритм “наиболее подходящего” (best fit). Данную задачу можно решить путем распространения нулей в x на $n-1$ позиций влево с последующим поиском кратчайшей строки единиц в модифицированном слове x (см. упражнения).

Упражнения

1. Закодируйте усовершенствованный алгоритм Хсиха [54], который находит как длину, так и позицию наидлиннейшей строки единичных битов в слове x . Вы можете использовать функцию `nlz`.
2. Закодируйте функцию поиска длины и позиции кратчайшей строки единичных битов в слове x , имеющей длину не менее заданного значения n .
3. Еще один способ поиска кратчайшей строки единичных битов в слове x заключается в последовательном обнулении крайней справа строки битов в слове x и фиксации изменения количества единичных битов в слове на каждом шаге. Закодируйте функцию для полного набора команд RISC, которая использует эту идею и находит не только длину кратчайшей строки единиц, но и ее позицию.
4. Чему для "полностью случайного" 32-битового слова x (каждый бит которого равен 0 или 1 с вероятностью 0.5 независимо от других битов) равно среднее количество строк единичных битов? Ответ определяет среднее время выполнения функции из упр. 3 для соответствующих входных данных.
5. Вновь обратимся к "полностью случайным" 32-битовым словам x . Чему равна средняя длина кратчайшей строки единичных битов в x ? Ответ определяет среднее время выполнения функции `fminstr1` из листинга 6.8 для соответствующих входных данных. Вычислите интересующее нас значение с помощью программы, работающей методом Монте-Карло или исчерпывающего перебора.
6. У скольких из 2^n бинарных слов длиной n кратчайшая строка единичных битов имеет длину 1? То есть сколько n -битовых слов начинаются с 10, заканчиваются 01 или содержат подпоследовательность 010? Найдите решение в аналитическом виде или в виде рекурсии, не прибегая к исчерпывающему перебору.
7. Аналогично у скольких из 2^n бинарных слов длиной n кратчайшая строка единичных битов имеет длину 2?

ГЛАВА 7

ПЕРЕСТАНОВКА БИТОВ И БАЙТОВ

7.1. Реверс битов и байтов

Под операцией реверса битов подразумевается отражение содержимого регистра относительно середины слова, т.е. размещение битов в слове в обратном порядке.

$$\text{rev}(0x01234567) = 0xE6A2C480$$

$$\text{rev}(00000001001000110100010101100111) = 11100110101000101100010010000000$$

Под операцией реверса байтов подразумевается аналогичное отображение четырех байтов регистра. Операция реверса байтов необходима при преобразовании данных между форматом, когда первым идет младший байт (прямой порядок, little-endian), который используется, например, DEC и Intel, и форматом, когда первым идет старший байт (обратный порядок, big-endian), использующийся большинством других производителей.

Реверс битов достаточно эффективно выполняется путем обмена соседних битов, затем — соседних двухбитовых полей и так далее, как показано ниже [6]. Все пять операторов присвоения можно выполнять в произвольном порядке. Это тот же алгоритм, что и первый алгоритм подсчета количества единичных битов слова в разделе 5.1, но сложение в нем заменено обменом.

```
x = (x & 0x55555555) << 1 | (x & 0xAAAAAAAA) >> 1;  
x = (x & 0x33333333) << 2 | (x & 0xCCCCCCCC) >> 2;  
x = (x & 0x0F0F0F0F) << 4 | (x & 0xF0F0F0F0) >> 4;  
x = (x & 0x00FF00FF) << 8 | (x & 0xFF00FF00) >> 8;  
x = (x & 0x0000FFFF) << 16 | (x & 0xFFFF0000) >> 16;
```

Для ряда машин можно воспользоваться усовершенствованием, приведенным в листинге 7.1, которое заключается в том, чтобы избежать больших непосредственно задаваемых значений. Этот код требует выполнения 30 базовых RISC-команд и не использует команд ветвления.

ЛИСТИНГ 7.1. Реверс битов

```
unsigned rev(unsigned x)  
{  
    x = (x & 0x55555555) << 1 | (x >> 1) & 0x55555555;  
    x = (x & 0x33333333) << 2 | (x >> 2) & 0x33333333;  
    x = (x & 0x0F0F0F0F) << 4 | (x >> 4) & 0x0F0F0F0F;  
    x = (x << 24) | ((x & 0xFF00) << 8) |  
        ((x >> 8) & 0xFF00) | (x >> 24);  
    return x;  
}
```

Последнее присваивание x в этом коде выполняет реверс байта за девять команд базового набора RISC. Если компьютер способен к выполнению циклического сдвига, вычислить x можно за семь команд.

$$x = \left((x \& 0x00FF00FF) \gg 8 \right) | \left((x \ll 8) \& 0x00FF00FF \right)$$

На PowerPC операцию реверса байта можно реализовать всего за три команды [47]: *циклический сдвиг влево* на 8 битов, за которым следуют две команды *rlwimi* (*rotate left word immediate then mask insert* — *циклический сдвиг влево непосредственного заданного слова с последующей вставкой маски*).

Следующий алгоритм, разработанный Кристофером Стречи (Christopher Strachey) [104] (1961), по компьютерным меркам уже стар, но очень поучителен. Он обращает правые 16 битов слова в предположении, что перед началом работы левые 16 битов слова сброшены, и помещает обращенную половину слова в левую половину регистра.

Его работа основана на количестве битовых позиций, на которые необходимо переместить каждый бит. 16 битов, взятых слева направо, необходимо переместить на 1,3,5,...,31 позиций. Биты, которые необходимо переместить на 16 или более позиций, перемещаются первыми; затем перемещаются те, которые необходимо переместить на восемь или более позиций, и т.д. Алгоритм проиллюстрирован ниже; здесь каждая буква обозначает один бит, а точка обозначает бит, значение которого не играет никакой роли.

0000 0000 0000 0000	abcd efgh ijkl mnop	Исходное слово
0000 0000	ijkl mnop abcd efgh	После shl 16
0000	mnop ijkl efgh abcd	После shl 8
00op	mnkl ijgh efcd ab..	После shl 4
0pon	mlkj ihgf edcb a...	После shl 2
ponm	lkji hgfe dcba	После shl 1

Непосредственная реализация требует выполнения 16 команд из базового набора RISC плюс 12 команд для загрузки констант.

```
x = x | ((x & 0x000000FF) << 16);
x = (x & 0xFOFOFOFO) | ((x & 0x0F0F0F0F) << 8);
x = (x & 0xCCCCCCCC) | ((x & 0x33333333) << 4);
x = (x & 0xAAAAAAAA) | ((x & 0x55555555) << 2);
x = x << 1;
```

Для снижения количества различных масок можно воспользоваться дополнением. Если использовать “менее красивые” маски, можно сохранить и 16 правых битов.

Если доступны циклические сдвиги, идея Стречи может быть применена для реверса 32-битовых слов. Эта идея заключается в выяснении, на сколько битовых позиций необходимо циклически переместить влево каждый бит так, чтобы он оказался в своей окончательной позиции. Рассматривая биты слева направо, получаем величины сдвигов, равные 1,3,5,...,31,1,3,5,...,31 (ни один бит не перемещается на четное число позиций). Алгоритм сначала циклически перемещает те биты, которые необходимо переместить на 16 или более позиций, затем те, которые необходимо переместить на восемь или более по-

зиций, и так далее до тех, которые необходимо переместить на одну позицию (это все биты, потому что все величины сдвигов нечетны). Описанная схема, примененная для реверса 32-битового слова x , показана ниже. Функция $\text{shlr}(x, y)$ выполняет циклический сдвиг x влево на y позиций.

```
x = shlr(x & 0x00FF00FF, 16) | x & ~0x00FF00FF;
x = shlr(x & 0x0F0F0F0F, 8) | x & ~0x0F0F0F0F;
x = shlr(x & 0x33333333, 4) | x & ~0x33333333;
x = shlr(x & 0x55555555, 2) | x & ~0x55555555;
x = shlr(x, 1);
```

Приведенный код использует команду *и с дополнением*, чтобы избежать загрузки некоторых масок. Если ваша машина такой командой не оснащена, ее можно избежать, переписав первую строку кода как мультиплексорную операцию

```
x = shlr(x, 16) & 0x00FF00FF | x & ~0x00FF00FF;
```

и использовать тождество

$$x \& m | y \& \sim m = ((x \oplus y) \& m) \oplus y$$

для получения

```
x = ((shlr(x, 16) ^ x) & 0x00FF00FF) ^ x;
```

Аналогично следует поступить и для прочих строк, в которых использована команда *и с дополнением*.

Для многих машин немного лучшим способом (в плане обеспечения небольшой параллельности на уровне команд) является применение тождества [63]

$$x \& \sim m = (x \& m) \oplus x$$

и использование общего подвыражения *и*. Это приводит к функции, код которой показан в листинге 7.2 (17 команд плюс 8 команд для загрузки констант; итого — 25).

ЛИСТИНГ 7.2. Реверс битов с помощью циклических сдвигов

```
unsigned rev(unsigned x)
{
    unsigned t;
    t = x & 0x00FF00FF; x = shlr(t, 16) | t ^ x;
    t = x & 0x0F0F0F0F; x = shlr(t, 8) | t ^ x;
    t = x & 0x33333333; x = shlr(t, 4) | t ^ x;
    t = x & 0x55555555; x = shlr(t, 2) | t ^ x;
    x = shlr(x, 1);
    return x;
}
```

Вероятно, следует отметить, что константы 0x00FF00FF, 0x0F0F0F0F и прочие могут быть сгенерированы одна из другой, как показано ниже. Применение этого способа не слишком полезно для 32-разрядных машин (оно может оказаться просто вредным из-за снижения степени параллельности), так как 32-разрядные RISC-машины в общем случае

могут загружать константы двумя командами. Но этот способ может быть полезен на 64-разрядных машинах, для которых он и проиллюстрирован.

$$C_0 \leftarrow 0x0000\ 0000\ FFFF\ FFFF$$

$$C_1 \leftarrow C_0 \oplus (C_0 \ll 16)$$

$$C_2 \leftarrow C_1 \oplus (C_1 \ll 8)$$

...

Еще один способ реверса битов состоит в разбиении слова на три группы битов и обмене крайней слева и крайней справа групп, оставляя центральную на месте [8]. Этот способ продемонстрирован ниже для 27-битового слова.

```
012345678 9abcdefgh ijklmnopq Исходное 27-битовое слово
ijklmnopq 9abcdefgh 012345678 Первый обмен третей
opqilmnijk fghcde9ab 678345012 Второй обмен третей
qponmlkji hgfedcba9 876543210 Третий обмен третей
```

Далее приведен простой код, выполняющий описанные действия. При запуске на 32-разрядной машине он выполняет реверс битов с 0 по 26 и сбрасывает биты с 27 по 31.

```
x = (x & 0x000001FF) << 18 | (x & 0x0003FE00) |
(x >> 18) & 0x000001FF;
x = (x & 0x001C0E07) << 6 | (x & 0x00E07038) |
(x >> 6) & 0x001C0E07;
x = (x & 0x01249249) << 2 | (x & 0x02492492) |
(x >> 2) & 0x01249249;
```

Этот код требует выполнения 21 базовой команды RISC плюс 10 команд для загрузки констант; итого — 31 команда. Сравним с количеством команд кода, приведенного в листинге 7.1: 24 базовые команды RISC плюс шесть команд для загрузки констант плюс сдвиг вправо на 5 для выравнивания результата вправо, или всего 31 команда. Таким образом, тернарный метод не хуже (или даже лучше) при реверсе 27 или менее битов.

Следующая функция, предложенная Дональдом Кнутом (Donald E. Knuth) [73], интересна тем, что выполняет реверс 32-битового слова всего лишь за четыре шага, причем сдвиги и маски оказываются неожиданно "некрасивыми". Метод работает следующим образом.

```
01234567 89abcdef ghijklmn opqrstuv Исходное слово
fghijklm nopqrstu v0123456 789abcde "Поворот" влево на 15
pqrstuvn ofghijk labcde56 78901234 Обмен через 10 битов
tuvspqrm nojklifg hebcda96 78541230 Обмен через 4 бита
vutsrqpo mnlkjihg fedcba98 76543210 Обмен через 2 бита
```

Вот как выглядит простейший код, реализующий этот метод.

```
x = shlr(x, 15); // Циклический сдвиг влево на 15
x = (x & 0x003F801F) << 10 | (x & 0x01C003E0) |
(x >> 10) & 0x003F801F;
x = (x & 0x0E038421) << 4 | (x & 0x11C439CE) |
(x >> 4) & 0x0E038421;
x = (x & 0x22488842) << 2 | (x & 0x549556B5) |
(x >> 2) & 0x22488842;
```

Улучшить код в смысле количества операций можно ценой снижения степени параллельности, переписывая

```
x = (x & M1) << s | (x & M2) | (x >> s) & M1;
```

Здесь M2 представляет собой $\sim(M1 | (M1 << s))$, как

```
t = (x ^ (x >> s)) & M1; x = (t | (t << s)) ^ x;
```

В результате получается код, показанный в листинге 7.3 (19 команд из полного набора RISC плюс 6 для загрузки констант; итого — 25 команд).

Листинг 7.3. Алгоритм Кнута для реверса битов

```
unsigned rev(unsigned x)
{
    unsigned t;
    x = shlr(x, 15); // Циклический сдвиг влево на 15
    t = (x ^ (x>>10)) & 0x003F801F; x = (t | (t<<10)) ^ x;
    t = (x ^ (x>> 4)) & 0x0E038421; x = (t | (t<< 4)) ^ x;
    t = (x ^ (x>> 2)) & 0x22488842; x = (t | (t<< 2)) ^ x;
    return x;
}
```

Хотя алгоритм Кнута и не превосходит алгоритма реверса 32-битового значения с применением правых сдвигов (17 команд плюс 8 для загрузки констант), показанного в листинге 7.2, в его коде применен только один циклический сдвиг. Если закодировать его как

```
x = (x << 15) | (x >> 17); // Циклический сдвиг влево на 15
```

то алгоритм Кнута состоит из 21 команды плюс 6 команд для загрузки констант — наилучший результат для циклического сдвига 32-битового слова с использованием только базовых команд RISC. Это заставляет задуматься, нет ли простого способа предсказать количество сдвигов и логических команд, необходимых для реверса битов в слове заданной длины.

Можно ли распространить алгоритм Кнута для реверса 64-битового слова на 64-разрядной машине? Да, есть простой способ и способ посложнее. Простой способ заключается в том, чтобы обменять местами половины 64-битового регистра, а затем параллельно применить к ним 32-разрядную версию алгоритма Кнута. Соответствующий код приведен в листинге 7.4. Он выполняется за 24 команды, если обмен (циклический сдвиг на 32 бита) считать одной командой.

Листинг 7.4. Алгоритм Кнута для реверса битов 64-битового слова

```
unsigned long long rev(unsigned long long x)
{
    unsigned long long t;

    x = (x << 32) | (x >> 32); // Обмен половин регистра
    x = (x & 0x0001FFFF0001FFFFLL) << 15 | // Циклический сдвиг
        (x & 0xFFFF0000FFFF0000LL) >> 17; // влево на 15 битов
    t = (x ^ (x >> 10)) & 0x003F801F003F801FLL;
```



```

x = (t | (t << 10)) ^ x;
t = (x ^ (x >> 4)) & 0x0E0384210E038421LL;
x = (t | (t << 4)) ^ x;
t = (x ^ (x >> 2)) & 0x2248884222488842LL;
x = (t | (t << 2)) ^ x;
return x;
}

```

Другой способ заключается в том, чтобы найти величины сдвигов и маски, аналогичные применяемым в 32-битовом алгоритме Кнута. Такой код показан ниже. Он выполняется за 25 команд, если считать циклический сдвиг влево на 31 бит одной командой.

```

unsigned long long rev(unsigned long long x)
{
    unsigned long long t;

    x = (x << 31) | (x >> 33); // T.e. shlr(x, 31)
    t = (x ^ (x >> 20)) & 0x00000FFF800007FFLL;
    x = (t | (t << 20)) ^ x;
    t = (x ^ (x >> 8)) & 0x00F8000F80700807LL;
    x = (t | (t << 8)) ^ x;
    t = (x ^ (x >> 4)) & 0x0808708080807008LL;
    x = (t | (t << 4)) ^ x;
    t = (x ^ (x >> 2)) & 0x1111111111111111LL;
    x = (t | (t << 2)) ^ x;
    return x;
}

```

При выполнении реверса битов можно воспользоваться поиском в таблице. Приведенный ниже код выполняет реверс побайтово, используя 256-байтовую таблицу и заполняя четырехбайтовый результат побайтово в обратном порядке. Если развернуть цикл, для работы потребуется 13 базовых команд RISC плюс четыре загрузки, так что для некоторых машин этот код является непревзойденным.

```

unsigned rev(unsigned x)
{
    static unsigned char table[256] = {0x00, 0x80, 0x40,
    0xC0, 0x20, 0xA0, 0x60, 0xE0, ..., 0xBF, 0x7F, 0xFF};
    int i;
    unsigned r;

    r = 0;
    for (i = 3; i >= 0; i--) {
        r = (r << 8) + table[x & 0xFF];
        x = x >> 8;
    }
    return r;
}

```

Обобщенный реверс битов

В [37] предложено следующее обобщение реверса битов, названное переворотом (flip) и являющееся хорошим кандидатом для добавления в набор команд компьютера.

```

if (k & 1) x = (x & 0x55555555) << 1 | (x & 0xAAAAAAAA) >> 1;
if (k & 2) x = (x & 0x33333333) << 2 | (x & 0xCCCCCCCC) >> 2;
if (k & 4) x = (x & 0x0F0F0F0F) << 4 | (x & 0xF0F0F0F0) >> 4;
if (k & 8) x = (x & 0x00FF00FF) << 8 | (x & 0xFF00FF00) >> 8;
if (k & 16) x = (x & 0x0000FFFF) << 16 | (x & 0xFFFF0000) >> 16;
    
```

(Последние две команды *и* можно опустить.) Если $k = 31$, этот код выполняет реверс всех битов слова. При $k = 24$ выполняется реверс байтов слова. Если $k = 7$, выполняется реверс битов в каждом байте, но расположение байтов в слове при этом не изменяется. При $k = 16$ выполняется перестановка левого и правого полуслов, и т.д. В общем случае бит, находящийся в позиции m , перемещается в позицию $m \oplus k$. На аппаратном уровне этот процесс реализуется подобно циклическому сдвигу (пять уровней мультиплексирования, каждым из которых управляет свой бит величины сдвига k).

Современные методы реверса битов

В [44, item 167] приводятся более сложные выражения для реверса 6-, 7- и 8-битовых целых чисел. Несмотря на то что все эти выражения разработаны для 36-битовых машин, выражение для реверса 6-битовых целых чисел работает и на 32-битовых машинах, а выражения для реверса 7- и 8-битовых чисел — на 64-битовых машинах. Приведем эти выражения.

```

6 бит:   getu((x * 0x00082082) & 0x01122408, 255)
7 бит:   getu((x * 0x40100401) & 0x442211008, 255)
8 бит:   getu((x * 0x202020202) & 0x10884422010, 1023)
    
```

В результате получаются “чистые” целые числа, т.е. числа, выровненные по правой границе, неиспользуемые старшие биты которых равны 0.

Во всех трех случаях вместо функции `getu` можно использовать функцию `get` или `mod`, поскольку аргументы функции `getu` положительны. Функция *остаток от деления* в данных формулах просто суммирует цифры числа в системе счисления по основанию 256 или 1024 подобно тому, как в десятичной системе счисления остаток от деления на 9 просто равен сумме цифр числа (если полученное число больше 9, процесс суммирования цифр повторяется). Следовательно, эту функцию можно заменить командами *умножения* и *сдвига вправо*. Например, реверс 6-битовых чисел на 32-битовом компьютере можно выполнить так (умножение выполняется по модулю 2^{32}).

$$\begin{aligned}
 t &\leftarrow (x * 0x00082082) \& 0x01122408 \\
 &((t * 0x01010101) \gg 24)
 \end{aligned}$$

Применение описанных формул ограничено, так как они включают получение остатка (требующее для выполнения 20 тактов и больше) и/или несколько команд умножения, а также команды загрузки больших констант. Последняя формула, например, требует выполнения 10 базовых RISC-команд, две из которых — команды *умножения*, что составляет примерно 20 тактов на современных RISC-машинах. С другой стороны, адаптация кода из

листинга 7.1 для реверса 6-битовых целых чисел потребует около 15 команд и примерно от 9 до 15 тактов, в зависимости от возможностей конкретной машины по распараллеливанию вычислений на уровне команд. Тем не менее подобные формулы дают очень компактный код. Ниже приведено еще несколько могущих пригодиться алгоритмов для 32-разрядной машины. При применении этого метода для 8- и 9-битовых целых чисел на 32-разрядной машине в алгоритмах дважды применяются идеи из [44].

Вот как выглядит методика реверса 8-битовых целых чисел.

$$\begin{aligned} s &\leftarrow (x * 0x02020202) \& 0x84422010 \\ t &\leftarrow (x * 8) \& 0x00000420 \\ \text{getu}(s + t, 1023) \end{aligned}$$

Здесь функция `getu` не может быть заменена командами *умножения* и *сдвига*. (Чтобы понять, почему такая замена в данном случае оказывается некорректной, выполните все описанные действия и посмотрите, какие при этом получаются битовые комбинации.)

Вот еще один метод реверса 8-битового целого числа, который интересен тем, что может быть несколько упрощен.

$$\begin{aligned} s &\leftarrow (x * 0x00020202) \& 0x01044010 \\ t &\leftarrow (x * 0x00080808) \& 0x02088020 \\ \text{getu}(s + t, 4095) \end{aligned}$$

Упрощение заключается в том, что второе произведение представляет собой *сдвиг влево* первого произведения, последняя маска может быть получена из второй маски с помощью единственной команды *сдвига*, а команду *остаток от деления* можно заменить командами *умножения* и *сдвига*. После всех упрощений код требует выполнения 14 базовых RISC-команд, две из которых — команды *умножения*.

$$\begin{aligned} u &\leftarrow x * 0x00020202 \\ m &\leftarrow 0x01044010 \\ s &\leftarrow u \& m \\ t &\leftarrow (u \ll 2) \& (m \ll 1) \\ &(0x01001001 * (s + t)) \gg 24 \end{aligned}$$

Реверс 9-битового целого числа выполняется следующим образом.

$$\begin{aligned} s &\leftarrow (x * 0x01001001) \& 0x84108010 \\ t &\leftarrow (x * 0x00040040) \& 0x00841080 \\ \text{getu}(s + t, 1023) \end{aligned}$$

Второго умножения можно избежать, так как оно равно первому произведению, сдвинутому вправо на 6 бит. Последняя маска равна второй маске, сдвинутой вправо на 8 бит. После соответствующих замен код будет требовать выполнения 12 базовых RISC-

команд, включая одну команду *умножения* и одну *остатка от деления*. Команда *остаток от деления* должна быть беззнаковой, и ее нельзя заменить командами *умножения* и *сдвига*.

Читатель, изучивший эти необычные методы, может создать аналогичный код и для других операций перестановки битов. Рассмотрим следующий простой (и достаточно искусственный) пример. Пусть имеется 8-битовая величина, из которой требуется извлечь каждый второй бит и упаковать четыре полученных бита в крайние справа разряды, т.е. выполнить следующее преобразование.

```
0000 0000 0000 0000 0000 0000 abcd efgh ⇒
0000 0000 0000 0000 0000 0000 0000 bdfh
```

Данное преобразование можно выполнить следующим образом.

$$i \leftarrow (x * 0x01010101) \& 0x40100401$$

$$(i * 0x08040201) \gg 27$$

Тем не менее на большинстве компьютеров наиболее практичным оказывается метод, основанный на индексировании таблицы из однобайтовых (или 9-битовых) целых чисел.

Увеличение обращенного целого числа

В алгоритме быстрого преобразования Фурье (БПФ) используются как целая переменная i в цикле, в котором происходит ее увеличение на 1, так и значение $\text{rev}(i)$ [94]. Простейший способ работы состоит, конечно, в вычислении в каждой итерации увеличенного значения i с последующим вычислением $\text{rev}(i)$. Безусловно, для малых значений i быстрее и практичнее проводить поиск значений $\text{rev}(i)$ в таблице, но что делать при больших i , когда поиск в таблице неприменим в силу непрактичности, а для вычисления $\text{rev}(i)$ требуется 29 команд?

Если использовать поиск в таблице невозможно, лучше хранить i как в нормальном, так и в реверсном виде, увеличивая оба значения при каждой итерации. При этом возникает вопрос: как наиболее эффективно выполнить приращение целого числа, которое хранится в регистре в реверсном виде? Например, на 4-битовой машине при таком увеличении необходимо последовательно пройти такие значения (в шестнадцатеричной записи).

0, 8, 4, C, 2, A, 6, E, 1, 9, 5, D, 3, B, 7, F

В алгоритме быстрого преобразования Фурье как i , так и реверсное значение имеют одинаковую длину, которая почти наверняка меньше 32, и оба значения в регистрах выровнены по правой границе. Однако мы считаем, что i — 32-битовое целое число. После прибавления 1 к реверсному 32-битовому значению выполняется *сдвиг вправо* определенного количества битов, и затем полученное значение используется в алгоритме БПФ для индексации массива в памяти.

Непосредственный метод приращения реверсного значения состоит в выполнении последовательного сканирования слова слева направо до тех пор, пока не будет обнару-

жен первый нулевой бит. Этот бит устанавливается равным 1, а все биты, расположенные слева от него (если таковые имеются), сбрасываются в 0. Вот код этого метода.

```
unsigned x, m;

m = 0x80000000;
x = x ^ m;
if ((int)x >= 0)
{
    do
    {
        m = m >> 1;
        x = x ^ m;
    } while (x < m);
}
```

Здесь происходит выполнение трех базовых RISC-команд, если x начинается с нулевого бита, плюс четыре дополнительные команды для каждой итерации цикла. Поскольку вероятность того, что x начинается с нулевого бита, равна $1/2$, что x начинается с 10 (в бинарной записи) равна $1/4$ и так далее, среднее число выполняемых команд приближенно равно

$$\begin{aligned} & 3 \cdot \frac{1}{2} + 7 \cdot \frac{1}{4} + 11 \cdot \frac{1}{8} + 15 \cdot \frac{1}{16} + \dots \\ &= 4 \cdot \frac{1}{2} + 8 \cdot \frac{1}{4} + 12 \cdot \frac{1}{8} + 16 \cdot \frac{1}{16} + \dots - 1 \\ &= 4 \left(\frac{1}{2} + \frac{2}{4} + \frac{3}{8} + \frac{4}{16} + \dots \right) - 1 \\ &= 7 \end{aligned}$$

(Во второй строке мы добавили и вычли 1, причем первая единица представлена в виде суммы $1/2 + 1/4 + 1/8 + \dots$. Полученный ряд похож на ряд, уже рассматривавшийся на с. 136.) Однако в худшем случае выполнять приходится достаточно большое количество команд — 131.

Если в компьютере есть команда вычисления количества ведущих нулевых битов в слове, то увеличить реверсное значение на 1 можно следующим образом.

Сначала вычисляется $s \leftarrow \text{nlz}(-x)$,

а затем либо $x \leftarrow x \oplus \left(0x80000000 \gg s \right)$,

либо $x \leftarrow \left((x \ll s) + 0x80000000 \right) \gg s$.

В любом случае выполняются пять команд из полного набора RISC-команд, а кроме того, для корректного перехода от $0xFFFFFFFF$ к 0 требуется, чтобы сдвиги выполнялись по модулю 64. (Из-за этого приведенные формулы не будут работать на компьютерах с процессорами Intel x86, поскольку сдвиги в них выполняются по модулю 32.)

Довольно загадочная однострочная формула из [86], приведенная ниже, увеличивает реверсное целое число за шесть базовых команд RISC. В ней нет ветвлений, но есть целочисленное деление. Формула работает для целых чисел длиной вплоть до размера машинного слова, уменьшенного на 1.

$$revi \leftarrow revi \oplus \left(m - \frac{m}{(i \oplus (i+1)) + 1} \right)$$

Чтобы воспользоваться этой формулой, должны быть доступны как само число i , так и реверсное число $revi$. Переменная m представляет собой модуль; при работе с n -битовыми числами $m = 2^n$. Применение указанной формулы дает следующее значение реверсного числа; исходное число i следует увеличивать независимо, т.е. это значение не сдвигается к старшему концу регистра, как в двух предыдущих методах.

Вариантом этого метода является следующая формула.

$$revi \leftarrow revi \oplus \left(m - \frac{m/2}{-i \& (i+1)} \right) \tag{1}$$

Она вычисляется за пять команд, если машина имеет команду *и-не* и если m представляет собой константу, так что вычисление $m/2$ не учитывается. Она работает для целых чисел длиной до размера машинного слова. (В случае целых чисел размером в полное слово вместо первого вхождения m в формулу используйте 0, а вместо $m/2$ — значение 2^{n-1} .)

7.2. Перемешивание битов

Другим важным видом перестановки битов в слове является операция идеального перемешивания или идеального тасования (*perfect shuffle*), которая используется в криптографии. Существуют две разновидности идеального перемешивания, которые называются внешним (*outer*) и внутренним (*inner*). Обе операции чередуют биты из двух половин слова. В целом процесс аналогичен тасованию колоды из 32 карт, разница лишь в том, какая карта должна быть первой. В результате внешнего идеального перемешивания внешние биты остаются во внешних позициях, а после внутреннего идеального перемешивания бит из позиции 15 перемещается в позицию 31. Если 32-битовое слово равно (здесь каждая буква означает один бит)

abcd efgh ijkl mnop ABCD EFGH IJKL MNOP,

то после выполнения внешнего идеального перемешивания расположение битов в слове окажется следующим.

aAbB cCdD eEfF gGhH iIjJ kKlL mMnN oOpP

После внутреннего полного перемешивания получаем следующее расположение битов в слове.

AaBb CcDd EeFf GgHh IiJj KkLl MmNn OoPp

Пусть длина слова W представляет собой степень 2. Тогда операцию внешнего идеального перемешивания можно выполнить с помощью базовых RISC-команд за $\log_2(W/2)$ шагов, причем на каждом шаге выполняется перестановка второй и третьей четвертей во все меньших частях слова [37], т.е. 32-битовое слово преобразуется следующим образом.

```
abcd efgh ijkl mnop ABCD EFGH IJKL MNOP
abcd efgh ABCD EFGH ijkl mnop IJKL MNOP
abcd ABCD efgh EFGH ijkl IJKL mnop MNOP
abAB cdCD eFEF ghGH iJIJ kLKL mNMN oPOP
aAbB cCdD eEeF gGhH iIjJ kKlL mMnN oOpP
```

Вот как выглядит код, непосредственно реализующий описанный алгоритм и требующий выполнения 42 базовых RISC-команд.

```
x = (x & 0x0000FF00) << 8 | (x >> 8) & 0x0000FF00 | x & 0xFF0000FF;
x = (x & 0x00F000F0) << 4 | (x >> 4) & 0x00F000F0 | x & 0xF00FF00F;
x = (x & 0x0C0C0C0C) << 2 | (x >> 2) & 0x0C0C0C0C | x & 0xC3C3C3C3;
x = (x & 0x22222222) << 1 | (x >> 1) & 0x22222222 | x & 0x99999999;
```

Этот код можно сократить до 30 команд (хотя при этом время выполнения возрастет с 17 до 21 такта на машинах с неограниченными возможностями распараллеливания вычислений на уровне команд), если использовать для обмена двух полей регистра метод *исключающего или* (который рассматривался в разделе “Обмен двух полей одного регистра” на с. 69). Все используемые величины беззнаковые.

```
t = (x ^ (x >> 8)) & 0x0000FF00; x = x ^ t ^ (t << 8);
t = (x ^ (x >> 4)) & 0x00F000F0; x = x ^ t ^ (t << 4);
t = (x ^ (x >> 2)) & 0x0C0C0C0C; x = x ^ t ^ (t << 2);
t = (x ^ (x >> 1)) & 0x22222222; x = x ^ t ^ (t << 1);
```

Для выполнения обратной операции достаточно изменить последовательность обменов на обратную.

```
t = (x ^ (x >> 1)) & 0x22222222; x = x ^ t ^ (t << 1);
t = (x ^ (x >> 2)) & 0x0C0C0C0C; x = x ^ t ^ (t << 2);
t = (x ^ (x >> 4)) & 0x00F000F0; x = x ^ t ^ (t << 4);
t = (x ^ (x >> 8)) & 0x0000FF00; x = x ^ t ^ (t << 8);
```

Использование только последних двух шагов обоих приведенных алгоритмов тасования перемешивает биты в пределах байтов; три последних шага приводят к перемешиванию битов в пределах полуслова и т.д. В двух последних строках тасования выполняется перемешивание битов в каждом отдельном байте; в трех последних строках — в каждом полуслове и т.д. То же замечание справедливо и для операции, обратной перемешиванию, только в этом случае вместо последних шагов рассматриваются *первые*.

Чтобы выполнить внутреннее идеальное перемешивание, перед выполнением указанных последовательностей следует выполнить обмен левой и правой половин регистра.

```
x = (x >> 16) | (x << 16);
```

(Можно также воспользоваться *циклическим сдвигом* на 16 битов.) Операция, обратная внутреннему идеальному перемешиванию, выполняется аналогично, но упомянутая строка *добавляется* в конец кода.

Если изменить описанный выше процесс перемешивания таким образом, чтобы на каждом шаге выполнялась перестановка *первой* и *четвертой* четвертей во все меньших частях слова, то в результате получим реверс слова, образованного внутренним идеальным перемешиванием.

Пожалуй, стоит упомянуть один частный случай, когда левая половина слова (левое полуслово) x состоит из одних нулевых битов; иначе говоря, требуется разместить биты из правой половины слова x так, чтобы они оказались в каждом втором разряде, т.е. преобразовать исходное 32-битовое слово

0000 0000 0000 0000 ABCD EFGH IJKL MNOP

в слово

0A0B 0C0D 0E0F 0G0H 0I0J 0K0L 0M0N 0O0P

Внешнее идеальное перемешивание может быть упрощено так, что будет требовать выполнения 22 базовых RISC-команд. Приведенный же ниже код для нашей частной задачи выполняется только за 19 базовых RISC-команд без увеличения времени работы на компьютере с неограниченными возможностями распараллеливания вычислений на уровне команд (12 тактов при любом методе). Для данного кода не требуется установка левой половины слова x равной 0.

```
x = (x & 0xFF00) << 8 | (x & 0x00FF);
x = ((x << 4) | x) & 0xF0F0F0F;
x = ((x << 2) | x) & 0x33333333;
x = ((x << 1) | x) & 0x55555555;
```

Аналогично код обратной операции (частный случай *упаковки*; см. раздел 7.4, “Сжатие, или Обобщенное извлечение”, на с. 176) может быть упрощен до 26 или 29 базовых RISC-команд в зависимости от того, требуется ли выполнение команды *и* для очистки битов в нечетных разрядах. Код, приведенный ниже, позволяет решить задачу обратного идеального перемешивания для частного случая ненулевых четных битов за 18 или 21 базовую RISC-команду и с меньшим временем работы на компьютере с неограниченными возможностями распараллеливания вычислений на уровне команд, составляющим 12 или 15 тактов.

```
x = x & 0x55555555; // (При необходимости)
x = ((x >> 1) | x) & 0x33333333;
x = ((x >> 2) | x) & 0xF0F0F0F;
x = ((x >> 4) | x) & 0x00FF00FF;
x = ((x >> 8) | x) & 0x0000FFFF;
```

7.3. Транспонирование битовой матрицы

При транспонировании матрицы A получается матрица, столбцы которой являются строками матрицы A , а строки — столбцами. В этом разделе рассматриваются методы транспонирования битовой матрицы, элементы которой являются отдельными битами, упакованными по 8 в одном байте, а строки и столбцы начинаются на границах байтов. Это несложное, на первый взгляд, преобразование оказывается чрезвычайно громоздким в смысле количества выполняемых команд.

На большинстве компьютеров загрузка и сохранение отдельных битов выполняется очень медленно, главным образом из-за кода, который должен извлекать и (что еще сложнее) сохранять отдельные биты. Лучше разделить исходную матрицу на подматрицы размером 8×8 , загрузить каждую подматрицу 8×8 в регистры, транспонировать каждую из подматриц в регистрах по отдельности и затем сохранить их в соответствующих местах результирующей матрицы. На рис. 7.1 проиллюстрировано транспонирование битовой матрицы размером 2×3 байт. A, B, \dots, F представляют собой подматрицы размером 8×8 бит, а A^T, B^T, \dots обозначает транспонированные подматрицы A, B, \dots .

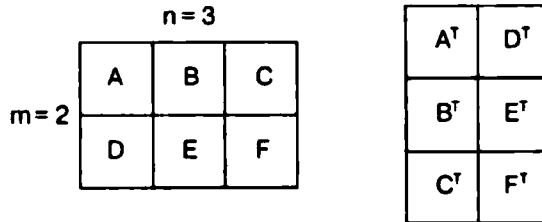


Рис. 7.1. Транспонирование матрицы 16×24

Как именно хранится матрица 8×8 — по строкам или столбцам, значения не имеет: в любом случае выполняются одни и те же действия. Предположим для определенности, что матрица хранится построчно. Тогда первый байт матрицы содержит верхнюю строку A , следующий байт содержит верхнюю строку B и т.д. Если обозначить через L адрес первого байта (верхней строки) подматрицы, то последовательные ее строки будут располагаться по адресам $L + n, L + 2n, \dots, L + 7n$.

В этой задаче мы откажемся от обычного предположения о 32-разрядной машине и будем считать, что компьютер имеет 64-битовые регистры общего назначения. Алгоритмы в этом случае оказываются проще (в том числе и для понимания), и их не так уж сложно преобразовать для выполнения на 32-разрядной машине. Фактически всю необходимую работу вместо вас в состоянии выполнить компилятор с поддержкой 64-битовых операций с целыми числами на 32-разрядной машине (хотя и, вероятно, не так эффективно, как это можно было бы сделать вручную).

Общая схема заключается в загрузке подматрицы восемью командами загрузки байта и упаковке байтов слева направо в 64-битовый регистр. Затем выполняется транспонирование содержимого регистра, после чего результат сохраняется в целевой области памяти с помощью восьми команд сохранения байта.

Ниже проиллюстрировано транспонирование битовой матрицы размером 8×8 , где каждый символ представляет отдельный бит.

0123 4567	08go wEMU
89ab cdef	19hp xFNV
ghij klmn	2aiq yGOW
opqr stuv	3bjr zHPX
wxyz ABCD	4cks AIQY
EFGH IJKL	5dlt BJRZ
MNOP QRST	6emu CKSS
UVWX YZS.	7fnv DLT.

В терминах двойных слов преобразование заключается в показанной далее замене первой строки второй.

```
01234567 89abcdef ghijklmn opqrstuv wxyzABCD EFGHIJKL MNOPQRST UVWXYZS.
08g0wEMU 19hpxFNV 2aiqyGOM 3bjrzHPX 4cksAIQY 5dltBJRZ 6emuCKSS 7fnvDLT.
```

Обратите внимание, что бит, обозначенный как 1, перемещается на 7 позиций вправо, бит, обозначенный как 2, — на 14 позиций вправо, а бит, обозначенный как 8, — на 7 позиций влево. Каждый бит перемещается на 0, 7, 14, 21, 28, 35, 42 или 49 позиций влево или вправо. Поскольку в двойном слове перемещаются 56 бит и имеется только 14 различных ненулевых перемещений, за один раз в среднем перемещается четыре бита, чего можно достичь соответствующими масками и сдвигами. Простейший непосредственно реализующий эту идею код имеет следующий вид.

```
y = x & 0x8040201008040201LL |
(x & 0x0080402010080402LL) << 7 |
(x & 0x0000804020100804LL) << 14 |
(x & 0x0000008040201008LL) << 21 |
(x & 0x0000000080402010LL) << 28 |
(x & 0x0000000000804020LL) << 35 |
(x & 0x0000000000008040LL) << 42 |
(x & 0x0000000000000080LL) << 49 |
(x >> 7) & 0x0080402010080402LL |
(x >> 14) & 0x0000804020100804LL |
(x >> 21) & 0x0000008040201008LL |
(x >> 28) & 0x0000000080402010LL |
(x >> 35) & 0x0000000000804020LL |
(x >> 42) & 0x0000000000008040LL |
(x >> 49) & 0x0000000000000080LL;
```

Этот код выполняется с помощью 43 команд из базового набора RISC, исключая генерацию масок (которая неважна в приложении, транспонирующем большую битовую матрицу, так как маски являются константами цикла). Циклические сдвиги в данном случае не помогают. Одни члены имеют вид $(x \& \text{mask}) \ll s$, а другие — $(x \gg s) \& \text{mask}$. Это уменьшает количество требуемых масок; последние семь из них являются повторением первых. Заметим также, что после генерации первой маски каждая последующая маска может быть получена из предыдущей с помощью команды *сдвига вправо*. Благодаря этому довольно просто написать более компактную версию этого кода с применением цикла `for`, выполняющегося семь раз.

Другой вариацией является применение метода Стила (Steele), который использует *исключающее или* для обмена битовых полей (описан на с. 70). Этот подход не слишком эффективен для применения в приложениях — в результате получается функция, выполняющая 42 команды, не считая генерации масок. Код начинается с

```
t = (x ^ (x >> 7)) & 0x0080402010080402LL;
x = x ^ t ^ (t << 7);
```

и содержит семь таких пар строк.

Хотя, похоже, не имеется *действительно эффективного* алгоритма для решения данной задачи, описанный далее метод превосходит простейший и его варианты при-

мерно в два раза при использовании базового набора команд RISC в части вычислений (без учета загрузки и сохранения подматриц или генерации масок). Мощность этого метода связана с высоким уровнем битовой параллельности. Если бы элементы матрицы представляли собой слова, он бы не был столь эффективен (в этой ситуации лучшее, что можно сделать — это загрузить каждое слово и сохранить его в нужном месте).

На первом шаге битовая матрица 8×8 интерпретируется как 16 битовых матриц размером 2×2 и выполняется транспонирование каждой из них. На втором шаге матрица интерпретируется как четыре подматрицы размером 2×2 (каждый элемент такой подматрицы является матрицей размером 2×2) и выполняется транспонирование всех четырех подматриц 2×2 . На последнем шаге матрица интерпретируется как матрица 2×2 , элементами которой являются матрицы размером 4×4 , и выполняется транспонирование этой матрицы 2×2 . Эти преобразования проиллюстрированы ниже.

0123 4567	082a 4c6e	08go 4cks	08go wEMU
89ab cdef	193b 5d7f	19hp 5dlt	19hp xFNV
ghij klmn	goiq ksmu	2aiq 6emu	2aiq yGOW
opqr stuv =>	hpjr ltnv =>	3bjr 7fnv =>	3bjr zHPX
wxyz ABCD	wEyG AICK	wEMU AIQY	4cks AIQY
EFGH IJKL	xFzH BJDL	xFNV BJRZ	5dlt BJRZ
MNOP QRST	MUOW QYSS	yGOW CKSS	6emu CKSS
UVWX YZ\$.	NVPX RZT.	zHPX DLT.	7fnv DLT.

Полностью данная процедура показана в листинге 7.5. Параметр A представляет собой первый байт подматрицы 8×8 исходной матрицы, а параметр B — первый байт подматрицы 8×8 целевой матрицы.

Листинг 7.5. Транспонирование матрицы размером 8×8 бит

```
void transpose8(unsigned char A[8], int m, int n,
               unsigned char B[8])
{
    unsigned long long x;
    int i;
    for (i = 0; i <= 7; i++) // Загрузка 8 байт из входного
        x = x << 8 | A[m*i]; // массива и их упаковка в x

    x = x & 0xAA55AA55AA55AA55LL |
        (x & 0x00AA00AA00AA00AALL) << 7 |
        (x >> 7) & 0x00AA00AA00AA00AALL;
    x = x & 0xCCCC3333CCCC3333LL |
        (x & 0x0000CCCC0000CCCCLL) << 14 |
        (x >> 14) & 0x0000CCCC0000CCCCLL;
    x = x & 0xF0F0F0F0F0F0F0FLL |
        (x & 0x00000000F0F0F0FOLL) << 28 |
        (x >> 28) & 0x00000000F0F0F0FOLL;

    for (i = 7; i >= 0; i--) // Сохранение результата в
    { // выходном массиве B
        B[n*i] = x; x = x >> 8;
    }
}
```

Вычислительная часть этой функции выполняется за 21 команду. Каждый из трех основных этапов представляет собой обмен битов, так что можно переписать их с применением метода Стила обмена битов *исключаящим или*. При этом первое присваивание переменной x в листинге 7.5 принимает следующий вид.

```
t = (x ^ (x >> 7)) & 0x00AA00AA00AA00AALL;
x = x ^ t ^ (t << 7);
```

Вычислительная часть модифицированной таким образом функции выполняется за 18 команд, но имеет малую параллельность на уровне команд.

Алгоритм в листинге 7.5 работает в направлении роста длины групп обмениваемых битов, но этот метод может быть применен и в обратном направлении. Для этого сначала матрица 8×8 должна рассматриваться как матрица 2×2 , элементы которой представляют собой матрицы размером 4×4 , и транспонировать матрицу 2×2 . Затем каждая из четырех подматриц размером 4×4 рассматривается как матрица размером 2×2 , элементами которой являются матрицы размером 2×2 , и т.д. Получающийся при этом код такой же, как приведенный в листинге 7.5, с тем отличием, что три присваивания, выполняющие перестановку битов, выполняются в обратном порядке.

Как уже упоминалось, эти функции могут быть модифицированы для выполнения на 32-разрядной машине с использованием двух регистров для каждой 64-битовой величины. Если это сделано и все вычисления, приводящие к нулевому результату, использованы для выполнения очевидных упрощений, то получится 32-битовая версия простого непосредственного метода, приведенного на с. 169. Он выполняется за 74 команды (сравните с 43 на 64-разрядной машине), а 32-битовая версия кода из листинга 7.5 выполняется за 36 команд (сравните с 21 командой на 64-разрядной машине). Применение методики обмена битов Стила приводит к снижению количества команд ценой снижения параллельности на уровне команд, как и в случае 64-разрядной машины.

Транспонирование битовой матрицы размером 32×32

Естественно, что при транспонировании матриц больших размеров можно использовать тот же рекурсивный метод, что и при транспонировании матриц 8×8 . Таким образом, транспонирование матрицы 32×32 выполняется в пять этапов.

Детали реализации транспонирования матрицы 32×32 существенно отличаются от процедуры, приведенной в листинге 7.5, так как мы полагаем, что матрицу 32×32 невозможно полностью разместить в регистрах общего назначения и поэтому необходимо разработать компактную процедуру, которая индексирует соответствующие слова битовой матрицы для выполнения перестановки битов. Описанный алгоритм работает лучше при движении от больших блоков к меньшим.

На первом этапе исходная матрица рассматривается как четыре матрицы размером 16×16 и преобразуется следующим образом.

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \Rightarrow \begin{bmatrix} A & C \\ B & D \end{bmatrix}$$

Здесь A обозначает левую половину первых 16 слов исходной матрицы, B — правую половину первых 16 слов и т.д. Очевидно, что такое преобразование может быть выполнено посредством следующих обменов.

Правая половина слова 0 с левой половиной слова 16,

правая половина слова 1 с левой половиной слова 17,

...

правая половина слова 15 с левой половиной слова 31.

При реализации этого алгоритма в коде используется индекс k , принимающий значения от 0 до 15. В цикле по k правая половина слова с номером k меняется местами с левой половиной слова $k + 16$.

На втором этапе матрица рассматривается как 16 матриц размером 8×8 бит и выполняется следующее преобразование.

$$\begin{bmatrix} A & B & C & D \\ E & F & G & H \\ I & J & K & L \\ M & N & O & P \end{bmatrix} \Rightarrow \begin{bmatrix} A & E & C & G \\ B & F & D & H \\ I & M & K & O \\ J & N & L & P \end{bmatrix}$$

Такое преобразование может быть выполнено с помощью следующих обменов.

Битов 0x00FF00FF слова 0 с битами 0xFF00FF00 слова 8,

битов 0x00FF00FF слова 1 с битами 0xFF00FF00 слова 9

и т.д.

Это означает, что биты 0–7 (младшие восемь битов) слова 0 меняются местами с битами 8–15 слова 8 и т.д. Индексы первого слова в этих перестановках принимают значения $k = 0, 1, 2, 3, 4, 5, 6, 7, 16, 17, 18, 19, 20, 21, 22, 23$. На каждом шаге очередное значение k можно вычислить по следующей формуле.

$$k' = (k + 9) \& -8$$

В цикле по k биты слова k меняются местами с битами слова $k + 8$.

Аналогично на третьем этапе меняются местами

биты 0x0F0F0F0F слова 0 с битами 0xF0F0F0F0 слова 4,

биты 0x0F0F0F0F слова 1 с битами 0xF0F0F0F0 слова 5

и т.д.

Индексы битов в первом слове в этих перестановках равны $k = 0, 1, 2, 3, 8, 9, 10, 11, 16, 17, 18, 19, 24, 25, 26, 27$. На каждом шаге очередное значение k вычисляется по следующей формуле.

$$k' = (k + 5) \& -4$$

В цикле по k биты слова k меняются местами с битами слова $k + 8$.

Все описанные действия в целях компактности представлены в виде отдельной функции на языке C в листинге 7.6 [37]. Внешний цикл выполняет описанные пять этапов; переменная *j* при этом принимает значения 16, 8, 4, 2 и 1. На каждом этапе формируется маска *m*, принимающая значения 0x0000FFFF, 0x00FF00FF, 0x0F0F0F0F, 0x33333333 и 0x55555555. (Код вычисления маски компактен и красив: $m = m \wedge (m \ll j)$). Обратного ему кода не существует, и в этом основная причина того, что преобразования выполняются от блоков больших размеров к блокам меньших размеров.) Во внутреннем цикле *k* принимает описанные выше значения. Здесь выполняется обмен битов *a[k]*, определяемых маской *m*, с битами *a[k+j]*, сдвинутыми вправо на *j* и выделенными с помощью маски *m*, что эквивалентно выделению битов *a[k+j]* с помощью дополнения к маске *m*. Код, выполняющий обмен, использует адаптированную методику трех *исключающих или*, описанную на с. 69 в столбце (в).

Листинг 7.6. Компактный код транспонирования битовой матрицы 32×32

```
void transpose32(unsigned A[32])
{
    int j, k;
    unsigned m, t;

    m = 0x0000FFFF;
    for (j = 16; j != 0; j = j >> 1, m = m ^ (m << j))
    {
        for (k = 0; k < 32; k = (k + j + 1) & ~j)
        {
            t = (A[k] ^ (A[k+j] >> j)) & m;
            A[k] = A[k] ^ t;
            A[k+j] = A[k+j] ^ (t << j);
        }
    }
}
```

Откомпилированный с помощью компилятора GNU C для компьютера, аналогичного RISC с базовым набором команд, этот код содержит 31 команду — 20 во внутреннем цикле и 7 во внешнем без учета внутреннего. Следовательно, всего выполняется $4 + 5(7 + 16 \cdot 20) = 1639$ команд. Если при транспонировании матрицы использовать 16 обращений к процедуре транспонирования матрицы 8×8 из листинга 7.5 (модифицированной для работы на 32-разрядной машине), то в этом случае транспонирование выполняется за $16(101 + 5) = 1696$ команд в предположении, что все 16 вызовов следуют один за другим. Расходы на вызов функции составляют, как показало рассмотрение скомпилированного кода, пять команд. Таким образом, с точки зрения времени выполнения оба метода практически эквивалентны.

С другой стороны, в случае 64-разрядной машины код из листинга 7.6 легко модифицировать для транспонирования битовой матрицы размером 64×64. В этом случае понадобится выполнить примерно $4 + 6(7 + 32 \cdot 20) = 3886$ команд. Транспонирование той же матрицы с вызовом 64 функций транспонирования матрицы 8×8 требует выполнения около $64(85 + 5) = 5760$ команд.

В алгоритме все преобразования выполняются без использования дополнительной памяти, следовательно, этот алгоритм может применяться для транспонирования матриц больших размеров с дополнительными действиями по пересылке подматриц размером 32×32 . Это же можно сделать, поместив результирующую матрицу в другой области памяти, отличной от исходной, и выполняя пересылку подматриц в первой либо последней итерации внешнего цикла.

Около половины команд, выполняемых в функции из листинга 7.6, используются для управления циклом, и функция пять раз загружает и сохраняет всю матрицу. Может быть, имеет смысл развертывание циклов, которое позволит убрать из кода команды управления? Да, имеет, если вас интересует, в первую очередь, скорость работы программы, если дополнительные расходы памяти не представляют для вас проблемы, если кеш процессора вашего компьютера имеет достаточный размер для хранения больших блоков кода, а главное — при медленном выполнении команд ветвления вашим компьютером. Тело программы будет состоять из обменов битов с помощью шести команд, повторенных 80 ($5 \cdot 16$) раз. Кроме того, в программе будет 32 команды загрузки исходной матрицы и 32 команды сохранения результата; итого — по меньшей мере 544 команды.

В листинге 7.7 представлена программа, в которой такое развертывание циклов выполнено вручную. Данная функция помещает транспонированную матрицу в другой массив, отличный от исходного; однако, если требуется работа “на месте”, без привлечения дополнительной памяти, достаточно вызвать функцию с одинаковыми аргументами. Количество строк `swap` в данном коде равно 80. После компиляции компилятором GNU C для RISC-компьютера с базовым набором команд код содержит, включая пролог и эпилог, 576 команд, среди которых нет команд ветвления (не считая команды возврата из функции). У целевого компьютера нет команд *сохранить несколько слов* и *загрузить несколько слов*, но он в состоянии сохранять и загружать по два регистра за такт с помощью команд *сохранить двойное слово* и *загрузить двойное слово*.

Листинг 7.7. Линейный код транспонирования матрицы 32×32

```
#define swap(a0, a1, j, m) t = (a0 ^ (a1 >> j)) & m; \
                          a0 = a0 ^ t; \
                          a1 = a1 ^ (t << j);

void transpose32(unsigned A[32], unsigned B[32])
{
    unsigned m, t;
    unsigned a0, a1, a2, a3, a4, a5, a6, a7,
             a8, a9, a10, a11, a12, a13, a14, a15,
             a16, a17, a18, a19, a20, a21, a22, a23,
             a24, a25, a26, a27, a28, a29, a30, a31;

    a0 = A[ 0]; a1 = A[ 1]; a2 = A[ 2]; a3 = A[ 3];
    a4 = A[ 4]; a5 = A[ 5]; a6 = A[ 6]; a7 = A[ 7];
    ...
    a28 = A[28]; a29 = A[29]; a30 = A[30]; a31 = A[31];

    m = 0x0000FFFF;
    swap(a0, a16, 16, m)
    swap(a1, a17, 16, m)
```

```

...
swap(a15, a31, 16, m)
m = 0x00FF00FF;
swap(a0, a8, 8, m)
swap(a1, a9, 8, m)
...
...
swap(a28, a29, 1, m)
swap(a30, a31, 1, m)

B[ 0] = a0; B[ 1] = a1; B[ 2] = a2; B[ 3] = a3;
B[ 4] = a4; B[ 5] = a5; B[ 6] = a6; B[ 7] = a7;
...
B[28] = a28; B[29] = a29; B[30] = a30; B[31] = a31;
)

```

Производительность кода можно немного увеличить при наличии в компьютере команды *циклического сдвига* (не важно — влево или вправо). Идея состоит в замене всех операций обмена битов из листинга 7.7, каждая из которых выполняется за шесть команд, более простыми обменами без сдвига, выполняющимися за четыре команды каждая (при этом используются имеющиеся макросы, но с опущенными сдвигами).

Сначала выполняется циклический сдвиг вправо на 16 разрядов слов $A[16..31]$ (т.е. $A[k]$ для $16 \leq k \leq 31$). Затем выполняется обмен местами правых половин слов: $A[0]$ с $A[16]$, $A[1]$ с $A[17]$ и так далее, аналогично коду из листинга 7.7. После этого выполняется циклический сдвиг вправо на восемь разрядов слов $A[0..8]$ и $A[24..31]$ и обмен битами, соответствующими маске $0x00FF00FF$, между словами $A[0]$ и $A[8]$, $A[1]$ и $A[9]$ и т.д. После выполнения пяти этапов транспонирование полностью не завершено. Необходимо выполнить циклический сдвиг слова $A[1]$ на один разряд, слова $A[2]$ — на два разряда и так далее (всего 31 команда). Код для этого метода не приводится, однако ниже показана последовательность шагов для битовой матрицы 4×4 .

```

abcd   abcd   abij   abij   aeim   aeim
efgh  => efgh => efmn => nefm => nbfj => bfjn
ijkl   klij   klcd   klcd   kocg   cgko
mnop   opmn   opgh   hopg   hlpd   dhlp

```

Часть программы из листинга 7.7, в которой выполняются перестановки битов, требует выполнения 480 команд (80 перестановок по шесть команд). В исправленной программе с использованием команд *циклического сдвига* выполняется 80 перестановок по четыре команды плюс 80 команд *циклического сдвига* ($16 \cdot 5$) на первых пяти этапах и завершающая 31 команда *циклического сдвига*; итого — 431 команда. Поскольку код пролога и эпилога остается неизменным, использование команд *циклического сдвига* позволяет сэкономить 49 команд.

Существует еще один метод транспонирования битовой матрицы, основанный на трех преобразованиях сдвига [37]. Если имеется матрица размером $n \times n$, то выполняются следующие шаги: 1) циклический сдвиг строки i вправо на i разрядов; 2) циклический сдвиг столбца j вверх на $(j+1) \bmod n$ разрядов; 3) циклический сдвиг строки i вправо на $(i+1) \bmod n$ разрядов; 4) отражение матрицы относительно горизонтальной оси, проходящей через ее середину. Проиллюстрируем сказанное на примере матрицы 4×4 .

abcd	abcd	hlpd	dhlp	aeim
efgh	-> hefg	-> kocg	-> cgko	-> bfjn
ijkl	kl ij	nb fj	bfjn	cgko
m nop	no pm	aeim	aeim	dhlp

Этот метод не может конкурировать с другими из-за высокой стоимости шага 2. (Чтобы выполнить этот шаг за разумное количество шагов, сначала выполняется циклический сдвиг вверх на $n/2$ позиции всех столбцов, которые должны быть сдвинуты на $n/2$ разряда или более (это столбцы от $n/2-1$ до $n-2$), затем выполняется циклический сдвиг вверх на $n/4$ позиции и т.д.) Шаги 1 и 3 требуют выполнения всего лишь по $n-1$ команд, а шаг 4 не требует никаких действий, если сохранять результаты вычислений в соответствующих позициях.

Если битовая матрица 8×8 сохраняется в 64-битовом слове как обычно (верхняя строка помещается в старшие восемь битов и т.д.), то операция транспонирования такой матрицы эквивалентна трем внешним идеальным перемешиваниям или обратным им процессам [37]. Если перемешивание реализовано на компьютере в виде отдельной команды, этим можно воспользоваться; однако для RISC-компьютеров с базовым набором команд этот способ не годится.

7.4. Сжатие, или Обобщенное извлечение

В языке программирования APL имеется операция, которая называется сжатием (compress) и записывается как B/V , где B — булев вектор, а V — вектор с произвольными элементами той же длины, что и B . Результатом данной операции является вектор, состоящий из тех элементов V , для которых соответствующий бит в B равен 1. Длина результирующего вектора равна количеству единичных битов в B .

Рассмотрим аналогичную операцию с битами в слове. Заданы маска m и слово x ; требуется выбрать те биты слова x , которые соответствуют единичным битам маски, и переместить их вправо (“сжать”). Например, если сжимается слово

abcd efgh ijkl mnopqrst uvwx yzAB CDEF,

где каждый символ представляет один бит, а маска имеет вид

0000 1111 0011 0011 1010 1010 0101 0101,

то в результате будет получено слово

0000 0000 0000 0000 efgh klop qsuw zBDF.

Эту операцию можно также назвать *обобщенным извлечением* (generalized extract) по аналогии с командой *извлечения*, имевшейся в ряде компьютеров.

Нас интересует код, осуществляющий данную операцию с минимальным временем работы в наихудшем случае. В качестве первого приближения, которое будет совершенствоваться, возьмем простой цикл из листинга 7.8. Этот код не содержит команд ветвления в теле цикла и в худшем случае требует выполнения 260 команд, включая пролог и эпилог функции.

Листинг 7.8. Операция сжатия на основе простого цикла

```

unsigned compress(unsigned x, unsigned m)
{
    unsigned r, s, b; // Результат, сдвиг, маскируемый бит

    r = 0;
    s = 0;
    do
    {
        b = m & 1;
        r = r | ((x & b) << s);
        s = s + b;
        x = x >> 1;
        m = m >> 1;
    } while (m != 0);
    return r;
}

```

Улучшить этот код можно путем многократного применения метода “параллельного суффикса” (см. с. 120) с операцией *исключающего или* [37] (далее будем обозначать эту операцию как PS-XOR). Основная идея состоит в том, чтобы прежде всего определить биты аргумента *x*, которые будут перемещены на нечетное количество позиций вправо, и переместить их (эта задача упрощается, если сначала применить операцию *и* к слову *x* и к маске, очистив биты, которые не имеют значения). Биты маски перемещаются точно таким же образом. Затем определяются биты *x*, перемещаемые на нечетное число, умноженное на 2 (2, 6, 10 и т.д.), после чего вправо перемещаются как эти биты, так и соответствующие биты маски. Далее работаем с битами, перемещаемыми на нечетные числа, умноженные на 4, затем — на 8 и 16.

Поскольку этот алгоритм достаточно сложен для понимания, рассмотрим его более детально. Предположим, что входными данными для нашего алгоритма являются следующие.

```

x = abcd efgh ijkl mnopqrst uvwx yzAB CDEF,
m = 1000 1000 1110 0000 0000 1111 0101 0101,
   1   1   111
   9   6   333           4444 3 2 1 0

```

Здесь каждая буква в *x* означает один бит (который может принимать значения 0 или 1). Числа под каждым единичным битом маски указывают, насколько далеко должен перемещаться вправо соответствующий бит *x* (это значение равно количеству нулевых битов в *m* справа от рассматриваемого бита). Как упоминалось ранее, сначала стоит сбросить все не имеющие значения биты *x*, что даст нам такое число.

```

x = a000 e000 ijk0 0000 0000 uvwx 0z0B 0D0F

```

Сначала нужно определить, какие биты будут перемещены вправо на нечетное количество позиций, и переместить их на одну позицию вправо. Вспомним, что операция PS-XOR дает нам единичные биты в позициях, для которых количество единичных битов в данной позиции и справа от нее нечетно. Нас же интересуют биты, для которых нечет-

но количество нулевых битов, расположенных строго справа от данной позиции. Найти их можно путем вычисления $mk = \sim m \ll 1$ и применения операции PS-XOR к полученному значению. Результат будет следующим.

```
mk = 1110 1110 0011 1111 1110 0001 0101 0100,
mp = 1010 0101 1110 1010 1010 0000 1100 1100.
```

Заметим, что mk указывает биты m , которые содержат нулевой бит рядом справа, а mp суммирует их справа по модулю 2. Таким образом, mp определяет биты m , которые имеют нечетное количество нулевых битов справа.

Биты, которые должны быть перемещены на одну позицию, находятся в позициях, имеющих нечетное количество нулевых битов справа (указываемые словом mp), и в этих позициях в исходной маске m находятся единичные биты, т.е. их можно получить с помощью простой операции $mv = mp \& m$.

```
mv = 1000 0000 1110 0000 0000 0000 0100 0100
```

Соответствующие биты m могут быть перемещены путем присваивания.

```
m = (m ^ mv) | (mv >> 1);
```

Биты x могут быть перемещены с помощью двух присваиваний.

```
t = x & mv;
x = (x ^ t) | (t >> 1);
```

(Перемещение битов m выполняется проще в силу того, что здесь все выбранные биты являются единичными.) Операция *исключающего или* сбрасывает биты, которые равны 1 в m и x , а операция *или* устанавливает биты, которые в m и x равны 0. Эти операции могут быть заменены операциями *вычитания* и *сложения* либо обе быть операциями *исключающего или*. После перемещения битов, выбранных с помощью mv , на одну позицию вправо получаем следующее.

```
m = 0100 1000 0111 0000 0000 1111 0011 0011
x = 0a00 e000 0ijk 0000 0000 uvwx 00zb 00DF
```

Теперь необходимо подготовить маску для второй итерации, где будут указаны биты, которые должны быть перемещены вправо на величину, представляющую собой нечетное число, умноженное на 2. Заметим, что величина $mk \& \sim mp$ определяет биты, соседом которых справа в исходной маске m является нулевой бит, и биты, которые в исходной маске имеют справа четное число нулевых битов. Эти свойства применяются совместно, хотя и не каждое в отдельности, к новому значению маски m (т.е. mk указывает *все* позиции в новой маске m , у которых справа расположен нулевой бит, а всего справа находится четное число нулей). Будучи просуммированной справа с помощью операции PS-XOR, эта величина указывает биты, которые перемещаются вправо на количество позиций, равное удвоенному нечетному числу (2, 6, 10 и т.д.). Таким образом, необходимо присвоить это значение переменной mk и выполнить вторую итерацию описанных выше шагов. Новое значение mk имеет следующий вид.

```
mk = 0100 1010 0001 0101 0100 0001 0001 0000
```

Полностью функция *C*, реализующая рассматриваемую операцию, приведена в листинге 7.9. Она требует выполнения 127 базовых RISC-команд (константное значение, не зависящее от входных данных)¹, включая пролог и эпилог. В листинге 7.10 показана последовательность значений, принимаемых различными переменными в ключевых точках вычислений (входные данные те же, которые использовались при рассмотрении данного алгоритма). Заметьте, что побочным результатом работы алгоритма является значение маски *m*, в котором все единичные биты перенесены вправо.

ЛИСТИНГ 7.9. Метод параллельного суффикса для операции сжатия

```
unsigned compress(unsigned x, unsigned m)
{
    unsigned mk, mp, mv, t;
    int i;

    x = x & m;           // Сброс незначащих битов
    mk = ~m << 1;       // Подсчет 0 справа

    for (i = 0; i < 5; i++)
    {
        mp = mk ^ (mk << 1); // Параллельный суффикс
        mp = mp ^ (mp << 2);
        mp = mp ^ (mp << 4);
        mp = mp ^ (mp << 8);
        mp = mp ^ (mp << 16);
        mv = mp & m;       // Перемешаемые биты
        m = m ^ mv | (mv >> (1 << i)); // Сжатие m
        t = x & mv;
        x = x ^ t | (t >> (1 << i)); // Сжатие x
        mk = mk & ~mp;
    }
    return x;
}
```

ЛИСТИНГ 7.10. Действия при использовании метода параллельного префикса для операции сжатия

```

x = abcd efgh ijkl mnopqrst uvwx yzAB CDEF
m = 1000 1000 1110 0000 0000 1111 0101 0101
x = a000 e000 ijk0 0000 0000 uvwx 0z0B 0D0F
i = 0,   mk = 1110 1110 0011 1111 1110 0001 0101 0100
После PS, mp = 1010 0101 1110 1010 1010 0000 1100 1100
          mv = 1000 0000 1110 0000 0000 0000 0100 0100
          m  = 0100 1000 0111 0000 0000 1111 0011 0011
x = 0a00 e000 0ijk 0000 0000 uvwx 0z0B 0D0F
i = 1,   mk = 0100 1010 0001 0101 0100 0001 0001 0000
После PS, mp = 1100 0110 0000 1100 1100 0000 1111 0000
          mv = 0100 0000 0000 0000 0000 0000 0011 0000
          m  = 0001 1000 0111 0000 0000 1111 0000 1111
x = 000a e000 0ijk 0000 0000 uvwx 0000 zBDF
i = 2,   mk = 0000 1000 0001 0001 0000 0001 0000 0000
```

¹ В действительности первый *сдвиг влево* может быть опущен, что снизит количество команд до 126. Значение *mv* получается верным в любом случае [21].

```

После PS,   mp = 0000 0111 1111 0000 1111 1111 0000 0000
             mv = 0000 0000 0111 0000 0000 1111 0000 0000
             m  = 0001 1000 0000 0111 0000 0000 1111 1111
             x  = 000a e000 0000 0ijk 0000 0000 uvwx zBDF
i = 3,
После PS,   mp = 0000 0111 1111 1111 0000 0000 0000 0000
             mv = 0000 0000 0000 0111 0000 0000 0000 0000
             m  = 0001 1000 0000 0000 0000 0111 1111 1111
             x  = 000a e000 0000 0000 0000 0ijk uvwx zBDF
i = 4,
После PS,   mp = 1111 1000 0000 0000 0000 0000 0000 0000
             mv = 0001 1000 0000 0000 0000 0000 0000 0000
             m  = 0000 0000 0000 0000 0000 0001 1111 1111
             x  = 0000 0000 0000 0000 000a eijk uvwx zBDF

```

Алгоритм из листинга 7.9 на 64-битовой машине с базовым RISC- набором выполняется за 169 команд, в то время как алгоритм из листинга 7.8 в наихудшем случае требует выполнения 516 команд.

Количество команд, выполнения которых требует алгоритм из листинга 7.9, может быть значительно снижено, если m представляет собой константу. Это может быть в двух ситуациях: 1) когда вызов `compress(x, m)` осуществляется в цикле, в котором значение m не известно заранее, но является константой цикла; 2) когда известно значение m и код функции `compress` генерируется заранее, возможно, компилятором.

Обратите внимание на то, что значение, присвоенное x в цикле в листинге 7.9, используется в этом цикле исключительно в строке, где выполняется присвоение x . Как видно из кода, x зависит только от себя самого и переменной mv . Следовательно, все обращения к x можно удалить, а пять вычисляемых значений mv сохранить в переменных $mv0, mv1, \dots, mv4$. Тогда в ситуации 1 вне цикла, в котором имеется обращение `compress(x, m)`, можно разместить функцию, которая не будет обращаться к x , а будет лишь вычислять значения $mv i$, а в самом цикле можно разместить только следующие строки.

```

x = x & m;
t = x & mv0; x = x ^ t | (t >> 1);
t = x & mv1; x = x ^ t | (t >> 2);
t = x & mv2; x = x ^ t | (t >> 4);
t = x & mv3; x = x ^ t | (t >> 8);
t = x & mv4; x = x ^ t | (t >> 16);

```

Таким образом, вычисления в цикле состоят из 21 команды (загрузка констант выносятся за пределы цикла), что является значительным улучшением по сравнению со 127 командами, необходимыми в случае полной подпрограммы из листинга 7.9.

В ситуации 2, в которой известно значение m , можно выполнить примерно те же действия; возможна и дальнейшая оптимизация. Так, может оказаться, что одна из пяти масок равна 0, и в этом случае можно опустить одну из приведенных выше строк. Например, если нет битов, которые должны быть перемещены на нечетное число позиций, маска $m1$ равна 0; если нет битов, которые перемещаются более чем на 15 позиций, равной 0 окажется маска $m4$.

Так, для

$m = 0101\ 0101\ 0101\ 0101\ 0101\ 0101\ 0101\ 0101$

вычисленные маски равны

```

mv0 = 0100 0100 0100 0100 0100 0100 0100 0100
mv1 = 0011 0000 0011 0000 0011 0000 0011 0000
mv2 = 0000 1111 0000 0000 0000 1111 0000 0000
mv3 = 0000 0000 1111 1111 0000 0000 0000 0000
mv4 = 0000 0000 0000 0000 0000 0000 0000 0000
    
```

Поскольку последняя маска равна 0, в скомпилированном коде операция сжатия выполняется за 17 команд (не считая команд для загрузки масок). Конечно, имеются еще более специализированные случаи, в частности сжатие чередующихся битов на с. 167 выполняется всего за 13 команд, не считая загрузки масок.

Использование команд вставки и извлечения

Если ваш компьютер имеет команду *вставки* (insert), причем предпочтительно с непосредственно задаваемыми значениями в качестве операндов, которые определяют битовое поле в целевом регистре, то она может использоваться для выполнения операции *сжатия* с меньшим числом команд, чем в описанном выше методе. Кроме того, при этом не накладываются ограничения на использование регистров, хранящих маски.

Целевой регистр инициализируется нулевым значением, после чего для каждой непрерывной группы единичных битов в маске m переменная x сдвигается вправо для выравнивания очередного поля и используется команда *вставки*, которая вставляет биты x в соответствующее место целевого регистра. Таким образом, операция сжатия выполняется с помощью $2n+1$ команд, где n — количество полей (групп единичных битов) в маске. В наихудшем случае требуется выполнить 33 команды, так как максимальное число полей равно 16 (случай чередования единичных и нулевых битов).

В качестве примера, когда метод с использованием команды вставки дает существенный выигрыш, рассмотрим $m=0x0010084A$. Сжатие с такой маской требует перемещения битов на 1, 2, 4, 8 и 16 позиций. Следовательно, при использовании метода параллельного суффикса требуется выполнение 21 команды, в то время как при использовании *вставки* достаточно 11 команд (в маске имеется пять полей). Еще более удачным примером может служить маска $m=0x80000000$. Здесь происходит перемещение одного бита на 31 позицию, что требует выполнения 21 команды в методе параллельного суффикса, только трех команд при использовании вставки и только одной команды (*сдвиг вправо на 31 бит*), если вы не ограничены определенной схемой вычислений.

Можно также воспользоваться командой *извлечения* (extract) для реализации операции сжатия посредством выполнения $3n-2$ команд, где n — количество полей в маске.

Совершенно ясно, что создание оптимального кода, реализующего сжатие для заранее известной маски, представляет собой далеко не тривиальную задачу.

Сжатие влево

Очевидным путем решения этой задачи является реверс аргумента x и маски m , выполнение сжатия вправо и реверс полученного результата. Еще один способ решения состоит в сжатии вправо с последующим сдвигом влево на $\text{pop}(\overline{m})$ позиций. Эти решения наиболее приемлемы, если ваш компьютер имеет команды реверса или подсчета количества единичных битов. Если же таких команд в наборе компьютера нет, можно легко адаптировать алгоритм из листинга 7.9, просто изменив направление всех сдвигов на обратное (за исключением двух сдвигов в выражениях $l \ll i$; всего — восемь изменений).

На машине БЭСМ-6 (ок. 1967) имелась команда для сжатия влево (“упаковка битов, маскированных X ”) и обратная к ней (“распаковка...”), работавшая с 48-битовыми регистрами этой машины. Эти команды не так просто реализовать. Эксперты в области криптографии предполагают, что их единственное применение — для взлома шифров США. Кроме того, БЭСМ-6 оснащена командой подсчета количества единичных битов, которая, как отмечалось, важна для секретных служб.

7.5. Расширение, или Обобщенная вставка

Обратная к *сжатию вправо* функция перемещает биты из младших разрядов в позиции, задаваемые маской, сохраняя их относительный порядок. Например, $\text{expand}(0000abcd, 10011010) = a00bc0d0$. Таким образом,

$$\text{compress}(\text{expand}(x, m), m) = x$$

Эта функция иногда называется *unpack* (распаковка), *scatter* (разброс) или *deposit* (отложение).

Ее можно получить путем выполнения кода из листинга 7.9 в обратном направлении [3]. Чтобы избежать перезаписывания битов в x , необходимо первыми перемещать влево биты, перемещаемые на большие расстояния, а последними — биты, перемещаемые на одну позицию. Это означает, что первые пять величин “перемещений” (mv в коде) должны быть вычислены, сохранены, а затем использованы в порядке, обратном порядку вычисления. Для многих приложений это не представляет проблемы, так как эти приложения применяют одну и ту же маску m к большому количеству данных, так что эти приложения могут вычислить величины перемещений заранее, а затем повторно их использовать.

Соответствующий код показан в листинге 7.11. Он выполняется примерно за 168 базовых команд RISC (константное значение, не зависящее от входных данных), включая пять сохранений и пять загрузок. 64-битовая версия для 64-разрядных машин будет выполняться примерно за 200 команд.

Листинг 7.11. Метод параллельного суффикса для операции *expand*

```
unsigned expand(unsigned x, unsigned m)
{
    unsigned m0, mk, mp, mv, t;
```

```

unsigned array[5];
int i;

m0 = m;          // Сохранение исходной маски
mk = ~m << 1;   // Подсчет нулей справа

for (i = 0; i < 5; i++)
{
    mp = mk ^ (mk << 1);    // Параллельный суффикс
    mp = mp ^ (mp << 2);
    mp = mp ^ (mp << 4);
    mp = mp ^ (mp << 8);
    mp = mp ^ (mp << 16);
    mv = mp & m;           // Перемешаемые биты
    array[i] = mv;
    m = (m ^ mv) | (mv >> (1 << i)); // Сжатие m
    mk = mk & ~mp;
}

for (i = 4; i >= 0; i--)
{
    mv = array[i];
    t = x << (1 << i);
    x = (x & ~mv) | (t & mv);
}

return x & m0;    // Очистка посторонних битов
}

```

В случае машины, не оснащенной командой *и-не*, мультиплексорная операция во втором цикле может быть закодирована следующим образом.

```
x = ((x ^ t) & mv) ^ x;
```

7.6. Аппаратные алгоритмы сжатия и расширения

В этом разделе рассматриваются ориентированные на аппаратную реализацию алгоритмы функции *сжатия вправо* и обратной к ней [114]. Подобно алгоритмам из предыдущих разделов время их работы пропорционально логарифму размера машинного слова. Они более подходят для аппаратной реализации, но не приводят к быстрому коду при реализации с применением базового набора команд RISC. Мы просто опишем их работу, не обращая ни к коду на языке программирования C, ни к машинному коду.

Сжатие

Чтобы проиллюстрировать работу алгоритма, представим каждый бит *x* буквой и рассмотрим конкретную маску *m*, приведенную ниже.

```

Вход  x = abcd efgh ijkl mnopqrst uvwx yzAB CDEF
Маска m = 0111 1110 0110 1100 1010 1111 0011 0010

```

Алгоритм работает в $\log_2(W)$ “этапов”, где *W* — размер машинного слова в битах. Каждый этап работает параллельно с “карманами” по 2^n бит, где *n* находится в диапазо-

не от 1 до $\log_2(W)$. В конце каждого этапа каждый карман x содержит исходный карман x , биты которого, выбранные маской m , сжаты вправо. Каждый карман m будет содержать целое число, представляющее собой количество нулевых битов в соответствующем кармане исходного значения m . Оно равно количеству битов x , не сжатых вправо. Они соответствуют старшим нулевым битам в кармане x .

На каждом этапе алгоритм выполняет следующие действия параллельно над каждым карманом x и m , где w — размер кармана в битах.

1. Установить L равным левой половине кармана x , дополненной $w/2$ нулевыми битами справа.
2. Сдвинуть L (все w бит) вправо на величину, заданную в правой половине соответствующего кармана m , вставив слева нулевые биты. Никакие единицы не выйдут за правую границу, поскольку максимальная величина сдвига равна $w/2$.
3. Установить $R = w/2$ нулевых битов, за которыми следует правая половина кармана x .
4. Заменить весь w -битовый карман x результатом операции или над R и сдвинутым L .
5. Сложить левую и правую половины кармана m и заменить весь карман суммой.

Чтобы применить эти шаги к первому этапу ($w = 2$), сначала нужно применить операцию *и* к x и m , чтобы сбросить не имеющие значения биты x , и дополнить m , так чтобы каждый бит m представлял собой число нулевых битов в каждом однобитовом полукармане. Проще всего сделать исключение для первого этапа и объединить эти шаги с первой операцией сжатия, применяя логику, показанную в таблице ниже, к каждому двухбитовому карману x и m .

Вход		Выход	
x	m	x	m
ab	00	00	10
ab	01	0b	01
ab	10	0a	01
ab	11	ab	00

Например, в третьей строке $m = 10$ (бинарное). Это означает, что левый бит x выбран как часть результата, а правый — нет. Таким образом, левый бит (a) сжимается вправо. Другой бит x сбрасывается, что гарантирует, что в конечном итоге все старшие (невывбранные) биты будут равны нулю.

Применение этой логики к исходным x и m дает следующее.

Пары бит, $x = 0bcd\ ef0g\ 0j0k\ mn00\ 0q0s\ uvwx\ 00AB\ 000E$
 $m = 0100\ 0001\ 0101\ 0010\ 0101\ 0000\ 1000\ 1001$

На втором этапе рассмотрим, например, второй полубайт ($ef0g$). Образуются значения $L = ef00$ и $R = 000g$. L сдвигается вправо на одну позицию (определяется правой половиной полубайта m), что дает нам $0ef0$. Затем выполняем операцию *или* этого зна-

чения с R , получая $0efg$ в качестве нового значения полубайта. Левая и правая половины m суммируются и дают 0001 (изменений нет).

Полубайты, $x = 0bcd\ 0efg\ 00jk\ 00mn\ 00qs\ uvwx\ 00AB\ 000E$
 $m = 0001\ 0001\ 0010\ 0010\ 0010\ 0000\ 0010\ 0011$

Аналогично третий, четвертый и пятый этапы сжимают байты, полуслова и слово x и обновляют значение m следующим образом.

Байты, $x = 00bc\ defg\ 0000\ jkmn\ 00qs\ uvwx\ 0000\ 0ABE$
 $m = 0000\ 0010\ 0000\ 0100\ 0000\ 0010\ 0000\ 0101$

Полуслова, $x = 0000\ 00bc\ defg\ jkmn\ 0000\ 000q\ suvw\ xABE$
 $m = 0000\ 0000\ 0000\ 0110\ 0000\ 0000\ 0000\ 0111$

Слова, $x = 0000\ 0000\ 0000\ 0bcd\ efgh\ kmnq\ suvw\ xABE$
 $m = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1101$

По завершении значение m дает целое число, указывающее количество ведущих нулевых битов x . Вычитание этого значения из размера слова дает количество сжатых битов в x , равное количеству единичных битов в исходном значении маски m .

Причина, по которой это не самый подходящий алгоритм для реализации с помощью базового набора команд RISC, заключается в трудности выполнения сдвигов полукарманов вправо на разные количества разрядов. С другой стороны, он может оказаться полезным на SIMD-машине, оснащенной командами для параллельной и независимой работы с карманами слова.

Расширение

Алгоритм аппаратного сжатия можно преобразовать в алгоритм расширения, по сути, запуская его сначала в прямом, а затем в обратном направлении. Как и в алгоритме, основанном на методе параллельного суффикса, выполняется вычисление пяти масок алгоритма аппаратного сжатия, они сохраняются, а затем используются в порядке, обратном порядку их вычисления. В действительности последняя маска не используется (она не используется и в алгоритме сжатия), но требуется дополнительная маска $m0$, которая представляет собой дополнение исходной маски. При прямом проходе следует выполнить только шаги, вычисляющие маски; шаги, в которых используется значение x , можно опустить.

Для иллюстрации предположим, что имеется следующее.

Вход $x = abcd\ efgh\ ijkl\ mnop\ qrst\ uvwx\ yzAB\ CDEF$
 Маска $m = 0111\ 1110\ 0110\ 1100\ 1010\ 1111\ 0011\ 0010$

В таком случае результат расширения имеет такой вид.

$0nop\ qrs0\ 0tu0\ vw00\ x0y0\ zABC\ 00DE\ 00F0.$

Вычисляемые маски показаны ниже.

$m0 = 1000\ 0001\ 1001\ 0011\ 0101\ 0000\ 1100\ 1101$
 $m1 = 0100\ 0001\ 0101\ 0010\ 0101\ 0000\ 1000\ 1001$
 $m2 = 0001\ 0001\ 0010\ 0010\ 0010\ 0000\ 0010\ 0011$

$m3 = 0000\ 0010\ 0000\ 0100\ 0000\ 0010\ 0000\ 0101$
 $m4 = 0000\ 0000\ 0000\ 0110\ 0000\ 0000\ 0000\ 0111$

Целочисленные значения в каждой половине $m4$ дают количества нулевых битов в соответствующей половине исходной маски m . В частности, в правой половине m имеется семь нулевых битов. Это означает, что семь старших битов правой половины x ей не принадлежат и должны находиться в левой половине x . Таким образом, биты x с 9 по 15 должны быть сдвинуты влево на расстояние, достаточное для размещения их в левой половине x , а чтобы принять их, старшие биты x должны быть соответствующим образом сдвинуты влево. Этого можно достичь, сдвигая влево все 32-битовое слово x на семь позиций и заменяя левую половину x левой половиной сдвинутой величины. Это даст следующее.

$x = hijk\ lmno\ pqrs\ tuvw\ qrst\ uvwx\ yzAB\ CDEF.$

В общем случае алгоритм работает с размерами карманов от 32 до 2, в пять этапов, используя маски с $m4$ по $m0$. Каждый карман (параллельно) сдвигается влево, при этом биты, выходящие за левую границу, отбрасываются, а пустые позиции справа заполняются нулями, так что величина после сдвига имеет ту же длину, что и исходный карман. Затем левая половина кармана заменяется левой половиной величины, полученной в результате сдвига. При этом в обеих половинах кармана остается "мусор". Он будет обнулен после последнего этапа путем операции \wedge с исходной маской.

Продолжая работу, рассматриваем $m3$ как два 16-битовых кармана. Левый карман имеет в правой половине значение 4, так что левый карман x сдвигается влево на четыре позиции (что дает нам $lmno\ pqrs\ tuvw\ 0000$), и левая половина получившейся величины заменяет левую половину левого кармана x . Выполнение тех же операций над правым 16-битовым карманом x дает следующее.

$x = lmno\ pqrs\ pqrs\ tuvw\ vwxy\ zABC\ yzAB\ CDEF$

На следующем этапе используется константа $m2$, состоящая из четырех 8-битовых карманов. Применив ее к x , получим

$x = mnop\ pqrs\ rstu\ tuvw\ vwxy\ zABC\ BCDE\ CDEF$

На следующем этапе используется константа $m1$, состоящая из восьми 4-битовых карманов. Применив ее к x , получим

$x = mnop\ qrrs\ sttu\ vwvw\ wxyx\ zABC\ BCDE\ DEEF$

На последнем этапе используется константа $m0$, состоящая из шестнадцати 2-битовых карманов. Применив ее к x , получим

$x = mnop\ qrss\ stuu\ vwvw\ xxyy\ zABC\ CCDE\ EEFF$

Последний шаг состоит в выполнении операции \wedge с исходной маской для сброса "ненужных" битов. Это приводит к следующему результату.

$x = 0nop\ qrs0\ 0tu0\ vw00\ x0y0\ zABC\ 00DE\ 00F0$

Полукарманы каждой вычисляемой маски содержат количество нулевых битов в соответствующем полукармане исходной маски m . Поэтому в качестве альтернативы вы-

числения и сохранения масок машина может использовать схемы для вычисления количества нулевых битов в полукарманах “на лету”.

7.7. Обобщенные перестановки

При выполнении обобщенной перестановки битов в слове (или где-либо еще) главной проблемой становится представление перестановок. Очень компактно представить перестановку невозможно. Поскольку в 32-битовом слове можно выполнить 32! перестановки, для указания одной из них требуется как минимум $\lceil \log_2(32!) \rceil = 118$ бит, или 3 слова и еще 22 бита.

Один интересный способ представления перестановок тесно связан с операцией сжатия, рассматриваемой в разделе 7.4 [37]. Начнем с непосредственного метода, в котором для каждого бита просто указывается позиция его перемещения. Например, для перестановки, которая выполняется путем циклического сдвига влево на четыре бита, бит в позиции 0 (младший бит) переместится в позицию 4, в позиции 1 — в позицию 5, ..., в позиции 31 — в позицию 3. Такая перестановка может быть представлена следующим вектором из 32 пятибитовых индексов.

```
00100
00101
...
11111
00000
00001
00010
00011
```

Рассматривая это представление как битовую матрицу, транспонируем ее так, чтобы при этом верхняя строка содержала младшие биты, а результат представим в формате с прямой нумерацией битов. Таким образом будет получен массив из пяти 32-битовых слов.

```
p[0] = 1010 1010 1010 1010 1010 1010 1010 1010
p[1] = 1100 1100 1100 1100 1100 1100 1100 1100
p[2] = 0000 1111 0000 1111 0000 1111 0000 1111
p[3] = 0000 1111 1111 0000 0000 1111 1111 0000
p[4] = 0000 1111 1111 1111 1111 0000 0000 0000
```

Каждый бит в $p[0]$ представляет собой младший бит позиции, в которую переносится соответствующий бит исходного слова x , биты в $p[1]$ представляют собой следующий значащий бит и т.д. Это очень похоже на то, как маска m в предыдущем разделе кодировалась значениями $m \cdot v$, только в алгоритме сжатия эти значения применялись к обновляемому значению маски, а не к исходному.

Операция сжатия, интересующая нас, должна сжимать влево все биты, помеченные в маске единицами, и вправо — биты, помеченные нулями². Иногда эту операцию называют операцией разделения “агнцев и козлиц” (sheep and goats — SAG) или “обобщен-

² Если используется формат с обратной нумерацией битов, в котором младший бит находится справа, влево сжимаются биты, помеченные нулями, а помеченные единицами сжимаются вправо.

ным упорядочением" (general unshuffle). Вычислить эту операцию можно по следующей формуле.

$$\text{SAG}(x, m) = \text{compress_left}(x, m) \mid \text{compress}(x, \sim m)$$

Используя SAG в качестве базовой операции, а перестановку p — как описано выше, биты слова x можно переставить в нужном порядке за 15 шагов.

```

x   = SAG(x,   p[0]);
p[1] = SAG(p[1], p[0]);
p[2] = SAG(p[2], p[0]);
p[3] = SAG(p[3], p[0]);
p[4] = SAG(p[4], p[0]);

x   = SAG(x,   p[1]);
p[2] = SAG(p[2], p[1]);
p[3] = SAG(p[3], p[1]);
p[4] = SAG(p[4], p[1]);

x   = SAG(x,   p[2]);
p[3] = SAG(p[3], p[2]);
p[4] = SAG(p[4], p[2]);

x   = SAG(x,   p[3]);
p[4] = SAG(p[4], p[3]);

x   = SAG(x,   p[4]);

```

Здесь операция SAG используется для выполнения устойчивой двоичной поразрядной сортировки (устойчивость означает сохранение порядка следования равных элементов при сортировке). Массив p используется как 32 пятибитовых ключа для сортировки битов x . На первом шаге все биты x , для которых $p[0] = 1$, перемещаются в левую половину результирующего слова, а биты, для которых $p[0] = 0$, — в правую. Порядок битов при этом не изменяется (т.е. сортировка устойчива). Затем аналогично сортируются все ключи, которые будут использоваться в следующих шагах алгоритма. Очередная, шестая, строка сортирует x на основе вторых младших битов и т.д.

Подобно сжатию, если некоторая перестановка p используется с несколькими словами x , можно значительно сэкономить выполняемые команды путем предварительного вычисления большинства перечисленных шагов алгоритма. Массив перестановки предварительно преобразуется в следующий.

```

p[1] = SAG(p[1], p[0]);
p[2] = SAG(SAG(p[2], p[0]), p[1]);
p[3] = SAG(SAG(SAG(p[3], p[0]), p[1]), p[2]);
p[4] = SAG(SAG(SAG(SAG(p[4], p[0]), p[1]), p[2]), p[3]);

```

После этого каждая перестановка осуществляется с помощью следующих действий.

```

x = SAG(x, p[0]);
x = SAG(x, p[1]);
x = SAG(x, p[2]);
x = SAG(x, p[3]);
x = SAG(x, p[4]);

```

Более прямой (хотя, пожалуй, менее интересный) способ осуществления обобщенных перестановок битов в слове состоит в представлении перестановки как последовательности 32 пятибитовых индексов. k -й индекс представляет собой номер бита в исходном слове, который будет помещен в k -ю позицию при перестановке. (Это список “откуда”, в то время как метод SAG использует список “куда”.) Эти индексы могут быть упакованы по 6 в 32-битовое слово; таким образом, для хранения всех 32 индексов требуется шесть слов. Аппаратно команда может быть реализована как

`bitgather Rt, Rx, Ri`

Здесь Rt — целевой регистр (а также источник), регистр Rx содержит переставляемые биты, а в регистре Ri содержится шесть пятибитовых индексов (и два неиспользуемых бита). Ниже приведена операция, выполняемая командой.

$$t \leftarrow (t \ll 6) | x_0 x_1 x_2 x_3 x_4 x_5$$

Говоря обычным языком, содержимое целевого регистра сдвигается влево на шесть позиций и из слова x выбираются шесть битов и помещаются в шесть освобожденных позиций t . Выбираемые биты определяются шестью пятибитовыми индексами в слове i , берущимися в порядке слева направо. Нумерация битов в индексах может быть как прямой, так и обратной и должна работать, как описано, на любом типе машин.

Для перестановки слова используется последовательность из шести таких команд, все с одинаковыми Rt и Rx , но с разными индексными регистрами. В первом индексном регистре последовательности значащими являются только индексы i_4 и i_5 , а биты, выбираемые остальными четырьмя индексами, просто выдвигаются за левую границу Rt .

Реализация этой команды, скорее всего, будет позволять индексным значениям повторяться, так что эта команда может использоваться не только для перестановки битов. Она может применяться для многократного повторения любого выбранного бита в целевом регистре. Операции SAG такая универсальность не свойственна.

Не слишком сложно реализовать эту команду как быструю (т.е. выполняемую за один такт). Модуль выбора бита состоит из шести 32:1 мультиплексоров. Если они построены на базе пяти сегодняшних мультиплексоров 2:1 (всего $6 \cdot 31 = 186$ мультиплексоров), то данная команда может выполняться быстрее, чем 32-битовая команда сложения [83].

Некоторые машины Intel имеют команды, работающие очень схоже с описанной командой перестановки, но переставляет байты, “слова” (16 бит) и “двойные слова” (32 бита). Это команды PSHUFB, PSHUFW и PSHUFD (Shuffle Packed Bytes/Words/Doublewords — тасование упакованных байтов/слов/двойных слов).

Перестановка битов применяется в криптографии, а для компьютерной графики характерна тесно связанная с перестановкой битов операция перестановки подслов (например, перестановка байтов в пределах слова). Оба эти применения работают чаще с 64-битовыми словами, иногда даже со 128-битовыми, чем с 32-битовыми. Методы SAG и *bitgather* легко модифицировать для работы со словами указанной длины.

Для шифрования и дешифрования сообщения с помощью алгоритма стандарта шифрования данных (Data Encryption Standard — DES) требуется большое количество отображений, аналогичных перестановкам. Сначала выполняется генерация ключа, которая содержит

17 перестановочных отображений. Первое отображает 64-битовую величину в 56-битовую (выбирая из ключа 56 битов, не являющихся битами четности, и переставляя их). Затем следуют 16 отображений 56 битов в 48 битов, использующих то же самое отображение.

Затем, после генерации ключа, каждый блок сообщения, состоящий из 64 бит, подвергается 34 операциям перестановки. Первая и последняя операции представляют собой 64-битовые перестановки, причем одна из них обратна другой. Кроме того, 16 перестановок с повторениями отображают 32-битовые величины на 48-битовые и еще 16 выполняются в пределах 32-битовых слов. Общее количество разных отображений равно шести. Все они являются константами, приведенными в [23].

Впрочем, в настоящее время DES вышел из употребления, после того как в 1998 году Electronic Frontier Foundation была доказана его низкая криптостойкость при наличии специального аппаратного обеспечения. Национальный институт стандартов и технологий (National Institute of Standards and Technology — NIST) в качестве временной замены предложил использовать Triple DES — метод шифрования, в котором к каждому 64-битовому блоку DES применяется трижды, каждый раз с другим ключом (таким образом, длина ключа равна 192 битам, включая 24 бита четности). Такой метод требует соответственно в три раза больше операций перестановки битов по сравнению с обычным DES.

Однако предлагаемая “постоянная” замена DES и Triple DES, стандарт Advanced Encryption Standard [1], совсем не использует перестановки на уровне битов. Ближайшая к перестановкам операция, используемая новым стандартом, — простой циклический сдвиг 32-битовых слов на расстояния, кратные 8. Другие предлагаемые или применяемые методы шифрования используют гораздо меньше битовых перестановок по сравнению с DES.

При сравнении двух рассмотренных методов перестановки следует указать достоинства каждого из них. Метод *bitgather* обладает следующими преимуществами: 1) простота подготовки индексных слов из данных, описывающих перестановку, 2) более простое аппаратное обеспечение, 3) более универсальное отображение. Метод SAG 1) выполняет перестановку за пять, а не шесть команд, 2) использует в команде только два регистра источника (что может играть существенную роль в ряде RISC-архитектур), 3) легче масштабируется для выполнения перестановки в двойных словах, 4) более эффективно выполняет перестановку подслов.

Пункт 3 подробно рассматривается в [81]. Команда SAG позволяет осуществить перестановку двойного слова путем выполнения двух команд SAG, нескольких базовых RISC-команд и двух полных перестановок единичных слов. Команда *bitgather* позволяет сделать то же самое путем выполнения *трех* полных перестановок единичных слов и нескольких базовых RISC-команд. Здесь не учитывается предварительная обработка перестановки, состоящая в получении значений, зависящих только от нее: этот подсчет остается читателю в качестве домашнего задания.

Что касается п. 4, то для перестановки, например, четырех байтов слова с помощью *bitgather* требуется выполнить шесть команд — столько же, сколько при обобщенной перестановке с помощью *bitgather*. Метод SAG позволяет выполнить те же действия за две команды, а не за пять, которые требуются для осуществления обобщенной перестановки с помощью этого метода. Выигрыш в эффективности достигается даже в том слу-

чае, когда размер переставляемых подслов не является степенью 2; для перестановки n подслов требуется $\lceil \log_2 n \rceil$ шагов (в n не входят возможные группы битов, расположенные на концах слова и не участвующие в перестановке).

Кроме команд *SAG* и *bitgather* (которые имеют имена *GRP* и *PPER* соответственно), в [81] рассматриваются и другие возможные команды перестановки, а также перестановки путем поиска в таблицах.

Есть красивый трюк для прибавления 1 к переставленному слову, т.е. для вычисления

$$SAG^{-1}(SAG(x, m) + 1, m)$$

без использования функции *SAG* или обратной к ней [73]. Здесь предполагается, что *SAG(x, m)* переносит все нули вправо и что суммирование не приводит к переполнению в левой части. Разработка этого трюка остается читателю в качестве упражнения.

7.8. Перегруппировки и преобразования индексов

Многие простые перегруппировки битов в компьютерном слове соответствуют еще более простым преобразованиям координат, или индексов битов [37]. Эти соответствия применимы к перегруппировкам элементов в любом одномерном массиве, в котором количество элементов представляет собой целую степень 2. В контексте практического программирования наиболее важен случай, когда элементы массива представляют собой слова (или превосходят их по размеру).

Рассмотрим в качестве примера внешнее идеальное перемешивание элементов массива A размером 8, дающее в результате массив B и состоящее из следующих шагов.

$$\begin{array}{cccc} A_0 \rightarrow B_0; & A_1 \rightarrow B_2; & A_2 \rightarrow B_4; & A_3 \rightarrow B_6; \\ A_4 \rightarrow B_1; & A_5 \rightarrow B_3; & A_6 \rightarrow B_5; & A_7 \rightarrow B_7; \end{array}$$

Каждый B -индекс представляет собой соответствующий A -индекс, циклически сдвинутый влево на одну позицию (при использовании циклического сдвига трех битов). Обратный внешнему идеальному перемешиванию процесс получается путем циклического сдвига каждого индекса *вправо*. Ряд подобных соответствий приведен в табл. 7.1. Здесь n — количество элементов массива, а циклические сдвиги производятся над $\log_2 n$ битами.

ТАБЛИЦА 7.1. ПЕРЕГРУППИРОВКИ И ПРЕОБРАЗОВАНИЯ ИНДЕКСОВ

Перегруппировка	Преобразование индексов	
	Индекс массива, или обратная нумерация битов	Прямая нумерация битов
Реверс	Дополнение	Дополнение
Обобщенный реверс (см. раздел "Обобщенный реверс битов" на с. 160)	Исключающее или с константой	Исключающее или с константой

Перегруппировка	Преобразование индексов	
	Индекс массива, или обратная нумерация битов	Прямая нумерация битов
Циклический сдвиг влево на k позиций	Вычитание $k \pmod n$	Прибавление $k \pmod n$
Циклический сдвиг вправо на k позиций	Прибавление $k \pmod n$	Вычитание $k \pmod n$
Внешнее идеальное перемешивание	Циклический сдвиг влево на одну позицию	Циклический сдвиг вправо на одну позицию
Операция, обратная внешнему идеальному перемешиванию	Циклический сдвиг вправо на одну позицию	Циклический сдвиг влево на одну позицию
Внутреннее идеальное перемешивание	Циклический сдвиг влево на одну позицию, затем дополнение младшего бита	Дополнение младшего бита, затем циклический сдвиг вправо на одну позицию
Операция, обратная внутреннему идеальному перемешиванию	Дополнение младшего бита, затем циклический сдвиг вправо на одну позицию	Циклический сдвиг влево на одну позицию, затем дополнение младшего бита
Транспонирование битовой матрицы 8×8 , хранящейся в 64-битовом слове	Циклический сдвиг (влево или вправо) на три позиции	Циклический сдвиг (влево или вправо) на три позиции
БПФ	Реверс битов	Реверс битов

7.9. Алгоритм LRU

Вас никогда не удивляло то, как компьютер отслеживает, какая строка кеша является последней использованной? Здесь мы опишем один такой алгоритм, известный как метод *матрицы обращений* (reference matrix). Этот алгоритм предназначен, в первую очередь, для аппаратной реализации, но может применяться и в программных приложениях.

Мы не будем вдаваться в длительное обсуждение интригующего мира кешей, заметим только, что мы имеем в виду высокоскоростной кеш, буферизующий данные между основной памятью и процессором компьютера. Такие кеши могут получать запросы слов на каждом такте процессора и обычно должны возвращать запрошенные данные в течение такта или двух, так что времени на сколь-нибудь сложный алгоритм просто не остается.

Кеш содержит копии подмножества данных из основной памяти, и рассматриваемая нами задача заключается в том, чтобы решить при обнаружении промаха кеша (т.е. когда запрошено слово по определенному адресу, но соответствующие данные в кеше отсут-

ствуют), какой именно блок (или *строку* на жаргоне кешей) заменить запрашиваемыми данными. В идеальном случае нужно перезаписать строку, к которой дольше всего не будет обращений. Но мы не можем знать будущее, поэтому будем просто угадывать. наилучшее решение для широкого спектра прикладных программ, пожалуй, заключается в применении стратегии *дальше всех неиспользовавшегося* (least recently used — LRU). Эта стратегия заключается в замене той строки, обращений к которой не было дольше всего.

Имеются три варианта кешей: с *прямым отображением* (direct-mapped), *полностью ассоциативные* (fully associative) и *множественно-ассоциативные* (set-associative). В случае кешей с прямым отображением определенные биты адреса загрузки или сохранения непосредственно адресуют конкретную строку кеша. При обнаружении промаха вопрос о том, какую строку заменять, не возникает, — это должна быть адресуемая строка. При этом не требуется ни LRU, ни какая-либо иная стратегия.

В полностью ассоциативном кеше блок из основной памяти может быть размещен в *любой* строке кеша. При выполнении загрузки или сохранения проверяется, не находится ли данный адрес в кеше. Если нет, необходимо заменить содержимое некоторой строки. Машине предоставлена полная свобода в решении вопроса о том, какую строку кеша заменить. При этом использовались различные стратегии (наиболее распространенные — очередь (FIFO), случайный выбор и LRU), и, как уже упоминалось, наиболее подходящей с точки зрения уменьшения количества промахов представляется стратегия LRU. К сожалению, при наличии большого количества строк-претендентов на замену она и наиболее дорогостояща в плане реализации.

Зачастую выбирается множественно-ассоциативная организация кеша, которая представляет собой компромисс между прямым отображением и полной ассоциативностью. Разработчик принимает решение о степени ассоциативности, которая обычно равна 2, 4, 8 или 16. Кеш разделяется на ряд “множеств”, каждое из которых содержит (обычно) 2, 4, 8 или 16 строк. У множества прямая адресация, для которой используются определенные биты адреса загрузки или сохранения, но внутри множества должен выполняться поиск строки. Этот поиск выполняется в основном так же, как и в полностью ассоциативном кеше. Теперь, когда требуется замена строки, алгоритм LRU должен определить только, какая из строк одного множества дольше всех не использовалась, и заменить ее.

Теперь, после этого краткого обзора, мы можем описать метод матрицы обращений. Для иллюстрации предположим, что кеш четырежды множественно-ассоциативен, т.е. имеются четыре строки кеша, среди которых мы должны отслеживать дольше всех неиспользовавшуюся. Кеш может быть полностью ассоциативным и состоять из четырех строк, а может быть множественно-ассоциативным с четырьмя строками в множестве.

Метод матрицы обращений использует квадратную битовую матрицу с размерностью, равной степени ассоциативности (в принципе; позже мы изменим это утверждение). Каждое ассоциативное множество имеет одну такую матрицу. Суть метода заключается в том, что, когда выполняется обращение к строке i , биты в i -й строке матрицы устанавливаются равными 1, после чего биты i -го столбца устанавливаются равными 0. На рис. 7.2 проиллюстрированы изменения в матрице от начального состояния до конфигурации после обращений к строкам 3, 1, 0, 2, 0, 3 и 2 в указанном порядке.

	Начало	3	1	0	2	0	3	2
	0123	0123	0123	0123	0123	0123	0123	0123
Строка 0	0111	0110	0010	0111	0101	0111	0110	0100
1	0011	0010	1011	0011	0001	0001	0000	0000
2	0001	0000	0000	0000	1101	0101	0100	1101
3	0000	1110	1010	0010	0000	0000	1110	1100

Рис. 7.2. Иллюстрация метода матрицы обращений

Каждая матрица имеет строку, содержащую три единицы, две единицы, одну единицу и ни одной единицы. Номер строки, в которой нет единичных битов, и есть номер дольше всех неиспользовавшейся строки. Следующая по времени обращения строка — с одной единицей, и т.д. Когда происходит промах кеша, машина находит строку матрицы со всеми нулями и заменяет соответствующую строку кеша. Затем она помечает эту строку как *последнюю* использованную, устанавливая в соответствующей строке все биты равными 1, а в столбце — равными 0.

Почему это работает? Если обозначить матрицу буквой M , то причина работоспособности данного алгоритма в том, что M_{ij} указывает, была ли строка i использована после строки j . Если $M_{ij} = 1$, то строка i использовалась после строки j , а если $M_{ij} = 0$, то строка i использовалась не позже строки j .

Рассмотрим произвольную матрицу размером 4×4 , в которой произошло обращение к строке 2. Затем матрица изменяется, как показано на рис. 7.3. Установка всех битов i -й строки равными 1 (за исключением элемента на главной диагонали) указывает, что строка i является использовавшейся позже строки j для всех $i \neq j$. Установка битов i -го столбца равными 0 указывает, что строка j использовалась не позже строки i , для всех j . Отношения между строками кеша, отличными от i -й, не изменяются. Когда обращение будет выполнено ко всем строкам кеша, все отношения “использован позже” будут установлены.

	Начало	2
	0123	0123
Строка 0	abcd	ab0d
1	efgh	ef0h
2	ijkl	1101
3	mnp	mn0p

Рис. 7.3. Один шаг метода матрицы обращений

Таким образом, матрица обращений антисимметрична, а все биты на главной диагонали — всегда нулевые. Значит, в кеше достаточно хранить либо элементы матрицы, находящиеся выше главной диагонали, либо элементы ниже нее. Так и делается на практике. Для ассоциативного множества с n строками требуется хранить $n(n-1)/2$ бит; для $n = 4$ требуется 6 бит, для $n = 8$ — 28. Двадцать восемь битов — достаточно большая величина, так что при использовании метода матрицы обращений и точной стратегии LRU степень ассоциативности редко превышает 8. Вместо этого применяют приближенные методы LRU и методы, не использующие стратегию LRU.

В программном обеспечении стратегия LRU, вероятно, имеет смысл реализовать с помощью списка номеров строк (простого вектора или связанного списка). При обращении к строке i в списке выполняется поиск i , после чего значение i помещается на вершину списка. Таким образом, номер дольше всего неиспользовавшейся строки перемещается в конец списка.

Этот метод достаточно медленный при выполнении обращения (из-за перестройки списка), но быстрый при выяснении, какая строка должна быть заменена. Другой метод, с противоположными скоростными характеристиками, заключается в поддержке вектора с длиной, равной степени ассоциативности, в которой в позиции i хранится как адрес, хранящийся в i -й строке кеша, так и ее "возраст" (или, точнее, "новизна"), закодированный целым числом. При обращении к строке i выполняется увеличение единственной переменной, хранящей возраст строки, и получающееся значение сохраняется в векторе в позиции i . Чтобы найти дольше всего неиспользовавшуюся строку, следует выполнить поиск в векторе наименьшего значения возраста. Метод сбрасывает некорректно при переполнении целочисленного значения возраста.

"Возраст" может быть как один для каждого ассоциативного множества, так и общий для всего кеша; можно также использовать аппаратный счетчик тактов.

Метод матрицы обращений может оказаться полезным в программном обеспечении при небольшой степени ассоциативности. Например, предположим, что приложение использует степень ассоциативности, равную 8, и работа выполняется на 64-разрядной машине. Тогда матрица обращений может храниться в одном 64-битовом регистре. Пусть младшие 8 бит регистра хранят строку 0 матрицы, следующие 8 бит — строку 1 и т.д. Тогда при обращении к строке i в байте i регистра все биты должны быть установлены равными 1, а биты $i, i+8, \dots, i+56$ должны быть сброшены в 0. Обозначая регистр через m , эту задачу можно решить следующим образом.

$$m \leftarrow m | (0xFF \ll (8 * i))$$

$$m \leftarrow m \& \sim (0x0101010101010101 \ll i)$$

Для указанных действий требуется выполнить пять или шесть команд плюс несколько команд для загрузки констант. Чтобы найти дольше всего неиспользовавшуюся строку, выполняем поиск нулевого байта (см. раздел 6.1). Преимущество этого метода над прочими программными методами, вкратце описанными выше, в том, что вся работа выполняется в пределах регистра.

Упражнения

1. Поясните работу второй формулы Мебиуса (формула (1) на с. 165).
2. Операция идеального внешнего перемешивания и обратная к ней используют следующие маски.

$$m_0 = 0x22222222$$

$$m_1 = 0x0C0C0C0C$$

$$m_2 = 0x00F000F0$$

$$m_3 = 0x0000FF00$$

Как выглядит формула для общего случая m_i ? Такая формула может быть полезна в ситуациях, когда верхняя граница длины перемешиваемых целых чисел заранее неизвестна, как, например, в приложениях, работающих с большими “многословными” числами.

3. Напишите код функции, подобной функции сжатия в листинге 7.8, но выполняющей операцию расширения.
4. Каково теоретически минимальное количество битов требуется для реализации стратегии LRU в множественно-ассоциативном кеше со степенью n ? Сравните это значение с количеством битов, требующихся для метода матрицы обращений для нескольких небольших значений n .

ГЛАВА 8

УМНОЖЕНИЕ

8.1. Умножение больших чисел

Умножение больших чисел может быть выполнено традиционным школьным способом — “в столбик”. Однако вместо использования массива промежуточных результатов гораздо эффективнее добавлять к произведению каждую новую строку сразу же после ее вычисления.

Если множимое состоит из m слов, множитель — из n слов, то произведение занимает не более $m + n$ слов независимо от того, какое умножение выполняется, — знаковое или беззнаковое.

При использовании метода “в столбик” каждое 32-битовое слово рассматривается как отдельная цифра. Такой метод вполне применим, если у нас есть команда, которая дает 64-битовое произведение двух 32-битовых слов. К сожалению, даже наличие такой команды в вашем компьютере не означает доступности к ней из языка высокого уровня. В действительности многие современные RISC-компьютеры не имеют такой команды именно потому, что она недоступна из языка высокого уровня и поэтому практически не используется. (Другая причина заключается в том, что эта команда — одна из немногих команд, возвращающих результат в двух регистрах.)

Код данной процедуры представлен в листинге 8.1. Здесь в качестве “цифр” используются полуслова. Параметр w содержит результат, а u и v — множимое и множитель соответственно. Каждый параметр представляет собой массив полуслов, в котором первое полуслово ($w[0]$, $u[0]$, $v[0]$) является младшей значащей цифрой (прямой порядок). Параметры m и n представляют собой количество полуслов в массивах u и v соответственно.

ЛИСТИНГ 8.1. Знаковое умножение больших чисел

```
void mulmns(unsigned short w[], unsigned short u[],
            unsigned short v[], int m, int n)
{
    unsigned int k, t, b;
    int i, j;

    for (i = 0; i < m; i++)
    {
        w[i] = 0;
    }
    for (j = 0; j < n; j++)
    {
        k = 0;
        for (i = 0; i < m; i++)
        {
            t = u[i]*v[j] + w[i + j] + k;
            w[i + j] = t; // (то есть, t & 0xFFFF).
        }
    }
}
```

```

        k = t >> 16;
    }
    w[j + m] = k;
}

// Теперь w[] содержит беззнаковое произведение.
// Корректируем вычитанием v*2**16m при u < 0 и
// вычитанием u*2**16n при v < 0
if ((short)u[m - 1] < 0)
{
    b = 0; // Инициализируем заем
    for (j = 0; j < n; j++)
    {
        t = w[j + m] - v[j] - b;
        w[j + m] = t;
        b = t >> 31;
    }
}

if ((short)v[n - 1] < 0)
{
    b = 0;

    for (i = 0; i < m; i++)
    {
        t = w[i + n] - u[i] - b;
        w[i + n] = t;
        b = t >> 31;
    }
}

return;
}

```

Понять код из листинга 8.1 поможет следующая иллюстрация (никакой связи между n и m нет, любое из них может быть больше другого).

$$\begin{array}{ccccccc}
 u_{n-1} & u_{n-2} & \dots & \dots & u_1 & u_0 & \\
 \hline
 & & & x & v_{n-1} & \dots & v_1 & v_0 \\
 \hline
 w_{n+n-1} & w_{n+n-2} & \dots & \dots & \dots & w_1 & w_0 &
 \end{array}$$

Данная процедура следует алгоритму М из [67, раздел 4.3.1], но в качестве рабочего языка использует С и модифицирована для выполнения знакового умножения. Заметим, что присвоение t в первой половине листинга 8.1 не может вызвать переполнения, поскольку максимальное значение, которое может быть присвоено t , равно $(2^{16} - 1)^2 + 2(2^{16} - 1) = 2^{32} - 1$.

Умножение больших чисел выполняется проще для беззнаковых операндов. По сути, код в листинге 8.1 выполняет беззнаковое умножение с коррекцией результата, так что, опустив код после трехстрочного комментария, можно получить процедуру беззнакового умножения больших чисел. Беззнаковую версию умножения можно распространить на знаковые операнды тремя способами.

1. Получить абсолютные значения каждого входного операнда, выполнить беззнаковое умножение, а затем изменить знак результата, если сомножители имеют разные знаки.
2. Выполнить умножение с использованием элементарных беззнаковых умножений, за исключением умножения одного из старших полуслов; в этом случае используется умножение "знаковое × беззнаковое" или "знаковое × знаковое".
3. Выполнить беззнаковое умножение, а затем скорректировать полученный результат.

Первый метод для вычисления абсолютного значения требует прохода по всем $m + n$ входным полусловам. Если один из операндов положителен, а второй отрицателен, то для получения дополнения отрицательного операнда и результата требуется проход по $\max(m, n) + m + n$ полусловам. Но гораздо неприятнее изменение значения входного операнда (который может быть передан в функцию по адресу), что может быть неприемлемо для ряда приложений. Как вариант обхода этой ситуации для размещения операндов можно временно выделять дополнительную память либо после их изменения и проведения вычислений восстанавливать их предыдущие значения. Впрочем, эти обходные пути также малопривлекательны.

Второй метод требует трех вариантов элементарных операций умножения (беззнаковое × беззнаковое, беззнаковое × знаковое и знаковое × знаковое), а также знакового расширения промежуточных результатов, что существенно увеличивает время вычислений.

Таким образом, лучше выбрать третий метод. Для того чтобы понять, как именно он работает, рассмотрим умножение двух знаковых целых значений — u и v размером M и N бит соответственно. Вычисления в первой части листинга 8.1 рассматривают u как беззнаковую величину, имеющую значение $u + 2^M u_{M-1}$, где u_{M-1} — знаковый бит u (т.е. $u_{M-1} = 1$, если u отрицательно, и $u_{M-1} = 0$ в противном случае). Аналогично v интерпретируется программой как имеющее значение $v + 2^N v_{N-1}$.

Программа вычисляет произведение этих беззнаковых величин.

$$(u + 2^M u_{M-1})(v + 2^N v_{N-1}) = uv + 2^M u_{M-1}v + 2^N v_{N-1}u + 2^{M+N} u_{M-1}v_{N-1}.$$

Для того чтобы получить требующийся результат (uv), необходимо вычесть из беззнакового произведения величину $2^M u_{M-1}v + 2^N v_{N-1}u$. Вычитать член $2^{M+N} u_{M-1}v_{N-1}$ нет необходимости, поскольку известно, что результат можно выразить с помощью $M + N$ битов, так что нет необходимости вычислять биты старше, чем бит в позиции $M + N - 1$. Описанные вычитания выполняются в листинге 8.1 после трехстрочкового комментария и требуют прохода максимум по $m + n$ полусловам.

Может оказаться соблазнительным использовать приведенную в листинге 8.1 программу, передавая ей в качестве параметров массивы из полных слов. Такая программа будет работать на машине с прямым порядком байтов, но не с обратным. Если хранить массивы в обратном порядке, когда $u[0]$ содержит старшее полуслово (и программа будет соответствующим образом изменена), то при передаче в качестве параметров массивов полных слов программа будет работоспособной на машине с обратным (но не с прямым) порядком байтов.

8.2. Старшее слово 64-битового произведения

Теперь рассмотрим задачу вычисления старших 32 бит произведения двух 32-битовых целых чисел. Это функции из нашего базового набора RISC-команд — *старшее слово знакового умножения* (multiply high signed, mulhs) и *старшее слово беззнакового умножения* (multiply high unsigned, mulhu).

В случае беззнакового умножения хорошо работает алгоритм из первой части листинга 8.1. Перепишем его для частного случая $m = n = 2$ с развернутыми циклами и выполнением очевидных упрощений (изменив параметры на 32-битовые целые числа).

В случае знакового умножения нет необходимости в корректирующем коде из второй половины листинга 8.1. Этот код можно опустить, если правильно подойти к промежуточным результатам, объявляя их как знаковые. Конечный алгоритм приведен в листинге 8.2. При использовании беззнаковой версии просто замените все объявления `int` объявлениями `unsigned`.

Листинг 8.2. Старшее слово результата знакового умножения

```
int mulhs(int u, int v)
{
    unsigned u0, v0, w0;
    int      u1, v1, w1, w2, t;

    u0 = u & 0xFFFF;    u1 = u >> 16;
    v0 = v & 0xFFFF;    v1 = v >> 16;
    w0 = u0*v0;
    t  = u1*v0 + (w0 >> 16);
    w1 = t & 0xFFFF;
    w2 = t >> 16;
    w1 = u0*v1 + w1;
    return u1*v1 + w2 + (w1 >> 16);
}
```

Данный алгоритм требует выполнения 16 базовых RISC-команд в каждой из версий (знаковой и беззнаковой), четыре из которых являются умножениями.

8.3. Преобразование знакового и беззнакового произведений одно в другое

Предположим, что наш компьютер позволяет легко вычислить старшую половину 64-битового произведения двух 32-битовых *беззнаковых* целых чисел, но нам требуется выполнение этой операции для *знаковых* чисел. Можно воспользоваться процедурой из листинга 8.2, но она требует выполнения четырех умножений. Более эффективную процедуру можно найти в [10].

Анализ данного алгоритма представляет собой частный случай анализа преобразованного для умножения знаковых чисел алгоритма Кнута (см. листинг 8.1). Пусть x и y означают два 32-битовых знаковых целых числа, которые требуется перемножить. Компьютер рассматривает x как *беззнаковое* целое число, имеющее значение $x + 2^{32} x_{31}$, где x_{31} — целая величина, равная 1, если x отрицательно, и 0 в противном случае. Аналогично y рассматривается как беззнаковое число, имеющее значение $y + 2^{32} y_{31}$.

Хотя нас интересуют только старшие 32 бита величины xy , компьютер выполняет следующее вычисление.

$$(x + 2^{32} x_{31})(y + 2^{32} y_{31}) = xy + 2^{32}(x_{31}y + y_{31}x) + 2^{64} x_{31}y_{31}$$

Для того чтобы получить интересующий нас результат, необходимо вычесть из полученной величины $2^{32}(x_{31}y + y_{31}x) + 2^{64} x_{31}y_{31}$. Поскольку известно, что результат может быть записан с помощью 64 бит, можно использовать арифметику по модулю 2^{64} . Это означает, что можно спокойно игнорировать последний член и вычислять знаковое произведение, как показано ниже (посредством семи базовых RISC-команд).

$$\begin{aligned} p &\leftarrow \text{mulhu}(x, y) \quad // \text{Команда беззнакового 64-битового умножения} \\ t_1 &\leftarrow (x \gg 31) \& y \quad // t_1 = x_{31}y \\ t_2 &\leftarrow (y \gg 31) \& x \quad // t_2 = y_{31}x \\ p &\leftarrow p - t_1 - t_2 \quad // p - \text{искомый результат} \end{aligned} \tag{1}$$

Беззнаковое произведение из знакового

Обратное преобразование — знакового произведения в беззнаковое — выполняется так же просто. Полученная в результате программа представляет собой (1) с тем отличием, что первая команда заменяется командой *знакового 64-битового умножения*, а последняя операция заменяется операцией $p \leftarrow p + t_1 + t_2$.

8.4. Умножение на константы

Совершенно очевидно, что умножение на константу можно преобразовать в последовательность *сдвигов влево и сложений*. Например, для умножения x на 13 (1101 в двоичном представлении) можно воспользоваться приведенным ниже способом (результат умножения — r).

$$\begin{aligned} t_1 &\leftarrow x \ll 2 \\ t_2 &\leftarrow x \ll 3 \\ r &\leftarrow t_1 + t_2 + x \end{aligned}$$

В этом разделе левые сдвиги обозначают умножения на степени двойки, так что приведенный план вычислений записывается как $r \leftarrow 8x + 4x + x$; это должно ясно показывать, что для вычислений на большинстве компьютеров нам потребуются четыре базовые RISC-команды.

Хотелось бы подчеркнуть, что тема умножения на константы сложнее, чем кажется на первый взгляд. Кроме прочего, при рассмотрении разных вариантов умножения следует учитывать не только количество *сдвигов* и *сложений*, необходимых для реализации того или иного умножения. Для иллюстрации сказанного рассмотрим два варианта умножения на 45 (101101 в двоичном представлении).

$t \leftarrow 4x$	$t_1 \leftarrow 4x$
$r \leftarrow x + t$	$t_2 \leftarrow 8x$
$t \leftarrow 2t$	$t_3 \leftarrow 32x$
$r \leftarrow r + t$	$r \leftarrow t_1 + x$
$t \leftarrow 4t$	$t_3 \leftarrow t_3 + t_2$
$r \leftarrow r + t$	$r \leftarrow r + t_3$

План вычислений, представленный слева, использует переменную t для хранения величины x , сдвинутой на количество позиций, соответствующих единичным битам множителя. Каждое сдвинутое значение получается из сдвинутого перед этим. Такой план обладает следующими преимуществами:

- для работы, кроме регистров для входного и выходного значений, требуется только один регистр;
- все команды, за исключением первых двух, двухадресные;
- величины сдвигов относительно невелики.

Эти же свойства остаются действительными и для других множителей.

В схеме справа сначала выполняются все *сдвиги* с x в качестве операнда. Главное преимущество такого решения — увеличение параллелизма вычислений. На компьютере с достаточными возможностями распараллеливания вычислений на уровне команд правая схема выполняется всего за три такта, в то время как левая схема на машине с неограниченными возможностями распараллеливания вычислений на уровне команд требует четыре такта.

Кроме того, задача поиска минимального количества команд (подразумеваемая под ними команды *сдвига* и *сложения*) для выполнения умножения на константу весьма нетривиальна. Далее будем считать, что множество допустимых команд состоит из *сложения*, *вычитания*, *сдвига влево* на любую константу и *изменения знака*. Кроме того, будем считать, что формат команд — трехадресный. Впрочем, если ограничиться только *сложением* (при этом сдвиг влево осуществляется путем сложения числа с самим собой, суммы с самой собой и т.д.), задача отнюдь не упрощается, как и в случае, если в набор команд добавить команду, комбинирующую *сдвиг* и *сложение* (т.е. такая команда вычисляет $z \leftarrow x + (y \ll n)$). Будем также считать, что нас интересуют только младшие 32 бита произведения.

Первое улучшение к предложенной выше базовой схеме бинарного разложения состоит в использовании *вычитания* для сокращения последовательности вычислений, если множитель содержит группу из трех или более последовательных единичных битов. Например, для умножения на 28 (11100 в двоичном представлении) можно вычислить $32x - 4x$ (три команды) вместо того, чтобы вычислять $16x + 8x + 4x$ за пять команд. На машинах с использованием дополнения до двойки конечный результат остается корректным, даже если промежуточное значение $32x$ вызывает переполнение.

Умножение на константу m по схеме бинарного разложения с использованием только команд *сдвига* и *сложения* требует выполнения

$$2 \operatorname{pop}(m) - 1 - \delta$$

команд, где $\delta = 1$, если m завершается единичным битом (нечетно), и равно 0 в противном случае. Если же в дополнение к перечисленным командам используются команды *вычитания*, то умножение требует выполнения

$$4g(m) + 2s(m) - 1 - \delta$$

команд, где $g(m)$ — количество групп из двух и более последовательных единичных битов в m , а $s(m)$ — количество “одиноких” единичных битов. Значение δ определяется так же, как и ранее.

Для группы размером 2 оба метода равноценны.

Следующее улучшение состоит в рассмотрении групп специализированного вида, состоящих из разделенных одним нулевым битом групп из последовательных единичных битов. Рассмотрим, например, $m = 55$ (110111 в двоичном представлении). Метод групп позволяет вычислить произведение путем выполнения шести команд: $(64x - 16x) + (8x - x)$. Если же вычислять произведение как $64x - 8x - x$, то потребуются выполнение только четырех команд. Аналогично умножение на двоичное число 110111011 требует выполнения шести команд, что иллюстрируется формулой $512x - 64x - 4x - x$.

Приведенные выше формулы дают верхнюю границу количества команд, требующихся для умножения переменной x на данное число m . Другая граница может быть получена исходя из размера m в битах: $n = \lfloor \log_2 m \rfloor + 1$.

ТЕОРЕМА. Умножение переменной x на n -битовую константу m , $m \geq 1$, может быть выполнено не более чем за n команд сложения, вычитания и сдвига влево на некоторое заданное значение.

Доказательство (по индукции по n). Умножение на 1 выполняется за 0 команд, так что для $n = 1$ теорема справедлива. Для $n > 1$ умножение на m (если m заканчивается нулевым битом) можно реализовать путем умножения на число, представляющее собой левые $n-1$ бит числа m (т.е. на $m/2$) за $n-1$ команду, и *сдвига влево* на одну позицию, т.е. всего за n команд.

Если m в двоичном исчислении заканчивается на 01, то mx можно вычислить путем умножения x на число, состоящее из левых $n-2$ бит m (для чего потребуется $n-2$ команд), с последующим *сдвигом* результата *влево* на две позиции и *сложения* с x . Всего для этого требуется выполнение n команд.

Если m в двоичном исчислении заканчивается на 11, то нужно рассмотреть случаи, когда m в двоичном представлении заканчивается на 0011, 0111, 1011 и 1111. Пусть t — результат умножения x на левые $n-4$ бит m . Если m заканчивается на 0011, то $mx = 16t + 2x + x$, для чего требуется выполнение $(n-4) + 4 = n$ команд. Если m закан-

чивается на 0111, то $mx = 16t + 8x - x$, а если на 1111, то $mx = 16t + 16x - x$; в обоих случаях для вычислений требуется выполнить n команд. Последний случай — когда двоичное представление m заканчивается на 1011.

При этом легко показать, что mx можно вычислить за n команд, если m заканчивается на 001011, 011011 или 111011. Остается рассмотреть случай 101011.

Эти рассуждения можно продолжать — и всякий раз “последний случай” будет представлять собой число вида 101010...10101011. В конце концов будет достигнут размер m и нам останется рассмотреть только n -битовое число 101010...10101011, которое содержит $n/2 + 1$ единичных битов. Но ранее было показано, что x на такое число можно умножить за $2(n/2 + 1) - 2 = n$ команд.

Таким образом, например, на 32-разрядном компьютере умножение на любую константу может быть выполнено максимум за 32 команды в соответствии с описанным выше методом.

Проверка показывает, что, когда n четно, n -битовое число 101010...101011 требует выполнения n команд, как и число 1010101...010110 в случае нечетного n , так что указанная граница точная.

Описанная методология не слишком сложна до тех пор, пока мы работаем с ней вручную или встраиваем в алгоритм для использования, например, в компиляторе. Однако такой алгоритм не всегда дает оптимальный код, поскольку иногда возможно дальнейшее его улучшение. Это улучшение может использовать, например, результаты разложения m или промежуточные результаты вычислений. Взглянем еще раз на случай $m = 45$ (101101 в двоичном представлении). Описанный выше метод требует выполнения шести команд, но разложение 45 на множители 5 и 9 дает нам решение с использованием всего четырех команд.

$$t \leftarrow 4x + x$$

$$r \leftarrow 8t + t$$

Разложение может быть скомбинировано с аддитивными методами. Например, умножение на 106 (1101010 в двоичном представлении) при использовании бинарного разложения требует выполнения семи команд, но запись 106 как $7 \cdot 15 + 1$ дает решение, состоящее из пяти команд. Для больших констант наименьшее количество команд, выполняющих умножение, может оказаться существенно меньшим, чем при использовании описанного метода бинарного разложения. Например, значение $m = 0xAAAAAAAB$ при бинарном разложении требует 32 команд, но если записать это значение как $2 \cdot 5 \cdot 17 \cdot 257 \cdot 65537 + 1$, то будет получено решение, состоящее из 10 команд (что, вероятно, нетипично для больших чисел; данное разложение отражает чередование единичных и нулевых битов в двоичном представлении числа).

Непохоже, чтобы существовала простая формула или процедура, определяющая минимальное количество команд *сдвига* и *сложения*, необходимых для умножения на заданную константу m . Практичная процедура поиска приведена в [9], но она не всегда находит минимум. Можно разработать для этой цели реализацию метода исчерпываю-

шего поиска, но он достаточно расточителен в смысле как используемой памяти, так и времени работы. (См., например, дерево на рис. 15 в [67, 4.6.3].)

В [67, раздел 4.6.3] рассматривается тесно связанная с данной задачей проблема вычисления a^m с использованием минимального количества умножений, аналогичная задаче умножения на m с использованием одних только команд сложения.

Упражнения

1. Покажите, что в случае умножения $32 \times 32 \Rightarrow 64$ бит младшие 32 бита произведения одинаковы как в случае, когда операнды представляют собой знаковые целые числа, так и в случае, когда они рассматриваются как беззнаковые целые числа.
2. Покажите, как модифицировать функцию `mulhs` (листинг 8.2) так, чтобы она вычисляла как младшую половину 64-битового произведения, так и старшую половину (покажите только изменения в вычислениях, а не в передаче параметров).
3. Умножение комплексных чисел определяется так.

$$(a + bi)(c + di) = ac - bd + (ad + bc)i$$

Это произведение может быть вычислено с помощью всего лишь трех умножений.¹ Пусть

$$\begin{aligned} p &= ac, \\ q &= bd, \\ r &= (a + b)(c + d). \end{aligned}$$

Тогда произведение комплексных чисел можно получить как

$$p - q + (r - p - q)i,$$

в чем легко может убедиться читатель.

Запрограммируйте аналогичный метод для получения 64-битового произведения двух 32-битовых беззнаковых чисел с помощью только трех команд умножения. Считаем, что команда *умножения* генерирует 32 младших бита произведения двух 32-битовых целых чисел (эти младшие биты одинаковы как при знаковом умножении, так и при беззнаковом).

¹ Утверждается, что это было известно еще Гауссу.

ГЛАВА 9 ЦЕЛОЧИСЛЕННОЕ ДЕЛЕНИЕ

9.1. Предварительные сведения

В этой и следующей главах приводится ряд приемов и алгоритмов “компьютерного деления” целых чисел. В математических формулах используется запись x/y для обозначения обычного деления, $x \div y$ — для знакового компьютерного деления целых чисел (с отсечением дробной части) и $x \dot{\div} y$ — для беззнакового компьютерного деления целых чисел. В коде на языке C x/y означает компьютерное деление, которое может быть как знаковым, так и беззнаковым — в зависимости от его операндов.

Деление представляет собой сложный процесс, и включающие его алгоритмы зачастую не слишком элегантны. Более того, серьезным вопросом является даже то, как именно должно быть определено знаковое деление. Большинство языков высокого уровня и наборов команд компьютера определяют результат целочисленного деления как результат рационального деления с отброшенной дробной частью по направлению к 0. Ниже проиллюстрированы данное и два других возможных определения целочисленного деления (*rem* означает остаток от деления).

	Отсечение	Модуль	Округление к меньшему значению
$7 \div 3$	= 2 rem 1	2 rem 1	2 rem 1
$(-7) \div 3$	= -2 rem -1	-3 rem 2	-3 rem 2
$7 \div (-3)$	= -2 rem 1	-2 rem 1	-3 rem -2
$(-7) \div (-3)$	= 2 rem -1	3 rem 2	2 rem -1

Соотношение *Делимое = Частное * Делитель + Остаток* выполняется при любом определении целочисленного деления. “Модульное” деление определяется как требующее, чтобы остаток от деления был неотрицательным¹, а деление с округлением к меньшему значению требует, чтобы результат представлял собой наибольшее целое число, не превосходящее результата рационального деления (функция “пол”). При положительном делителе эти два метода дают одинаковые результаты. Есть и четвертый, редко используемый вариант целочисленного деления, при котором результат рационального деления округляется до ближайшего целого числа.

¹ Я знаю, что меня будут критиковать за такую классификацию, поскольку из понятия “модуль” не следует понятие “неотрицательный”. Оператор *mod* у Кнута [66] означает остаток от деления с округлением к меньшему значению (*floor*), который имеет отрицательное значение (или равен 0), если делитель отрицателен. Ряд языков программирования определяет функцию *mod* как остаток при отсекающем делении. Однако, например, в теории комплексных чисел модуль есть величина неотрицательная, как и в ряде других разделов математики.

Одно из достоинств модульного деления и деления с округлением к меньшему значению заключается в том, что они упрощают ряд используемых приемов. Например, деление на 2^n можно заменить *знаковым сдвигом вправо* на n позиций, а остаток от деления x на 2^n получается в результате применения операции $\&$ к x и $2^n - 1$. Я также подозреваю, что модульное деление и деление с округлением к меньшему значению чаще дают тот результат, который вас интересует. Например, предположим, что вы пишете программу для вывода графика целочисленной функции, значения которой находятся в диапазоне от $imin$ до $imax$, и хотите, чтобы крайними значениями по оси ординат были наименьшие кратные 10 числа, включающие $imin$ и $imax$. Эти крайние значения равны $(imin + 10) * 10$ и $((imax + 9) + 10) * 10$, если воспользоваться “модульным” делением или делением с округлением к меньшему значению. Если же использовать обычное деление, вам придется писать код наподобие следующего.

```
if (imin >= 0) gmin = (imin/10) * 10;
else          gmin = ((imin - 9)/10) * 10;
if (imax >= 0) gmax = ((imax + 9)/10) * 10;
else          gmax = (imax/10) * 10;
```

Помимо того что “модульное” деление и деление с округлением к меньшему значению практичнее деления с отсечением, следует отметить, что неотрицательный остаток от деления требуется, пожалуй, чаще, чем остаток, который может принимать отрицательные значения.

Сделать выбор между “модульным” делением и делением с округлением к меньшему значению непросто, поскольку они отличаются только в случае отрицательного делителя, что встречается довольно редко. Обращение к существующим языкам программирования высокого уровня в этой ситуации не поможет, поскольку практически везде деление знаковых целых чисел x/y означает деление с отсечением. Некоторые языки программирования при делении целых чисел возвращают число с плавающей точкой, т.е. рациональное. Ничуть не проще ситуация с остатками от деления. В Fortran 90 функция MOD дает остаток от деления с отсечением, а функция MODULO — остаток от деления с округлением к меньшему значению (который может быть отрицательным). Аналогично в Common Lisp и ADA функция REM дает остаток от деления с отсечением, а MOD — остаток от деления с округлением к меньшему значению. В PL/I значение MOD всегда неотрицательно (остаток от модульного деления). В Pascal операция $A \bmod B$ определена только для $B > 0$, а ее результат всегда неотрицателен (т.е. это остаток от модульного деления либо от деления с округлением к меньшему значению).

Как бы то ни было, мы не можем изменить мир, даже если бы и хотели это сделать², так что далее, следуя за большинством, будем использовать обычное (отсекающее) определение деления $x + y$.

Самое ценное свойство отсекающего деления заключается в том, что при $d \neq 0$ выполняются следующие соотношения.

² Впрочем, некоторые пытаются. Язык IBM PL.8 использует модульное деление, а команда деления в машине Кнута MMIX использует деление с округлением к меньшему значению [72].

$$(-n) \div d = n + (-d) = -(n + d).$$

Однако к подобным преобразованиям в программах следует подходить осторожно, так как если n или d представляют собой наибольшее (по модулю) отрицательное число, то $-n$ или $-d$ не могут быть представлены 32 битами. Операция $(-2^{32}) + (-1)$ вызывает переполнение (результат не может быть представлен в виде знаковой величины в дополнительном коде), и на большинстве компьютеров результат окажется неопределенным либо операция — невыполненной.

Знаковое целое (отсекающее) деление связано с обычным рациональным делением следующим соотношением.

$$n \div d = \begin{cases} \lfloor n/d \rfloor, & \text{если } d \neq 0, nd \geq 0 \\ \lceil n/d \rceil, & \text{если } d \neq 0, nd < 0 \end{cases} \quad (1)$$

Беззнаковое целое деление, т.е. деление, при котором и n , и d рассматриваются как беззнаковые целые числа, удовлетворяет верхней части (1).

В дальнейшем рассмотрении темы используется ряд элементарных арифметических свойств, которые приводятся здесь без доказательства. Функции, возвращающие наибольшее целое число, меньшее данного (floor — “пол”), и наименьшее, большее данного (ceil — “потолок”), подробно описаны в [66, 36].

ТЕОРЕМА D1. Для действительного x и целого k справедливы следующие соотношения.

$\lfloor x \rfloor = -\lceil -x \rceil$	$\lceil x \rceil = -\lfloor -x \rfloor$
$x - 1 < \lfloor x \rfloor \leq x$	$x \leq \lceil x \rceil < x + 1$
$\lfloor x \rfloor \leq x < \lfloor x \rfloor + 1$	$\lceil x \rceil - 1 < x \leq \lceil x \rceil$
$x \geq k \Leftrightarrow \lfloor x \rfloor \geq k$	$x \leq k \Leftrightarrow \lceil x \rceil \leq k$
$x > k \Rightarrow \lfloor x \rfloor \geq k$	$x < k \Rightarrow \lceil x \rceil \leq k$
$x \leq k \Rightarrow \lfloor x \rfloor \leq k \Rightarrow x < k + 1$	$x \geq k \Rightarrow \lceil x \rceil \geq k \Rightarrow x > k - 1$
$x < k \Leftrightarrow \lfloor x \rfloor < k$	$x > k \Leftrightarrow \lceil x \rceil > k$

ТЕОРЕМА D2. Для целых n и d ($d > 0$) справедливы следующие соотношения.

$$\left\lfloor \frac{n}{d} \right\rfloor = \left\lfloor \frac{n - d + 1}{d} \right\rfloor \quad \text{и} \quad \left\lceil \frac{n}{d} \right\rceil = \left\lceil \frac{n + d - 1}{d} \right\rceil$$

При $d < 0$ справедливы следующие соотношения.

$$\left\lfloor \frac{n}{d} \right\rfloor = \left\lfloor \frac{n - d - 1}{d} \right\rfloor \quad \text{и} \quad \left\lceil \frac{n}{d} \right\rceil = \left\lceil \frac{n + d + 1}{d} \right\rceil$$

ТЕОРЕМА D3. Для действительного x и целого $d > 0$ справедливы следующие соотношения.

$$\lfloor \lfloor x \rfloor / d \rfloor = \lfloor x/d \rfloor \quad \text{и} \quad \lceil \lceil x \rceil / d \rceil = \lceil x/d \rceil$$

СЛЕДСТВИЕ. Для действительных чисел a и b ($b \neq 0$) и целого $d > 0$ справедливы следующие соотношения.

$$\left\lfloor \left\lfloor \frac{a}{b} \right\rfloor / d \right\rfloor = \left\lfloor \frac{a}{bd} \right\rfloor \quad \text{и} \quad \left\lceil \left\lceil \frac{a}{b} \right\rceil / d \right\rceil = \left\lceil \frac{a}{bd} \right\rceil$$

ТЕОРЕМА D4. Для целых n и d ($d \neq 0$) и действительного x справедливы следующие соотношения.

$$\left\lfloor \frac{n}{d} + x \right\rfloor = \left\lfloor \frac{n}{d} \right\rfloor, \quad \text{если} \quad 0 \leq x < \frac{1}{|d|}, \quad \text{и} \quad \left\lceil \frac{n}{d} + x \right\rceil = \left\lceil \frac{n}{d} \right\rceil, \quad \text{если} \quad -\frac{1}{|d|} < x \leq 0$$

В приведенной ниже теореме $\text{rem}(n, d)$ обозначает остаток от деления n на d . Для отрицательных значений d действует следующее определение: $\text{rem}(n, -d) = \text{rem}(n, d)$, как в случае деления с отсечением и модульного деления. При $n < 0$ функция $\text{rem}(n, d)$ не используется (таким образом, при нашем использовании данной функции ее значения всегда неотрицательны).

ТЕОРЕМА D5. Для $n \geq 0$ и $d \neq 0$ справедливы следующие соотношения.

$$\text{rem}(2n, d) = \begin{cases} 2 \text{rem}(n, d) & \text{или} \\ 2 \text{rem}(n, d) - |d| \end{cases} \quad \text{и} \quad \text{rem}(2n+1, d) = \begin{cases} 2 \text{rem}(n, d) + 1 & \text{или} \\ 2 \text{rem}(n, d) - |d| + 1 \end{cases}$$

(причем каждое из значений больше или равно 0 и меньше $|d|$)

ТЕОРЕМА D6. Для $n \geq 0$ и $d \neq 0$ справедливо следующее соотношение.

$$\text{rem}(2n, 2d) = 2 \text{rem}(n, d)$$

Теоремы D5 и D6 легко доказываются исходя из базового определения остатка, т.е. что для некоторого целого q он удовлетворяет следующему соотношению.

$$n = qd + \text{rem}(n, d), \quad 0 \leq \text{rem}(n, d) < |d|$$

Здесь $n \geq 0$ и $d \neq 0$ (n и d могут быть не целыми, однако мы будем применять наши теоремы только к целым числам).

9.2. Деление больших чисел

Как и в случае умножения, деление больших (“многословных”) чисел можно осуществить традиционным школьным способом “в столбик”. Однако детали реализации такого метода на удивление сложны. В листинге 9.1 приведен алгоритм Кнута [67, раздел 4.3.1], реа-

лизованный на языке С. В его основе лежит беззнаковое деление типа $32^m \div 16 \Rightarrow 32$ (на самом деле частное при такой операции деления имеет длину не более 17 бит).

ЛИСТИНГ 9.1. Беззнаковое деление больших целых чисел

```
int divmnu(unsigned short q[], unsigned short r[],
          const unsigned short u[],
          const unsigned short v[],
          int m, int n)
{
    const unsigned b = 65536; // Основание чисел (16 бит)
    unsigned short *un, *vn; // Нормализованный вид u и v
    unsigned qhat; // Предполагаемая цифра частного
    unsigned rhat; // Остаток
    unsigned p; // Произведение двух цифр
    int s, i, j, t, k;

    if (m < n || n <= 0 || v[n-1] == 0)
        return 1; // Возвращается при некорректном параметре

    if (n == 1) // Частный случай делителя из одной цифры
    {
        k = 0;
        for (j = m - 1; j >= 0; j--)
        {
            q[j] = (k*b + u[j])/v[0];
            k = (k*b + u[j]) - q[j]*v[0];
        }
        if (r != NULL) r[0] = k;
        return 0;
    }

    // Нормализация путем сдвига v влево, такого, что
    // старший бит становится единичным, и сдвига u влево
    // на ту же величину. Нам может потребоваться добавить
    // старшую цифру к частному; мы делаем это безусловно
    s = nlz(v[n-1]) - 16; // 0 <= s <= 16.
    vn = (unsigned short *)alloca(2*n);
    for (i = n - 1; i > 0; i--)
        vn[i] = (v[i] << s) | (v[i-1] >> 16-s);
    vn[0] = v[0] << s;

    un = (unsigned short *)alloca(2*(m + 1));
    un[m] = u[m-1] >> 16-s;
    for (i = m - 1; i > 0; i--)
        un[i] = (u[i] << s) | (u[i-1] >> 16-s);
    un[0] = u[0] << s;
    for (j = m - n; j >= 0; j--) // Главный цикл
    {
        // Вычисляем оценку q[j]
        qhat = (un[j+n]*b + un[j+n-1])/vn[n-1];
        rhat = (un[j+n]*b + un[j+n-1]) - qhat*vn[n-1];
    again:
        if (qhat >= b || qhat*vn[n-2] > b*rhat + un[j+n-2])
        {
            qhat = qhat - 1;
            rhat = rhat + vn[n-1];
            if (rhat < b) goto again;
        }
    }
}
```

```

    }

    // Умножение и вычитание
    k = 0;
    for (i = 0; i < n; i++)
    {
        p = qhat*vn[i];
        t = un[i+j] - k - (p & 0xFFFF);
        un[i+j] = t;
        k = (p >> 16) - (t >> 16);
    }
    t = un[j+n] - k;
    un[j+n] = t;

    q[j] = qhat;      // Сохранение цифры частного
    if (t < 0)        // Если мы вычли слишком много,
    {                 // вернем назад
        q[j] = q[j] - 1;
        k = 0;
        for (i = 0; i < n; i++)
        {
            t = un[i+j] + vn[i] + k;
            un[i+j] = t;
            k = t >> 16;
        }
        un[j+n] = un[j+n] + k;
    }
} // for j

// Если вызывающей функции нужно значение остатка,
// денормализуем и возвращаем его
if (r != NULL)
{
    for (i = 0; i < n; i++)
        r[i] = (un[i] >> s) | (un[i+1] << 16-s);
}
return 0;
}

```

Алгоритм обрабатывает входные и выходные данные по полслова. Конечно, было бы предпочтительнее обрабатывать данные по слову, но такой алгоритм требует наличия команды деления $64 \div 32 \Rightarrow 32$. Однако предполагается, что такой команды в компьютере нет (или она недоступна из языка высокого уровня). Хотя в общем случае предполагается наличие команды деления $32 \div 32 \Rightarrow 32$, для решения нашей задачи достаточно команды $32 \div 16 \Rightarrow 16$.

Таким образом, в данной реализации алгоритма Кнута основание счисления b равно 65536. В [67] вы найдете подробное пояснение данного алгоритма.

Делимое u и делитель v используют прямой порядок, т.е. $u[0]$ и $v[0]$ представляют собой младшие цифры числа (впрочем, данный код корректно работает при использовании как прямого, так обратного порядка). Параметры m и n представляют собой количество полуслов в u и v соответственно (Кнут определяет m как длину частного). Вызывающая функция должна обеспечить память для хранения частного q и (необязательно) для остатка r .

Пространство для частного должно иметь размер как минимум $m-p+1$ полуслов и p полуслов для остатка. Если значение остатка нас не интересует, вместо адреса для его размещения можно передать параметр NULL.

Алгоритм требует, чтобы старшая цифра делителя $v[p-1]$ была ненулевой. Это упрощает нормализацию и помогает убедиться, что вызывающая функция выделила достаточное количество памяти для частного. Код проверяет, является ли значение $v[p-1]$ ненулевым, а также выполнение условий $p \geq 1$ и $m \geq p$. Если любое из этих условий нарушено, возвращается код ошибки (значение 1).

После этих проверок код выполняет деление для частного случая, когда делитель имеет длину 1. Это выделение частного случая предназначено не для ускорения работы: просто оставшаяся часть алгоритма требует, чтобы длина делителя была не менее 2.

Если делитель имеет длину 2 или большую, алгоритм нормализует делитель, сдвигая его влево на величину, достаточную, чтобы старший бит делителя был равен 1. На ту же величину сдвигается и делимое, так что эти сдвиги никак не влияют на величину частного. Как поясняется у Кнута, эти шаги необходимы для того, чтобы облегчить оценку каждой цифры частного с хорошей точностью. Для определения величины сдвига используется функция *количества ведущих нулевых битов* $nlz(x)$.

При выполнении нормализации для делимого и делителя выделяется новая память. Это делается в связи с тем, что с точки зрения вызывающей функции крайне нежелательно изменение входных данных, которые вообще могут оказаться константами, расположенными в памяти, доступной только для чтения. Кроме того, из-за сдвига делимому может понадобиться дополнительное полуслово для хранения старшей цифры. Для выделения этой памяти в C идеально подходит функция `alloca()`, которая обычно выделяет память с помощью двух-трех команд и не требует явного ее освобождения. Эта память выделяется в программном стеке, так что она автоматически освобождается при возврате из функции.

В основном цикле происходит быстрое вычисление цифр частного, по одной цифре за итерацию: делимое при этом уменьшается и в конечном счете становится равным остатку. Оценка значения очередной цифры частного `qhat` после уточнения в цикле, помеченном меткой `again`, всегда оказывается либо точной, либо больше точного значения на 1.

Следующие шаги состоят в умножении `qhat` на делитель и вычитании произведения из полученного остатка, как и в методе деления столбиком. Если остаток оказывается отрицательным, необходимо уменьшить цифру частного на 1 и либо заново вычислить произведение и вычесть его из остатка, либо, что гораздо проще, прибавить к отрицательному значению остатка делитель. Такой шаг делается не более одного раза, поскольку цифра частного либо точна, либо больше точного значения на 1.

Последнее действие состоит в возврате остатка вызывающей функции, если переданный адрес памяти для хранения остатка не NULL. Остаток при этом должен быть сдвинут вправо на ту же величину, на которую в целях нормализации сдвигались влево делитель и делимое.

Шаги, связанные с неверной оценкой цифры частного, выполняются достаточно редко. Чтобы убедиться в этом, заметим, что первое вычисление оценки каждой цифры частного $qhat$ выполняется путем деления двух старших цифр текущего остатка на старшую цифру делителя. Шаги цикла *again* уточняют оценку путем деления *трех* старших цифр текущего остатка на *две* старшие цифры делителя (доказательство опущено; убедиться в корректности этих действий можно, проведя несколько тестовых вычислений с основанием системы счисления $b = 10$). Заметим, что из-за выполненной нормализации делитель имеет значение, не меньшее $b/2$, а делимое не более чем в b раз превышает делитель (поскольку любой остаток меньше делителя).

Насколько точна оценка частного при использовании только трех цифр делимого и двух — делителя? Можно показать, что в связи с выполненной нормализацией такая оценка довольно точна. Чтобы увидеть это в какой-то мере интуитивно (без формального доказательства), рассмотрим оценку u/v в арифметике с основанием счисления 10. Можно показать, что такая оценка всегда несколько завышена (или точна). Следовательно, наихудшим случаем является тот, при котором усечение делителя до двух цифр максимально уменьшает его, а усечение делимого до трех цифр не уменьшает его величину, которая должна быть максимально возможной. Это происходит в случае $49900\dots 0/5099\dots 9$, что приводит к оценке $499/50 = 9.98$. Точный результат составляет примерно $499/51 \approx 9.7843$. Разница в 0.1957 показывает, что оценка цифры частного отличается от истинного значения не более чем на 1, причем это происходит примерно в 20% случаев (в предположении, что цифры частного распределены равномерно). Это, в свою очередь, означает, что дополнительные действия, связанные с неточной оценкой цифры частного, будут выполняться примерно в 20% случаев.

Проведение этого (нестроого) анализа для общего случая системы счисления b приводит к выводу о том, что оценочное и истинное значения цифры частного различаются не более чем на $2/b$. В случае $b = 65536$ мы получим, что оценочное и истинное значения цифры частного отличаются не более чем на 1 и происходит это с вероятностью примерно $2/65536 \approx 0.00003$. Следовательно, дополнительные действия выполняются только для примерно 0.003% цифр частного.

В качестве конкретного примера, когда требуются дополнительные корректирующие действия, в десятичной системе счисления можно привести деление $4500/501$. Аналогичный пример для системы счисления по основанию 65536 — $0x7FFF\ 8000\ 0000\ 0000/0x8000\ 0000\ 0001$.

Не будем пытаться оценить время работы программ в целом, ограничимся лишь замечанием, что для больших m и n время выполнения определяется циклом умножения/вычитания. Хороший компилятор в состоянии использовать для этого цикла всего 16 базовых RISC-команд, одна из которых — *умножение*. Цикл `for j` выполняется $m - n + 1$ раз, а цикл умножения/вычитания — n раз, что дает время выполнения этой части программы, равное $(15 + mul)n(m - n + 1)$ тактов, где *mul* — время умножения двух 16-битовых переменных. Кроме того, в программе выполняется $m - n + 1$ команд *деления* и одна — вычисления количества ведущих нулевых битов.

Знаковое деление больших чисел

Приводить специализированный алгоритм для знакового деления больших чисел нет необходимости, просто укажем, как можно адаптировать для этого алгоритм беззнакового деления больших чисел.

1. Изменить знак у делимого, если оно отрицательно. Сделать то же для делителя.
2. Преобразовать делимое и делитель в беззнаковое представление.
3. Использовать алгоритм беззнакового деления больших чисел.
4. Преобразовать частное и остаток в знаковое представление.
5. Изменить знак частного, если знаки делимого и делителя были различны.
6. Если делимое отрицательно, изменить знак остатка.

Иногда выполнение этих шагов требует добавления или удаления старшего бита. Например, предположим для простоты, что числа представлены в системе счисления с основанием 256 (один байт на цифру) и что при знаковом представлении числа старший бит последовательности цифр является знаковым. Такое представление наиболее похоже на обычное представление чисел в дополнительном коде. В этом случае делитель 255, который в знаковом представлении имеет вид 0x00FF, при выполнении шага 2 должен быть сокращен до 0xFF. Аналогично, если частное из шага 3 начинается с единичного бита, для корректного представления в виде знакового числа нужно добавить к нему ведущий нулевой байт.

9.3. Беззнаковое короткое деление на основе знакового

Под "коротким делением" подразумевается деление одного слова на другое (т.е. деление типа $32 + 32 \Rightarrow 32$). Это обычное деление, которое в языке C и многих других языках программирования высокого уровня выполняется оператором "/" с целыми операндами. В C имеются и знаковое, и беззнаковое короткое деления, но в ряде компьютеров в набор команд входит только знаковое деление. Каким образом можно реализовать беззнаковое деление на такой машине? Непохоже, чтобы такая задача решалась гладко и просто, тем не менее рассмотрим некоторые возможные варианты.

Использование знакового длинного деления

Даже если компьютер оснащен знаковым длинным делением ($64 + 32 \Rightarrow 32$), осуществление беззнакового короткого деления — не такая простая задача, как может показаться на первый взгляд. В компиляторе XLC для IBM RS/6000 $q \leftarrow \left(n \div d \right)$ реализуется следующим образом (запись с использованием псевдокода).

```

if  $n < d$  then  $q \leftarrow 0$ 
else if  $d = 1$  then  $q \leftarrow n$ 
else if  $d \leq 1$  then  $q \leftarrow 1$ 
else  $q \leftarrow (0 \parallel n) \div d$ 

```


Третья строка в действительности проверяет выполнение условия $d \stackrel{n}{\geq} 2^{31}$. Если d в этом месте алгебраически меньше или равно 1, то, поскольку оно не равно 1 (эта проверка выполнялась во второй строке), алгебраически оно должно не превышать 0. Нас не беспокоит случай $d = 0$, поэтому, если проверяемое в третьей строке условие истинно, это означает, что установлен знаковый бит d , т.е. $d \stackrel{n}{\geq} 2^{31}$. Поскольку из первой строки известно, что $n \stackrel{n}{\geq} d$, и n не может превышать $2^{32} - 1$, то $n \stackrel{n}{+} d = 1$.

Обозначение в четвертой строке означает формирование целого числа двойной длины, состоящего из 32 нулевых битов, за которыми следует 32-битовая величина n , и деление его на d . Проверка $d = 1$ во второй строке необходима для того, чтобы гарантировать, что это деление не вызовет переполнения (переполнение может возникнуть при $n \stackrel{n}{\geq} 2^{31}$, и частное в этом случае оказывается неопределенным).

Если объединить сравнения во второй и третьей строках³, то описанный выше алгоритм может быть реализован с помощью 11 команд, три из которых — команды ветвления. Если необходимо выполнение деления при $d = 0$ (завершающееся генерацией прерывания), то третью строку можно заменить строкой “else if $d < 0$ then $q \leftarrow 1$ ”, что приводит к 12-командной реализации на машине RS/6000.

Не так сложно изменить код таким образом, чтобы наиболее вероятные значения ($2 \stackrel{n}{\leq} d < 2^{31}$) не требовали такого большого количества проверок (начав код с проверки if $d \leq 1$), но это приведет к некоторому увеличению объема получаемого кода.

Использование знакового короткого деления

Этот раздел написан для 32-разрядной машины, но применим и для 64-разрядной машины (т.е. для получения беззнакового деления $64 + 64 \Rightarrow 64$ из аналогичного знакового деления) путем замены всех встречающихся 31 на 63. Этот способ может быть применен в языке программирования Java, в котором беззнаковые целые числа отсутствуют.

Если знаковое длинное деление недоступно, но есть знаковое короткое деление, то $n \stackrel{n}{+} d$ можно реализовать приведением задачи к случаю $n, d < 2^{31}$ и использованием машинной команды *деления*. Если $d \stackrel{n}{\geq} 2^{31}$, то $n \stackrel{n}{+} d$ может быть только 0 или 1, так что от этого условия легко освободиться. Далее можно воспользоваться тем фактом, что выражение

$\left(\left(\left(n \stackrel{n}{+} 2 \right) + d \right) \times 2 \right)$ приближенно равно значению $n \stackrel{n}{+} d$ с ошибкой, равной 0 или 1.

Это приводит нас к следующему методу (записан с использованием псевдокода).

1. if $d < 0$ then if $n \stackrel{n}{+} d$ then $q \leftarrow 0$
2. else $q \leftarrow 1$
3. else do

³ Выполнение команды *сравнения* на RS/6000 устанавливает несколько битов состояния, указывающих выполнение отношений “меньше”, “больше” и “равно”.

4. $q \leftarrow \left(\left(\frac{n}{2} + d \right) \times 2 \right)$
5. $r = n - qd$
6. if $r \geq d$ then $q \leftarrow q + 1$
7. end

Проверка $d < 0$ в первой строке на самом деле выясняет, не выполняется ли условие $d \geq 2^{31}$. Если это условие истинно, то наибольшее частное может быть равно только $(2^{31} - 1) \div 2^{31} = 1$, так что первые две строки представляют собой вычисление частного для случая $d \geq 2^{31}$.

Четвертая строка представляет собой код, выполняющий команды *беззнакового сдвига вправо на 1 бит, деления, сдвига влево на 1 бит*. Понятно, что $\frac{n}{2} + 2 < 2^{31}$, и к этому моменту мы убедились, что и $d < 2^{31}$, а потому эти величины могут использоваться в машинной команде знакового деления (если $d = 0$, машина сообщит о переполнении).

Вычисления в четвертой строке дают следующую оценку частного (с использованием следствия из теоремы D3).

$$q = \left\lfloor \left\lfloor \frac{n}{2} \right\rfloor / d \right\rfloor \cdot 2 = \left\lfloor n / (2d) \right\rfloor \cdot 2 = \frac{n - \text{rem}(n, 2d)}{d}$$

В строке 5 вычисляется соответствующий этому частному остаток.

$$r = n - \frac{n - \text{rem}(n, 2d)}{d} \cdot 2 = \text{rem}(n, 2d)$$

Таким образом, $0 \leq r < 2d$. Если $r < d$, то q представляет собой корректный остаток. Если же $r \geq d$, то верное значение частного получается добавлением 1 к вычисленному в строке 4 значению (при выполнении проверки значения остатка программа должна использовать беззнаковое сравнение, так как возможна ситуация, когда $r \geq 2^{31}$).

Загружая командой *загрузки непосредственного значения 0* в q до выполнения сравнения $n < d$ и кодируя присвоение $q \leftarrow 1$ в строке 2 как переход к присвоению $q \leftarrow q + 1$ в строке 6, получим 14 команд на большинстве машин, четыре из которых будут представлять собой команды ветвления. Достаточно просто дополнить данный код для вычисления не только частного, но и остатка от деления. Для этого необходимо добавить к строке 1 присвоение $r \leftarrow n$, к строке 2 — $r \leftarrow n - d$, а к части "then" строки 6 — $r \leftarrow r - d$ (либо, если пойти на использование команды *умножения*, можно обойтись единственным добавлением присвоения $r \leftarrow n - qd$ после выполнения всех вычислений).

Альтернативой строкам 1 и 2 являются следующие строки.

- if $n < d$ then $q \leftarrow 0$
- else if $d < 0$ then $q \leftarrow 1$

Это приводит к более компактному коду, состоящему из 13 команд, три из которых являются командами ветвления. Однако в этом случае для наиболее вероятной ситуации (небольшие величины, $n > d$) будет выполняться большее количество команд.

Используя выражения предикатов, программу можно переписать следующим образом.

1. if $d < 0$ then $q \leftarrow \left(n \overset{\cdot}{\geq} d \right)$
2. else do
3. $q \leftarrow \left(\left(n \overset{\cdot}{+} 2 \right) + d \right) \times 2$
4. $r \leftarrow n - qd$
5. $q \leftarrow q + \left(r \overset{\cdot}{\geq} d \right)$
6. end

Это экономит две команды ветвления, если в компьютере имеется возможность вычисления указанных предикатов без использования ветвления. В случае базового набора команд RISC они могут быть вычислены одной командой (CMPGEU); на машине MIPS требуется две команды (SLTU, XORI). На большинстве компьютеров вычисление каждого из этих предикатов требует выполнения четырех команд (трех при наличии полного набора логических команд) путем использования выражения для $x \overset{\cdot}{\leq} y$ из раздела "Предикаты сравнения" на с. 43 с упрощениями, основанными на том, что в строке 1 известно, что $d_{31} = 1$, а в строке 5 — что $d_{31} = 0$. Это позволяет упростить выражения.

$$\text{В строке 1: } n \overset{\cdot}{\geq} d = (n \& \neg(n-d)) \gg 31$$

$$\text{В строке 5: } r \overset{\cdot}{\geq} d = (r | \neg(r-d)) \gg 31$$

Можно получить код без ветвлений, если считать, что при $d \overset{\cdot}{\geq} 2^{31}$ частное должно быть равно 0. В этом случае делитель можно использовать в команде знакового деления, поскольку при неверном его распознавании как отрицательной величины результат устанавливается равным 0 (что находится в корректных пределах частного — до 1). Случай большого делимого обрабатывается, как и раньше, сдвигом перед делением на одну позицию вправо и сдвигом частного на одну позицию влево. Это дает нам следующую программу (требующую выполнения 10 базовых RISC-команд).

1. $t \leftarrow d \gg 31$
2. $n' \leftarrow n \& \neg t$
3. $q \leftarrow \left(\left(n' \overset{\cdot}{+} 2 \right) + d \right) \times 2$

$$4. \quad r \leftarrow n - qd$$

$$5. \quad q \leftarrow q + \left(r \overset{\circ}{\geq} d \right)$$

9.4. Беззнаковое длинное деление

Под длинным делением подразумевается деление двойного слова на одинарное. В случае 32-битовой машины это деление $64 \overset{\circ}{\div} 32 \Rightarrow 32$ с неопределенным результатом случае переполнения (в частности, при делении на 0).

Некоторые 32-разрядные компьютеры оснащены командами беззнакового длинного деления, однако эти команды малоупотребимы в связи с тем, что из большинства языков программирования высокого уровня доступно только деление вида $32 \overset{\circ}{\div} 32 \Rightarrow 32$. Таким образом, разработчики компьютеров могут предпочесть обеспечить наличие только команды деления формата $32 \overset{\circ}{\div} 32$. Далее будут приведены два алгоритма, которые обеспечивают реализацию отсутствующего в компьютере беззнакового длинного деления.

Аппаратный алгоритм сдвига и вычитания

В качестве первой попытки выполнения длинного деления рассмотрим, как это деление реализуется аппаратно. Существует два широко распространенных алгоритма, называемых *восстанавливающим* (restoring) и *невосстанавливающим* (nonrestoring) делением [56, sec. A-2; 29]. Оба алгоритма представляют собой алгоритмы "сдвига и вычитания". В восстанавливающей версии, показанной ниже, восстанавливающий шаг состоит в прибавлении делителя, когда вычитание дает отрицательный результат. Величины x , y и z хранятся в 32-битовых регистрах. Изначально $x \parallel y$ представляет собой делимое размером в два слова, а делитель располагается в z . Для хранения переполнения при вычитании нам потребуется однобитовый регистр c . Вот как выглядит этот алгоритм.

```
do  $i \leftarrow 1$  to 32
   $c \parallel x \parallel y \leftarrow 2(x \parallel y)$  // Сдвиг влево на один бит
   $c \parallel x \leftarrow (c \parallel x) - (0b0 \parallel z)$  // Вычитание (33 бита)
   $y_0 \leftarrow \neg c$  // Установка одного бита частного
  if  $c$  then  $c \parallel x \leftarrow (c \parallel x) + (0b0 \parallel z)$  // Восстановление
end
```

По окончании вычислений частное находится в регистре y , а остаток — в регистре x .

В случае переполнения этот алгоритм *не* дает никаких полезных результатов. При делении величины $x \parallel y$ на 0 частное представляет собой поразрядное дополнение x до единиц, а остаток — величину y . В частности, $0 \overset{\circ}{\div} 0 \Rightarrow 2^{32} - 1 \text{ gem } 0$. Охарактеризовать прочие случаи переполнения гораздо сложнее.

Было бы неплохо, если бы при ненулевом делителе алгоритм давал корректное значение частного по модулю 2^{32} и корректное значение остатка. Однако, похоже, единственный способ добиться этого — создать регистр длиной свыше 97 бит для представ-

ления $c \parallel x \parallel y$ и выполнить тело цикла 64 раза, что обеспечит выполнение деления вида $62 \div 32 \Rightarrow 64$. Вычитание при этом остается 33-битовой операцией; однако дополнительные аппаратное обеспечение и время выполнения делают такое усовершенствование слишком дорогостоящим.

Реализовать этот алгоритм программно достаточно сложно, поскольку у большинства компьютеров нет 33-битового регистра, который бы хранил значение $c \parallel x$. Тем не менее в листинге 9.2 приведен вариант кода, который в определенной степени следует рассмотренному аппаратному алгоритму.

Листинг 9.2. Беззнаковое длинное деление, алгоритм сдвига и вычитания

```
unsigned divlu(unsigned x, unsigned y, unsigned z)
{
    // Деление (x || y) на z
    int i;
    unsigned t;

    for (i = 1; i <= 32; i++)
    {
        t = (int)x >> 31;           // Все биты = 1,
                                   // если x(31) = 1
        x = (x << 1) | (y >> 31); // Сдвиг x||y влево
        y = y << 1;                // на один бит
        if ((x | t) >= z)
        {
            x = x - z;
            y = y + 1;
        }
    }
    return y;                       // Остаток хранится в x
}
```

Переменная t используется для того, чтобы обеспечить правильное сравнение. После сдвига $x \parallel y$ необходимо выполнить 33-битовое сравнение. Если первый бит x равен 1 (до выполнения сдвига), то совершенно очевидно, что интересующая нас 33-битовая величина больше, чем 32-битовый делитель. В этом случае все биты $x \parallel t$ равны единице, так что сравнение с z дает правильный результат (`true`). Если же первый бит x равен 0, то 32-битового сравнения достаточно.

Код в листинге 9.2 требует выполнения от 321 до 385 базовых RISC-команд, в зависимости от того, насколько часто результат сравнения оказывается истинным. Если компьютер имеет команду *сдвига влево двойного слова*, операция сдвига может быть выполнена с помощью одной команды вместо четырех. Это позволяет снизить общее количество выполняемых команд до 225–289 (в предположении, что для управления циклом требуется выполнение двух команд на итерацию).

Если принять $x = 0$, то алгоритм из листинга 9.2 может быть использован для реализации деления $32 \div 32 \Rightarrow 32$. Единственное упрощение при этом состоит в том, что переменную t можно не использовать, поскольку ее значение всегда равно 0.

Ниже приведен невосстанавливающий аппаратный алгоритм беззнакового деления. Основная идея состоит в том, что после вычитания делителя z из 33-битовой величины,

которую мы обозначаем как $c \parallel x$, нет необходимости добавлять z , если результат оказался отрицательным. Вместо этого в следующей итерации достаточно выполнить сложение вместо вычитания. Это связано с тем, что добавление z (для коррекции ошибки, связанной с вычитанием z в предыдущей итерации), сдвиг влево и вычитание z эквивалентно его добавлению ($2(u+z) - z = 2u + z$). Преимущество данного метода в контексте аппаратной реализации состоит в том, что в каждой итерации цикла выполняется только одна команда сложения или вычитания и сумматор оказывается самой медленной схемой в цикле⁴. В конце алгоритма требуется корректировка остатка, если он отрицателен (соответствующая корректировка частного не требуется).

Делимое представляет собой двойное слово $x \parallel y$, а делитель — слово z . По окончании вычислений частное находится в регистре y , а остаток — в регистре x .

```

c = 0
do i ← 1 to 32
  if c = 0 then do
    c ∥ x ∥ y ← 2(x ∥ y)      // Сдвиг влево на один бит
    c ∥ x ← (c ∥ x) - (0b0 ∥ z) // Вычитание делителя
  end
  else do
    c ∥ x ∥ y ← 2(x ∥ y)      // Сдвиг влево на один бит
    c ∥ x ← (c ∥ x) + (0b0 ∥ z) // Прибавление делителя
  end
  y0 ← ¬c                      // Установка одного бита частного
end
if c = 1 then x ← x + z         // Коррекция отрицательного остатка

```

Похоже, что хорошей адаптации к 32-битовому алгоритму не существует.

Миникомпьютер 801 (ранняя экспериментальная RISC-разработка IBM) имеет команду *пошагового деления*, которая, по сути, выполняет приведенные выше в теле цикла действия. Эта команда использует бит переноса для хранения c и MQ (32-битовый регистр) для хранения y . Для реализации данной команды требуются 33-битовый сумматор и 33-битовое вычитающее устройство. Команда *пошагового деления* миникомпьютера 801 несколько сложнее, чем действия, выполняемые в теле цикла, поскольку выполняет знаковое деление и оснащена проверкой переполнения. При ее использовании подпрограмма деления может быть записана с помощью 32 последовательных команд *пошагового деления*, за которыми следует коррекция частного и остатка (для того, чтобы обеспечить верный знак остатка).

⁴ В действительности восстанавливающий алгоритм деления может избежать восстанавливающего шага, помещая результат вычитания в отдельный регистр и записывая этот регистр в x , только если результат 33-битового вычитания неотрицателен. Однако в некоторых реализациях для этого требуется дополнительный регистр и, возможно, большее время работы.

Использование короткого деления

Алгоритм для деления $64 \div 32 \Rightarrow 32$ может быть получен из алгоритма беззнакового деления больших целых чисел в листинге 9.1 на с. 211 для частного случая $m = 4$ и $n = 2$. Кроме того, требуется внесение ряда других изменений. Параметры в данном случае должны представлять собой полные слова, передаваемые по значению, а не массивы полуслов. Условие переполнения также отличается от исходного; в данном случае переполнение происходит, если частное не помещается в полном слове. Все это дает возможность внести определенные упрощения в программу. Можно показать, что оценка `qhat` в данном случае всегда точная, так как делитель состоит только из двух цифр по полслова. Это означает, что корректирующий шаг с прибавлением делителя может быть опущен. Если развернуть “главный цикл” листинга 9.1 и цикл внутри него, становятся возможными еще некоторые небольшие упрощения кода.

Результат этих преобразований показан в листинге 9.3. Делимое находится в словах `u1` и `u0`; в `u1` содержится старшее слово. Параметр `v` содержит делитель. Функция по завершении вычислений возвращает значение частного. Если вызывающая функция передает в качестве `r` ненулевой указатель, то остаток будет возвращен в слове, на которое указывает `r`.

ЛИСТИНГ 9.3. Длинное беззнаковое деление с использованием команды деления полных слов

```
unsigned divlu(unsigned u1, unsigned u0, unsigned v,
              unsigned *r)
{
    const unsigned b
        = 65536; // Основание счисления (16 битов)
    unsigned un1, un0, // Слова нормализованного делимого
            vn1, vn0, // Цифры нормализованного делителя
            q1, q0, // Цифры частного
            un32, un21, // Пары цифр делимого
            un10,
            rhat; // Остаток
    int s; // Величина сдвига нормализации

    if (u1 >= v) // При переполнении
    { // устанавливается невозможное
        if (r != NULL) // значение остатка и
            *r = 0xFFFFFFFF; // возвращается максимально
        return 0xFFFFFFFF; // возможное число
    }

    s = nlz(v); // 0 <= s <= 31.
    v = v << s; // Нормализация делителя
    vn1 = v >> 16; // Разбиение делителя
    vn0 = v & 0xFFFF; // на две 16-битовые цифры

    un32 = (u1 << s) | (u0 >> 32 - s) & (~s >> 31);
    un10 = u0 << s; // Сдвиг делимого влево

    un1 = un10 >> 16; // Разбиение правой половины
    un0 = un10 & 0xFFFF; // делимого на две цифры
```

```

    q1 = un32/vn1;           // Вычисление первой цифры
    rhat = un32 - q1*vn1;    // частного q1
again1:
    if (q1 >= b || q1*vn0 > b*rhat + un1)
    {
        q1 = q1 - 1;
        rhat = rhat + vn1;
        if (rhat < b) goto again1;
    }

    un21 = un32*b + un1 - q1*v; // Умножение и вычитание

    q0 = un21/vn1;           // Вычисление второй цифры
    rhat = un21 - q0*vn1;    // частного q0
again2:
    if (q0 >= b || q0*vn0 > b*rhat + un0)
    {
        q0 = q0 - 1;
        rhat = rhat + vn1;
        if (rhat < b) goto again2;
    }

    if (r != NULL)           // Если требуется,
        *r = (un21*b + un0   // вычисляем остаток
              - q0*v) >> s;
    return q1*b + q0;
}

```

В качестве индикатора переполнения программа возвращает значение остатка, равное максимально возможному целому числу. Такой остаток невозможен ни при каком корректном делении, поскольку остаток всегда меньше делителя. Кроме того, при переполнении функция возвращает частное, равное максимально возможному целому числу (что также может служить индикатором переполнения в случае, когда вызывающая функция не запрашивает значение остатка).

Странное выражение `(-s>>31)` в присвоении `u32` предназначено для того, чтобы программа работала в случае `s = 0` на машинах, сдвиг у которых выполняется по модулю 32 (например, в Intel x86).

Эксперименты с равномерно распределенными случайными числами показывают, что тело цикла `again` выполняется примерно около 0.38 раза на каждый вызов функции. Это дает нам среднее количество выполняемых команд, равное 52. Среди них одна команда подсчета *количества ведущих нулевых битов*, две — *деления*, а также 6.5 — *умножения* (не считая умножений на `b`, которые выполняются с помощью *сдвигов*). Если требуется значение остатка, добавляется шесть команд (включая команду сохранения `r`), одна из которых — *умножение*.

Что касается знаковой версии `divlu`, то, вероятно, пошаговая модификация приведенного в листинге 9.3 кода для получения знакового варианта представляет собой непростую задачу. Однако данный алгоритм можно использовать для реализации знакового деления, если взять абсолютное значение аргументов, вызвать `divlu` и обратить знак результата, если знаки аргументов были различны. При этом не возникает никаких проблем с экстремальными значениями типа наибольшего по абсолютной величине отрица-

тельного значения, так как абсолютное значение любого знакового числа корректно представимо беззнаковым числом. Такой алгоритм показан в листинге 9.4.

Листинг 9.4. Длинное знаковое деление с использованием длинного беззнакового деления

```
int divls(int u1, unsigned u0, int v, int *r)
{
    int q, uneg, vneg, diff, borrow;
    uneg = u1 >> 31; // -1, если u < 0.
    if (uneg)
    {
        // Вычисление абсолютного
        u0 = -u0; // значения делимого u
        borrow = (u0 != 0);
        u1 = -u1 - borrow;
    }

    vneg = v >> 31; // -1, если v < 0.
    v = (v ^ vneg) - vneg; // Абсолютное значение v

    if ((unsigned)u1 >= (unsigned)v) goto overflow;

    q = divlu(u1, u0, v, (unsigned *)r);

    diff = uneg ^ vneg; // Изменяем знак q, если
    q = (q ^ diff) - diff; // знаки u и v различны
    if (uneg && r != NULL)
        *r = -*r;

    if ((diff ^ q) < 0 && q != 0)
    {
        // При переполнении даем
        overflow: // остатку невозможное
        if (r != NULL) // значение и возвращаем
            *r = 0x80000000; // отрицательное частное
        q = 0x80000000; // с максимально возможным
    } // абсолютным значением
    return q;
}
```

В случае знакового деления достаточно трудно разработать действительно хороший код для обнаружения переполнения. Алгоритм, приведенный в листинге 9.4, выполняет предварительное обнаружение переполнения, идентичное случаю длинного беззнакового деления, которое гарантирует, что $|u/v| < 2^{32}$. После этого нам остается только убедиться, что остаток имеет корректный знак или равен 0.

9.5. Деление двойных слов из длинного деления

В этом разделе рассматривается выполнение деления $64 + 64 \Rightarrow 64$ на основе деления $64 + 32 \Rightarrow 32$ как для знакового, так и для беззнакового случаев. Рассматриваемые далее алгоритмы лучше всего подходят для машин, оснащенных командой длинного деления ($64 + 32$), как минимум для беззнакового случая. Очень хорошо, если машина имеет команду подсчета количества ведущих нулевых битов. Машина может быть оснащена как 32-, так и 64-битовыми регистрами, но мы предполагаем, что если регистры машины

32-битовые, то компилятор в состоянии реализовать такие базовые операции, как сложение и сдвиг, над 64-битовыми операндами (тип `long long` языка программирования C).

В мире GNU C эти функции известны как `__udivdi3` и `__divdi3`, и здесь также использованы аналогичные имена.

Беззнаковое деление двойных слов

Процедура для этой операции приведена в листинге 9.5.

Листинг 9.5. Беззнаковое деление двойных слов на основе длинного деления

```
unsigned long long udivdi3(unsigned long long u,
                          unsigned long long v)
{
    unsigned long long u0, u1, v1, q0, q1, k, n;

    if (v >> 32 == 0)          // Если v < 2**32:
    {
        if (u >> 32 < v)      // Если u/v не может вызвать
            return DIVU(u, v) // переполнения, выполняем
                & 0xFFFFFFFF; // одно деление
        else                  // Если u/v вызывает
        {                      // переполнение:
            u1 = u >> 32;     // Разбиваем u на две
            u0 = u & 0xFFFFFFFF; // половины.
            q1 = DIVU(u1, v)  // Первая цифра частного.
                & 0xFFFFFFFF;
            k = u1 - q1*v;    // Первая цифра остатка, < v.
            q0 = DIVU((k << 32) + u0, v) // Вторая цифра
                & 0xFFFFFFFF; // частного.
            return (q1 << 32) + q0;
        }
    }

    // Здесь v >= 2**32.
    n = nlz64(v);             // 0 <= n <= 31.
    v1 = (v << n) >> 32;    // Нормализация делителя так, что
    // его младший значащий бит равен 1.
    u1 = u >> 1;             // Чтобы обеспечить
    // отсутствие переполнения.
    q1 = DIVU(u1, v1)        // Получаем частное из
    & 0xFFFFFFFF;           // беззнакового деления.
    q0 = (q1 << n) >> 31;    // Откат нормализации
    // и деление u на 2.
    if (q0 != 0)             // Делаем q0 правильным
    {                          // или меньшим на 1.
        q0 = q0 - 1;
    }

    if ((u - q0*v) >= v)
    {
        q0 = q0 + 1;        // Теперь значение q0 верное.
    }

    return q0;
}
```

В этом коде различаются три случая: (1) когда используется одно выполнение машинного беззнакового длинного деления DIVU, (2) когда неприменим случай (1), но делитель является 32-битовой величиной, и (3) когда делитель нельзя представить 32-битовой величиной. Нетрудно увидеть, что приведенный код корректно работает в случаях (1) и (2). В случае (2) вспомните о выполнении длинного деления "в столбик".

Случай (3), однако, заслуживает доказательства, потому что в некоторых случаях он очень близок к тому, чтобы быть неработоспособным. Обратите внимание, что в этом случае требуется только одно выполнение DIVU, но требуется выполнение команд *количества ведущих нулевых битов и умножения*.

Доказательство будет основываться на следующих тождествах.

$$\lfloor \lfloor a/b \rfloor / d \rfloor = \lfloor a / (bd) \rfloor \quad (2)$$

$$b \lfloor a/b \rfloor = a - \text{rem}(a, b) \quad (3)$$

Из первой строки интересующей нас процедуры (мы полагаем, что $v \neq 0$) получаем

$$0 \leq n \leq 31.$$

Очевидно, что при вычислении v_1 левый сдвиг не может вызвать переполнения. Следовательно,

$$v_1 = \lfloor v / 2^{32-n} \rfloor \quad \text{и}$$

$$u_1 = \lfloor u / 2 \rfloor.$$

При вычислении q_1 значения u_1 и v_1 находятся в диапазоне команды DIVU и не могут вызвать переполнения. Следовательно,

$$q_1 = \lfloor u_1 / v_1 \rfloor.$$

При первом вычислении q_0 левый сдвиг не может вызвать переполнения, поскольку $q_1 < 2^{32}$ (так как максимальное значение u_1 равно $2^{63} - 1$, а минимальное значение v_1 равно 2^{31}). Таким образом,

$$q_0 = \lfloor q_1 / 2^{31-n} \rfloor.$$

Теперь в качестве основной части доказательства мы хотим показать, что

$$\lfloor u/v \rfloor \leq q_0 \leq \lfloor u/v \rfloor + 1,$$

т.е. что первое вычисление q_0 дает требуемый результат или требуемый результат плюс 1.

Дважды используя (2), получаем

$$q_0 = \left\lfloor \frac{u}{2^{32-n} v_1} \right\rfloor = \left\lfloor \frac{u}{2^{32-n} \left\lfloor \frac{v}{2^{32-n}} \right\rfloor} \right\rfloor.$$

Применение (3) дает

$$q_0 = \left\lfloor \frac{u}{v - \text{rem}(v, 2^{32-n})} \right\rfloor.$$

Выполнение алгебраических преобразований для приведения полученной формулы к виду $u/v + \text{нечто}$ даст

$$q_0 = \left\lfloor \frac{u}{v} + \frac{u \text{rem}(v, 2^{32-n})}{v(v - \text{rem}(v, 2^{32-n}))} \right\rfloor.$$

Эта формула имеет вид

$$\left\lfloor \frac{u}{v} + \delta \right\rfloor,$$

и мы покажем, что $\delta < 1$.

δ имеет наибольшее значение, когда $\text{rem}(v, 2^{32-n})$ велико настолько, насколько это возможно, и при этом v настолько возможно мало. Максимальное значение $\text{rem}(v, 2^{32-n})$ равно $2^{32-n} - 1$. Исходя из способа определения n через v , $v \geq 2^{63-n}$. Таким образом, наименьшее значение v , дающее указанный остаток, равно $2^{63-n} + 2^{32-n} - 1$. Следовательно,

$$\delta \leq \frac{u(2^{32-n} - 1)}{(2^{63-n} + 2^{32-n} - 1)2^{63-n}} < \frac{u(2^{32-n} - 1)}{(2^{63-n})^2}.$$

Методом подбора для n в диапазоне от 0 до 31

$$\delta < \frac{u}{2^{64}}.$$

Поскольку u не превышает $2^{64} - 1$, $\delta < 1$. Так как $q_0 = \lfloor u/v + \delta \rfloor$, и $\delta < 1$ (и, очевидно, $\delta \geq 0$)

$$\left\lfloor \frac{u}{v} \right\rfloor \leq q_0 \leq \left\lfloor \frac{u}{v} \right\rfloor + 1.$$

Для исправления результата путем при необходимости вычитания 1 нам требуется код

if ($u < q_0 * v$) $q_0 = q_0 - 1$;

(Иначе говоря, если остаток $u - q_0 v$ отрицателен, вычитаем 1 из q_0 .) Однако это не все, поскольку $q_0 v$ может вызвать переполнение (например, при $u = 2^{64} - 1$ и $v = 2^{32} + 3$). Так что вместо описанного мы вычитаем 1 из q_0 , так что результат либо корректен, либо меньше на 1, и $q_0 v$ не может вызвать переполнение. Нам нужно избежать вычитания 1, если $q_0 = 0$ (при $q_0 = 0$ это уже верное частное).

Таким образом, окончательная коррекция имеет следующий вид.

if ($(u - q_0 * v) >= v$) $q_0 = q_0 - 1$;

Убедимся в корректности вычислений. Мы уже видели, что q_0v не приводит к переполнению. Легко показать, что

$$0 \leq u - q_0v < 2v.$$

Не может ли вычитание привести к переполнению при очень большом v ($\geq 2^{63}$) в попытке получения результата, большего, чем v ? Нет, поскольку $u < 2^{64}$, а $q_0v \geq 0$.

Кстати, строкам

```
if (q0 != 0)           // Делаем q0 правильным
{                       // или меньшим на 1.
    q0 = q0 - 1;
}
```

имеются альтернативы, которые на некоторых машинах могут оказаться предпочтительными. Одной из них является замена на

```
if (q0 == 0) return 0;
```

Еще одна альтернатива заключается в размещении в начале этого раздела процедуры (или в начале всей процедуры) строки

```
if (u < v) return 0; // Избегаем последующих проблем
```

Эти альтернативы предпочтительны, если ветвление не является дорогостоящей операцией. Код, показанный в листинге 9.5, хорошо работает в случае, когда команда сравнения машины дает целочисленный результат 0/1 в регистре общего назначения. Тогда компилятор может заменить указанный код вычислением

```
q0 = q0 - (q0 != 0);
```

(либо вы можете записать этот код указанным образом, если ваш компилятор на такую оптимизацию не способен). На таких машинах это просто две команды — *сравнение* и *вычитание*.

Знаковое деление двойных слов

В случае знакового деления, похоже, нет лучшего способа, чем выполнить деление двойных слов, представляющих собой абсолютные значения операндов, и обратить знак частного, если операнды имеют разные знаки. Если машина имеет команду знакового длинного деления, которая обозначается здесь как DIVS, то может иметь смысл отдельный случай, когда можно применить DIVS вместо вызова `udivdi3` (предполагается, что такие случаи встречаются достаточно часто). Такая функция показана в листинге 9.6.

Листинг 9.6. Знаковое деление двойных слов на основе беззнакового

```
#define llabs(x) \
    ((unsigned long long t = (x) >> 63; ((x) ^ t) - t;))

long long divdi3(long long u, long long v)
{
    unsigned long long au, av;
    long long q, t;

    au = llabs(u);
    av = llabs(v);
```

```

if (av >> 31 == 0)      // Если |v| < 2**31 и
{                       // |u|/|v| не приводит
    if (au < av << 31) // к переполнению,
    {
        q = DIVS(u, v); // используем DIVS.
        return (q << 32) >> 32;
    }
}

q = au/av;              // Используем udivd13.
t = (u ^ v) >> 63;     // Если знаки u и v различны,
return (q ^ t) - t;    // меняем знак q.
}

```

Директива `#define` в коде в листинге 9.6 использует возможность GCC размещения составной инструкции в скобках для построения выражения — возможность, которой большинство компиляторов C не имеют. В некоторых компиляторах `llabs(x)` представляет собой встроенную функцию.

Проверка диапазона v не является точной; она дает сбой при $v = -2^{31}$. Если использование команды `DIVS` в этом случае критично, вместо третьей выполнимой строки в листинге 9.6 можно использовать проверку

```
if ((v << 32) >> 32 == v)
```

ценой одной лишней команды. Аналогично проверка того, что $|u|/|v|$ не может привести к переполнению, упрощена, и несколько важных случаев оказываются пропущенными; код, учитывающий использование $\delta = 0$ при проверке переполнения знакового деления приведен в разделе “Деление” на с. 55.

Упражнения

1. Покажите, что для действительных чисел $\lfloor x \rfloor = -\lceil -x \rceil$.
2. Разработайте код без ветвлений для вычисления частного и остатка при модульном делении на машине с базовым набором команд RISC, в который входят команды деления и получения остатка при делении с отсечением.
3. Аналогично разработайте код без ветвлений для вычисления частного и остатка при делении с округлением к меньшему значению на машине с базовым набором команд RISC, в который входят команды деления и получения остатка при делении с отсечением.
4. Как бы вы вычисляли $\lceil n/d \rceil$ в случае беззнаковых целых чисел n и d , $0 \leq n \leq 2^{32} - 1$ и $1 \leq d \leq 2^{32} - 1$? Считайте, что ваша машина имеет команду для беззнакового деления, вычисляющую $\lfloor n/d \rfloor$.
5. Теорема Д3 гласит, что для действительного x и целого d $\lfloor \lfloor x \rfloor / d \rfloor = \lfloor x/d \rfloor$. Покажите, что справедливо более общее утверждение: если $f(x)$ представляет собой (а) непрерывную, (б) монотонно неубывающую и (в) такую, что $f(x)$ при целом x является целым числом, то $\lfloor f(\lfloor x \rfloor) \rfloor = \lfloor f(x) \rfloor$ [36].

ГЛАВА 10

ЦЕЛОЕ ДЕЛЕНИЕ НА КОНСТАНТЫ

На многих компьютерах операция деления выполняется довольно медленно, и по возможности ее использования стараются избежать. Время выполнения команды деления, как правило, превышает время выполнения элементарного сложения в 20 и более раз, причем обычно оно одинаково велико как для больших, так и для малых операндов. В этой главе приводится ряд методов, позволяющих избежать команды деления в случаях, когда делитель представляет собой константу.

10.1. Знаковое деление на известную степень 2

Похоже, многие ошибочно считают, что *знаковый сдвиг вправо на k позиций* делит число на 2^k с использованием обычного отскакивающего деления [38]. Однако на самом деле не все так просто. Приведенный ниже код выполняет деление $q = n + 2^k$, где $1 \leq k \leq 31$ [52].

shrsi	t, n, k-1	Формируем целое число
shri	t, t, 32-k	$2^{**k} - 1$, если $n < 0$; иначе 0
add	t, n, t	Добавляем его к n
shrsi	q, t, k	и выполняем знаковый сдвиг вправо

Этот код не содержит команд ветвления. Он может быть упрощен до трех команд при делении на 2 ($k = 1$). Однако данный код имеет смысл только на компьютерах, которые способны выполнить сдвиг на большое значение за малое время. Случай $k = 31$ не имеет особого смысла, поскольку число 2^{31} не представимо на компьютере. Тем не менее приведенный код дает правильный результат и в этом случае ($q = -1$, если $n = -2^{31}$, и $q = 0$ для прочих значений n).

Для деления на -2^k после приведенного выше кода должна следовать команда изменения знака. Лучшего варианта для данного деления пока не придумано.

Вот еще один, более очевидный и простой код для деления на 2^k .

bge	n, label	Ветвление при $n \geq 0$
addi	n, n, $2^{**k}-1$	Прибавление $2^{**k} - 1$ к n
label	shrsi	n, n, k и знаковый сдвиг вправо

Этот код предпочтительнее использовать на машине с медленным сдвигом и быстрыми командами ветвления.

Компьютер PowerPC имеет необычное устройство для ускорения деления на степень 2 [34]. Команда *знакового сдвига вправо* устанавливает бит переноса, если сдвигаемое число отрицательно и был сдвиг одного или нескольких единичных битов. Кроме того, данный компьютер имеет команду *addze* для сложения бита переноса с регистром. Все это позволяет реализовать деление на любую (положительную) степень 2 с помощью двух команд.


```
shrsi q,n,k
addze q,q
```

Единственная команда сдвига `shrsi` на k позиций выполняет знаковое деление на 2^k , которое совпадает как с модульным делением, так и с делением с округлением к меньшему значению. Это наводит на мысль о том, что такое деление может быть предпочтительнее для использования в языках высокого уровня, чем деление с отсечением, так как позволяет компилятору транслировать выражение $n/2$ в единственную команду `shrsi`. Кроме того, команда `shrsi` с последующей за ней командой `neg` выполняет модульное деление на -2^k , что, в свою очередь, указывает на преимущества использования модульного деления (впрочем, в основном это вопрос эстетического восприятия, в силу того что задача деления на отрицательную константу встречается достаточно редко).

10.2. Знаковый остаток от деления на известную степень 2

Если нам требуется вычислить и частное, и остаток от деления $n + 2^k$, простейший путь получения значения остатка — вычислить его по формуле $r = n - q * 2^k$, для чего после вычисления частного требуется выполнение двух команд.

```
shli r,q,k
sub r,n,n
```

Для вычисления одного лишь остатка, похоже, необходимо выполнить около четырех или пяти команд. Один из способов вычисления состоит в использовании рассмотренной ранее последовательности из четырех команд для знакового деления на 2^k , за которой следуют две команды вычисления остатка. При этом две последовательные команды сдвига можно заменить командой `u`, что дает нам решение из пяти команд (четыре при $k = 1$).

<code>shrsi</code>	<code>t,n,k-1</code>	Формируем целое число
<code>shri</code>	<code>t,t,32-k</code>	$2^{**k} - 1$, если $n < 0$; иначе 0
<code>add</code>	<code>t,n,t</code>	Добавляем его к n ,
<code>andi</code>	<code>t,t,-2**k</code>	сбрасываем k правых битов
<code>sub</code>	<code>r,n,t</code>	и вычитаем его из n

Еще один метод основан на том, что

$$\text{rem}(n, 2^k) = \begin{cases} n \& (2^k - 1), & n \geq 0, \\ -((-n) \& (2^k - 1)), & n < 0. \end{cases}$$

Для того чтобы воспользоваться этой формулой, сначала вычисляем $t \leftarrow n \dot{\gg} 31$, а затем

$$r \leftarrow ((\text{abs}(n) \& (2^k - 1)) \oplus t) - t$$

(требуется пять команд) или в случае $k = 1$, так как $(-n) \& 1 = n \& 1$, вычисляем

$$r \leftarrow ((n \& 1) \oplus t) - t$$

(требуется четыре команды). Этот метод не очень хорош при $k > 1$, если компьютер не оснащен командой вычисления *абсолютного значения* (в таком случае вычисление остатка будет требовать выполнения семи команд).

Еще один метод основан на том, что

$$\text{rem}(n, 2^k) = \begin{cases} n \& (2^k - 1), & n \geq 0, \\ (((n + 2^k - 1) \& (2^k - 1)) - (2^k - 1)), & n < 0. \end{cases}$$

Эта формула приводит к вычислениям

$$t \leftarrow (n \gg k - 1) \gg 32 - k,$$

$$r \leftarrow (((n + t) \& (2^k - 1)) - t)$$

(при $k > 1$ требуется выполнение пяти команд, при $k = 1$ — четырех).

Все перечисленные методы работают при $1 \leq k \leq 31$.

Кстати, если в компьютере нет команды *знакового сдвига вправо*, то значение, которое равно $2^k - 1$ для $n < 0$ и 0 для $n \geq 0$, может быть построено следующим образом.

$$t_1 \leftarrow n \gg 31$$

$$r \leftarrow (t_1 \ll k) - t_1$$

Это приводит к добавлению только одной команды.

10.3. Знаковое деление и вычисление остатка для других случаев

Основной прием в этом случае состоит в умножении на некоторое соответствующее делителю d число, примерно равное $2^{32}/d$, с последующим выделением 32 левых битов произведения. Однако реальные вычисления для ряда делителей, в частности для 7, оказываются существенно сложнее.

Рассмотрим сначала несколько конкретных примеров, которые пояснят код, генерируемый обобщенным методом. Обозначим регистры следующим образом.

- n — входное целое число (числитель)
- M — в него загружается "магическое число"
- t — временный регистр
- q — в нем будет размещено частное
- r — в нем будет размещен остаток

Деление на 3

```
li    M, 0x55555556 Загрузка "магического числа" (2**32+2)/3
mulhs q, M, n      q = floor(M*n/2**32)
shri  t, n, 31     Прибавляем 1 к q, если
add   q, q, t      n отрицательно
```

```

mul  t, q, 3      Вычисляем остаток как
sub   r, n, t     r = n - q*3

```

Доказательство. Операция *старшее слово знакового умножения* (`mulhs`) не может вызвать переполнения, поскольку произведение двух 32-битовых чисел всегда может быть представлено 64-битовым числом, а команда `mulhs` возвращает старшие 32 бит 64-битового произведения. Это действие эквивалентно делению 64-битового произведения на 2^{32} и вычислению функции `floor()` от полученного результата, причем это замечание справедливо независимо от того, отрицательно рассматриваемое произведение или положительно. Таким образом, при $n \geq 0$ приведенный выше код вычисляет

$$q = \left\lfloor \frac{2^{32} + 2n}{3 \cdot 2^{32}} \right\rfloor = \left\lfloor \frac{n}{3} + \frac{2n}{3 \cdot 2^{32}} \right\rfloor.$$

Далее, $n < 2^{31}$, поскольку $2^{31} - 1$ является наибольшим представимым целым числом. Следовательно, член-ошибка $2n/(3 \cdot 2^{32})$ меньше $1/3$ (и это значение неотрицательно), так что в соответствии с теоремой D4 (с. 210) получаем $q = \lfloor n/3 \rfloor$, что и требовалось получить (формула (1) на с. 209).

В случае $n < 0$ к частному добавляется 1. Следовательно, приведенный выше код вычисляет в этом случае величину

$$q = \left\lfloor \frac{2^{32} + 2n}{3 \cdot 2^{32}} \right\rfloor + 1 = \left\lfloor \frac{2^{32}n + 2n + 3 \cdot 2^{32}}{3 \cdot 2^{32}} \right\rfloor = \left\lfloor \frac{2^{32}n + 2n + 1}{3 \cdot 2^{32}} \right\rfloor.$$

Здесь была использована теорема D2. Следовательно,

$$q = \left\lfloor \frac{n}{3} + \frac{2n + 1}{3 \cdot 2^{32}} \right\rfloor.$$

При $-2^{31} \leq n \leq -1$ получаем

$$-\frac{1}{3} + \frac{1}{3 \cdot 2^{32}} \leq \frac{2n + 1}{3 \cdot 2^{32}} \leq -\frac{1}{3 \cdot 2^{32}}.$$

Таким образом, ошибка отрицательна и больше, чем $-1/3$, а значит, в соответствии с теоремой D4 $q = \lceil n/3 \rceil$, что и является искомым результатом (формула (1) на с. 209).

Итак, установлено, что вычисленное значение частного корректно. Тогда корректно и значение остатка, поскольку остаток должен удовлетворять соотношению

$$n = qd + r.$$

Умножение на 3 не может вызвать переполнения (так как $-2^{31}/3 \leq q \leq (2^{31}-1)/3$), а *вычитание* не может привести к переполнению в связи с тем, что результат должен находиться в диапазоне от -2 до $+2$.

Команда *умножения на непосредственно заданное значение* может быть выполнена с помощью двух сложений или сдвига и сложения в том случае, если такая замена дает выигрыш во времени вычислений.

На многих современных RISC-компьютерах частное может быть вычислено так, как показано выше, за девять-десять тактов, в то время как команда *деления* может потребовать 20 тактов или около того.

Деление на 5

Для деления на 5 хотелось бы использовать код того же типа, что и для деления на 3, но, понятно, с множителем $(2^{32} + 4)/5$. К сожалению, при этом ошибка оказывается слишком большой и результат отличается на единицу от точного примерно для пятой части значений $|n| \geq 2^{30}$. Но оказывается, можно использовать множитель $(2^{33} + 3)/5$ и добавить к коду команду *знакового сдвига вправо*. В результате получается следующий код.

```
li      M, 0x66666667   Загрузка "магического числа" (2**33+3)/5
mulhs  q, M, n         q = floor(M*n/2**32)
shrsi  q, q, 1
shri   t, n, 31        Прибавляем 1 к q, если
add     q, q, t         n отрицательно

mul1   t, q, 5         Вычисляем остаток как
sub     r, n, t        r = n - q*5
```

Доказательство. Команда `mulhs` дает 32 старших бита 64-битового произведения, после чего код знаково сдвигает полученное значение вправо на одну позицию. Это действие эквивалентно делению 64-битового произведения на 2^{33} и вычислению функции `floor()` от полученного результата. Таким образом, для $n \geq 0$ приведенный код вычисляет

$$q = \left\lfloor \frac{2^{33} + 3}{5} \frac{n}{2^{33}} \right\rfloor = \left\lfloor \frac{n}{5} + \frac{3n}{5 \cdot 2^{33}} \right\rfloor.$$

Для $0 \leq n < 2^{31}$ ошибка составляет $3n/(5 \cdot 2^{33})$; это значение неотрицательно и меньше, чем $1/5$, так что в соответствии с теоремой D4 $q = \lfloor n/5 \rfloor$.

При $n < 0$ приведенный выше код вычисляет значение

$$q = \left\lfloor \frac{2^{33} + 3}{5} \frac{n}{2^{33}} \right\rfloor + 1 = \left\lfloor \frac{n}{5} + \frac{3n+1}{5 \cdot 2^{33}} \right\rfloor.$$

Здесь член-ошибка отрицателен и превышает $-1/5$, так что $q = \lceil n/5 \rceil$. Корректность вычисленного остатка доказывается так же, как и в случае деления на 3. Как и ранее, команда *умножения на непосредственно заданное значение* может быть заменена — в данном случае *сдвигом влево* на две позиции и *сложением*.

Деление на 7

При делении на 7 возникают новые проблемы. Множители $(2^{32} + 3)/7$ и $(2^{33} + 6)/7$ дают слишком большие значения ошибки. Множитель $(2^{34} + 5)/7$ подходит, но он

слишком велик для размещения в 32-битовом знаковом слове. Однако умножение на такое большое число можно выполнить путем умножения на отрицательное число $(2^{34} + 5)/7 - 2^{32}$ с последующей коррекцией произведения путем сложения. В результате получаем следующий код для деления на 7.

```
li      M, 0x92492493  "Магическое число" (2**34+5)/7 - 2**32
mulhs  q, M, n        q = floor(M*n/2**32).
add    q, q, n        q = floor(M*n/2**32) + n
shrsi  q, q, 2        q = floor(q/4)
shri   t, n, 31      Прибавляем 1 к q, если
add    q, q, t        n отрицательно

mul    t, q, 7        Вычисляем остаток как
sub    r, n, t        r = n - q*7
```

Доказательство. Важно обратить внимание на то, что команда "add q, q, n" не может вызвать переполнения. Дело в том, что q и n имеют противоположные знаки; это связано с умножением на отрицательное число. Таким образом, это "компьютерное сложение" выполняется так же, как и обычное арифметическое сложение, и для $n \geq 0$ приведенный выше код вычисляет

$$q = \left\lfloor \left(\left(\left(\frac{2^{34} + 5}{7} - 2^{32} \right) \frac{n}{2^{32}} + n \right) / 4 \right) \right\rfloor = \left\lfloor \left(\frac{2^{34}n + 5n - 7 \cdot 2^{32}n + 7 \cdot 2^{32}n}{7 \cdot 2^{32}} \right) / 4 \right\rfloor = \left\lfloor \frac{n}{7} + \frac{5n}{7 \cdot 2^{34}} \right\rfloor$$

(здесь при преобразованиях использовано следствие из теоремы D3).

Для $0 \leq n < 2^{31}$ ошибка, составляющая $5n/(7 \cdot 2^{34})$, неотрицательна и меньше $1/7$, так что $q = \lfloor n/7 \rfloor$.

Для $n < 0$ приведенный выше код вычисляет значение

$$q = \left\lfloor \left(\left(\left(\frac{2^{34} + 5}{7} - 2^{32} \right) \frac{n}{2^{32}} + n \right) / 4 \right) + 1 \right\rfloor = \left\lfloor \frac{n}{7} + \frac{5n+1}{7 \cdot 2^{34}} \right\rfloor.$$

В этом случае ошибка не положительна и больше $-1/7$, так что $q = \lceil n/7 \rceil$.

Команда умножения на непосредственно заданное значение может быть заменена сдвигом влево на три позиции и вычитанием.

10.4. Знаковое деление на делитель, не меньший 2

После рассмотренного материала вы можете заинтересоваться, не возникают ли новые проблемы при работе с другими делителями. В этом разделе вы узнаете, что не возникают: все проблемы при делителе $d \geq 2$ исчерпываются тремя рассмотренными случаями.

Некоторые из приводимых доказательств весьма сложны, так что читайте дальнейший материал внимательно и не забывайте о том, что работаете со словом размером W .

Итак, пусть даны размер слова $W \geq 3$ и делитель $2 \leq d < 2^{W-1}$ и требуется найти наименьшее целое $0 \leq m < 2^W$ и целое $p \geq W$, такие, что

$$\left\lfloor \frac{mn}{2^p} \right\rfloor = \left\lfloor \frac{n}{d} \right\rfloor \text{ для } 0 \leq n < 2^{W-1} \text{ и} \quad (1,а)$$

$$\left\lfloor \frac{mn}{2^p} \right\rfloor + 1 = \left\lfloor \frac{n}{d} \right\rfloor \text{ для } -2^{W-1} \leq n \leq -1. \quad (1,б)$$

Найти *наименьшее* целое m необходимо потому, что меньший множитель может потребовать меньшую величину сдвига (возможно, 0) либо привести к коду наподобие кода из примера “деление на 5”, но не “деление на 7”. Требование $m \leq 2^W - 1$ обеспечивает не большее количество команд, чем в примере деления на 7 (т.е. с множителем в диапазоне от 2^{W-1} до $2^W - 1$ можно справиться с помощью дополнительной команды add, как это было сделано в примере деления на 7, но предпочтительнее обойтись меньшими множителями). Требование $p \geq W$ нужно постольку, поскольку генерируемый код выделяет левую половину произведения mn , что эквивалентно сдвигу вправо на W позиций. Таким образом, общий сдвиг вправо составляет не менее W позиций.

Есть определенное различие между множителем m и “магическим числом”, обозначаемым как M . Магическое число — это значение, которое используется в команде *умножения* и задается следующим образом:

$$M = \begin{cases} m, & \text{если } 0 \leq m < 2^{W-1}, \\ m - 2^W, & \text{если } 2^{W-1} \leq m < 2^W. \end{cases}$$

В силу того, что соотношение (1,б) должно выполняться при $n = -d$, т.е. $\left\lfloor -md/2^p \right\rfloor + 1 = -1$, получаем

$$\frac{md}{2^p} > 1. \quad (2)$$

Пусть n_c — наибольшее (положительное) значение n , такое, что $\text{rem}(n_c, d) = d - 1$. Оно существует, поскольку имеется как минимум одно возможное значение $n_c = d - 1$. Его можно вычислить как $n_c = \left\lfloor 2^{W-1}/d \right\rfloor d - 1 = 2^{W-1} - \text{rem}(2^{W-1}, d) - 1$. Поскольку n_c представляет собой одно из d наибольших допустимых значений n , получаем

$$2^{W-1} - d \leq n_c \leq 2^{W-1} - 1 \quad (3,а)$$

и, очевидно,

$$n_c \geq d - 1. \quad (3,б)$$

Поскольку (1,а) должно выполняться при $n = n_c$,

$$\left\lfloor \frac{mn_c}{2^p} \right\rfloor = \left\lfloor \frac{n_c}{d} \right\rfloor = \frac{n_c - (d - 1)}{d}$$

или

$$\frac{mn_c}{2^p} < \frac{n_c + 1}{d}.$$

Объединяя полученный результат с (2), получим

$$\frac{2^p}{d} < m < \frac{2^p}{d} \frac{n_c + 1}{n_c}. \quad (4)$$

Поскольку m должно быть наименьшим целым, удовлетворяющим (4), оно представляет собой целое число, следующее за $2^p/d$, т.е.

$$m = \frac{2^p + d - \text{rem}(2^p, d)}{d}. \quad (5)$$

Комбинируя полученную формулу с правой частью (4), получим

$$2^p > n_c (d - \text{rem}(2^p, d)). \quad (6)$$

Алгоритм

Таким образом, алгоритм поиска магического числа M и величины сдвига s для данного делителя d начинается с вычисления n_c , а затем неравенство (6) решается путем подстановки последовательных возрастающих значений. Если $p < W$, то устанавливаем $p = W$ (теорема ниже показывает, что данное значение также удовлетворяет неравенству (6)). Когда найдено наименьшее $p \geq W$, удовлетворяющее неравенству (6), из (5) вычисляется значение m . Это наименьшее возможное значение m , поскольку нами найдено наименьшее приемлемое p , а из (4) понятно, что, чем меньше значение p , тем меньше и значение m . И наконец $s = p - W$ и M просто являются интерпретацией m как знаково-го целого числа (как оно рассматривается командой `mulhs`).

То, что при $p < W$ мы устанавливаем $p = W$, обосновывается следующей теоремой.

ТЕОРЕМА DC1. Если для некоторого значения p справедливо неравенство (6), то оно справедливо и для больших значений p .

Доказательство. Предположим, что неравенство (6) справедливо для $p = p_0$. Умножение (6) на 2 дает

$$2^{p_0+1} > n_c (2d - 2\text{rem}(2^{p_0}, d)).$$

Из теоремы D5 $\text{rem}(2^{p_0+1}, d) \geq 2\text{rem}(2^{p_0}, d) - d$. Объединяя эти выражения, получим

$$2^{p_0+1} > n_c (2d - (\text{rem}(2^{p_0+1}, d) + d)), \text{ или } 2^{p_0+1} > n_c (d - \text{rem}(2^{p_0+1}, d)).$$

Следовательно, неравенство (6) справедливо для $p = p_0 + 1$, а потому и для всех больших значений.

Таким образом, можно решить (6) путем бинарного поиска, хотя, по-видимому, предпочтительнее простой линейный поиск, начинающийся со значения $p = W$, поскольку d обычно мало, а малые значения d приводят к малым значениям p .

Доказательство пригодности алгоритма

Покажем, что (6) всегда имеет решение и что $0 \leq m < 2^W$ (нет необходимости показывать, что $p \geq W$, так как это условие выполняется принудительно).

Покажем, что (6) всегда имеет решение, получив верхнюю границу p . Кроме того, в познавательных целях получим также нижнюю границу в предположении, что p не обязано быть равным как минимум W . Для получения указанных границ p заметим, что для любого положительного целого x существует степень 2, большая x и не превосходящая $2x$. Следовательно, из (6)

$$n_c (d - \text{rem}(2^p, d)) < 2^p \leq 2n_c (d - \text{rem}(2^p, d)).$$

Поскольку $0 \leq \text{rem}(2^p, d) \leq d - 1$,

$$n_c + 1 \leq 2^p \leq 2n_c d. \tag{7}$$

Из (3,а) и (3,б) получаем $n_c \geq \max(2^{p-1} - d, d - 1)$. Графики функций $f_1(d) = 2^{p-1} - d$ и $f_2(d) = d - 1$ пересекаются в точке $d = (2^{p-1} + 1)/2$. Следовательно, $n_c \geq (2^{p-1} - 1)/2$. Поскольку n_c — целое число, $n_c \geq 2^{p-2}$. С учетом того, что $n_c d \leq 2^{p-1} - 1$, (7) превращается в

$$2^{p-2} + 1 \leq 2^p \leq 2(2^{p-1} - 1)^2$$

или

$$W - 1 \leq p \leq 2W - 2. \tag{8}$$

Нижняя граница $p = W - 1$ вполне может быть достигнута (например, для $W = 32$, $d = 3$), но в этом случае мы устанавливаем $p = W$.

Если не делать p равным W принудительно, то из (4) и (7) получаем

$$\frac{n_c + 1}{d} < m < \frac{2n_c d n_c + 1}{d n_c}.$$

Применение (3,б) дает

$$\frac{d-1+1}{d} < m < 2(n_c + 1).$$

Так как $n_c \leq 2^{p-1} - 1$ (3,а),

$$2 \leq m \leq 2^W - 1.$$

Если p принудительно делается равным W , то из (4) получаем

$$\frac{2^W}{d} < m < \frac{2^W n_c + 1}{d n_c}.$$

Поскольку $2 \leq d \leq 2^{w-1} - 1$ и $n_c \geq 2^{w-2}$,

$$\frac{2^w}{2^{w-1} - 1} < m < \frac{2^w}{2} \frac{2^{w-2} + 1}{2^{w-2}} \text{ или}$$

$$3 \leq m \leq 2^{w-1} + 1.$$

Следовательно, в любом случае m находится в пределах границ для схемы, проиллюстрированной примером деления на 7.

Доказательство корректности произведения

Покажем, что если p и m вычисляются из (6) и (5), то выполняются уравнения (1,а) и (1,б).

Уравнение (5) и неравенство (6), как легко видеть, вытекают из (4). (В случае, когда p принудительно приравнивается к W , как показывает теорема DC1, неравенство (6) остается справедливым.) Далее рассмотрим отдельно пять диапазонов значений n :

$$\begin{aligned} 0 &\leq n \leq n_c, \\ n_c + 1 &\leq n \leq n_c + d - 1, \\ -n_c &\leq n \leq -1, \\ -n_c - d + 1 &\leq n \leq -n_c - 1 \text{ и} \\ n &= -n_c - d. \end{aligned}$$

Из (4), поскольку m — целое число, следует

$$\frac{2^p}{d} < m \leq \frac{2^p (n_c + 1) - 1}{dn_c}.$$

При умножении на $n/2^p$ для $n \geq 0$ это соотношение превращается в

$$\frac{n}{d} \leq \frac{mn}{2^p} \leq \frac{2^p n(n_c + 1) - n}{2^p dn_c}, \text{ так что}$$

$$\left\lfloor \frac{n}{d} \right\rfloor \leq \left\lfloor \frac{mn}{2^p} \right\rfloor \leq \left\lfloor \frac{n}{d} + \frac{(2^p - 1)n}{2^p dn_c} \right\rfloor.$$

Если $0 \leq n \leq n_c$, то $0 \leq (2^p - 1)n / (2^p dn_c) < 1/d$, так что по теореме D4

$$\left\lfloor \frac{n}{d} + \frac{(2^p - 1)n}{2^p dn_c} \right\rfloor = \left\lfloor \frac{n}{d} \right\rfloor.$$

Следовательно, в случае $0 \leq n \leq n_c$ уравнение (1,а) выполняется.

Для $n > n_c$ значение n ограничено диапазоном

$$n_c + 1 \leq n \leq n_c + d - 1, \tag{9}$$

так как $n \geq n_c + d$ противоречит выбору n_c как наибольшего значения n , такого, что $\text{rem}(n_c, d) = d - 1$ (кроме того, из (3,а) видно, что $n \geq n_c + d$ влечет за собой $n \geq 2^{p-1}$). Из (4) для $n \geq 0$ следует

$$\frac{n}{d} < \frac{mn}{2^p} < \frac{n n_c + 1}{d n_c}.$$

Путем элементарных алгебраических преобразований это неравенство может быть записано следующим образом:

$$\frac{n}{d} < \frac{mn}{2^p} < \frac{n_c + 1}{d} + \frac{(n - n_c)(n_c + 1)}{d n_c}. \tag{10}$$

Из (9) $1 \leq n - n_c \leq d - 1$, так что

$$0 < \frac{(n - n_c)(n_c + 1)}{d n_c} \leq \frac{d - 1}{d} \frac{n_c + 1}{n_c}.$$

Поскольку (из (3,б)) $n_c \geq d - 1$ и $(n_c + 1)/n_c$ максимально при минимальном n_c ,

$$0 < \frac{(n - n_c)(n_c + 1)}{d n_c} \leq \frac{d - 1}{d} \frac{d - 1 + 1}{d - 1} = 1.$$

В (10) член $(n_c + 1)/d$ является целым числом. Член $(n - n_c)(n_c + 1)/d n_c$ меньше или равен 1. Следовательно, (10) превращается в

$$\left\lfloor \frac{n}{d} \right\rfloor \leq \left\lfloor \frac{mn}{2^p} \right\rfloor \leq \frac{n_c + 1}{d}.$$

Для всех n из диапазона (9) $\lfloor n/d \rfloor = (n_c + 1)/d$. Следовательно, (1,а) выполняется и в этом случае ($n_c + 1 \leq n \leq n_c + d - 1$).

Для $n < 0$ из (4), так как m — целое число, имеем

$$\frac{2^p + 1}{d} \leq m < \frac{2^p n_c + 1}{d n_c}.$$

При умножении на $n/2^p$ с учетом того, что $n < 0$, это выражение преобразуется:

$$\frac{n n_c + 1}{d n_c} < \frac{mn}{2^p} \leq \frac{n 2^p + 1}{d 2^p}$$

или

$$\left\lfloor \frac{n n_c + 1}{d n_c} \right\rfloor + 1 \leq \left\lfloor \frac{mn}{2^p} \right\rfloor + 1 \leq \left\lfloor \frac{n 2^p + 1}{d 2^p} \right\rfloor + 1.$$

Применив теорему D2, получим

$$\left\lceil \frac{n_c(n_c+1) - dn_c + 1}{dn_c} \right\rceil + 1 \leq \left\lfloor \frac{mn}{2^p} \right\rfloor + 1 \leq \left\lfloor \frac{n(2^p+1) - 2^p d + 1}{2^p d} \right\rfloor + 1,$$

$$\left\lceil \frac{n(n_c+1)+1}{dn_c} \right\rceil \leq \left\lfloor \frac{mn}{2^p} \right\rfloor + 1 \leq \left\lfloor \frac{n(2^p+1)+1}{2^p d} \right\rfloor.$$

Поскольку $n+1 \leq 0$, правое неравенство может быть ослаблено, что дает нам

$$\left\lceil \frac{n}{d} + \frac{n+1}{dn_c} \right\rceil \leq \left\lfloor \frac{mn}{2^p} \right\rfloor + 1 \leq \left\lfloor \frac{n}{d} \right\rfloor. \quad (11)$$

Для $-n_c \leq n \leq -1$

$$\frac{-n_c+1}{dn_c} \leq \frac{n+1}{dn_c} \leq 0 \text{ или}$$

$$-\frac{1}{d} < \frac{n+1}{dn_c} \leq 0.$$

Следовательно, в силу теоремы D4

$$\left\lceil \frac{n}{d} + \frac{n+1}{dn_c} \right\rceil = \left\lfloor \frac{n}{d} \right\rfloor,$$

так что в данном случае ($-n_c \leq n \leq -1$) уравнение (1,6) выполняется.

При $n < -n_c$ величина n ограничена диапазоном

$$-n_c - d \leq n \leq -n_c - 1. \quad (12)$$

(Исходя из (3,а), из $n < -n_c - d$ вытекает $n < -2^{p-1}$, что невозможно.) Выполняя элементарные алгебраические преобразования в левой части (11), получаем

$$\left\lceil \frac{-n_c-1}{d} + \frac{(n+n_c)(n_c+1)+1}{dn_c} \right\rceil \leq \left\lfloor \frac{mn}{2^p} \right\rfloor + 1 \leq \left\lfloor \frac{n}{d} \right\rfloor. \quad (13)$$

Для $-n_c - d + 1 \leq n \leq -n_c - 1$

$$\frac{(-d+1)(n_c+1)}{dn_c} + \frac{1}{dn_c} \leq \frac{(n+n_c)(n_c+1)+1}{dn_c} \leq \frac{-(n_c+1)+1}{dn_c} = -\frac{1}{d}.$$

Отношение $(n_c+1)/n_c$ максимально при минимальном значении n_c ; таким образом, $n_c = d-1$. Следовательно,

$$\frac{(-d+1)(d-1+1)}{d(d-1)} + \frac{1}{dn_c} \leq \frac{(n+n_c)(n_c+1)+1}{dn_c} < 0, \text{ или}$$

$$-1 < \frac{(n+n_c)(n_c+1)+1}{dn_c} < 0.$$

Из (13), поскольку $(-n_c - 1)/d$ — целое число, а прибавляемая величина находится между 0 и -1 , получим

$$\frac{-n_c - 1}{d} \leq \left\lfloor \frac{mn}{2^p} \right\rfloor + 1 \leq \left\lceil \frac{n}{d} \right\rceil.$$

Для n из диапазона $-n_c - d + 1 \leq n \leq -n_c - 1$

$$\left\lceil \frac{n}{d} \right\rceil = \frac{-n_c - 1}{d}.$$

Следовательно, $\lfloor mn/2^p \rfloor + 1 = \lceil n/d \rceil$, т.е. выполняется (1,6).

Последний случай, $n = -n_c - d$, может осуществиться только для некоторых значений d . Из (3,а) вытекает, что $-n_c - d \leq -2^{W-1}$, поэтому если n принимает указанное значение, то $n = -n_c - d = -2^{W-1}$ и, следовательно, $n_c = 2^{W-1} - d$. Таким образом, $\text{rem}(2^{W-1}, d) = \text{rem}(n_c + d, d) = d - 1$ (т.е. d является делителем $2^{W-1} + 1$).

Для рассматриваемого случая $n = -n_c - d$ формула (6) имеет решение $p = W - 1$ (наименьшее возможное значение p), поскольку при этом

$$n_c (d - \text{rem}(2^p, d)) = (2^{W-1} - d)(d - \text{rem}(2^{W-1}, d)) =$$

$$= (2^{W-1} - d)(d - (d - 1)) = 2^{W-1} - d < 2^{W-1} = 2^p.$$

Тогда из (5)

$$m = \frac{2^{W-1} + d - \text{rem}(2^{W-1}, d)}{d} = \frac{2^{W-1} + d - (d - 1)}{d} = \frac{2^{W-1} + 1}{d}.$$

Таким образом,

$$\left\lfloor \frac{mn}{2^p} \right\rfloor + 1 = \left\lfloor \frac{2^{W-1} + 1}{d} \frac{-2^{W-1}}{2^{W-1}} \right\rfloor + 1 = \left\lfloor \frac{-2^{W-1} - 1}{d} \right\rfloor + 1 =$$

$$= \left\lfloor \frac{-2^{W-1} - d}{d} \right\rfloor + 1 = \left\lceil \frac{-2^{W-1}}{d} \right\rceil = \left\lceil \frac{n}{d} \right\rceil,$$

так что уравнение (1,6) справедливо.

Этим завершается доказательство того, что если m и p вычислены по формулам (5) и (6), то уравнения (1,а) и (1,б) выполняются для всех возможных значений n .

10.5. Знаковое деление на делитель, не превышающий -2

Поскольку знаковое целое деление удовлетворяет соотношению $n + (-d) = -(n + d)$, вполне приемлемой реализацией будет генерация кода для деления $n + |d|$ с последую-

щей командой обращения знака частного. (В этом случае мы не получим верный результат при делении на $d = -2^{n-1}$, но для этой и других отрицательных степеней 2 можно воспользоваться кодом из раздела 10.1, "Знаковое деление на известную степень 2", на с. 231, после которого просто изменить знак полученного результата.) Менять знак делимого при этом не следует из-за того, что делимое может оказаться максимальным по абсолютному значению отрицательным числом.

Однако можно избежать команды обращения знака. Схема такого вычисления имеет следующий вид.

$$q = \left\lfloor \frac{mn}{2^p} \right\rfloor \quad \text{при } n \leq 0$$

$$q = \left\lfloor \frac{mn}{2^p} \right\rfloor + 1 \quad \text{при } n > 0$$

Однако добавление 1 при положительных значениях n неудобно (так как невозможно просто использовать знаковый бит n), так что вместо этого будет выполняться прибавление 1, если значение q отрицательно. Эти действия эквивалентны, поскольку сомножитель m отрицателен (как будет показано далее).

Генерируемый для случая $W = 32$, $d = -7$ код показан ниже.

```
li      M, 0x6DB6DB6D      Магическое число  $-(2^{34}+5)/7 + 2^{32}$ 
mulhs  q, M, n              $q = \text{floor}(M \cdot n / 2^{32})$ 
sub     q, q, n              $q = \text{floor}(M \cdot n / 2^{32}) - n$ 
shlwi  q, q, 2              $q = \text{floor}(q/4)$ 
shri   t, q, 31           Прибавляем 1 к q, если
add     q, q, t             q отрицательно (n положительно)

mulwi  t, q, -7           Вычисляем остаток по формуле
sub     r, n, t             $r = n - q \cdot (-7)$ 
```

Этот код очень напоминает код для деления на +7 с тем отличием, что он использует множитель для деления на +7, но с обратным знаком, команду sub вместо команды add после умножения, а кроме того, команда сдвига shri на 31 позицию использует q , а не n (о чем говорилось ранее). (Заметим, что в случае деления на +7 также можно использовать в качестве операнда q , но при этом код будет иметь меньшую степень параллелизма.) Команда *вычитания* не может привести к переполнению, поскольку операнды имеют один и тот же знак. Тем не менее эта схема работает не всегда! Хотя приведенный код для $W = 32$, $d = -7$ вполне корректен, аналогичные изменения кода деления на 3 для получения кода для деления на -3 приведут к неверному результату при $W = 32$, $n = -2^{31}$.

Рассмотрим ситуацию подробнее.

Пусть даны размер слова $W \geq 3$ и делитель $-2^{n-1} \leq d \leq -2$ и требуется найти наименьшее (по абсолютному значению) целое $-2^n \leq m \leq 0$ и целое $p \geq W$, такие, что

$$\left\lfloor \frac{mn}{2^p} \right\rfloor = \left\lfloor \frac{n}{d} \right\rfloor \quad \text{для } -2^{n-1} \leq n \leq 0 \quad \text{и} \quad (14.a)$$

$$\left\lfloor \frac{mn}{2^p} \right\rfloor + 1 = \left\lfloor \frac{n}{d} \right\rfloor \text{ для } 1 \leq n < 2^{p-1}. \quad (14,6)$$

Поступим так же, как и в случае деления на положительный делитель. Пусть n_c — наибольшее по абсолютной величине отрицательное значение n , такое, что $n_c = kd + 1$ для некоторого целого k . Такая величина n_c существует в силу того, что имеется по крайней мере одно значение $n_c = d + 1$. Это значение можно вычислить, исходя из того, что $n_c = \left\lfloor \frac{-2^{p-1} - 1}{d} \right\rfloor d + 1 = -2^{p-1} + \text{rem}(2^{p-1} + 1, d)$. Величина n_c представляет собой одно из $|d|$ наименьших допустимых значений n , так что

$$-2^{p-1} \leq n_c \leq -2^{p-1} - d - 1 \quad (15,a)$$

и, очевидно,

$$n_c \leq d + 1. \quad (15,b)$$

Поскольку (14,6) должно выполняться при $n = -d$, а (14,a) — для $n = n_c$, то аналогично (4) получим

$$\frac{2^p n_c - 1}{d n_c} < m < \frac{2^p}{d}. \quad (16)$$

Так как m является наибольшим целым числом, удовлетворяющим (16), оно представляет собой ближайшее целое, меньшее $2^p/d$, т.е.

$$m = \frac{2^p - d - \text{rem}(2^p, d)}{d}. \quad (17)$$

Объединяя это выражение с левой частью (16) и упрощая, получим

$$2^p > n_c (d + \text{rem}(2^p, d)). \quad (18)$$

Доказательство пригодности предложенного в (17) и (18) алгоритма и корректности произведения выполняется аналогично доказательству для положительного делителя и здесь не приводится. Трудности возникают только при попытках доказать, что $-2^p \leq m \leq 0$. Для доказательства рассмотрите по отдельности случаи, когда d представляет собой отрицательное число, абсолютное значение которого является степенью двойки, и прочие числа. Для $d = -2^k$ легко показать, что $n_c = -2^{p-1} + 1$, $p = W + k - 1$ и $m = -2^{p-1} - 1$ (которое находится в изначально определяемом диапазоне). Для тех d , которые имеют отличный от -2^k вид, достаточно просто изменить доказательство, приведенное ранее, при рассмотрении положительного делителя.

Для каких делителей $m(-d) \neq -m(d)$?

Под $m(d)$ подразумевается множитель, соответствующий делителю d . Если $m(-d) = -m(d)$, код для деления на отрицательный делитель может быть сгенерирован

путем вычисления множителя для $|d|$, изменения его знака и генерации кода, аналогичного коду из примера “деления на -7 ”, проиллюстрированному ранее.

Сравнивая (18) с (6), а (17) с (5), можно увидеть, что если значение n_c для $-d$ представляет собой значение n_c для d , но с обратным знаком, то $m(-d) = -m(d)$. Следовательно, $m(-d) \neq -m(d)$ может быть только тогда, когда значение n_c , вычисленное для отрицательного делителя, представляет собой наибольшее по абсолютному значению отрицательное число -2^{w-1} . Такие делители представляют собой отрицательные множители $2^{w-1} + 1$. Эти числа встречаются очень редко, что иллюстрируется приведенными ниже разложениями.

$$2^{15} + 1 = 3^2 \cdot 11 \cdot 331$$

$$2^{31} + 1 = 3 \cdot 715827883$$

$$2^{63} + 1 = 3^3 \cdot 19 \cdot 43 \cdot 5419 \cdot 77158673929$$

Для всех этих множителей $m(-d) \neq -m(d)$. Вот набросок доказательства. Для $d > 0$ мы имеем $n_c = 2^{w-1} - d$. Поскольку $\text{rem}(2^{w-1}, d) = d - 1$, значение $p = W - 1$ удовлетворяет (6), а следовательно, этому неравенству удовлетворяет и $p = W$. Однако для $d < 0$ мы имеем $n_c = -2^{w-1}$ и $\text{rem}(2^{w-1}, d) = |d| - 1$. Следовательно, ни $p = W - 1$, ни $p = W$ не удовлетворяют (18), так что $p > W$.

10.6. Встраивание в компилятор

Для того чтобы компилятор мог заменить деление на константу произведением, он должен вычислить для данного делителя d магическое число M и величину сдвига s . Простейший способ состоит в вычислении (6) или (18) для $p = W, W + 1, \dots$ до тех пор, пока условие будет выполняться. Затем из (5) или (17) вычисляется значение m . Магическое число M представляет собой не что иное, как интерпретацию m как знакового целого числа, а $s = p - W$.

Описанная ниже схема обрабатывает положительные и отрицательные d с помощью небольшого дополнительного кода и позволяет избежать арифметических выражений с двойными словами.

Вспомним, что n_c задается следующим образом.

$$n_c = \begin{cases} 2^{w-1} - \text{rem}(2^{w-1}, d) - 1 & \text{при } d > 0 \\ -2^{w-1} + \text{rem}(2^{w-1} + 1, d) & \text{при } d < 0 \end{cases}$$

Следовательно, $|n_c|$ можно вычислить так.

$$t = 2^{w-1} + \begin{cases} 0, & d > 0 \\ 1, & d < 0 \end{cases}$$

$$|n_c| = t - 1 - \text{rem}(t, |d|)$$

Остаток должен вычисляться с использованием беззнакового деления, в соответствии со значениями аргументов. Здесь используется запись $\text{rem}(t, |d|)$, а не эквивалентная ей

$\text{rem}(t, d)$ именно для того, чтобы подчеркнуть, что программа должна работать с двумя положительными (и беззнаковыми) аргументами.

Исходя из (6) и (18), значение p может быть вычислено из соотношения

$$2^p > |n_c|(|d| - \text{rem}(2^p, |d|)), \quad (19)$$

после чего $|m|$ можно вычислить следующим образом (используя (5) и (17)).

$$|m| = \frac{2^p + |d| - \text{rem}(2^p, |d|)}{|d|} \quad (20)$$

Непосредственное вычисление $\text{rem}(2^p, |d|)$ в (19) требует применения “длинного деления” (деление $2W$ -битового делимого на W -битовый делитель и получение W -битовых частного и остатка), причем это длинное деление должно быть *беззнаковым*. Однако имеется путь решения (19), который позволяет избежать использования длинного деления и может быть легко реализован в обычных языках программирования высокого уровня с использованием только W -битовой арифметики. Тем не менее нам потребуются беззнаковое деление и беззнаковое сравнение.

Вычислить $\text{rem}(2^p, |d|)$ можно инкрементно, инициализируя переменные q и r значениями частного и остатка от деления 2^p на $|d|$ при $p = W - 1$, а затем обновляя значения q и r по мере роста p .

В процессе поиска при увеличении значения p на 1 значения q и r изменяются следующим образом (см. теорему D5 (первая формула)).

```
q = 2*q;
r = 2*r;
if (r >= abs(d))
{
    q = q + 1;
    r = r - abs(d);
}
```

Из левой половины неравенства (4) и правой половины (16) вместе с доказанными границами для m следует, что $q = \lfloor 2^p / |d| \rfloor < 2^w$, так что q можно представить W -битовым беззнаковым числом. Кроме того, $0 \leq r < |d|$, так что r можно представить W -битовым знаковым или беззнаковым числом. (*Внимание:* промежуточный результат $2r$ может превысить $2^{w-1} - 1$, поэтому r должно быть беззнаковым числом, и выполняющиеся выше сравнения также должны быть беззнаковыми.)

Далее вычисляется $\delta = |d| - r$. Оба члена разности представимы в виде W -битовых беззнаковых целых чисел (так же, как и разность $1 \leq \delta \leq |d|$), так что данное вычисление не представляет никаких трудностей.

Для того чтобы избежать длинного умножения в (19), перепишем его.

$$\frac{2^p}{|n_c|} > \delta$$

Величина $2^p / |n_c|$ представима в виде W -битового беззнакового целого числа (аналогично (7) из (19) может быть показано, что $2^p \leq 2|n_c| \cdot |d|$, и для $d = -2^{W-1}$, $n_c = -2^{W-1} + 1$ и $p = 2W - 2$, так что $2^p / |n_c| = 2^{2W-2} / (2^{W-1} - 1) < 2^W$ при $W \geq 3$). Эти вычисления так же, как и в случае $\text{rem}(2^p, |d|)$, могут быть инкрементными (при увеличении p). Сравнение для случая $2^p / |n_c| \geq 2^{W-1}$ (что может произойти при больших d) должно быть беззнаковым.

Для вычисления m не требуется непосредственное вычисление (20), при котором может понадобиться длинное деление. Заметим, что

$$\frac{2^p + |d| - \text{rem}(2^p, |d|)}{|d|} = \left\lfloor \frac{2^p}{|d|} \right\rfloor + 1 = q + 1.$$

Проверка завершения цикла $2^p / |n_c| > \delta$ вычисляется достаточно сложно. Величина $2^p / |n_c|$ доступна только в виде частного q_i и остатка r_i . Величина $2^p / |n_c|$ может быть (а может и не быть) целой (эта величина является целой только для $d = 2^{W-2} + 1$ и некоторых отрицательных значений d). Проверка $2^p / |n_c| \leq \delta$ может быть закодирована следующим образом.

$$q_i < \delta \mid (q_i = \delta \ \& \ r_i = 0)$$

Полностью процедура для вычисления M и s для данного d показана в листинге 10.1, в котором представлена программа на языке C для $W = 32$. Здесь есть несколько мест, где может произойти переполнение, однако если его игнорировать, то полученный результат оказывается корректным.

ЛИСТИНГ 10.1. Вычисление магического числа для знакового деления

```
struct ms
{
    int M;      // Магическое число
    int s;      // Величина сдвига
};

struct ms magic(int d)
{
    // d должно удовлетворять условию
    // 2 <= d <= 2**31-1 или -2**31 <= d <= -2
    int p;
    unsigned ad, anc, delta, q1, r1, q2, r2, t;
    const unsigned two31 = 0x80000000; // 2**31
    struct ms mag;
    ad = abs(d);
    t = two31 + ((unsigned)d >> 31);
    anc = t - 1 - t*ad; // Абсолютное значение nc
    p = 31;           // Инициализация p
    q1 = two31/anc;   // Инициализация q1 = 2**p/|nc|
    r1 = two31 - q1*anc; // Инициализация r1 = rem(2**p, |nc|)
    q2 = two31/ad;    // Инициализация q2 = 2**p/|d|
    r2 = two31 - q2*ad; // Инициализация r2 = rem(2**p, |d|)
    do {
        p = p + 1;
        q1 = 2*q1;      // Обновление q1 = 2**p/|nc|.
        r1 = 2*r1;      // Обновление r1 = rem(2**p, |nc|)
```

```

if (r1 >= anc) // Здесь требуется беззнаковое
{
    // сравнение
    q1 = q1 + 1;
    r1 = r1 - anc;
}
q2 = 2*q2;      // Обновление q2 = 2**p/|d|
r2 = 2*r2;      // Обновление r2 = rem(2**p, |d|)
if (r2 >= ad) // Здесь требуется беззнаковое
{
    // сравнение
    q2 = q2 + 1;
    r2 = r2 - ad;
}
delta = ad - r2;
} while(q1 < delta || (q1 == delta && r1 == 0));
mag.M = q2 + 1;
if (d < 0)
    mag.M = -mag.M; // Магическое число и
mag.s = p - 32;    // величина сдвига
return mag;
}

```

Для использования результата этой программы компилятор должен сгенерировать команды `li` и `mulhs`, команду `add` при $d > 0$ и $M < 0$ или `sub` при $d < 0$ и $M > 0$, а также команду `shrsi` при $s > 0$. После этого генерируются команды `shri` и `add`.

В случае $W = 32$ можно избежать обработки отрицательного делителя, если использовать предвычисленный результат для $d = 3$ и $d = 715\,827\,883$, а для остальных отрицательных делителей использовать формулу $m(-d) = -m(d)$. Однако при этом программа из листинга 10.1 не станет значительно короче (если вообще сократится).

10.7. Дополнительные вопросы

ТЕОРЕМА DC2. *Наименьший множитель m нечетен, если p не равно W в принудительном порядке.*

Доказательство. Предположим, что уравнения (1,а) и (1,б) выполняются при наименьшем (не установленном принудительно) значении p и четном m . Очевидно, что тогда m можно разделить на 2, а p уменьшить на 1 и уравнения останутся удовлетворены. Но это противоречит предположению о том, что p — наименьшее значение, при котором выполняются данные уравнения.

Единственность

Магическое число для данного делителя может быть единственным (как, например, в случае $W = 32$, $d = 7$), но зачастую это не так. Более того, эксперименты показывают, что как раз обычно имеется несколько магических чисел для одного делителя. Например, для $W = 32$, $d = 6$ имеется четыре магических числа.

$$M = 715827833 \quad \left((2^{32} + 2) / 6 \right), \quad s = 0$$

$$M = 1431655766 \quad \left((2^{32} + 2) / 3 \right), \quad s = 1$$

$$M = -1431655765 \quad \left((2^{33} + 1) / 3 - 2^{32} \right), \quad s = 2$$

$$M = -1431655764 \quad \left((2^{33} + 4) / 3 - 2^{32} \right), \quad s = 2$$

Тем не менее имеется следующее свойство единственности.

ТЕОРЕМА DC3. Для данного делителя d существует только один множитель m , имеющий минимальное значение p , если p не приравнивается к W в принудительном порядке.

Доказательство. Рассмотрим сначала случай $d > 0$. Разность между верхней и нижней границами в неравенстве (4) составляет $2^p / dn_c$. Мы уже доказали (7), т.е. что если p минимально, то $2^p / dn_c \leq 2$. Таким образом, может быть не более двух значений m , которые удовлетворяют неравенству (4). Пусть m равно минимальному из этих значений, задаваемому соотношением (5); тогда второе допустимое значение — $m + 1$.

Пусть p_0 — наименьшее значение p , для которого $m + 1$ удовлетворяет правой части неравенства (4) (значение p_0 не приравнено к W в принудительном порядке). Тогда

$$\frac{2^{p_0} + d - \text{rem}(2^{p_0}, d)}{d} + 1 < \frac{2^{p_0}}{d} \frac{n_c + 1}{n_c}.$$

Это выражение упрощается до

$$2^{p_0} > n_c (2d - \text{rem}(2^{p_0}, d)).$$

Деление на 2 дает

$$2^{p_0-1} > n_c \left(d - \frac{1}{2} \text{rem}(2^{p_0}, d) \right).$$

Поскольку в соответствии с теоремой D5 на с. 210 $\text{rem}(2^{p_0}, d) \leq 2 \text{rem}(2^{p_0-1}, d)$, получаем

$$2^{p_0-1} > n_c (d - \text{rem}(2^{p_0-1}, d)),$$

что противоречит предположению о минимальности p_0 .

Доказательство для случая $d < 0$ аналогично и здесь не приводится.

Делители с лучшими программами

Программа для $d = 3$, $W = 32$ оказывается короче программы для общего случая, поскольку в ней не требуется выполнение команд `add` и `shrsi` после команды `mulhs`. Возникает вопрос: для каких делителей программа также оказывается укороченной?

Рассмотрим только положительные делители. Итак, требуется найти целые числа m и p , которые удовлетворяют уравнениям (1,а) и (1,б) и для которых выполняются соотношения $p = W$ и $0 \leq m \leq 2^{W-1}$. Поскольку любые целые числа m и p , которые удовлетворяют уравнениям (1,а) и (1,б), должны также удовлетворять неравенству (4), достаточно найти те делители d , для которых (4) имеет решение при $p = W$ и $0 \leq m \leq 2^{W-1}$. Все решения (4) при $p = W$ задаются уравнением

$$m = \frac{2^W + kd - \text{rem}(2^W, d)}{d}, \quad k = 1, 2, 3, \dots$$

Объединяя его с правой частью (4) и упрощая, получим

$$\text{rem}(2^W, d) > kd - \frac{2^W}{n_c}. \quad (21)$$

Наименьшим ограничением на $\text{rem}(2^W, d)$ является ограничение при $k = 1$ и n_c , равном своему минимальному значению 2^{W-2} , так что

$$\text{rem}(2^W, d) > d - 4,$$

т.е. d является делителем $2^W + 1$, $2^W + 2$ или $2^W + 3$.

Посмотрим теперь, какие из данных делителей в действительности имеют оптимальные программы.

Если d является делителем $2^W + 1$, то $\text{rem}(2^W, d) = d - 1$. Тогда решением (6) является $p = W$, так как неравенство превращается в очевидное:

$$2^W > n_c(d - (d - 1)) = n_c,$$

поскольку $n_c < 2^{W-1}$. При вычислении m имеем

$$m = \frac{2^W + d - (d - 1)}{d} = \frac{2^W + 1}{d},$$

и при $d \geq 3$ это значение оказывается меньше, чем 2^{W-1} (d не может быть равно 2, поскольку является делителем $2^W + 1$). Следовательно, все делители $2^W + 1$ имеют оптимальные программы.

Аналогично, если d является делителем $2^W + 2$, то $\text{rem}(2^W, d) = d - 2$. Здесь также решением (6) является $p = W$, поскольку неравенство также превращается в

$$2^W > n_c(d - (d - 2)) = 2n_c,$$

которое, очевидно, является истинным. Вычисление m дает значение

$$m = \frac{2^W + d - (d - 2)}{d} = \frac{2^W + 2}{d},$$

превышающее $2^{W-1} - 1$ для $d = 2$, но не превышающее $2^{W-1} - 1$ при $W \geq 3$, $d \geq 3$ (случай $W = 3$ и $d = 3$ невозможен, поскольку 3 не является делителем $2^3 + 2 = 10$). Следовательно, все делители $2^W + 2$, за исключением числа 2 и дополнительного к нему (дополнительным к 2 делителем является $(2^W + 2)/2$, т.е. число, которое нельзя представить как W -битовое знаковое целое).

Если d является делителем $2^W + 3$, то показать, что оптимальной программы для него нет, можно следующим образом. Поскольку $\text{rem}(2^W, d) = d - 3$, из неравенства (21) имеем

$$n_c < \frac{2^W}{kd-d+3} (2^W + 2)$$

для некоторого $k=1, 2, 3, \dots$. Самое слабое ограничение — при $k=1$, так что должно выполняться $n_c < 2^W/3$.

Из (3,а) $n_c \geq 2^{W-1} - d$ или $d \geq 2^{W-1} - n_c$, следовательно, получаем

$$d > 2^{W-1} - \frac{2^W}{3} = \frac{2^W}{6}.$$

Кроме того, поскольку 2, 3 и 4 не могут быть делителями $2^W + 3$, наименьшим возможным делителем $2^W + 3$ является 5 и, следовательно, наибольший возможный делитель равен $(2^W + 3)/5$. Таким образом, если d является делителем $2^W + 3$ и d имеет оптимальную программу, то

$$\frac{2^W}{6} < d \leq \frac{2^W + 3}{5}.$$

Преобразуя данное неравенство, для дополнительного к d делителя $(2^W + 3)/d$ получаем следующие границы:

$$5 \leq \frac{2^W + 3}{d} < \frac{(2^W + 3) \cdot 6}{2^W} = 6 + \frac{18}{2^W}.$$

Эти границы указывают, что при $W \geq 5$ единственными возможными дополнительными делителями являются 5 и 6. Легко убедиться, что при $W < 5$ делителей $2^W + 3$ не существует. Поскольку 6 не может быть делителем $2^W + 3$, единственным возможным делителем этого числа может быть 5. Таким образом, единственным возможным делителем $2^W + 3$, который может иметь оптимальную программу, является $(2^W + 3)/5$.

Для $d = (2^W + 3)/5$

$$n_c = \left\lfloor \frac{2^{W-1}}{(2^W + 3)/5} \right\rfloor \left(\frac{2^W + 3}{5} \right) - 1,$$

а так как для $W \geq 4$

$$2 < \frac{2^{W-1}}{(2^W + 3)/5} < 2.5,$$

получаем

$$n_c = 2 \left(\frac{2^W + 3}{5} \right) - 1.$$

Эта величина превышает $2^W/3$, так что $d = (2^W + 3)/5$ оптимальной программы не имеет. Поскольку для $W < 4$ величина $2^W + 3$ делителей не имеет, можно сделать вывод о том, что делителей $2^W + 3$ с оптимальной программой не существует.

Резюмируем: все делители $2^W + 1$ и $2^W + 2$, за исключением 2 и $(2^W + 2)/2$, имеют оптимальные программы, и никакие другие числа оптимальных программ не имеют. Более того, приведенное выше доказательство показывает, что алгоритм из листинга 10.1 всегда дает оптимальную программу, если таковая существует.

Рассмотрим частные случаи значений W , равные 16, 32 и 64. Соответствующие разложения на множители показаны ниже.

$$\begin{aligned}
 2^{16} + 1 &= 65537 \text{ (простое)} \\
 2^{16} + 2 &= 2 \cdot 3^2 \cdot 11 \cdot 331 \\
 2^{32} + 1 &= 641 \cdot 6\,700\,417 \\
 2^{32} + 2 &= 2 \cdot 3 \cdot 715\,827\,883 \\
 2^{64} + 1 &= 247\,177 \cdot 67\,280\,421\,310\,721 \\
 2^{64} + 2 &= 2 \cdot 3^3 \cdot 19 \cdot 43 \cdot 5419 \cdot 77\,158\,673\,929
 \end{aligned}$$

Следовательно, для $W = 16$ имеется 20 делителей, для которых существует оптимальная программа. Из них меньше 100 делители 3, 6, 9, 11, 18, 22, 33, 66 и 99.

Для $W = 32$ есть шесть таких делителей: 3, 6, 641, 6 700 417, 715 827 883, 1 431 655 766.

Для $W = 64$ есть 126 таких делителей. Из них меньше 100 делители 3, 6, 9, 18, 19, 27, 38, 43, 54, 57 и 86.

10.8. Беззнаковое деление

Беззнаковое деление на величину, представляющую собой степень двойки, само собой, реализуется с помощью единственной команды *сдвига вправо*, а остаток — с помощью команды *и с непосредственно задаваемым операндом*.

Может показаться, что обработка других делителей должна быть несложной: нужно просто использовать результаты для знакового деления с $d > 0$, опуская две команды, которые добавляют 1 при отрицательном частном. Однако вы увидите, что некоторые детали беззнакового деления в действительности несколько сложнее.

Беззнаковое деление на 3

Начнем рассмотрение беззнакового деления на другие числа с беззнакового деления на 3 на 32-разрядной машине. Поскольку делимое n теперь может достигать по величине $2^{32} - 1$, множитель $(2^{32} + 2)/3$ оказывается неадекватен (член-ошибка $2n/(3 \cdot 2^{32})$ может превысить $1/3$). Однако можно использовать множитель $(2^{33} + 1)/3$. Соответствующий код имеет следующий вид.

```

11      M, 0хААААААВ      Загрузка магического числа (2**33+1)/3
mulhu  q, M, n           q = floor(M*n/2**32)
shri   q, q, 1
    
```

mul	t, q, 3	Вычисление остатка по формуле
sub	r, n, t	$r = n - q * 3$

В этом случае нам требуется команда mulhu, которая возвращает старшие 32 бита беззнакового 64-битового произведения.

Чтобы убедиться в корректности приведенного кода, заметим, что он вычисляет следующее значение.

$$q = \left\lfloor \frac{2^{33} + 1}{3} \frac{n}{2^{33}} \right\rfloor = \left\lfloor \frac{n}{3} + \frac{n}{3 \cdot 2^{33}} \right\rfloor$$

При $0 \leq n < 2^{32}$ выполняется неравенство $0 \leq n / (3 \cdot 2^{33}) < 1/3$, так что в соответствии с теоремой D4 $q = \lfloor n/3 \rfloor$.

При вычислении остатка может произойти переполнение при выполнении команды *умножения на непосредственно заданное значение*, если операнды рассматриваются как знаковые целые числа, но если считать операнды и результат беззнаковыми, то переполнение при выполнении данной команды невозможно. Переполнение невозможно и при выполнении команды *вычитания*, так как результат находится в диапазоне от 0 до 2, поэтому получаемый таким образом остаток корректен.

Беззнаковое деление на 7

При беззнаковом делении на 7 на 32-разрядном компьютере множители $(2^{32} + 3)/7$, $(2^{33} + 6)/7$ и $(2^{34} + 5)/7$ оказываются неадекватными, поскольку дают слишком большой член-ошибку. Можно использовать множитель $(2^{35} + 3)/7$, но он слишком велик для представления 32-битовым беззнаковым числом. Умножение на это число можно выполнить посредством умножения на $(2^{35} + 3)/7 - 2^{32}$ с последующей коррекцией путем сложения. Соответствующий код имеет следующий вид.

li	M, 0x24924925	Магическое число $(2^{35} + 3)/7 - 2^{32}$
mulhu	q, M, n	$q = \text{floor}(M * n / 2^{32})$
add	q, q, n	Может вызвать переполнение (перенос)
shrx	q, q, 3	Сдвиг вправо с битом переноса

mul	t, q, 7	Вычисление остатка по формуле
sub	r, n, t	$r = n - q * 7$

Здесь возникает небольшая проблема: команда add может вызвать переполнение. Для того чтобы обойти эту ситуацию, была придумана новая команда *расширенного сдвига вправо на непосредственно заданную величину* (shrx), которая рассматривает бит переноса команды add и 32 бита регистра q как единую 33-битовую величину и выполняет ее сдвиг вправо с заполнением освободившихся разрядов нулевыми битами. В семействе процессоров Motorola 68000 эти действия можно реализовать с помощью двух команд: *расширенного циклического сдвига* вправо на один бит с последующим *беззнаковым сдвигом вправо* на 3 бита (команда rotx на самом деле использует X-бит, но команда add присваивает ему то же значение, что и биту переноса). На большинстве

компьютеров для выполнения этих действий требуется большее количество команд. Например, на PowerPC требуется выполнение трех команд: сброс правых трех битов q , прибавление бита переноса к q и циклический сдвиг вправо на три позиции.

При той или иной реализации команды `shrx i` приведенный выше код вычисляет

$$q = \left\lfloor \left(\left(\left(\frac{2^{35} + 3}{7} - 2^{32} \right) \frac{n}{2^{32}} \right) + n \right) / 2^3 \right\rfloor = \left\lfloor \frac{n}{3} + \frac{3n}{7 \cdot 2^{35}} \right\rfloor.$$

При $0 \leq n < 2^{32}$ выполняется неравенство $0 \leq 3n/(7 \cdot 2^{35}) < 1/7$, так что в соответствии с теоремой D4 $q = \lfloor n/7 \rfloor$.

Гранлунд (Granlund) и Монтгомери (Montgomery) [39] предложили остроумную схему, позволяющую избежать применения команды `shrx i`. Она требует того же количества команд, что и рассмотренная выше схема, но использует только элементарные команды, которые есть практически у каждого компьютера: переполнения при этом оказываются просто невозможны. Схема использует следующее тождество.

$$\left\lfloor \frac{q+n}{2^p} \right\rfloor = \left\lfloor \left(\left\lfloor \frac{n-q}{2} \right\rfloor + q \right) / 2^{p-1} \right\rfloor, \quad p \geq 1$$

Применим данное тождество к нашей задаче, полагая $q = \lfloor Mn/2^{32} \rfloor$, где $0 \leq M < 2^{32}$. Вычитание не может вызвать переполнения, поскольку

$$0 \leq q = \left\lfloor \frac{Mn}{2^{32}} \right\rfloor \leq n,$$

поэтому очевидно, что $0 \leq n - q < 2^{32}$. Сложение также не может вызвать переполнения, поскольку

$$\left\lfloor \frac{n-q}{2} \right\rfloor + q = \left\lfloor \frac{n-q}{2} + q \right\rfloor = \left\lfloor \frac{n+q}{2} \right\rfloor$$

и $0 \leq n, q < 2^{32}$.

Использование предложенных идей приводит к следующему коду для беззнакового деления на 7.

```
li    M,0x24924925    Магическое число (2**35+3)/7 - 2**32
mulhu q,M,n          q = floor(M*n/2**32)
sub   t,n,q          t = n - q
shri  t,t,1         t = (n - q)/2
add   t,t,q         t = (n - q)/2 + q = (n + q)/2
shri  q,t,2         q = (n+Mn/2**32)/8 = floor(n/7)

mulu  t,q,7         Вычисление остатка по формуле
sub   r,n,t         r = n - q*7
```

Чтобы эта схема работала, величина сдвига в гипотетической команде `shrx i` должна быть больше 0. Можно показать, что если $d > 1$ и множитель $m \geq 2^{32}$ (т.е. необходима команда `shrx i`), то величина сдвига больше 0.

10.9. Беззнаковое деление на делитель, не меньший 1

Для данного размера слова $W \geq 1$ и делителя $1 \leq d < 2^W$ требуется найти наименьшее целое m и целое p , такие, что

$$\left\lfloor \frac{mn}{2^p} \right\rfloor = \left\lfloor \frac{n}{d} \right\rfloor \text{ для } 0 \leq n < 2^W, \quad (22)$$

при этом $0 \leq m < 2^{W+1}$ и $p \geq W$.

В случае беззнакового деления магическое число M определяется как

$$M = \begin{cases} m, & 0 \leq m < 2^W, \\ m - 2^W, & 2^W \leq m < 2^{W+1}. \end{cases}$$

Поскольку (22) должно выполняться для $n = d$, $\lfloor md/2^p \rfloor = 1$ или

$$\frac{md}{2^p} \geq 1. \quad (23)$$

Пусть, как и в случае знакового деления, n_c представляет собой наибольшее значение n , такое, что $\text{rem}(n_c, d) = d - 1$. Его можно вычислить как $n_c = \lfloor 2^W/d \rfloor d - 1 = 2^W - \text{rem}(2^W, d) - 1$. Тогда

$$2^W - d \leq n_c \leq 2^W - 1 \quad (24,а)$$

и

$$n_c \geq d - 1. \quad (24,б)$$

Это означает, что $n_c \geq 2^{W-1}$.

Поскольку (22) должно выполняться при $n = n_c$,

$$\left\lfloor \frac{mn_c}{2^p} \right\rfloor = \left\lfloor \frac{n_c}{d} \right\rfloor = \frac{n_c - (d-1)}{d}$$

или

$$\frac{mn_c}{2^p} < \frac{n_c + 1}{d}.$$

Объединение этого неравенства с (23) дает

$$\frac{2^p}{d} \leq m < \frac{2^p}{d} \frac{n_c + 1}{n_c}. \quad (25)$$

Поскольку m представляет собой наименьшее целое, удовлетворяющее неравенству (25), оно должно быть ближайшим целым числом, большим или равным $2^p/d$, т.е.

$$m = \left\lceil \frac{2^p + d - 1 - \text{rem}(2^p - 1, d)}{d} \right\rceil. \quad (26)$$

Комбинируя это уравнение с правой частью (25) и упрощая, получим

$$2^p > n_c \left(d - 1 - \text{rem}(2^p - 1, d) \right). \quad (27)$$

Беззнаковый алгоритм

Таким образом, алгоритм состоит в поиске методом проб и ошибок минимального $p \geq W$, удовлетворяющего неравенству (27), после чего из уравнения (26) вычисляется m . Это наименьшее возможное значение m , удовлетворяющее (22) при $p \geq W$. Как и в случае знакового деления, если неравенство (27) выполняется для некоторого значения p , то оно выполняется и для всех больших значений p . Доказательство этого факта, по сути, ничем не отличается от доказательства теоремы DC1, только вместо первой формулы из теоремы D5 в доказательстве используется вторая.

Доказательство пригодности беззнакового алгоритма

Покажем, что (27) всегда имеет решение и что $0 \leq m < 2^{W+1}$.

Поскольку для любого неотрицательного числа x имеется степень 2, большая x , но не превышающая $2x + 1$, из (27) получаем

$$n_c \left(d - 1 - \text{rem}(2^p - 1, d) \right) < 2^p \leq 2n_c \left(d - 1 - \text{rem}(2^p - 1, d) \right) + 1.$$

Так как $0 \leq \text{rem}(2^p - 1, d) \leq d - 1$,

$$1 \leq 2^p \leq 2n_c(d - 1) + 1. \quad (28)$$

С учетом того, что $n_c, d \leq 2^W - 1$, это неравенство превращается в

$$1 \leq 2^p \leq 2(2^W - 1)(2^W - 2) + 1$$

или

$$0 \leq p \leq 2W. \quad (29)$$

Таким образом, неравенство (27) всегда имеет решение.

Если p не приравнивается принудительно к W , то из (25) и (28) следует

$$\frac{1}{d} \leq m < \frac{2n_c(d - 1) + 1}{d} \frac{n_c + 1}{n_c},$$

$$1 \leq m < \frac{2d - 2 + 1/n_c}{d} (n_c + 1),$$

$$1 \leq m < 2(n_c + 1) \leq 2^{W+1}.$$

Если же p в принудительном порядке приравнено к W , то из (25) вытекает

$$\frac{2^W}{d} \leq m < \frac{2^W}{d} \frac{n_c + 1}{n_c}.$$

В силу того, что $1 \leq d \leq 2^W - 1$ и $n_c \geq 2^{W-1}$,

$$\frac{2^W}{2^W - 1} \leq m < \frac{2^W}{1} \frac{2^{W-1} + 1}{2^{W-1}},$$

$$2 \leq m \leq 2^W + 1.$$

Следовательно, в любом случае m находится в пределах схемы, проиллюстрированной примером беззнакового деления на 7.

Доказательство корректности произведения в случае беззнакового деления

Покажем, что если p и m вычислены на основании формул (27) и (26), то выполняется уравнение (22).

Легко показать, что уравнение (26) и неравенство (27) приводят к неравенству (25), которое практически тождественно неравенству (4); остальная часть доказательства, по сути, идентична доказательству для знакового деления при $n \geq 0$.

10.10. Встраивание в компилятор при беззнаковом делении

В реализации алгоритма, основанного на непосредственных вычислениях, использованных в доказательстве, имеются трудности. Хотя, как доказано выше, $p \leq 2W$, случай $p = 2W$ не исключается (например, для $d = 2^W - 2$, $W \geq 4$). Если $p = 2W$, вычислить m становится сложно в силу того, что делимое в (26) не размещается в $2W$ -битовом слове.

Однако реализация этих действий возможна при использовании методики "инкрементного деления и взятия остатка", уже рассматривавшейся нами ранее и примененной для данного случая в алгоритме, который приведен в листинге 10.2 (для случая $W = 32$). В этом алгоритме, помимо магического числа и величины сдвига, возвращается индикатор необходимости генерации команды `add` (в случае знакового деления необходимость добавления этой команды распознавалась по тому, что M и d имели разные знаки).

Листинг 10.2. Вычисление магического числа для случая беззнакового деления

```
struct mu
{
    unsigned M; // Магическое число
    int a; // Индикатор "add"
    int s; // Величина сдвига
};

struct mu magicu(unsigned d)
{
    // d должно быть в границах 1 <= d <= 2**32-1
    int p;
    unsigned nc, delta, q1, r1, q2, r2;
    struct mu magu;
    magu.a = 0; // Инициализация индикатора "add"
    nc = -1 - (-d)&d; // Используется беззнаковая арифметика
    p = 31; // Инициализация p.
    q1 = 0x80000000/nc; // Инициализация q1=2**p/nc
    r1 = 0x80000000 - q1*nc; // Инициализация r1=rem(2**p,nc)
```

```

q2 = 0x7FFFFFFF/d; // Инициализация q2=(2**p-1)/d
r2 = 0x7FFFFFFF - q2*d; // Инициализация r2=rem(2**p-1,d)
do {
    p = p + 1;
    if (r1 >= nc - r1)
    {
        q1 = 2*q1 + 1; // Коррекция q1
        r1 = 2*r1 - nc; // Коррекция r1
    }
    else
    {
        q1 = 2*q1;
        r1 = 2*r1;
    }
    if (r2 + 1 >= d - r2)
    {
        if (q2 >= 0x7FFFFFFF) magu.a = 1;
        q2 = 2*q2 + 1; // Коррекция q2
        r2 = 2*r2 + 1 - d; // Коррекция r2
    }
    else
    {
        if (q2 >= 0x80000000) magu.a = 1;
        q2 = 2*q2;
        r2 = 2*r2 + 1;
    }
    delta = d - 1 - r2;
} while (p < 64 &&
        (q1 < delta ||
         (q1 == delta && r1 == 0)));
magu.M = q2 + 1; // Возвращаемое магическое число
magu.s = p - 32; // и величина сдвига
return magu; // (magu.a установлено ранее)
}

```

Приведем ряд ключевых моментов в понимании этого алгоритма.

- В ряде мест алгоритма может произойти беззнаковое переполнение, которое должно быть проигнорировано.
- $n_c = 2^p - \text{rem}(2^p, d) - 1 = (2^p - 1) - \text{rem}(2^p - d, d)$.
- Частное и остаток от деления 2^p на n_c не могут быть подкорректированы так же, как и в алгоритме из листинга 10.1, поскольку здесь величина $2*r1$ может вызвать переполнение. Следовательно, алгоритм должен выполнить проверку “if (r1>=nc-r1)” вместо более естественной проверки “if (2*r1>=nc)”. Такое же замечание применимо и к вычислению частного и остатка от деления $2^p - 1$ на d .
- $0 \leq \delta \leq d - 1$, так что δ представимо в виде 32-битового беззнакового целого числа.
- $m = (2^p + d - 1 - \text{rem}(2^p - 1, d)) / d = \lfloor (2^p - 1) / d \rfloor + 1 = q_1 + 1$.
- В программе не выполняется явное вычитание 2^p для случая, когда магическое число M превышает $2^p - 1$; это вычитание происходит в случае, когда вычисление q_2 вызывает переполнение.

- Индикатор `magu.a` необходимости команды `add` из-за переполнения не может быть установлен путем непосредственного сравнения M с 2^{32} или $q2$ с $2^{32} - 1$. Вместо этого программа проверяет $q2$ до того, как может произойти переполнение. Если $q2$ достигнет величины $2^{32} - 1$, так что M станет не менее 2^{32} , то индикатору `magu.a` будет присвоено значение 1; в противном случае это значение останется равным 0.
- Неравенство (27) эквивалентно неравенству $2^p/n_c > \delta$.
- Проверка $p < 64$ в условии цикла необходима постольку, поскольку без нее переполнение $q1$ может привести к тому, что будет выполнено слишком большое количество итераций и будет получен неверный результат.

Чтобы воспользоваться результатами этой программы, компилятор должен генерировать команды `li` и `mulhu` и, в случае, если индикатор `a` применения команды `add` равен 0, генерировать команду `shri` на s бит (если $s > 0$), как проиллюстрировано в примере из раздела “Беззнаковое деление на 3” на с. 253. Если $a = 1$ и машина оснащена командой `shrx`, то компилятор должен генерировать `add` и `shri` на s бит, как проиллюстрировано в примере из раздела “Беззнаковое деление на 7” на с. 254. Если $a = 1$, но машина не имеет команды `shrx`, используется пример, приведенный на с. 255: генерируется `sub`, `shri` на 1 бит, `add` и наконец `shri` на $s-1$ бит (если $s-1 > 0$; в этом случае s не может быть равным 0, за исключением тривиального случая деления на 1, который, как мы предполагаем, компилятор просто удалит).

10.11. Дополнительные вопросы (беззнаковое деление)

ТЕОРЕМА DC2U. *Наименьший множитель t нечетен, если p принудительно не присваивается значение W .*

ТЕОРЕМА DC3U. *Для данного делителя d существует только один множитель t , имеющий минимальное значение p , если p не приравнивается к W в принудительном порядке.*

Доказательства этих теорем практически идентичны доказательствам соответствующих теорем для знакового деления.

Делители с лучшими программами (беззнаковое деление)

Для того чтобы найти в случае беззнакового деления делители (если таковые имеются) с оптимальными программами поиска частного, состоящими из двух команд (`li`, `mulhu`), можно выполнить почти такой же анализ, как и для случая знакового деления (см. выше раздел “Делители с лучшими программами” на с. 244). В результате анализа выясняется, что такими делителями являются делители чисел 2^n и $2^n + 1$ (за исключением $d = 1$). Для распространенных размеров слов таких нетривиальных делителей оказывается очень мало. При использовании 16-битовых слов таких делителей нет вовсе;

для 32-битовых слов их всего два — 641 и 6700417. Два таких делителя (274177 и 67280421310721) есть и при работе с 64-битовыми словами.

Случай $d = 2^k$, $k = 1, 2, \dots$ заслуживает отдельного упоминания. В этом случае рассмотренный алгоритм *magic* из листинга 10.2 дает нам $p = W$ (принудительное присвоение) и $m = 2^{32-k}$. Это минимальное значение m , но не минимальное значение M . Наилучший код получается при использовании $p = W + k$. Тогда $m = 2^k$, M равно 0, $a = 1$, $s = k$. Сгенерированный код включает умножение на 0 и может быть упрощен до одной команды *сдвига вправо на величину k*. На практике делители, представляющие собой степень 2, стоит рассматривать как особый случай и не привлекать для их обработки программу *magic*. (Этого не происходит при знаковом делении, поскольку в этом случае m не может быть степенью 2. *Доказательство*: при $d > 0$ неравенство (4), будучи скомбинированным с неравенством (3,б), дает $d - 1 < 2^p / m < d$. Следовательно, $2^p / m$ не может быть целым числом. При $d < 0$ аналогичный результат получается при комбинировании неравенств (16) и (15,б).)

При беззнаковом делении, когда компьютер не имеет команды *shrx*, код для случая $m \geq 2^k$ оказывается значительно хуже, чем код для $m < 2^k$. Поэтому возникает вопрос, как часто приходится иметь дело с большими множителями. Для размера слова 32 бита среди целых чисел, не превышающих 100, имеется 31 “плохой” делитель: 1, 7, 14, 19, 21, 27, 28, 31, 35, 37, 38, 39, 42, 45, 53, 54, 55, 56, 57, 62, 63, 70, 73, 74, 76, 78, 84, 90, 91, 95 и 97.

Использование знакового деления вместо беззнакового и наоборот

Если ваш компьютер не оснащен командой *mulhu*, но имеет команду *mulhs* (или команду знакового длинного умножения), то можно воспользоваться приемом, о котором говорилось в разделе “Преобразование знакового и беззнакового произведений одно в другое” на с. 200.

В этом разделе была приведена последовательность из семи команд, которая позволяет получить результат выполнения команды *mulhu* из команды *mulhs*. Однако в нашей ситуации этот прием упрощается, поскольку магическое число M известно заранее, так что компилятор может протестировать старший бит заранее и сгенерировать в соответствии с его значением ту или иную последовательность действий после выполнения команды “*mulhu q, M, n*”. Здесь t обозначает временный регистр.

	$M_{31} = 0$	$M_{31} = 1$
	<i>mulhs q, M, n</i>	<i>mulhs q, M, n</i>
	<i>shrsi t, n, 31</i>	<i>shrsi t, n, 31</i>
	<i>and t, t, M</i>	<i>and t, t, M</i>
	<i>add q, q, t</i>	<i>add t, t, n</i>
		<i>add q, q, t</i>

С учетом других команд, используемых вместе с *mulhu*, всего для получения частного и остатка от беззнакового деления на константу на компьютере, не оснащенном беззнаковым умножением, требуется от шести до восьми команд.

Этот прием можно обратить для получения команды `mulhs` из команды `mulhu`. При этом получается практически такой же код, только команда `mulhs` заменяется командой `mulhu`, а последние команды `add` в обоих столбцах заменяются командами `sub`.

Более простой беззнаковый алгоритм

Отказ от требования минимальности магического числа приводит к более простому алгоритму. Вместо (27) можно использовать

$$2^p \geq 2^m (d - 1 - \text{rem}(2^p - 1, d)), \quad (30)$$

а затем вычислить m , как и ранее, по формуле (26).

Должно быть очевидно, что формально этот алгоритм корректен (т.е. вычисленное значение m удовлетворяет уравнению (22)), поскольку единственное его отличие от предыдущего алгоритма состоит в том, что он вычисляет значение p , которое для некоторых значений d не обязательно будет большим. Можно доказать, что значение m , вычисленное по (30) и (26), меньше 2^{m+1} . Здесь доказательство опускается и просто приводится конечный результат — алгоритм, представленный в листинге 10.3.

ЛИСТИНГ 10.3. Упрощенный алгоритм вычисления магического числа для беззнакового деления

```
struct mu
{
    unsigned m; // Магическое число
    int a;      // Индикатор "add"
    int s;      // Величина сдвига
};

struct mu magicu2(unsigned d)
{
    // d должно быть в границах 1 <= d <= 2**32-1
    int p;
    unsigned p32, q, r, delta;
    struct mu magu;
    magu.a = 0; // Инициализация индикатора "add"
    p = 31; // Инициализация p
    q = 0x7FFFFFFF/d; // Инициализация q = (2**p-1)/d
    r = 0x7FFFFFFF - q*d; // Инициализация r = rem(2**p-1,d)
    do {
        p = p + 1;
        if (p == 32) p32 = 1; // p32 = 2**(p-32).
        else p32 = 2*p32;
        if (r + 1 >= d - r)
        {
            if (q >= 0x7FFFFFFF) magu.a = 1;
            q = 2*q + 1; // Коррекция q
            r = 2*r + 1 - d; // Коррекция r
        }
        else
        {
            if (q >= 0x80000000) magu.a = 1;
            q = 2*q;
        }
    } while (1);
    return magu;
}
```

```

    r = 2*r + 1;
  }
  delta = d - 1 - r;
} while (p < 64 && p32 < delta);
magu.M = q + 1; // Возвращаемое магическое число
magu.s = p - 32; // и величина сдвига
return magu;    // (magu.a установлено ранее)
}

```

Алверсон (Alverson) в [4] предложил гораздо более простой алгоритм, который рассматривается в следующем разделе, но он иногда дает очень большие значения m . Главное в алгоритме *magicu2* в том, что он почти всегда дает минимальное значение m при $d \leq 2^{n-1}$. Если размер слова равен 32 битам, наименьший делитель, для которого *magicu2* не дает минимального множителя, равен 102 807 (в этом случае *magicu* дает значение m , равное 2737896999, а *magicu2* — значение 5475793997).

Существует аналог алгоритма *magicu2* и для знакового деления на положительные делители, но он плохо работает для знакового деления на произвольные делители.

10.12. Применение к модульному делению и делению с округлением к меньшему значению

Может показаться, что преобразовать в умножение модульное деление на константу или деление с округлением к меньшему значению — задача еще более простая, чем при делении с отсечением, и для этого достаточно всего лишь убрать шаг “добавить 1, если делимое отрицательно”. Однако это не так. Предложенные выше методы не применимы к другим типам деления путем простых и очевидных преобразований. Вероятно, разработанное преобразование такого типа будет включать изменение множителя m в зависимости от знака делимого.

10.13. Другие похожие методы

Вместо кодирования алгоритма *magic* можно воспользоваться таблицей с магическими числами и величинами сдвигов для некоторых небольших делителей. Делители, равные табличным значениям, умноженным на степень 2, легко обрабатываются следующим образом.

1. Подсчитывается количество завершающих нулевых битов в d , которое далее обозначается как k .
2. В качестве аргумента поиска в таблице используется значение $d/2^k$ (сдвиг вправо на k позиций).
3. Используем магическое число, найденное в таблице.
4. Используем найденную в таблице величину сдвига, увеличенную на k .

Таким образом, если в таблице имеются делители 3, 5, 25, то могут быть обработаны делители 6, 10, 100 и т.д.

Такая процедура обычно дает наименьшее магическое число, хотя и не всегда. Если размер слова равен 32 бита, то наименьший положительный делитель, когда процедура дает наименьшее магическое число, равен 334972. Данный метод при этом дает значения $m = 3361176179$ и $s = 18$; минимальное же магическое число для данного делителя — 840294045, а величина сдвига — 16. Процедура также дает неминимальный результат при $d = -6$. В обоих приведенных случаях снижается качество генерируемого кода деления.

Алверсон [4] является первым известным автором, который указал на корректность работы описываемого далее метода для всех делителей. Используя наши обозначения, можно сказать, что его метод для беззнакового целого деления на d состоит в использовании величины сдвига $p = W + \lceil \log_2 d \rceil$ и множителя $m = \lceil 2^p / d \rceil$ с последующим выполнением деления $n \div d = \lfloor mn / 2^p \rfloor$ (т.е. путем умножения и сдвига вправо). Он доказал, что множитель m меньше 2^{W+1} и что этот метод даст точное значение частного для всех n , выражаемых с помощью W -битового числа.

Метод Алверсона представляет собой более простой вариант нашего, в котором для определения p не используется поиск методом проб и ошибок, а следовательно, он в большей степени подходит для аппаратной реализации, в чем и состоит его главное назначение. Однако его множитель m всегда оказывается не меньшим, чем 2^W , а потому при программной реализации всегда требуется код, проиллюстрированный в примере деления на 7 (т.е. всегда требуются либо команды `add` и `shrx1`, либо четыре альтернативные команды). Поскольку большинство небольших делителей могут быть обработаны с помощью множителей, меньших 2^W , представляется разумным рассмотреть эти случаи.

В случае знакового деления Алверсон предложил искать множитель для $|d|$ и длины слова $W - 1$ (тогда $2^{W-1} \leq m < 2^W$), умножать на него делимое n , если операнды имеют противоположные знаки, изменять знак полученного результата. (Множитель должен быть таким, чтобы давать корректный результат для делимого 2^{W-1} , которое представляет собой абсолютное значение максимального отрицательного числа). Похоже, что такое предложение может дать код, лучший по сравнению с кодом для множителя $m \geq 2^W$. Применяя данный прием к знаковому делению на 7, получим следующий код (для того чтобы избежать ветвления, в нем использовано соотношение $-x = \bar{x} + 1$).

```
abs    an, n
li     M, 0x92492493    Магическое число (2**34+5)/7
mulhu  q, M, an         q = floor(M*an/2**32)
shri   q, q, 2
shrsi  t, n, 31        Эти три команды
xor     q, q, t         изменяют знак q,
sub     q, q, t         если n отрицательно
```

Этот код не столь хорош, как тот, который был предложен для знакового деления на 7 ранее (7 и 6 команд соответственно), но он может оказаться полезным на компьютере, на котором есть команды `abs` и `mulhu` и нет команды `mulhs`.

10.14. Некоторые магические числа

ТАБЛИЦА 10.1. НЕКОТОРЫЕ МАГИЧЕСКИЕ ЧИСЛА ДЛЯ $W = 32$

d	Знаковое число		Беззнаковое число		
	M (шестнадцатеричное)	s	M (шестнадцатеричное)	a	s
-5	99999999	1			
-3	55555555	1			
-2 ^k	7FFFFFFF	k-1			
1	-	-	0	1	0
2 ^k	80000001	k-1	2 ^{32-k}	0	0
3	55555556	0	AAAAAAAB	0	1
5	66666667	1	CCCCCCCD	0	2
6	2AAAAAAB	0	AAAAAAAB	0	2
7	92492493	2	24924925	1	3
9	38E38E39	1	38E38E39	0	1
10	66666667	2	CCCCCCCD	0	3
11	2E8BA2E9	1	BA2E8BA3	0	3
12	2AAAAAAB	1	AAAAAAAB	0	3
25	51EB851F	3	51EB851F	0	3
125	10624DD3	3	10624DD3	0	3
625	68DB8BAD	8	D1B71759	0	9

ТАБЛИЦА 10.2. НЕКОТОРЫЕ МАГИЧЕСКИЕ ЧИСЛА ДЛЯ $W = 64$

d	Знаковое число		Беззнаковое число		
	M (шестнадцатеричное)	s	M (шестнадцатеричное)	a	s
-5	99999999 99999999	1			
-3	55555555 55555555	1			
-2 ^k	7FFFFFFF FFFFFFFF	k-1			
1	-	-	0	1	0
2 ^k	80000000 00000001	k-1	2 ^{64-k}	0	0
3	55555555 55555556	0	AAAAAAAAB	0	1
5	66666666 66666667	1	CCCCCCCC CCCCCCD	0	2
6	2AAAAAAA AAAAAAAB	0	AAAAAAAAB	0	2
7	49249249 24924925	1	24924924 92492493	1	3
9	1C71C71C 71C71C72	0	E38E38E3 8E38E38F	0	3
10	66666666 66666667	2	CCCCCCCC CCCCCCD	0	3
11	2E8BA2E8 BA2E8BA3	1	2E8BA2E8 BA2E8BA3	0	1
12	2AAAAAAA AAAAAAAB	1	AAAAAAAAB	0	3
25	A3D70A3D 70A3D70B	4	47AE147A E147AE15	1	5
125	20C49BA5 E353F7CF	4	0624DD2F 1A9FBE77	1	7
625	346DC5D6 3886594B	7	346DC5D6 3886594B	0	7

10.15. Простой код на языке программирования Python

Вычисление магического числа значительно упрощается, если вычисления не ограничены тем же размером слова, что и среда, в которой вычисленное магическое число будет использоваться. Например, в случае беззнакового деления в Python достаточно просто вычислить n_c , а затем (27) и (26), как описано в разделе 10.9. Соответствующая функция показана в листинге 10.4.

Листинг 10.4. Код на языке программирования Python вычисления магического числа для беззнакового деления

```
def magicgu(nmax, d):
    nc = (nmax//d)*d - 1
    nbits = int(log(nmax, 2)) + 1
    for p in range(0, 2*nbits + 1):
        if 2**p > nc*(d - 1 - (2**p - 1)%d):
            m = (2**p + d - 1 - (2**p - 1)%d)//d
            return (m, p)
    print "Не могу найти p, что-то не так".
    sys.exit(1)
```

Эта функция получает максимальное значение делимого n_{\max} и делителя d и возвращает пару целых чисел: магическое число m и величину сдвига p . Для деления делимого x на d следует умножить x на m , а затем сдвинуть (полную длину) произведения вправо на p бит.

Эта программа является более обобщенной, чем прочие программы из данной главы, по двум направлениям: (1) в ней определено максимальное значение делимого (n_{\max}), а не количество битов, требуемых для делимого, и (2) программа может использоваться для произвольно больших делимого и делителя. Преимущество указания максимального значения делимого заключается в возможности получить меньшее магическое число, чем при использовании в качестве этого максимального значения очередной степени двойки, уменьшенной на 1. Предположим, например, что максимальное значение делимого равно 90, а делитель равен 7. Тогда функция `magicgu` вернет пару (37, 8), означающую, что магическое число равно 37 (6-битовое число), а величина сдвига составляет 8 бит. Но если запросить, какое магическое число обеспечивает деление до делимого, равного 127, то результатом будет пара (147, 10), а 147 представляет собой 8-битовое число.

10.16. Точное деление на константу

Под "точным делением" подразумевается деление, о котором заранее известно, что его остаток равен 0. Хотя такая ситуация и не очень распространена, она имеет место, например, при вычитании двух указателей в языке C. В нем результат вычитания $p - q$, где p и q — указатели, определен и переносим, только если p и q указывают на объекты в одном и том же массиве [57, раздел 7.6.2]. Если размер элемента массива равен s , то код для вычисления разности двух указателей реально вычисляет $(p - q)/s$.

Материал этого раздела основан на [39, раздел 9].

Рассматриваемый здесь метод применим как к знаковому, так и к беззнаковому точному делению и базируется на следующей теореме.

ТЕОРЕМА M1. *Если a и m — взаимно простые целые числа, то существует целое число $1 \leq \bar{a} < m$, такое, что $a\bar{a} \equiv 1 \pmod{m}$.*

Таким образом, \bar{a} представляет собой обратный мультипликативный по модулю m элемент по отношению к a . Есть несколько способов доказательства этой теоремы, три из них описаны в [90, с. 52]. Доказательство, приведенное далее, требует только знания основ теории сравнений.

Доказательство. Докажем несколько более общее по сравнению с данной теоремой утверждение. Если a и m взаимно просты (и следовательно, имеют ненулевые значения), то для множества x , состоящего из m различных по модулю m значений, значения ax также будут различны по модулю m . Например, при $a = 3$ и $m = 8$ и значениях x из диапазона от 0 до 7 получим, что значения ax представляют собой множество 0, 3, 6, 9, 12, 15, 18, 21, или, по модулю 8, числа 0, 3, 6, 1, 4, 7, 2, 5. Заметим, что в полученной последовательности вновь представлены все числа от 0 до 7.

Чтобы доказать это в общем случае, воспользуемся методом доказательства от противного. Предположим, что существуют два различных по модулю m целых числа, которые отображаются в одно и то же число по модулю m при умножении на a , т.е. существуют такие x и y ($x \not\equiv y \pmod{m}$), что

$$ax \equiv ay \pmod{m}.$$

Но в таком случае существует целое число k , такое, что

$$ax - ay = km \quad \text{или} \\ a(x - y) = km.$$

Поскольку a не имеет общих множителей с m , разность $x - y$ должна быть кратной m , т.е.

$$x \equiv y \pmod{m}.$$

Однако этот вывод противоречит исходному предположению.

Теперь, поскольку значения ax принимают все m возможных различных значений по модулю m , среди значений x найдется такое, для которого ax по модулю m равно 1.

Приведенное доказательство показывает также, что имеется только одно значение x (по модулю m), такое, что $ax \equiv 1 \pmod{m}$, т.е. что мультипликативное обращение однозначно. Можно также показать, что имеется единственное (по модулю m) целое число x , такое, что $ax \equiv b \pmod{m}$, где b — любое целое число.

В качестве примера рассмотрим случай $m = 16$. Тогда $\bar{3} = 11$, так как $3 \cdot 11 = 33 \equiv 1 \pmod{16}$. Можно также считать, что $\bar{3} = -5$, поскольку $3 \cdot (-5) = -15 \equiv 1 \pmod{16}$. Аналогично $\bar{-3} = 5$, поскольку $(-3) \cdot 5 = -15 \equiv 1 \pmod{16}$.

Важность этих наблюдений заключается в том, что они показывают применимость рассматриваемых концепций как к знаковым, так и к беззнаковым числам. Если вы бу-

дете работать с беззнаковыми числами на 4-битовой машине, то получите $\bar{3} = 11$; если же будете иметь дело со знаковыми числами, то получите $\bar{3} = -5$. Однако и 11, и -5 имеют одно и то же представление в дополнительном коде (поскольку их разность равна 16), так что в обоих случаях в качестве мультипликативного обратного элемента используется одно и то же содержимое машинного слова.

Рассмотренная теорема непосредственно применима к задаче деления (знакового и беззнакового) на нечетное число d на W -битовом компьютере. Поскольку любое нечетное число является взаимно простым с 2^W , теорема гласит, что если d нечетно, то существует целое \bar{d} (единственное в диапазоне от 0 до $2^W - 1$ или от -2^{W-1} до $2^{W-1} - 1$), такое, что

$$d\bar{d} \equiv 1 \pmod{2^W}.$$

Следовательно, для любого n , кратного d , справедливо соотношение

$$\frac{n}{d} \equiv \frac{n}{d}(d\bar{d}) \equiv n\bar{d} \pmod{2^W}.$$

Другими словами, n/d можно вычислить путем умножения n на \bar{d} , после чего взять правые W бит полученного произведения.

Если делитель d представляет собой четное число, то пусть $d = d_0 \cdot 2^k$, где d_0 нечетно, а $k \geq 1$. Тогда просто выполняется сдвиг числа n вправо на k позиций (устраняя нулевые биты), после чего выполняется умножение на \bar{d}_0 (сдвиг можно выполнить и после операции умножения).

Ниже приведен код для деления на 7 числа n , кратного 7. Этот код дает корректный результат как при знаковом, так и при беззнаковом делении.

```
li m, 0xB6DB6DB7    Мультипликативное обратное, (5*2**32+1)/7
mul q, M, n          q = n/7
```

Вычисление мультипликативного обратного по алгоритму Евклида

Каким же образом можно вычислить мультипликативное обратное число? Стандартный метод — использование расширенного алгоритма Евклида. Этот метод вкратце рассматривается далее, в приложении к нашей основной задаче, а заинтересовавшимся читателям можно посоветовать обратиться к [90, с. 13] и [67, раздел 4.5.2].

Пусть задан нечетный делитель d и нам требуется найти число x , такое, что

$$dx \equiv 1 \pmod{m},$$

причем в нашей задаче $m = 2^W$ (где W — размер машинного слова). Это может быть выполнено, если уравнение

$$dx + my = 1$$

будет решено в целых числах x и y (положительных, отрицательных, равных 0).

Начнем с того, что сделаем d положительным, добавляя к нему достаточное количество раз число m (d и $d + km$ имеют одно и то же мультипликативное обратное). Затем запишем следующие уравнения (в которых $d, m > 0$):

$$d(-1) + m(1) = m - d, \tag{i}$$

$$d(1) + m(0) = d. \tag{ii}$$

Если $d = 1$, поставленная задача решена, так как (ii) показывает, что $x = 1$. В противном случае вычислим

$$q = \left\lfloor \frac{m-d}{d} \right\rfloor.$$

Далее умножим уравнение (ii) на q и вычтем его из (i). Это даст нам уравнение

$$d(-1-q) + m(1) = m - d - qd = \text{rem}(m-d, d).$$

Это уравнение справедливо, поскольку мы просто умножили одно из уравнений на константу и вычли его из другого. Если $\text{rem}(m-d, d) = 1$, то задача решена (последнее уравнение и является решением) и $x = -1 - q$.

Повторим описанный процесс с последними двумя уравнениями, получая новое уравнение. Будем продолжать эти действия до тех пор, пока в правой части уравнения не будет получена 1. Тогда множитель при d , приведенный по модулю m , и будет представлять собой искомое мультипликативное обратное к d .

Кстати, если $m - d < d$, так что первое частное равно 0, то третья строка будет точной копией первой, так что второе частное будет ненулевым. Кроме того, в разных источниках зачастую в качестве первой строки используется следующая:

$$d(0) + m(1) = m,$$

но в нашем приложении $m = 2^8$ невозможно представить в компьютере.

Лучше всего проиллюстрировать описываемый процесс на конкретном примере. Пусть $m = 256$ и $d = 7$. Тогда процесс вычислений выглядит как показано ниже (для получения третьей строки используется $q = \lfloor 249/7 \rfloor = 35$).

$$\begin{aligned} 7(-1) + 256(1) &= 249 \\ 7(1) + 256(0) &= 7 \\ 7(-36) + 256(1) &= 4 \\ 7(37) + 256(-1) &= 3 \\ 7(-73) + 256(2) &= 1 \end{aligned}$$

Таким образом, мультипликативным обратным по модулю 256 числа 7 является -73 (или, используя числа из диапазона от 0 до 255, число 183). Проверяем: $7 \cdot 183 = 1281 \equiv 1 \pmod{256}$.

Начиная с третьей строки числа в правом столбце представляют собой остатки от деления чисел, расположенных выше, так что это строго убывающая последовательность

неотрицательных чисел, которая, таким образом, должна завершиться 0 (в примере выше для этого требуется еще один шаг). Кроме того, число перед 0 должно быть 1 по следующей причине. Предположим, что последовательность остатков завершается некоторым числом b , не равным 1, после которого идет число 0. Тогда целое, предшествующее b , должно быть кратным b (скажем, числом k_1b) с тем, чтобы следующим остатком в последовательности было число 0. Такие же рассуждения приводят к выводу, что для того, чтобы получить остаток b , предшествующим k_1b числом должно быть $k_1k_2b + b$. Продолжая эту последовательность, вы увидите, что каждое из этих чисел должно быть кратно b , включая числа из первых двух строк, которые равны $m-d$ и d и являются взаимно простыми.

Приведенное выше рассуждение является неформальным доказательством того, что рассматриваемый процесс завершается, причем наличием 1 в правом столбце, а следовательно, находит мультипликативное обратное d число.

Для проведения данных вычислений на компьютере заметим, что если $d < 0$, то к нему необходимо прибавить 2^m . Однако в случае арифметики, использующей дополнительный к 2 код, это делать необязательно; достаточно просто рассматривать d как беззнаковое число, независимо от того, как именно его рассматривает приложение.

Вычисление q должно использовать беззнаковое деление.

Заметим, что все вычисления можно выполнять по модулю m , так как это никак не влияет на правый столбец (его значения все равно остаются в диапазоне от 0 до $m-1$). Это важно, так как позволяет нам выполнять вычисления с "одинарной точностью", используя беззнаковую компьютерную арифметику по модулю 2^m .

Большинство величин в рассматриваемой таблице не обязательно должны быть представлены при вычислениях, например столбец множителей при 256, поскольку при решении уравнения $dx + my = 1$ нас не интересует значение y . Нет необходимости и в представлении d в первом столбце. Оставляя в таблице только необходимое, получим ее сокращенный вид.

255	249
1	7
220	4
37	3
183	1

Программа на языке C для проведения этих вычислений представлена в листинге 10.5.

Листинг 10.5. Поиск мультипликативного обратного по модулю 2^{31} с помощью алгоритма Евклида

```
unsigned mulinv(unsigned d) // d должно быть нечетно
{
    unsigned x1, v1, x2, v2, x3, v3, q;
    x1 = 0xFFFFFFFF;    v1 = -d;
    x2 = 1;              v2 = d;
    while (v2 > 1)
    {
        q = v1/v2;
```

```

x3 = x1 - q*x2;   v3 = v1 - q*v2;
x1 = x2;         v1 = v2;
x2 = x3;         v2 = v3;
}
return x2;
}

```

Причина использования условия цикла $v2 > 1$ вместо более естественного $v2 != 1$ заключается в том, что если при использовании последнего условия функции будет ошибочно передано четное значение, то цикл может никогда не завершиться (если аргумент d четен, $v2$ никогда не получит значение 1, но в конечном итоге примет значение 0).

Что же вычисляет данная программа, если передать ей четный аргумент? Как и следует из формул, она вычисляет число x , такое, что $dx \equiv 0 \pmod{2^{32}}$, которое вряд ли может оказаться полезным. Однако минимальные изменения в программе (изменение условия цикла на $v2 != 0$ и возврат значения $x1$ вместо $x2$) приводят к вычислению ею числа x , такого, что $dx \equiv g \pmod{2^{32}}$, где g — наибольший общий делитель d и 2^{32} , т.е. наибольшей степени 2, являющейся делителем d . Такая модифицированная программа, как и ранее, вычисляет для нечетного d мультипликативное обратное число, хотя и требует для этого на одну итерацию больше.

Что касается количества итераций, выполняемых данной программой, то для нечетного d , не превышающего 20, требуется максимум три итерации; среднее их число равно 1.7. Для d порядка 1000 требуется максимум 11, а в среднем — шесть итераций.

Вычисление мультипликативного обратного по методу Ньютона

Хорошо известно, что при работе с действительными числами значение $1/d$, где $d \neq 0$, может быть вычислено с возрастающей точностью с помощью итеративного вычисления

$$x_{n+1} = x_n (2 - dx_n) \tag{31}$$

при начальном приближении x_0 , достаточно близком к $1/d$. Количество точных цифр в результате при каждой итерации примерно удваивается.

Однако гораздо менее известно, что данная формула может применяться для поиска мультипликативного обратного числа по модулю любой степени 2! Например, для поиска мультипликативного обратного к 3 по модулю 256 начнем с $x_0 = 1$ (можно использовать любое нечетное число). Тогда

$$\begin{aligned}
 x_1 &= 1(2 - 3 \cdot 1) = -1, \\
 x_2 &= -1(2 - 3(-1)) = -5, \\
 x_3 &= -5(2 - 3(-5)) = -85, \\
 x_4 &= -85(2 - 3(-85)) = -21845 \equiv -85 \pmod{256}.
 \end{aligned}$$

Итерации достигли неподвижной точки по модулю 256, так что -85 , или 171, и есть мультипликативное обратное числа 3 по модулю 256. Все приведенные вычисления выполнялись по модулю 256.

Почему работает этот метод? Поскольку, если x_n удовлетворяет условию

$$dx_n \equiv 1 \pmod{m},$$

а x_{n+1} определяется формулой (31), то

$$dx_{n+1} \equiv 1 \pmod{m^2}.$$

Чтобы увидеть это, положим $dx_n = 1 + km$. Тогда

$$\begin{aligned} dx_{n+1} &= dx_n(2 - dx_n) \\ &= (1 + km)(2 - (1 + km)) \\ &= (1 + km)(1 - km) \\ &= 1 - k^2 m^2 \\ &= 1 \pmod{m^2}. \end{aligned}$$

В нашем приложении m является степенью 2, скажем, 2^h . В этом случае, если

$$\begin{aligned} dx_n &\equiv 1 \pmod{2^N}, \text{ то} \\ dx_{n+1} &\equiv 1 \pmod{2^{2N}}. \end{aligned}$$

В этом смысле, если x_n рассматривается как некоторое приближение \bar{d} , то каждая итерация (31) удваивает количество битов "точности" приближения.

Так случилось, что по модулю 8 мультипликативным обратным любого нечетного числа d является само это число. Таким образом, в качестве начального приближения можно взять $x_0 = d$. Тогда по формуле (31) получим значения x_1, x_2, \dots , такие, что

$$\begin{aligned} dx_1 &\equiv 1 \pmod{2^6}, \\ dx_2 &\equiv 1 \pmod{2^{12}}, \\ dx_3 &\equiv 1 \pmod{2^{24}}, \\ dx_4 &\equiv 1 \pmod{2^{48}} \text{ и т.д.} \end{aligned}$$

Таким образом, четырех итераций достаточно, чтобы найти мультипликативное обратное по модулю 2^{32} (если $x \equiv 1 \pmod{2^{48}}$, то $x \equiv 1 \pmod{2^n}$ для $n \leq 48$). Это приводит нас к программе на языке C (листинг 10.6), в которой все вычисления выполняются по модулю 2^{32} .

Листинг 10.6. Вычисление мультипликативного обратного по модулю 2^{32}
методом Ньютона

```
unsigned mulinv(unsigned d) // d должно быть нечетно
{
    unsigned xn, t;

    xn = d;
loop:
    t = d*xn;
    if (t == 1) return xn;
    xn = xn*(2 - t);
    goto loop;
}
```

Примерно для половины значений d этой программе требуется 4.5 итерации, т.е. девять умножений. Для другой половины (у которой “верны четыре бита” начального значения x_n , т.е. $d^2 \equiv 1 \pmod{16}$) требуется, как правило, семь (или меньше) умножений. Таким образом, можно считать, что этот метод требует в среднем восьми умножений.

Один из вариантов этой программы состоит в простом выполнении цикла четыре раза независимо от значения d , что позволяет обойтись восемью умножениями и без команд управления циклом. Другой вариант предусматривает выбор в качестве начального приближения мультипликативного обратного “с точностью 4 бита”, т.е. перед началом вычислений нужно найти x_0 , которое удовлетворяет условию $dx_0 \equiv 1 \pmod{16}$. Тогда для вычислений потребуется выполнить только три итерации цикла. Найти подходящее начальное приближение можно, например, следующими способами.

$$x_0 \leftarrow d + 2((d+1) \& 4)$$

$$x_0 \leftarrow d^2 + d - 1$$

Здесь умножение на 2 выполняется с помощью сдвига влево, а все вычисления производятся по модулю 2^{32} (игнорируя переполнения). Поскольку во второй формуле используется умножение, ее применение позволяет сэкономить только одну такую команду.

Такое пристальное рассмотрение вопроса времени вычислений, конечно, не имеет смысла, если метод применяется в компиляторе. В этом случае данная программа используется настолько редко, что при ее кодировании, в первую очередь, должен интересоваться ее размер. Тем не менее могут существовать приложения, в которых поиск мультипликативного обратного должен выполняться максимально быстро.

Описанный здесь метод Ньютона применяется, только когда (1) модуль представляет собой целую степень некоторого числа a и (2) известно мультипликативное обратное d по модулю a . В частности, метод хорошо работает при $a = 2$, поскольку тогда тут же известно мультипликативное обратное любого нечетного числа d по модулю 2 — это просто 1.

Некоторые мультипликативные обратные

В завершение раздела перечислим в табл. 10.3 некоторые мультипликативные обратные числа.

ТАБЛИЦА 10.3. НЕКОТОРЫЕ МУЛЬТИПЛИКАТИВНЫЕ ОБРАТНЫЕ ЧИСЛА

d	\bar{d}		
	mod 16	mod 2^{32}	mod 2^{64}
(десятичное)	(десятичное)	(шестнадцатеричное)	(шестнадцатеричное)
-7	-7	49249249	92492492 49249249
-5	3	33333333	33333333 33333333
-3	5	55555555	55555555 55555555
-1	-1	FFFFFFFF	FFFFFFFF FFFFFFFF
1	1	1	1
3	11	AAAAAAAAB	AAAAAAAA AAAAAAAB
5	13	CCCCCCCD	CCCCCCCC CCCCCCCD
7	7	B6DB6DB7	6DB6DB6D B6DB6DB7
9	9	38E38E39	8E38E38E 38E38E39
11	3	BA2E8BA3	2E8BA2E8 BA2E8BA3
13	5	C4EC4EC5	4EC4EC4E C4EC4EC5
15	15	EEEEEEEF	EEEEEEEE EEEEEEEF
25		C28F5C29	8F5C28F5 C28F5C29
125		26E978D5	1CAC0831 26E978D5
625		3AFB7E91	D288CE70 3AFB7E91

Вы можете заметить, что в ряде случаев ($d = 3, 5, 9, 11$) мультипликативное обратное к d совпадает с магическим числом для беззнакового деления на d (см. раздел 10.14, "Некоторые магические числа" на с. 265). Это в большей или меньшей мере случайность. Подобное происходит для тех чисел, для которых магическое число M равно множителю m , и эти величины имеют вид $(2^p + 1)/d$ при $p \geq 32$. В этом случае

$$Md = \left(\frac{2^p + 1}{d} \right) d \equiv 1 \pmod{2^{32}},$$

так что $M \equiv \bar{d} \pmod{2^{32}}$.

10.17. Проверка нулевого остатка при делении на константу

Мультипликативное обратное делителя d может использоваться для проверки равенства 0 остатка после деления на d [39].

Беззнаковое деление

Сначала рассмотрим беззнаковое деление на нечетный делитель d . Обозначим через \bar{d} мультипликативное обратное к d число. Тогда, поскольку $d\bar{d} \equiv 1 \pmod{2^w}$, где w — размер машинного слова в битах, \bar{d} также нечетно. Следовательно, \bar{d} является взаимно простым с 2^w и, как показано в доказательстве теоремы М1 в предыдущем разделе, если n представляет собой множество всех 2^w различных чисел по модулю 2^w , то $n\bar{d}$ также представляет собой множество всех 2^w различных чисел по модулю 2^w .

В предыдущем разделе было показано, что если n кратно d , то

$$\frac{n}{d} \equiv n\bar{d} \pmod{2^w}.$$

Следовательно, $n\bar{d} \equiv 0, 1, 2, \dots, \lfloor (2^w - 1)/d \rfloor \pmod{2^w}$ при $n = 0, d, 2d, \dots, \lfloor (2^w - 1)/d \rfloor d$. Таким образом, для n , не являющегося кратным d , значение $n\bar{d}$, приведенное по модулю 2^w к диапазону от 0 до $2^w - 1$, должно превышать величину $\lfloor (2^w - 1)/d \rfloor$.

Этот вывод можно использовать для проверки на равенство нулю остатка от деления. Например, чтобы проверить, что число n кратно 25, умножим n на $\bar{25}$ и сравним правые w бит со значением $\lfloor (2^w - 1)/25 \rfloor$. Использование базового набора RISC-команд дает следующий код.

```
li    m, 0xC28F5C29    Загружаем мультипликативное обратное 25
mul   q, m, n          q = правая половина M*n
li    c, 0x0A3D70A3    c = floor((2**32-1)/25)
cmpleu t, q, c        Сравниваем q и c; переход,
bt    t, is_mult      если n кратно 25
```

Для того чтобы распространить этот метод на четные делители, представим делитель в следующем виде: $d = d_0 \cdot 2^k$, где d_0 нечетное, а $k \geq 1$. Тогда, поскольку целое число делится на d тогда и только тогда, когда оно делится на d_0 и 2^k , и поскольку n и $n\bar{d}_0$ имеют одинаковое количество завершающих нулевых битов (\bar{d}_0 нечетно), проверка того, что n кратно d , состоит в следующем:

установить $q = \text{mod}(n\bar{d}_0, 2^k)$;

$q \leq \lfloor (2^k - 1)/d_0 \rfloor$ и q завершается не менее чем k нулевыми битами,

где под функцией “mod” подразумевается приведение $n\bar{d}_0$ к интервалу $[0, 2^k - 1]$.

Непосредственная реализация описанного метода требует двух проверок и условных переходов, однако если компьютер имеет команду *циклического сдвига*, то метод можно реализовать с помощью одного *сравнения с ветвлением*. Это следует из приведенной далее теоремы, в которой запись $a \ggg k$ означает компьютерное слово a , циклически сдвинутое вправо на k бит ($0 \leq k \leq 32$).

ТЕОРЕМА ZRU. $x \stackrel{\text{н}}{\leq} a$ и x заканчивается k нулевыми битами тогда и только тогда, когда $x \stackrel{\text{н}}{\gg} k \stackrel{\text{н}}{\leq} \lfloor a/2^t \rfloor$.

Доказательство. (Будем считать, что имеем дело с 32-разрядным компьютером.)

Предположим, что $x \stackrel{\text{н}}{\leq} a$ и x заканчивается k нулевыми битами. Тогда, так как $x \stackrel{\text{н}}{\leq} a$, то $\lfloor x/2^t \rfloor \stackrel{\text{н}}{\leq} \lfloor a/2^t \rfloor$. Однако $\lfloor x/2^t \rfloor = x \stackrel{\text{н}}{\gg} k$, и, таким образом, $x \stackrel{\text{н}}{\gg} k \stackrel{\text{н}}{\leq} \lfloor a/2^t \rfloor$. Если x не заканчивается k нулевыми битами, то $x \stackrel{\text{н}}{\gg} k$ не начинается с k нулевых битов, в то время как $\lfloor a/2^t \rfloor$ начинается с них, так что $x \stackrel{\text{н}}{\gg} k \stackrel{\text{н}}{>} \lfloor a/2^t \rfloor$. И наконец, если $x \stackrel{\text{н}}{>} a$ и x заканчиваются k нулевыми битами, то целое число, образованное первыми $32-k$ битами x должно превышать число, образованное первыми $32-k$ битами a , так что $\lfloor x/2^t \rfloor \stackrel{\text{н}}{>} \lfloor a/2^t \rfloor$.

При использовании этой теоремы проверка того, что n кратно d , где n и d — ненулевые беззнаковые целые числа, причем $d = d_0 \cdot 2^t$, где d_0 — нечетно, выполняется следующим образом.

$$\boxed{\begin{array}{l} q \leftarrow \text{mod}(n\bar{d}_0, 2^n) \\ q \stackrel{\text{н}}{\gg} k \stackrel{\text{н}}{\leq} \lfloor (2^n - 1)/d \rfloor \end{array}}$$

Здесь мы воспользовались тем, что

$$\lfloor \lfloor (2^n - 1)/d_0 \rfloor / 2^t \rfloor = \lfloor (2^n - 1)/(d_0 \cdot 2^t) \rfloor = \lfloor (2^n - 1)/d \rfloor.$$

Далее в качестве примера приведен код, проверяющий, является ли беззнаковое целое число n кратным 100.

```
li      M,0x28F5C29    Мультипликативное обратное 25
mul     q,M,n          q = правая половина M*n
shrrl  q,q,2          Циклический сдвиг вправо на два бита
li      c,0x028F5C28  c = floor((2**32-1)/100)
cmpleu t,q,c          Сравнение q и c и переход, если
bt      t,is_mult     n кратно 100
```

Знаковое деление, делитель ≥ 2

В случае знакового деления, как было показано в предыдущем разделе, если n кратно d и d нечетно, получаем

$$\frac{n}{d} \equiv n\bar{d} \pmod{2^n}.$$

Таким образом, для $n = \lceil -2^{n-1}/d \rceil \cdot d, \dots, -d, 0, d, \dots, \lfloor (2^{n-1}-1)/d \rfloor \cdot d$ имеем $n\bar{d} \equiv \lceil -2^{n-1}/d \rceil, \dots, -1, 0, 1, \dots, \lfloor (2^{n-1}-1)/d \rfloor \pmod{2^n}$. Кроме того, поскольку \bar{d} взаимно про-

стое с 2^n , если n принимает 2^n различных значений по модулю 2^n , то $n\bar{d}$ также принимает 2^n различных значений по модулю 2^n . Следовательно, n кратно d тогда и только тогда, когда

$$\lceil -2^{n-1}/d \rceil \leq \text{mod}(n\bar{d}, 2^n) \leq \lfloor (2^{n-1} - 1)/d \rfloor,$$

где под функцией "mod" подразумевается приведение $n\bar{d}$ к интервалу $[-2^{n-1}, 2^{n-1} - 1]$.

Этот способ можно немного упростить, заметив, что, так как d нечетно, а из наших начальных предположений оно положительно и не равно единице, d не является делителем 2^{n-1} . Следовательно,

$$\lceil -2^{n-1}/d \rceil = \lceil (-2^{n-1} + 1)/d \rceil = -\lfloor (2^{n-1} - 1)/d \rfloor.$$

Для знаковых чисел проверка того, что n является кратным d , где $d = d_0 \cdot 2^t$, а d_0 нечетно, выполняется следующим образом:

установить $q = \text{mod}(n\bar{d}_0, 2^n)$;

$-\lfloor (2^{n-1} - 1)/d_0 \rfloor \leq q \leq \lfloor (2^n - 1)/d_0 \rfloor$ и q завершается не менее чем k нулевыми битами.

На первый взгляд, при таком методе требуется три проверки и перехода. Однако, как и в случае беззнакового деления, все можно свести к одной команде *сравнения с ветвлением*, если воспользоваться следующей теоремой.

ТЕОРЕМА ZRS. Если $a \geq 0$, то следующие утверждения эквивалентны:

(1) $-a \leq x \leq a$ и x заканчивается k или большим количеством нулевых битов,

(2) $\text{abs}(x) \gg k \leq \lfloor a/2^t \rfloor$,

(3) $x + a' \gg k \leq \lfloor 2a'/2^t \rfloor$,

где a' представляет собой a с установленными равными 0 правыми k битами (т.е. $a' = a \& -2^k$).

Доказательство. (Будем считать, что имеем дело с 32-разрядным компьютером.) Чтобы убедиться, что утверждение (1) эквивалентно утверждению (2), достаточно заметить, что выражения $-a \leq x \leq a$ и $\text{abs}(x) \leq a$ эквивалентны. После этого эквивалентность утверждений (1) и (2) следует из теоремы ZRU.

Чтобы убедиться в эквивалентности утверждений (1) и (3), заметим, что утверждение (1) справедливо и при замене a на a' . Тогда, в соответствии с "теоремой границ" на с. 90, оно, в свою очередь, эквивалентно

$$x + a' \leq 2a'.$$

Поскольку $x + a'$ заканчивается k нулевыми битами тогда и только тогда, когда k нулевыми битами заканчивается x , можно применить теорему ZRU, которая и дает требующийся нам результат.

Используя утверждение (3) рассмотренной теоремы, проверку кратности n числу d , где n и d — знаковые целые числа, не меньшие 2, и $d = d_0 \cdot 2^t$, где d_0 нечетно, можно провести следующим образом.

$$\begin{aligned} q &\leftarrow \text{mod}(n\bar{d}_0, 2^n) \\ a' &\leftarrow \lfloor (2^{n-1} - 1) / d_0 \rfloor \& -2^t \\ q + a' \gg k &\leq \lfloor (2a') / 2^t \rfloor \end{aligned}$$

(Поскольку d — константа, a' можно вычислить во время компиляции.)

Далее в качестве примера приведен код, проверяющий, является ли знаковое целое число n кратным 100. Обратите внимание, что константа $\lfloor 2a' / 2^t \rfloor$ в любой момент может быть получена из константы a' путем сдвига на $k-1$ позиций, что позволяет сохранить одну команду либо загрузку из памяти для получения операнда сравнения.

```
li      M, 0x28F5C29   Загрузка мультипликативного обратного 25
mul     q, M, n        q = правая половина M*n
li      c, 0x051EB850  c = floor((2**31 - 1)/25) & -4
add     q, q, c        Прибавляем c
shrrl  q, q, 2        Циклический сдвиг вправо на 2 позиции
shrl   c, c, 1        Вычисление константы для сравнения
cmpleu t, q, c        Сравнение q и c и переход, если
bt      t, is_mult    n кратно 100
```

10.18. Методы, не использующие команды умножения со старшим словом

В этом разделе мы рассмотрим некоторые методы деления на константы, в которых не используется команда *умножения со старшим словом* (*старшего слова умножения*), т.е. умножения, дающего в качестве результата двойное слово. Мы покажем, как заметить деление на константу последовательностью команд *сдвига* и *сложения* или *сдвига*, *сложения* и *умножения* для получения более компактного кода.

Беззнаковое деление

В этом случае беззнаковое деление оказывается проще знакового, так что мы начнем с него. Один из методов заключается в применении методов на основе “двухсловного” умножения, но с заменой последнего кодом, показанным в листинге 8.2 на с. 200. В листинге 10.7 показано, как такая технология работает для случая (беззнакового) деления на 3. Это комбинация кода, приведенного на с. 253, и кода из листинга 8.2 с заменой `int` на `unsigned`. Код выполняется с помощью 15 команд, включая четыре умножения. Умножения на большие константы при преобразовании к *сдвигам* и *умножениям* приводят к достаточно небольшому количеству команд. Очень похожий код может быть раз-

работан и в случае знакового деления. В целом этот метод не так уж хорош, и далее нами не рассматривается.

ЛИСТИНГ 10.7. Беззнаковое деление на 3 с использованием имитации команды старшего слова беззнакового умножения

```
unsigned divu3(unsigned n)
(
    unsigned n0, n1, w0, w1, w2, t, q;

    n0 = n & 0xFFFF;
    n1 = n >> 16;
    w0 = n0*0xAAAAB;
    t = n1*0xAAAAB + (w0 >> 16);
    w1 = t & 0xFFFF;
    w2 = t >> 16;
    w1 = n0*0xAAAAA + w1;
    q = n1*0xAAAAA + w2 + (w1 >> 16);
    return q >> 1;
)
```

Другой метод [37] заключается в вычислении обратного к делителю и умножении на него делимого с помощью ряда *сдвигов вправо* и *сложений*. Так мы получим приближенное значение частного. Это просто приближение, поскольку обратное к делителю (который, как предполагается, не является точной степенью двойки) не может быть точно выражен с помощью 32 бит, а также потому что каждый *сдвиг вправо* отбрасывает биты делимого. Затем вычисляется остаток по отношению к приближенному частному и делится на делитель для получения поправки, которая добавляется к приближенному частному и дает его точное значение. Остаток обычно мал по сравнению с делителем (в несколько раз меньше), так что зачастую имеется простой способ вычисления поправки без применения команды *деления*.

Для иллюстрации этого метода рассмотрим деление на 3, т.е. вычисление $\lfloor n/3 \rfloor$, где $0 \leq n < 2^{32}$. Значение, обратное к 3, в бинарном представлении приближенно равно

$$0.0101\ 0101\ 0101\ 0101\ 0101\ 0101\ 0101\ 0101.$$

Для вычисления приближенного произведения этого числа и n можно воспользоваться формулой

$$q \leftarrow \left(n \overset{\cdot}{\gg} 2 \right) + \left(n \overset{\cdot}{\gg} 4 \right) + \left(n \overset{\cdot}{\gg} 6 \right) + \dots + \left(n \overset{\cdot}{\gg} 30 \right) \tag{32}$$

(всего 29 команд; последняя 1 в обратном значении игнорируется, поскольку соответствует добавлению члена $n \overset{\cdot}{\gg} 32$, который, очевидно, равен 0). Однако простое повторение нулей и единиц в обратном значении позволяет разработать метод, оказывающийся и более быстрым (девять команд), и более точным.

$$\begin{aligned}
 q &\leftarrow \left(n \gg 2 \right) + \left(n \gg 4 \right) \\
 q &\leftarrow q + \left(q \gg 4 \right) \\
 q &\leftarrow q + \left(q \gg 8 \right) \\
 q &\leftarrow q + \left(q \gg 16 \right)
 \end{aligned}
 \tag{33}$$

Для сравнения точности этих методов рассмотрим биты, удаляемые сдвигом каждого из членов (32), когда n состоит из одних единичных битов. Первый член убирает два единичных бита, следующий — четыре и т.д. Каждый сдвиг вносит ошибку, почти равную 1 младшего разряда. Поскольку всего таких членов 16 (учитывая игнорируемый член), сдвиги дают ошибку, почти равную 16. Это ошибка, дополнительная к возникающей из-за ограниченного 32-битового представления обратного к делителю значения; оказывается, максимальная общая ошибка равна 16.

В случае процедуры (33) каждый правый сдвиг также вносит ошибку, почти равную единице в младшем разряде. Однако здесь имеется только 5 операций сдвига. Они вносят ошибку, близкую к 5; имеется также ошибка, связанная с обрезкой обратного значения 32 битами. Как оказывается, максимальная общая ошибка равна 5.

После вычисления оценочного значения частного q остаток r вычисляется как

$$r \leftarrow n - q * 3.$$

Остаток не может быть отрицательным, поскольку значение q никогда не превышает точное частное. Чтобы разработать наиболее простой метод вычисления $r + 3$, нам надо знать, насколько большим может быть значение r . В общем случае для делителя d и оценочного значения частного q , меньше точного значения частного на k , остаток находится в диапазоне от $k * d$ до $k * d + d - 1$ (верхняя граница консервативна и в действительности может не быть достигнута). Таким образом, используя (33), где q может оказаться меньше, чем следует, не более чем на 5, получаем, что остаток не превышает величину $5 * 3 + 2 = 17$. Эксперименты показывают, что в действительности он не превышает 15. Таким образом, поправка, которую мы должны вычислить, составляет (в точности)

$$r + 3 \text{ при } 0 \leq r \leq 15.$$

Поскольку r мало по сравнению с наибольшим значением, которое может храниться в регистре, это вычисление может быть получено приближенно путем умножения r на некоторое приближение $1/3$ вида a/b , где b представляет собой точную степень 2. Такое умножение легко вычислить, поскольку деление выполняется путем простого сдвига. Значение a/b должно быть немного больше $1/3$, чтобы после сдвига результат соответствовал делению с отсечением. Последовательность таких приближений имеет вид

$$1/2, 2/4, 3/8, 6/16, 11/32, 22/64, 43/128, 86/256, 171/512, 342/1024, \dots$$

Обычно чем меньше числитель и знаменатель дроби в последовательности, тем легче ее вычислить, так что выберем наименьшую подходящую. В нашем случае это дробь 11/32. Таким образом, окончательное точное значение частного вычисляется как

$$q \leftarrow q + \left(11 * r \gg 5 \right).$$

Решение включает два умножения на небольшие числа (3 и 11), которые можно заменить командами *сдвига* и *сложения*.

В листинге 10.8 показано окончательное решение на языке программирования C. Как видно из листинга, оно включает 14 команд, в том числе два умножения. Если эти умножения заменить командами *сдвига* и *сложения*, то получится 18 элементарных команд. Однако если желательно избежать умножений, то любой из указанных альтернативных операторов `return` дает решение из 17 элементарных команд. Вторая альтернатива дает несколько меньшую степень параллельности, но, по правде говоря, очень небольшой степенью параллельности обладает сам метод.

Листинг 10.8. Беззнаковое деление на 3

```
unsigned divu3(unsigned n)
{
    unsigned q, r;

    q = (n >> 2) + (n >> 4);    // q = n*0.0101 (приблизленно)
    q = q + (q >> 4);          // q = n*0.01010101.
    q = q + (q >> 8);
    q = q + (q >> 16);
    r = n - q*3;                // 0 <= r <= 15.
    return q + (11*r >> 5);    // Возврат q + r/3.
// return q + (5*(r + 1) >> 4); // Альтернатива 1.
// return q + ((r + 5 + (r << 2)) >> 4); // Альтернатива 2.
}
```

Более точная оценка получается при замене первой выполнимой строки строкой

```
q = (n >> 1) + (n >> 3);
```

(это делает `q` слишком большим (в два раза), но при этом добавляется один бит точности) и вставке непосредственно перед присваиванием `r` строки

```
q = q >> 1;
```

В этом варианте остаток не превышает 9. Однако не похоже, чтобы ограниченность `r` величиной 9 давала лучший код вычисления $r \div 3$ по сравнению с границей 15 (четыре элементарные команды в любом случае). Таким образом, стоимость этой идеи — одна команда. Рассмотренная возможность упомянута здесь, поскольку она *улучшает* код для большинства делителей.

В листинге 10.9 показаны два варианта этого метода для деления на 5. Бинарное представление обратного к 5 имеет вид

```
0.0011 0011 0011 0011 0011 0011 0011 0011.
```

Как и в случае деления на 3 простое повторение шаблона из единиц и нулей позволяет достаточно эффективно и точно вычислить оценку частного. Оценка частного, вычисленная кодом слева, может отличаться от точного значения не более чем на 5, а остаток — не более чем на 25. Код справа сохраняет два бита точности при оценке частного, которая в этом случае отличается от точного значения не более чем на 2, а остаток — не более чем на 10. В этом случае наименьшим допустимым приближением $1/5$ является $7/32$, а не $13/64$, что приводит к несколько более эффективной программе при замене умножения *сдвигами* и *сложениями*. Количество команд в коде слева составляет 14 команд, включая два умножения, или 18 элементарных команд; в коде справа — 15 команд, включая 2 умножения, или 17 элементарных команд. Альтернативный код в операторе `return` используется только в том случае, когда машина оснащена командами предикатов. Он не уменьшает количество команд, а просто обладает небольшой степенью параллельности на уровне команд.

Листинг 10.9. Беззнаковое деление на 5

<pre> unsigned divu5a(unsigned n) { unsigned q, r; q = (n >> 3) + (n >> 4); q = q + (q >> 4); q = q + (q >> 8); q = q + (q >> 16); r = n - q*5; return q + (13*r >> 6); } </pre>	<pre> unsigned divu5b(unsigned n) { unsigned q, r; q = (n >> 1) + (n >> 2); q = q + (q >> 4); q = q + (q >> 8); q = q + (q >> 16); q = q >> 2; r = n - q*5; return q + (7*r >> 5); // return q + (r>4) + (r>9); } </pre>
---	---

Деление на 6 можно реализовать с помощью деления на 3 с последующим *сдвигом вправо* на 1. Однако дополнительную команду можно сэкономить при непосредственном вычислении, используя бинарное приближение

$$4/6 \approx 0.1010\ 1010\ 1010\ 1010\ 1010\ 1010\ 1010.$$

Соответствующий код показан в листинге 10.10. Версия слева выполняет умножение на приближение $1/6$, а затем выполняет коррекцию путем умножения на $11/64$. Версия справа использует преимущество того факта, что при умножении на приближение $4/6$ оценка частного отличается от точного значения не более чем на 1. Это позволяет получить более простой код коррекции — простое прибавление 1 к q при $r \geq 6$. Код второго оператора `return` предназначен для машин, оснащенных командами предикатов. Функция `divu6b` включает, как показано, 15 команд, в том числе одно *умножение*, или 17 элементарных команд, если умножение на 6 заменено *сдвигами* и *сложениями*.

ЛИСТИНГ 10.10. Беззнаковое деление на 6

```

unsigned divu6a(unsigned n)
{
    unsigned q, r;

    q = (n >> 3) + (n >> 5);
    q = q + (q >> 4);
    q = q + (q >> 8);
    q = q + (q >> 16);
    r = n - q*6;
    return q + (11*r >> 6);
}

unsigned divu6b(unsigned n)
{
    unsigned q, r;

    q = (n >> 1) + (n >> 3);
    q = q + (q >> 4);
    q = q + (q >> 8);
    q = q + (q >> 16);
    q = q >> 2;
    r = n - q*6;
    return q + ((r + 2) >> 3);
}
// return q + (r > 5);
}
    
```

Похоже, что для больших делителей наилучшим оказывается использование приближения к $1/d$, сдвинутого влево так, чтобы старший бит был равен 1. Обычно при этом оценка частного не отличается от точного значения более чем на 1 (возможно, это верно всегда — этого автор не знает), так что шаг коррекции реализуется весьма эффективно. В листинге 10.11 показан код для деления на 7 и 9, использующий бинарные приближения

$$4/7 \approx 0.1001\ 0010\ 0100\ 1001\ 0010\ 0100\ 1001\ 0010\ \text{и}$$

$$8/9 \approx 0.1110\ 0011\ 1000\ 1110\ 0011\ 1000\ 1110\ 0011.$$

При замене умножений на 7 и 9 сдвигами и сложениями эти функции состоят из 16 и 15 элементарных команд соответственно.

ЛИСТИНГ 10.11. Беззнаковое деление на 7 и 9

```

unsigned divu7(unsigned n)
{
    unsigned q, r;

    q = (n >> 1) + (n >> 4);
    q = q + (q >> 6);
    q = q + (q >> 12) + (q >> 24);
    q = q >> 2;
    r = n - q*7;
    return q + ((r + 1) >> 3);
}
// return q + (r > 6);
}

unsigned divu9(unsigned n)
{
    unsigned q, r;

    q = n - (n >> 3);
    q = q + (q >> 6);
    q = q + (q >> 12) + (q >> 24);
    q = q >> 3;
    r = n - q*9;
    return q + ((r + 7) >> 4);
}
// return q + (r > 8);
}
    
```

В листингах 10.12 и 10.13 показан код для деления на 10, 11, 12 и 13. Приведенные функции основаны на следующих бинарных приближениях.

$$8/10 \approx 0.1100\ 1100\ 1100\ 1100\ 1100\ 1100\ 1100\ 1100$$

$$8/11 \approx 0.1011\ 1010\ 0010\ 1110\ 1000\ 1011\ 1010\ 0010$$

$$8/12 \approx 0.1010\ 1010\ 1010\ 1010\ 1010\ 1010\ 1010\ 1010$$

$$8/13 \approx 0.1001\ 1101\ 1000\ 1001\ 1101\ 1000\ 1001\ 1101$$

При замене умножений *сдвигами* и *сложениями* эти функции состоят из 17, 20, 17 и 20 элементарных команд соответственно.

ЛИСТИНГ 10.12. Беззнаковое деление на 10 и 11

```

unsigned divu10(unsigned n)
{
    unsigned q, r;

    q = (n >> 1) + (n >> 2);
    q = q + (q >> 4);
    q = q + (q >> 8);
    q = q + (q >> 16);
    q = q >> 3;
    r = n - q*10;
    return q + ((r + 6) >> 4);
// return q + (r > 9);
}

unsigned divu11(unsigned n)
{
    unsigned q, r;

    q = (n >> 1) + (n >> 2) -
        (n >> 5) + (n >> 7);
    q = q + (q >> 10);
    q = q + (q >> 20);
    q = q >> 3;
    r = n - q*11;
    return q + ((r + 5) >> 4);
// return q + (r > 10);
}

```

ЛИСТИНГ 10.13. Беззнаковое деление на 12 и 13

```

unsigned divu12(unsigned n)
{
    unsigned q, r;

    q = (n >> 1) + (n >> 3);
    q = q + (q >> 4);
    q = q + (q >> 8);
    q = q + (q >> 16);
    q = q >> 3;
    r = n - q*12;
    return q + ((r + 4) >> 4);
// return q + (r > 11);
}

unsigned divu13(unsigned n)
{
    unsigned q, r;

    q = (n>>1) + (n>>4);
    q = q + (q>>4) + (q>>5);
    q = q + (q>>12) + (q>>24);
    q = q >> 3;
    r = n - q*13;
    return q + ((r + 3) >> 4);
// return q + (r > 12);
}

```

Случай деления на 13 весьма поучителен, поскольку показывает, как нужно искать повторяющиеся строки в бинарном представлении обратного к делителю. Первое присваивание устанавливает q равным $n * 0.1001$. Второе присваивание q добавляет $n * 0.00001001$ и $n * 0.000001001$. В этот момент q (приблизленно) равно $n * 0.100111011$. Третье присваивание q добавляет повторения этого шаблона. Иногда помогает применение вычитания, как в функции `divu9` ранее. Однако при использовании вычитания вы должны быть особенно осторожны, так как оно может привести к слишком большой оценке частного, так что остаток окажется отрицательным, и рассматриваемый метод будет неработоспособен. Достаточно сложно получить оптимальный код, и у нас нет какого-то универсального метода, который можно было бы вставить в компилятор для обработки любого делителя.

Приведенные выше примеры обеспечивали экономное использование команд благодаря простому повторяющемуся шаблону бинарного представления обратного к делителю, а также тому, что умножение при вычислении остатка r выполнялось на небольшую константу, так что его можно было заменить только несколькими командами *сдвига* и *сложения*. Может вызвать интерес вопрос о том, насколько эффективен рассматриваемый метод для больших делителей. Для получения грубого представления по этому во-

просу в листингах 10.14 и 10.15 показан код для деления на (десятичные) 100 и 1000. Соответствующие обратные значения имеют следующий вид.

$$64/100 \approx 0.1010\ 0011\ 1101\ 0111\ 0000\ 1010\ 0011\ 1101$$

$$512/1000 \approx 0.1000\ 0011\ 0001\ 0010\ 0110\ 1110\ 1001\ 0111$$

Если умножения заменяются *сдвигами* и *сложениями*, то эти функции выполняются с помощью 25 и 23 элементарных команд соответственно.

ЛИСТИНГ 10.14. Беззнаковое деление на 100

```
unsigned divu100(unsigned n)
{
    unsigned q, r;

    q = (n >> 1) + (n >> 3) + (n >> 6) - (n >> 10) +
        (n >> 12) + (n >> 13) - (n >> 16);
    q = q + (q >> 20);
    q = q >> 6;
    r = n - q*100;
    return q + ((r + 28) >> 7);
// return q + (r > 99);
}
```

ЛИСТИНГ 10.15. Беззнаковое деление на 1000

```
unsigned divu1000(unsigned n)
{
    unsigned q, r, t;

    t = (n >> 7) + (n >> 8) + (n >> 12);
    q = (n >> 1) + t + (n >> 15) + (t >> 11) + (t >> 14);
    q = q >> 9;
    r = n - q*1000;
    return q + ((r + 24) >> 10);
// return q + (r > 999);
}
```

В случае деления на 1000 младшие восемь битов оценки обратного значения, по сути, игнорируются. Код в листинге 10.15 заменяет биты 10010111 битами 01000000, но оценка частного при этом остается в пределах единицы от истинного частного. Таким образом, похоже, что, хотя бинарное представление обратных больших делителей может иметь мало повторений, как минимум некоторые биты можно игнорировать, чтобы снизить количество необходимых *сдвигов* и *сложений* при вычислении оценки частного.

В этом разделе было показано, пусть и не совсем строгим способом, как беззнаковое деление на константу можно свести к последовательности, как правило, состоящей из около 20 элементарных команд. Получение алгоритма, генерирующего такие последовательности, годного для встраивания в компилятор, — задача нетривиальная из-за трех сложностей в получении оптимального кода.

1. Требуется выполнить поиск повторяющегося шаблона в битовой строке оценки обратного к делителю.

2. Иногда могут использоваться отрицательные члены (как в `divu10` и `divu100`), но соответствующий анализ для определения таких ситуаций весьма сложен.
3. Иногда можно игнорировать некоторые из младших битов оценки обратного к делителю (сколько именно?).

Еще одной сложностью для некоторых целевых машин является то, что имеется много вариантов примеров кода, в которых оказывается больше команд, но которые могут выполняться быстрее на машинах с несколькими модулями суммирования или сдвига.

Коды, представленные в листингах 10.7–10.15, протестированы для всех 2^{32} значений делителей.

Знаковое деление

Методы, рассмотренные выше, можно сделать применимыми для знакового деления. Команда *сдвига вправо* при вычислении оценки частного становится командой *знакового сдвига вправо*, которая вычисляет результат деления с округлением к меньшему значению на степени 2. Таким образом, оценка частного оказывается слишком малой (алгебраически), так что остаток оказывается неотрицательным, как и в случае беззнакового деления.

Код наиболее естественным образом вычисляет результат деления с округлением к меньшему значению, так что нам нужна поправка, которая приведет его к обычному результату с отсечением по направлению к 0. Это можно сделать с помощью трех вычислительных команд, добавляя $d-1$ к делимому, если оно отрицательно. Например, если делитель равен 6, то код начинается с

```
n = n + (n >> 31 & 5);
```

(здесь сдвиг является *знаковым сдвигом*).

В остальном код очень похож на код в случае беззнакового деления. Количество требуемых элементарных операций обычно на три больше, чем в соответствующей функции беззнакового деления. Несколько примеров приведены в листингах 10.16–10.22. Все коды прошли исчерпывающее тестирование.

Листинг 10.16. Знаковое деление на 3

```
int divs3(int n)
{
    int q, r;

    n = n + (n >> 31 & 2);           // Прибавить 2, если n < 0
    q = (n >> 2) + (n >> 4);        // q = n*0.0101 (приблизленно)
    q = q + (q >> 4);               // q = n*0.01010101.
    q = q + (q >> 8);
    q = q + (q >> 16);
    r = n - q*3;                    // 0 <= r <= 14
    return q + (11*r >> 5);         // Возврат q + r/3
// return q + (5*(r + 1) >> 4);    // Альтернатива 1
// return q + ((r + 5 + (r << 2)) >> 4); // Альтернатива 2
}
```

ЛИСТИНГ 10.17. Знаковое деление на 5 и 6

```
int divs5(int n)
{
    int q, r;

    n = n + (n>>31 & 4);
    q = (n >> 1) + (n >> 2);
    q = q + (q >> 4);
    q = q + (q >> 8);
    q = q + (q >> 16);
    q = q >> 2;
    r = n - q*5;
    return q + (7*r >> 5);
// return q + (r>4) + (r>9);
}
```

```
int divs6(int n)
{
    int q, r;

    n = n + (n>>31 & 5);
    q = (n >> 1) + (n >> 3);
    q = q + (q >> 4);
    q = q + (q >> 8);
    q = q + (q >> 16);
    q = q >> 2;
    r = n - q*6;
    return q + ((r + 2) >> 3);
// return q + (r > 5);
}
```

ЛИСТИНГ 10.18. Знаковое деление на 7 и 9

```
int divs7(int n)
{
    int q, r;

    n = n + (n>>31 & 6);
    q = (n >> 1) + (n >> 4);
    q = q + (q >> 6);
    q = q + (q>>12) + (q>>24);
    q = q >> 2;
    r = n - q*7;
    return q + ((r + 1) >> 3);
// return q + (r > 6);
}
```

```
int divs9(int n)
{
    int q, r;

    n = n + (n>>31 & 8);
    q = (n >> 1) + (n >> 2) +
        (n >> 3);
    q = q + (q >> 6);
    q = q + (q>>12) + (q>>24);
    q = q >> 3;
    r = n - q*9;
    return q + ((r + 7) >> 4);
// return q + (r > 8);
}
```

ЛИСТИНГ 10.19. Знаковое деление на 10 и 11

```
int divs10(int n)
{
    int q, r;

    n = n + (n>>31 & 9);
    q = (n >> 1) + (n >> 2);
    q = q + (q >> 4);
    q = q + (q >> 8);
    q = q + (q >> 16);
    q = q >> 3;
    r = n - q*10;
    return q + ((r + 6) >> 4);
// return q + (r > 9);
}
```

```
int divs11(int n)
{
    int q, r;

    n = n + (n>>31 & 10);
    q = (n >> 1) + (n >> 2) -
        (n >> 5) + (n >> 7);
    q = q + (q >> 10);
    q = q + (q >> 20);
    q = q >> 3;
    r = n - q*11;
    return q + ((r + 5) >> 4);
// return q + (r > 10);
}
```


Листинг 10.20. Знаковое деление на 12 и 13

```

int divs12(int n)
{
    int q, r;

    n = n + (n>>31 & 11);
    q = (n >> 1) + (n >> 3);
    q = q + (q >> 4);
    q = q + (q >> 8);
    q = q + (q >> 16);
    q = q >> 3;
    r = n - q*12;
    return q + ((r + 4) >> 4);
// return q + (r > 11);
}

int divs13(int n)
{
    int q, r;

    n = n + (n>>31 & 12);
    q = (n>>1) + (n>>4);
    q = q + (q>>4) + (q>>5);
    q = q + (q>>12) + (q>>24);
    q = q >> 3;
    r = n - q*13;
    return q + ((r + 3) >> 4);
// return q + (r > 12);
}

```

Листинг 10.21. Знаковое деление на 100

```

int divs100(int n)
{
    int q, r;

    n = n + (n>>31 & 99);
    q = (n >> 1) + (n >> 3) + (n >> 6) - (n >> 10) +
        (n >> 12) + (n >> 13) - (n >> 16);
    q = q + (q >> 20);
    q = q >> 6;
    r = n - q*100;
    return q + ((r + 28) >> 7);
// return q + (r > 99);
}

```

Листинг 10.22. Знаковое деление на 1000

```

int divs1000(int n)
{
    int q, r, t;

    n = n + (n>>31 & 999);
    t = (n >> 7) + (n >> 8) + (n >> 12);
    q = (n >> 1) + t + (n >> 15) + (t >> 11) + (t >> 14) +
        (n >> 26) + (t >> 21);
    q = q >> 9;
    r = n - q*1000;
    return q + ((r + 24) >> 10);
// return q + (r > 999);
}

```

10.19. Получение остатка суммированием цифр

В этом разделе рассматривается задача вычисления остатка от деления на константу без вычисления частного. Методы из этого раздела применимы только для делителей вида $2^k \pm 1$, где k — целое число, не меньшее 2, и в большинстве случаев код прибегает к получению значения из таблицы (команда индексированной *загрузки*) после достаточно короткого вычисления.

Мы будем часто использовать следующее элементарное свойство сравнимости по модулю.

ТЕОРЕМА С. Если $a \equiv b \pmod{m}$ и $c \equiv d \pmod{m}$, то

$$a + c \equiv b + d \pmod{m} \quad \text{и} \\ ac \equiv bd \pmod{m}.$$

Беззнаковый случай проще и рассматривается первым.

Беззнаковый остаток

При делении на 3 многократное умножение тривиального сравнения $1 \equiv 1 \pmod{3}$ на сравнение $2 \equiv -1 \pmod{3}$ в соответствии с теоремой С дает

$$2^k \equiv \begin{cases} 1 \pmod{3}, & k \text{ четно.} \\ -1 \pmod{3}, & k \text{ нечетно.} \end{cases}$$

Следовательно, число n , записываемое в бинарном виде как $\dots b_3 b_2 b_1 b_0$, удовлетворяет соотношению

$$n = \dots + b_3 \cdot 2^3 + b_2 \cdot 2^2 + b_1 \cdot 2 + b_0 \equiv \dots - b_3 + b_2 - b_1 + b_0 \pmod{3},$$

которое получается при многократном применении теоремы С. Таким образом, можно чередовать сложение и вычитание битов в бинарном представлении числа, чтобы получить меньшее число, которое имеет тот же остаток при делении на 3. Если сумма отрицательна, можно добавить кратное трем, чтобы сделать ее неотрицательной. Затем этот процесс может повторяться до тех пор, пока результат не окажется в диапазоне от 0 до 2.

Тот же метод работает и при поиске остатка десятичного числа на 11.

Таким образом, если машина оснащена командой подсчета количества единичных битов, то функция, которая вычисляет остаток по модулю 3 беззнакового целого числа n , может начинаться с

$$n = \text{pop}(n \& 0x55555555) - \text{pop}(n \& 0xAAAAAAAA);$$

Это выражение можно упростить, используя следующее удивительное тождество, открытое Паоло Бонзини (Paolo Bonzini) [12].

$$\text{pop}(x \& \bar{m}) - \text{pop}(x \& m) = \text{pop}(x \oplus m) - \text{pop}(m) \tag{34}$$

Доказательство

$$\begin{aligned} \text{pop}(x \& \bar{m}) - \text{pop}(x \& m) &= \\ &= \text{pop}(x \& \bar{m}) - (32 - \text{pop}(\overline{x \& m})) = && \text{pop}(a) = 32 - \text{pop}(\bar{a}) \\ &= \text{pop}(x \& \bar{m}) + \text{pop}(\bar{x} \mid \bar{m}) - 32 = && \text{Законы де Моргана} \\ &= \text{pop}(x \& \bar{m}) + \text{pop}(\bar{x} \& m) + \text{pop}(\bar{m}) - 32 = && \text{pop}(a \mid b) = \text{pop}(a \& \bar{b}) + \text{pop}(b) \\ &= \text{pop}((x \& \bar{m}) \mid (\bar{x} \& m)) + \text{pop}(\bar{m}) - 32 = && \text{Непересекающиеся множества} \\ &= \text{pop}(x \oplus m) - \text{pop}(m) \end{aligned}$$

Поскольку ссылки на размер слова 32 сокращаются, результат справедлив для любого размера слова. Другой способ доказательства (34) состоит в том, чтобы заметить, что оно справедливо при $x = 0$, и если нулевой бит в x заменяется единичным там, где бит m равен 1, то обе стороны (34) уменьшаются на 1, а если нулевой бит x заменяется единичным там, где бит m равен 0, то обе части (34) увеличиваются на 1.

Применение (34) к приведенной выше строке кода на языке программирования C дает следующее.

```
n = pop(n ^ 0xAAAAAAAA) - 16;
```

Мы хотим применять это преобразование и далее, пока n не окажется в диапазоне от 0 до 2, если это возможно. Лучше избежать получения отрицательного значения n , чтобы не столкнуться с некорректной трактовкой знакового бита на следующей итерации. Избежать отрицательного значения можно, добавляя достаточно большое кратное 3 к n . Код Бончини, показанный в листинге 10.23, увеличивает константу на 39. Это более чем необходимо, чтобы сделать n неотрицательным, но при этом после второго этапа приведения n приводится к диапазону от -3 до 2 (вместо диапазона от -3 до 3). Это упрощает код оператора `return`, который прибавляет 3, если n отрицательно. Эта функция требует для выполнения 11 команд, с учетом двух команд для загрузки большой константы.

Листинг 10.23. Остаток при беззнаковом делении на 3 с использованием команды подсчета количества единичных битов

```
int remu3(unsigned n)
{
    n = pop(n ^ 0xAAAAAAAA) + 23; // Теперь 23 <= n <= 55
    n = pop(n ^ 0x2A) - 3;      // Теперь -3 <= n <= 2
    return n + (((int)n >> 31) & 3);
}
```

В листинге 10.24 показана версия кода, выполняющаяся за 4 команды, плюс простая команда поиска в таблице (например, команда индексированной *загрузки байта*).

Листинг 10.24. Остаток при беззнаковом делении на 3 с использованием команды подсчета количества единичных битов и поиска в таблице

```
int remu3(unsigned n)
{
    static char table[33] = {2, 0,1,2, 0,1,2, 0,1,2,
                            0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2,
                            0,1,2, 0,1};

    n = pop(n ^ 0xAAAAAAAA);
    return table[n];
}
```

Чтобы избежать команды *подсчета количества единичных битов*, заметим, что поскольку $4 \equiv 1 \pmod{3}$, постольку $4^k \equiv 1 \pmod{3}$. Бинарное число можно рассматривать как число в системе счисления 4, беря его биты парами и рассматривая биты от 00 до 11

как цифры в системе счисления с основанием 4 в диапазоне от 0 до 3. Пары битов можно суммировать с использованием кода из листинга 5.1 на с. 104, пропуская первую выполнимую строку (при сложении переполнение не возникает). Конечная сумма находится в диапазоне от 0 до 48, и для приведения ее к диапазону от 0 до 2 можно использовать поиск в таблице. Результирующая функция выполняется за 16 элементарных команд, плюс индексированная загрузка.

Имеется подобный, но несколько лучший путь. В качестве первого шага n можно привести к меньшему значению, находящемуся в том же классе сравнимости по модулю 3 с помощью

```
n = (n >> 16) + (n & 0xFFFF);
```

Этот код разбивает число на две 16-битовые части и суммирует их. Вклад левых 16 бит по модулю 3 не изменяется при сдвиге на 16 разрядов, поскольку сдвинутое число, умноженное на 2^{16} , представляет собой исходное число, а $2^{16} \equiv 1 \pmod{3}$. В общем случае $2^k \equiv 1 \pmod{3}$, если k четно. Это свойство неоднократно (пять раз) использовано в коде в листинге 10.25. Код состоит из 19 команд. Это количество может быть уменьшено, если прервать суммирование раньше и прибегнуть к поиску в таблице, как показано в листинге 10.26 (9 команд плюс индексированная загрузка). Количество команд можно уменьшить до шести (плюс индексированная загрузка), если использовать таблицу размером $0x2FE = 766$ байт.

Листинг 10.25. Остаток при беззнаковом делении на 3 с использованием суммирования цифр и поиска в регистре

```
int remu3(unsigned n)
{
    n = (n >> 16) + (n & 0xFFFF); // Максимум 0x1FFFE
    n = (n >> 8) + (n & 0x00FF); // Максимум 0x2FD
    n = (n >> 4) + (n & 0x000F); // Максимум 0x3D
    n = (n >> 2) + (n & 0x0003); // Максимум 0x11
    n = (n >> 2) + (n & 0x0003); // Максимум 0x6
    return (0x0924 >> (n << 1)) & 3;
}
```

Листинг 10.26. Остаток при беззнаковом делении на 3 с использованием суммирования цифр и поиска в памяти

```
int remu3(unsigned n)
{
    static char table[62] = {0,1,2, 0,1,2, 0,1,2, 0,1,2,
        0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2,
        0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2,
        0,1,2, 0,1,2, 0,1};
    n = (n >> 16) + (n & 0xFFFF); // Максимум 0x1FFFE
    n = (n >> 8) + (n & 0x00FF); // Максимум 0x2FD
    n = (n >> 4) + (n & 0x000F); // Максимум 0x3D
    return table[n];
}
```

Для вычисления остатка от беззнакового деления на 5 код в листинге 10.27 использует соотношения $16^4 \equiv 1 \pmod{5}$ и $4 \equiv -1 \pmod{5}$. Он состоит из 21 элементарной команды в предположении, что умножение на 3 выполняется с помощью *сдвига и сложения*.

ЛИСТИНГ 10.27. Остаток при беззнаковом делении на 5 с использованием суммирования цифр

```
int remu5(unsigned n)
{
    n = (n >> 16) + (n & 0xFFFF);           // Максимум 0x1FFFE
    n = (n >> 8) + (n & 0x00FF);            // Максимум 0x2FD
    n = (n >> 4) + (n & 0x000F);           // Максимум 0x3D
    n = (n>>4) - ((n>>2) & 3) + (n & 3);    // От -3 до 6
    return (01043210432
           >> 3*(n + 3)) & 7;
}
```

Количество команд можно уменьшить, прибегнув к таблице, подобной использованной в листинге 10.26. Код оказывается идентичным, за исключением таблицы.

```
static char table[62] = {0,1,2,3,4, 0,1,2,3,4,
    0,1,2,3,4, 0,1,2,3,4, 0,1,2,3,4, 0,1,2,3,4,
    0,1,2,3,4, 0,1,2,3,4, 0,1,2,3,4, 0,1,2,3,4,
    0,1,2,3,4, 0,1,2,3,4, 0,1};
```

Для поиска остатка от беззнакового деления на 7 код в листинге 10.28 использует соотношение $8^4 \equiv 1 \pmod{7}$ (9 элементарных команд плюс индексированная *загрузка*).

В качестве последнего примера код из листинга 10.29 вычисляет остаток от беззнакового деления на 9. Он основан на соотношении $8 \equiv -1 \pmod{9}$. Как видно из листинга, для вычисления требуются 9 элементарных команд плюс индексированная *загрузка*. Количество элементарных команд может быть уменьшено до шести, если использовать таблицу из 831 элемента.

ЛИСТИНГ 10.28. Остаток при беззнаковом делении на 7 с использованием суммирования цифр

```
int remu7(unsigned n)
{
    static char table[75] = {0,1,2,3,4,5,6, 0,1,2,3,4,5,6,
        0,1,2,3,4,5,6, 0,1,2,3,4,5,6, 0,1,2,3,4,5,6,
        0,1,2,3,4,5,6, 0,1,2,3,4,5,6, 0,1,2,3,4,5,6,
        0,1,2,3,4,5,6, 0,1,2,3,4,5,6, 0,1,2,3,4};

    n = (n >> 15) + (n & 0x7FFF); // Максимум 0x27FFE
    n = (n >> 9) + (n & 0x001FF); // Максимум 0x33D
    n = (n >> 6) + (n & 0x0003F); // Максимум 0x4A
    return table[n];
}
```

Листинг 10.29. Остаток при беззнаковом делении на 9 с использованием суммирования цифр

```
int remu9(unsigned n)
{
    int r;
    static char table[75] = {0,1,2,3,4,5,6,7,8,
        0,1,2,3,4,5,6,7,8, 0,1,2,3,4,5,6,7,8,
        0,1,2,3,4,5,6,7,8, 0,1,2,3,4,5,6,7,8,
        0,1,2,3,4,5,6,7,8, 0,1,2,3,4,5,6,7,8,
        0,1,2,3,4,5,6,7,8, 0,1,2};

    r = (n & 0x7FFF) - (n >> 15); // От FFFE0001 до 7FFF
    r = (r & 0x01FF) - (r >> 9); // От FFFFFFFC1 до 2FF
    r = (r & 0x003F) + (r >> 6); // От 0 до 4A
    return table[r];
}
```

Знаковый остаток

Метод суммирования цифр можно адаптировать для вычисления остатка, получающегося при знаковом делении. Похоже, лучшие всего это делать путем добавления нескольких шагов для коррекции результата, полученного методом, применяемым для беззнакового деления. Необходимы две коррекции: (1) поправка на разную интерпретацию знакового бита и (2) добавление или вычитание кратного делителю d для получения результата в диапазоне от 0 до $-(d-1)$.

Для деления на 3 код получения остатка в беззнаковом случае рассматривает знаковый бит делимого n как вклад 2 в остаток (поскольку $2^{31} \bmod 3 = 2$). В случае знакового деления знаковый бит дает вклад, равный 1 (поскольку $(-2^{31}) \bmod 3 = 1$). Следовательно, можно использовать код для беззнакового остатка и исправить полученный результат, вычитая 1. Далее, результат должен находиться в диапазоне от 0 до -2 , т.е. результат вычисления беззнакового остатка должен быть отображен следующим образом.

$$(0,1,2) \Rightarrow (-1,0,1) \Rightarrow (-1,0,-2)$$

Достаточно эффективно этого можно добиться путем вычитания 1 из беззнакового остатка, если он равен 0 или 1, и 4, если остаток равен 2 (когда делимое отрицательно). Код не должен менять значение делимого n , поскольку оно требуется на последнем шаге.

Эту процедуру легко применить к любой из функций для вычисления остатка от беззнакового деления на 3. Например, применение к коду из листинга 10.26 на с. 291 дает функцию, показанную в листинге 10.30. Она состоит из 13 элементарных команд, плюс индексированная *загрузка*. Количество команд может быть уменьшено ценой увеличения таблицы.

Листинг 10.30. Остаток при знаковом делении на 3 с использованием суммирования цифр

```
int rem3(int n)
{
    unsigned r;
```

```

static char table[62] = {0,1,2, 0,1,2, 0,1,2, 0,1,2,
    0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2,
    0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2, 0,1,2,
    0,1,2, 0,1,2, 0,1};

r = n;
r = (r >> 16) + (r & 0xFFFF); // Максимум 0x1FFFE
r = (r >> 8) + (r & 0x00FF); // Максимум 0x2FD
r = (r >> 4) + (r & 0x000F); // Максимум 0x3D
r = table[r];
return r - (((unsigned)n >> 31) << (r & 2));
}

```

В листингах 10.31–10.33 показан аналогичный код для вычисления остатка от знакового деления на 5, 7 и 9. Все эти функции состоят из 15 элементарных операций, плюс индексированная *загрузка*. Они используют знаковые сдвиги вправо, а последняя коррекция состоит в вычитании модуля, если делимое отрицательно, а остаток ненулевой. Количество команд может быть уменьшено ценой увеличения таблицы.

ЛИСТИНГ 10.31. Остаток при знаковом делении на 5 с использованием суммирования цифр

```

int rem5(int n)
{
    int r;
    static char table[62] = {2,3,4, 0,1,2,3,4, 0,1,2,3,4,
        0,1,2,3,4, 0,1,2,3,4, 0,1,2,3,4, 0,1,2,3,4,
        0,1,2,3,4, 0,1,2,3,4, 0,1,2,3,4, 0,1,2,3,4,
        0,1,2,3,4, 0,1,2,3};

    r = (n >> 16) + (n & 0xFFFF); // От FFFF8000 до 17FFE
    r = (r >> 8) + (r & 0x00FF); // От FFFFFFF80 до 27D
    r = (r >> 4) + (r & 0x000F); // От -8 до 53 (десятичных)
    r = table[r + 8];
    return r - (((int)(n & -r) >> 31) & 5);
}

```

ЛИСТИНГ 10.32. Остаток при знаковом делении на 7 с использованием суммирования цифр

```

int rem7(int n)
{
    int r;
    static char table[75] = {5,6, 0,1,2,3,4,5,6,
        0,1,2,3,4,5,6, 0,1,2,3,4,5,6, 0,1,2,3,4,5,6,
        0,1,2,3,4,5,6, 0,1,2,3,4,5,6, 0,1,2,3,4,5,6,
        0,1,2,3,4,5,6, 0,1,2,3,4,5,6, 0,1,2,3,4,5,6,
        0,1,2};

    r = (n >> 15) + (n & 0x7FFF); // От FFFF0000 до 17FFE
    r = (r >> 9) + (r & 0x001FF); // От FFFFFFF80 до 2BD
    r = (r >> 6) + (r & 0x0003F); // От -2 до 72 (десятичных)
    r = table[r + 2];
    return r - (((int)(n & -r) >> 31) & 7);
}

```

Листинг 10.33. Остаток при знаковом делении на 9 с использованием суммирования цифр

```
int rem9(int n)
{
    int r;
    static char table[75] = {7,8, 0,1,2,3,4,5,6,7,8,
        0,1,2,3,4,5,6,7,8, 0,1,2,3,4,5,6,7,8,
        0,1,2,3,4,5,6,7,8, 0,1,2,3,4,5,6,7,8,
        0,1,2,3,4,5,6,7,8, 0,1,2,3,4,5,6,7,8,
        0,1,2,3,4,5,6,7,8, 0};

    r = (n & 0x7FFF) - (n >> 15); // От FFFF7001 до 17FFF
    r = (r & 0x01FF) - (r >> 9); // От FFFFFFF41 до 0x27F
    r = (r & 0x003F) + (r >> 6); // От -2 до 72 (десятичных)
    r = table[r + 2];
    return r - (((int)(n & -r) >> 31) & 9);
}
```

10.20. Получение остатка путем умножения и сдвига вправо

Описанный в этом разделе метод применим, в принципе, ко всем целым делителям, большим 2, но с точки зрения практического применения он годится только для малых делителей и делителей вида $2^t - 1$. Как и в предыдущем разделе, в большинстве случаев код выполняет поиск в таблице после достаточно короткого вычисления.

Беззнаковый остаток

В этом разделе используется математическое (а не компьютерно-алгебраическое) обозначение $a \bmod b$, где a и b — целые числа и $b > 0$, и которое используется для указания числа $0 \leq x < b$, удовлетворяющего условию $x \equiv a \pmod{b}$.

Чтобы вычислить $x \bmod 3$, заметим, что

$$n \bmod 3 = \left\lfloor \frac{4}{3}n \right\rfloor \bmod 4. \quad (35)$$

Доказательство. Пусть $n = 3k + \delta$, где δ и k — целые числа и $0 \leq \delta \leq 2$. Тогда

$$\left\lfloor \frac{4}{3}(3k + \delta) \right\rfloor \bmod 4 = \left\lfloor 4k + \frac{4\delta}{3} \right\rfloor \bmod 4 = \left\lfloor \frac{4\delta}{3} \right\rfloor \bmod 4.$$

Ясно, что значение последнего выражения равно 0, 1 или 2 при $\delta = 0, 1$ или 2 соответственно. Это позволяет заменить задачу вычисления остатка по модулю 3 задачей вычисления остатка по модулю 4, что, конечно, на бинарном компьютере выполнить проще.

Соотношения наподобие (35) выполняются не для всех модулей, но аналогичные соотношения имеются для модулей вида $2^t - 1$ при целом k , большем 1. Например, легко показать, что

$$n \bmod 7 = \left\lfloor \frac{8}{7}n \right\rfloor \bmod 8.$$

Для чисел вида, отличного от $2^t - 1$, таких простых соотношений нет, но имеется определенное свойство единственности, которое может быть использовано для вычисления

остатка для других делителей. Например, если делитель представляет собой десятичное число 10, рассмотрим выражение

$$\left\lfloor \frac{16}{10} n \right\rfloor \bmod 16. \quad (36)$$

Пусть $n = 10k + \delta$, где $0 \leq \delta \leq 9$. Тогда

$$\left\lfloor \frac{16}{10} n \right\rfloor \bmod 16 = \left\lfloor \frac{16}{10} (10k + \delta) \right\rfloor \bmod 16 = \left\lfloor \frac{16\delta}{10} \right\rfloor \bmod 16.$$

Для $\delta = 0, 1, 2, 3, 4, 5, 6, 7, 8$ и 9 последнее выражение принимает значения $0, 1, 3, 4, 6, 8, 9, 11, 12$ и 14 соответственно. Все числа в последнем множестве различны. Следовательно, если имеется достаточно простой способ вычисления (36), можно преобразовать 0 в 0 , 1 в 1 , 3 в 2 , 4 в 3 и так далее, получая таким образом остатки от деления на 10 . В общем случае этот метод требует применения таблицы преобразования с размером, равным ближайшей степени 2 , большей делителя, так что этот метод практичен только для достаточно малых делителей (и для делителей вида $2^t - 1$, для которых поиск в таблице не требуется).

Приведенный далее код получен при скромном применении описанной выше теории ценой огромного количества проб и ошибок.

Рассмотрим остаток от беззнакового деления на 3 . Следуя (35), мы хотим вычислить два крайних справа бита целой части $4n/3$. Это можно сделать приближенно, выполняя умножение на $\lfloor 2^{32}/3 \rfloor$, а затем деление на 2^{30} с помощью команды *сдвига вправо*. После того как умножение на $\lfloor 2^{32}/3 \rfloor$ выполнено (с помощью команды *умножения*, которая дает младшие 32 бита результата), старшие биты оказываются потерянными. Но это не имеет значения и в действительности даже полезно, поскольку нас интересует результат по модулю 4 . Следовательно, так как $\lfloor 2^{32}/3 \rfloor = 0x55555555$, возможный план заключается в вычислении

$$r \leftarrow (0x55555555 * n) \gg 30.$$

Эксперимент показывает, что этот способ работает для n из диапазона от 0 до $2^{30} + 2$. Я бы сказал, что он почти работает. — если n не равно 0 и кратно 3 , в результате мы получаем 3 . Следовательно, необходим шаг, преобразующий $(0, 1, 2, 3)$ соответственно в $(0, 1, 2, 0)$.

Чтобы расширить диапазон применимости, умножение должно выполняться более точно. Двух битов точности оказывается достаточно (т.е. умножения на число $0x55555555.4$). Приведенное далее вычисление с последующим шагом преобразования работает для всех n , представимых с помощью беззнакового 32-битового числа.

$$r \leftarrow \left(0x55555555 * n + \left(n \gg 2 \right) \right) \gg 30$$

Конечно, можно дать формальное доказательство этого факта, но соответствующие выкладки очень длинные, и в них легко ошибиться.

Шаг преобразования результата можно выполнить с помощью трех или четырех команд на большинстве машин, но есть способ избежать его ценой двух команд. Приведенное выше выражение для вычисления r дает заниженную оценку. Если ее немного повысить, результат всегда оказывается равным 0, 1 или 2. Это приводит нас к функции на языке программирования C, показанной в листинге 10.34 (восемь команд, включая *умножение*).

ЛИСТИНГ 10.34. Беззнаковый остаток по модулю 3, метод умножения

```
int remu3(unsigned n)
{
    return (0x55555555*n + (n >> 1) - (n >> 3)) >> 30;
}
```

Умножение можно заменить *сдвигами* и *сложениями*, так что получится код, приведенный в листинге 10.35 и использующий только 13 таких команд.

ЛИСТИНГ 10.35. Беззнаковый остаток по модулю 3, метод умножения на основе сдвигов и сложений

```
int remu3(unsigned n)
{
    unsigned r;

    r = n + (n << 2);
    r = r + (r << 4);
    r = r + (r << 8);
    r = r + (r << 16);
    r = r + (n >> 1);
    r = r - (n >> 3);
    return r >> 30;
}
```

Остаток от беззнакового деления на 5 можно вычислить очень схожим с вычисление остатка от деления на 3 способом. Пусть $n = 5k + r$, где $0 \leq r \leq 4$. Тогда $(8/5)n \bmod 8 = (8/5)(5k + r) \bmod 8 = (8/5)r \bmod 8$. При $r = 0, 1, 2, 3$ и 4 это дает нам значения 0, 1, 3, 4 и 6 соответственно. Поскольку $\lfloor 2^{32}/5 \rfloor = 0x33333333$, это приводит к функции, показанной в листинге 10.36 (11 команд, включая *умножение*). Последний шаг (код оператора return) отображает (0,1,3,4,6,7) на (0,1,2,3,4,0) соответственно, используя метод работы с регистром вместо индексированной *загрузки* из памяти. Отображение, кроме того, 2 на 2 и 5 на 4 снижает требуемую точность при умножении на $2^{32}/5$, что позволяет обойтись одним лишь членом $n \gg 3$ для приближения отсутствующей части множителя (шестнадцатеричное значение 0.333...). Если опустить корректирующий член $n \gg 3$, код останется работоспособным для значений n от 0 до 0x60000004.

Листинг 10.36. Беззнаковый остаток по модулю 5, метод умножения

```
int remu5(unsigned n)
{
    n = (0x33333333*n + (n >> 3)) >> 29;
    return (0x04432210 >> (n << 2)) & 7;
}
```

Код для вычисления остатка от деления на 7 похож на приведенный выше, но в нем проще шаг отображения — требуется только преобразование 7 в 0. Один из способов решения этой задачи показан в листинге 10.37 (11 команд, включая *умножение*). Если опустить корректирующий член $n \gg 4$, код останется работоспособным для значений n до $0x40000006$. При удалении обоих корректирующих членов код работает для значений n до $0x08000006$.

Листинг 10.37. Беззнаковый остаток по модулю 7, метод умножения

```
int remu7(unsigned n)
{
    n = (0x24924924*n + (n >> 1) + (n >> 4)) >> 29;
    return n & ((int)(n - 7) >> 31);
}
```

Код для вычисления остатка от деления на 9 приведен в листинге 10.38. Он состоит из шести команд, включая *умножение*, плюс индексированная *загрузка*. Если опустить корректирующий член $n \gg 1$, а множитель заменить на $0x1C71C71D$, то функция работает для значений n до $0x1999999E$.

Листинг 10.38. Беззнаковый остаток по модулю 9, метод умножения

```
int remu9(unsigned n)
{
    static char table[16] = {0, 1, 1, 2, 2, 3, 3, 4,
                           5, 5, 6, 6, 7, 7, 8, 8};

    n = (0x1C71C71C*n + (n >> 1)) >> 28;
    return table[n];
}
```

В листинге 10.39 показан способ вычисления беззнакового остатка по модулю 10. Он состоит из восьми команд, включая *умножение*, плюс индексированная *загрузка*. Если опустить корректирующий член $n \gg 3$, то код работает для значений n до $0x40000004$. При удалении обоих корректирующих членов код работает для значений n до $0x0AAAAAAD$.

Листинг 10.39. Беззнаковый остаток по модулю 10, метод умножения

```
int remu10(unsigned n)
{
    static char table[16] = {0, 1, 2, 2, 3, 3, 4, 5,
                           5, 6, 7, 7, 8, 8, 9, 0};

    n = (0x19999999*n + (n >> 1) + (n >> 3)) >> 28;
    return table[n];
}
```

В качестве последнего примера рассмотрим вычисление остатка по модулю 63. Эта функция используется в программе вычисления степени заполнения на с. 106. Джо Кин (Joe Keane) [64] предложил совершенно удивительный код, показанный в листинге 10.40. Он состоит из 12 элементарных команд из базового набора RISC.

Листинг 10.40. Беззнаковый остаток по модулю 63, метод Кина

```
int remu63(unsigned n)
{
    unsigned t;

    t = (((n >> 12) + n) >> 10) + (n << 2);
    t = ((t >> 6) + t + 3) & 0xFF;
    return (t - (t >> 6)) >> 2;
}
```

Метод умножения и сдвига вправо приводит к коду, показанному в листинге 10.41. Здесь выполняется 11 элементарных команд из базового набора RISC, одна из которых — *умножение*. Этот метод не так быстр, как метод Кина, если только машина не выполняет умножение очень быстро, а загрузку константы 0x04104104 можно вынести из цикла.

Листинг 10.41. Беззнаковый остаток по модулю 63, метод умножения

```
int remu63(unsigned n)
{
    n = (0x04104104*n + (n >> 4) + (n >> 10)) >> 26;
    return n & ((n - 63) >> 6); // Замена 63 нулевым
                                // значением
}
```

На некоторых машинах можно ускорить код, заменив умножение сдвигами и сложениями следующим образом (15 элементарных команд для всей функции).

```
r = (n << 2) + (n << 8); // r = 0x104*n
r = r + (r << 12);       // r = 0x104104*n
r = r + (n << 26);       // r = 0x04104104*n
```

Знаковый остаток

Как и в случае метода суммирования цифр, метод умножения и сдвига вправо можно адаптировать для вычисления остатка от знакового деления. Похоже, что и в этом случае нет более подходящего пути, чем добавление нескольких корректирующих шагов к методу, применяющемуся для беззнакового деления. Например, в листинге 10.42 показан код, полученный из листинга 10.34 на с. 297 (12 команд, включая *умножение*).

Листинг 10.42. Знаковый остаток по модулю 3, метод умножения

```
int rems3(int n)
{
    unsigned r;

    r = n;
    r = (0x55555555*r + (r >> 1) - (r >> 3)) >> 30;
    return r - (((unsigned)n >> 31) << (r & 2));
}
```



```

r = n;
r = (0x19999999*r + (r >> 1) + (r >> 3)) >> 28;
return table[r + (((unsigned)n >> 31) << 4)];
}

```

10.21. Преобразование в точное деление

Поскольку остаток оказывается возможным вычислить без вычисления частного, можно попробовать вычислять частное $q = \lfloor n/d \rfloor$, сначала вычисляя остаток, вычитая его из делимого n , а затем деля разность на делитель d . Последнее деление является точным, так что его можно выполнять путем умножения на мультипликативное обратное к d (см. раздел 10.16, “Точное деление на константу”, на с. 266). Этот метод оказывается особенно привлекательным, когда требуется найти и частное, и остаток.

Давайте попробуем применить наши рассуждения к случаю беззнакового деления на 3. Вычисление остатка методом умножения (листинг 10.34 на с. 297) приводит к функции, показанной в листинге 10.47.

Листинг 10.47. Беззнаковый остаток и частное для делителя 3 с использованием точного деления

```

unsigned divu3(unsigned n)
{
    unsigned r;

    r = (0x55555555*n + (n >> 1) - (n >> 3)) >> 30;
    return (n - r)*0xAAAAAAAB;
}

```

Эта функция включает 11 команд, в том числе два умножения на большие числа (константа 0x55555555 получается из константы 0xAAAAAAAB сдвигом вправо на одну позицию). Более прямолинейный метод вычисления частного q с использованием, например, кода из листинга 10.8 на с. 281 требует 14 команд, включая два умножения на небольшие числа, или 17 команд при замене умножений *сдвигами* и *сложениями*. Если требуется и значение остатка, которое вычисляется как $r = n - q * 3$, прямолинейный метод требует выполнения 16 команд, включая три умножения на небольшие числа или 20 команд при замене умножений *сдвигами* и *сложениями*.

Код из листинга 10.47 при замене умножений *сдвигами* и *сложениями* привлекательнее не становится: результат состоит из 24 элементарных команд. Таким образом, метод точного деления может быть применим, в первую очередь, на машинах, на которых нет команды вычисления *старшего слова умножения*, но есть быстрое *умножение по модулю 2^{32}* и медленное *деление*, в особенности если они легко работают с большими константами.

Для знакового деления на 3 метод точного деления может быть закодирован так, как показано в листинге 10.48. Этот код состоит из 15 команд, включая два умножения на большие константы.

**Листинг 10.48. Знаковый остаток и частное для делителя 3
с использованием точного деления**

```
int divs3(int n)
{
    unsigned r;

    r = n;
    r = (0x55555555*r + (r >> 1) - (r >> 3)) >> 30;
    r = r - (((unsigned)n >> 31) << (r & 2));
    return (n - r)*0xAAAAAAAB;
}
```

В качестве последнего примера в листинге 10.49 приведен код для вычисления частного и остатка от беззнакового деления на 10. В нем 12 команд, включая два умножения на большие константы, плюс команда индексированной загрузки.

**Листинг 10.49. Беззнаковый остаток и частное для делителя 10
с использованием точного деления**

```
unsigned divu10(unsigned n)
{
    unsigned r;
    static char table[16] = {0, 1, 2, 2, 3, 3, 4, 5,
                           5, 6, 7, 7, 8, 8, 9, 0};

    r = (0x19999999*n + (n >> 1) + (n >> 3)) >> 28;
    r = table[r];
    return ((n - r) >> 1)*0xCCCCCCCD;
}
```

10.22. Проверка времени выполнения

На многих машинах имеется команда *умножения* $32 \times 32 \Rightarrow 64$, так что можно ожидать, что при делении на константу, такую, как 3, код, показанный на с. 253, окажется самым быстрым. Если такой команды *умножения* на машине нет, но есть быстрая команда *умножения* $32 \times 32 \Rightarrow 32$, то лучшим для машин с быстрым умножением и медленным делением может оказаться метод точного деления. Чтобы проверить это утверждение, для сравнения четырех методов деления на 3 была написана программа на ассемблере, результаты работы которой подытожены в табл. 10.4. Использовалась машина Pentium III 667 МГц (ок. 2000 г.); следует ожидать сходных результатов и на других машинах.

Таблица 10.4. Беззнаковое деление на 3 на Pentium III

Метод деления	Количество тактов
Использование машинной команды деления (divl)	41.08
Использование <i>умножения</i> $32 \times 32 \Rightarrow 64$ (код на с. 253)	4.28
Все элементарные команды (листинг 10.8 на с. 281)	14.10
Приведение к точному делению (листинг 10.47 на с. 301)	6.68

Первая строка дает время в тактах для двух команд (`xorl` для сброса левой половины 64-битового регистра источника и для команды `divl`), которые требуют для выполнения 40 тактов. Вторая строка также дает время выполнения двух команд: *умножения* и *сдвига вправо* на 1 бит (`mull` и `shrl`). Третья строка указывает время выполнения 21 элементарной команды. Это код из листинга 10.8 на с. 281, использующий вторую альтернативу, и с умножением на 3, осуществляемым одной командой (`leal`). Требуются также несколько команд *перемещения* в силу того, что машина (в основном) двухадресная. Последняя строка содержит время выполнения последовательности из 10 команд, включая два умножения (`imull`). Эти две команды `imull` используют для больших констант четырехбайтовые непосредственные поля. (Использована команда знакового *умножения* `imull`, а не ее беззнаковый двойник `mull`, так как они дают один и тот же результат в младших 32 битах и при этом команда `imull` имеет большее количество режимов адресации).

Метод точного деления даже предпочтительнее второго и третьего методов, если требуются как частное, так и остаток, поскольку в таком случае к этим методам следует добавить код вычисления $r \leftarrow n - q * 3$ (команда `divl` дает как частное, так и остаток).

10.23. Аппаратная схема для деления на 3

Имеется простая схема деления на 3 — примерно того же уровня сложности, что и сумматор. Она может быть построена методом, очень схожим с тем, который применяется для построения n -битового сумматора из n 1-битовых “полных сумматоров”. Однако в делителе сигналы идут от старшего бита к младшему.

Рассмотрим деление на 3 школьным методом “в столбик”, но в бинарной системе счисления. Для получения каждого бита частного мы делим на 3 очередной бит, которому предшествует остаток 0, 1 или 2 от предыдущего шага. Соответствующая логика показана в табл. 10.5. Здесь остаток представлен двумя битами, r_i и s_i , где r_i представляет собой старший бит. Остаток никогда не равен 3, так что последние две строки таблицы представляют собой случаи “не имеет значения”.

ТАБЛИЦА 10.5. ЛОГИКА ДЕЛЕНИЯ НА 3

r_{i+1}	s_{i+1}	x_i	y_i	r_i	s_i
0	0	0	0	0	0
0	0	1	0	0	1
0	1	0	0	1	0
0	1	1	1	0	0
1	0	0	1	0	1
1	0	1	1	1	0
1	1	0	-	-	-
1	1	1	-	-	-

Схема для 32-битового деления на 3 показана на рис. 10.1. Частное представляет собой слово, состоящее из битов с y_{31} по y_0 , а остаток равен $2r_0 + s_0$.

Еще один способ аппаратной реализации операции деления на 3 — использование умножителя для того, чтобы умножить делимое на обратное к 3 (бинарное $0.010101\dots$) с соответствующим округлением и масштабированием. Эта методика показана на с. 233 и 253.

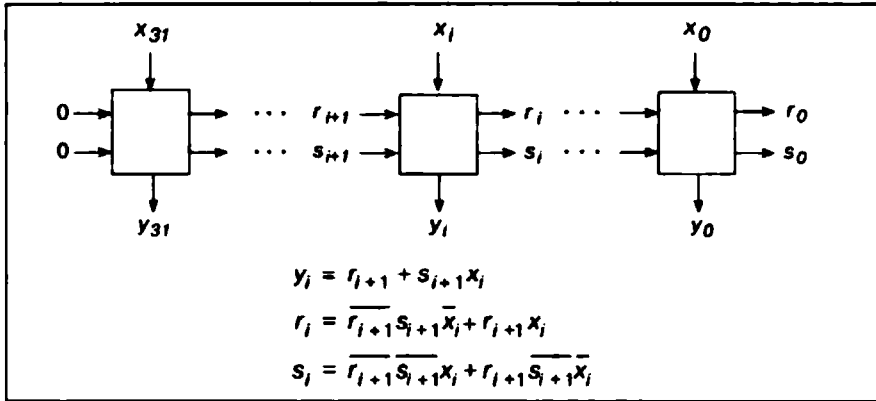


Рис. 10.1. Логическая схема для деления на 3

Упражнения

1. Покажите, что в случае беззнакового деления на четное число команды `shrci` (или эквивалентного кода) можно избежать путем (а) сброса младшего бита делимого (операцией `и`) [14] или (б) деления делимого на 2 (команда `сдвига вправо на 1 бит`) с последующим делением на половину делителя.
2. Запишите на языке программирования Python код функции, аналогичной функции из листинга 10.4 на с. 266, но для вычисления магического числа для знакового деления. Рассматривайте только положительные делители.
3. Покажите, как использовать метод Ньютона для вычисления мультипликативного обратного целого числа d по модулю 81. Покажите ход вычислений для $d = 146$.

ГЛАВА 11

НЕКОТОРЫЕ ЭЛЕМЕНТАРНЫЕ ФУНКЦИИ

11.1. Целочисленный квадратный корень

Под целочисленным квадратным корнем подразумевается функция $\lfloor \sqrt{x} \rfloor$. Чтобы расширить область применения и избежать решения вопроса о том, как поступать с отрицательными аргументами функции, будем считать x беззнаковой величиной: $0 \leq x \leq 2^{22} - 1$.

Метод Ньютона

Для чисел с плавающей точкой, по сути, универсальным методом вычисления квадратного корня является метод Ньютона. Этот метод начинается с некоего значения g_0 , которое является начальной оценкой \sqrt{a} . Затем выполняется серия уточнений значения квадратного корня по формуле

$$g_{n+1} = \left(g_n + \frac{a}{g_n} \right) / 2.$$

Приведенный метод имеет квадратичную сходимость, т.е. если в некоторый момент g_n имеет точность n бит, то g_{n+1} имеет точность $2n$ бит. Рабочая программа должна иметь средства для выявления достижения необходимой точности и прекращения вычислений.

Приятной неожиданностью оказывается то, что метод Ньютона хорошо работает и в области целых чисел. Для того чтобы убедиться в этом, рассмотрим следующую теорему.

ТЕОРЕМА. Пусть $g_{n+1} = \lfloor (g_n + \lfloor a/g_n \rfloor) / 2 \rfloor$, где g_n — целое число, большее 0. Тогда

- а) если $g_n > \lfloor \sqrt{a} \rfloor$, то $\lfloor \sqrt{a} \rfloor \leq g_{n+1} < g_n$, и
- б) если $g_n = \lfloor \sqrt{a} \rfloor$, то $\lfloor \sqrt{a} \rfloor \leq g_{n+1} \leq \lfloor \sqrt{a} \rfloor + 1$.

Иными словами, если есть завышенная целочисленная оценка g_n величины $\lfloor \sqrt{a} \rfloor$, то следующее приближение будет строго меньше предыдущего, но не меньше $\lfloor \sqrt{a} \rfloor$. Таким образом, если начать со слишком большого приближения, то будет получена монотонно убывающая последовательность приближений. Если же начальное приближение $g_n = \lfloor \sqrt{a} \rfloor$, то следующее приближение оказывается либо равным предыдущему, либо на 1 большим. Таким образом, мы получаем способ определения того, что последовательность приближений сошлась: если начать с $g_0 \geq \lfloor \sqrt{a} \rfloor$, то точное решение g_n будет достигнуто в тот момент, когда $g_{n+1} \geq g_n$, и результат оказывается в точности равным g_n .

Случай $a = 0$ следует рассматривать отдельно в связи с тем, что при этом может возникнуть ситуация деления 0 на 0.

Доказательство. а) Поскольку g_n — целое число,

$$g_{n+1} = \left\lfloor \left(g_n + \left\lfloor \frac{a}{g_n} \right\rfloor \right) / 2 \right\rfloor = \left\lfloor \left\lfloor g_n + \frac{a}{g_n} \right\rfloor / 2 \right\rfloor = \left\lfloor \left(g_n + \frac{a}{g_n} \right) / 2 \right\rfloor = \left\lfloor \frac{g_n^2 + a}{2g_n} \right\rfloor.$$

Так как $g_n > \lfloor \sqrt{a} \rfloor$ и g_n — целое число, $g_n > \sqrt{a}$. Определим ε следующим образом: $g_n = (1 + \varepsilon)\sqrt{a}$. Тогда $\varepsilon > 0$ и

$$\begin{aligned} \left\lfloor \frac{g_n^2 + a}{2g_n} \right\rfloor &= g_{n+1} \leq \frac{g_n^2 + a}{2g_n}, \\ \left\lfloor \frac{(1 + \varepsilon)^2 a + a}{2(1 + \varepsilon)\sqrt{a}} \right\rfloor &= g_{n+1} < \frac{g_n^2 + g_n^2}{2g_n}, \\ \left\lfloor \frac{2 + 2\varepsilon + \varepsilon^2}{2(1 + \varepsilon)} \sqrt{a} \right\rfloor &= g_{n+1} < g_n, \\ \left\lfloor \frac{2 + 2\varepsilon}{2(1 + \varepsilon)} \sqrt{a} \right\rfloor &\leq g_{n+1} < g_n, \\ \lfloor \sqrt{a} \rfloor &\leq g_{n+1} < g_n. \end{aligned}$$

б) Поскольку $g_n = \lfloor \sqrt{a} \rfloor$, $\sqrt{a} - 1 < g_n \leq \sqrt{a}$, так что $g_n^2 \leq a < (g_n + 1)^2$. Следовательно,

$$\begin{aligned} \left\lfloor \frac{g_n^2 + g_n^2}{2g_n} \right\rfloor &\leq g_{n+1} \leq \left\lfloor \frac{g_n^2 + (g_n + 1)^2}{2g_n} \right\rfloor, \\ \lfloor g_n \rfloor &\leq g_{n+1} \leq \left\lfloor g_n + 1 + \frac{1}{2g_n} \right\rfloor, \\ \lfloor \sqrt{a} \rfloor &\leq g_{n+1} \leq \lfloor g_n + 1 \rfloor \quad (\text{так как } g_n \text{ — целое и } \frac{1}{2g_n} < 1), \\ \lfloor \sqrt{a} \rfloor &\leq g_{n+1} \leq \lfloor g_n \rfloor + 1 = \lfloor \sqrt{a} \rfloor + 1. \end{aligned}$$

Сложной частью при использовании метода Ньютона для вычисления $\lfloor \sqrt{x} \rfloor$ является получение первого приближения. Процедура в листинге 11.1 устанавливает первое приближение g_0 равным наименьшей степени 2, которая больше или равна \sqrt{x} . Например, $g_0 = 2$ для $x = 4$, а для $x = 5$ в качестве первого приближения принимается $g_0 = 4$.

Листинг 11.1. Метод Ньютона поиска целочисленного квадратного корня

```

int isqrt(unsigned x)
{
    unsigned x1;
    int s, g0, g1;

    if (x <= 1) return x;
    s = 1;
    x1 = x - 1;
    if (x1 > 65535) {s = s + 8; x1 = x1 >> 16;}
    if (x1 > 255)  {s = s + 4; x1 = x1 >> 8;}
    if (x1 > 15)   {s = s + 2; x1 = x1 >> 4;}
    if (x1 > 3)    {s = s + 1;}

    g0 = 1 << s; // g0 = 2**s
    g1 = (g0 + (x >> s)) >> 1; // g1 = (g0 + x/g0)/2
    while (g1 < g0) { // Повторяем, пока приближения
        g0 = g1; // строго уменьшаются
        g1 = (g0 + (x/g0)) >> 1;
    }
    return g0;
}

```

Поскольку первое приближение g_0 представляет собой степень 2, для получения g_1 нет необходимости в реальном делении — можно обойтись *сдвигом вправо*.

Поскольку точность первого приближения составляет около 1 бита, а метод Ньютона обеспечивает квадратичную сходимость (число точных битов удваивается при каждой итерации), следует ожидать, что процедура поиска квадратного корня завершится за пять итераций (на 32-разрядном компьютере), для чего потребуется четыре деления (первое деление заменяется сдвигом вправо). Эксперименты показывают, что максимальное количество делений равно пяти (четыре для аргументов, меньших 16785407).

Если в компьютере имеется команда вычисления *количества ведущих нулевых битов*, то получить первое приближение очень легко: заменить первые семь выполнимых строк в приведенной выше процедуре строками

```

if (x <= 1) return x;
s = 16 - nlz(x - 1)/2;

```

Еще один вариант в случае отсутствия команды вычисления *количества ведущих нулевых битов* — вычисление s посредством бинарного дерева поиска. Этот метод позволяет получить несколько лучшее приближение g_0 : наименьшую степень 2, которая больше или равна $\lfloor \sqrt{x} \rfloor$. Для некоторых значений x это дает меньшее значение g_0 , но достаточно большое, чтобы выполнялся критерий сходимости из рассмотренной ранее теоремы. Различие этих схем показано ниже.

Диапазон x		Первое приближение
для листинга 11.1	для листинга 11.2	
0	0	0
1	От 1 до 3	1
От 2 до 4	От 4 до 8	2
От 5 до 16	От 9 до 24	4
От 17 до 64	От 25 до 80	8
От 65 до 256	От 81 до 288	16
...
От $2^{2^k} + 1$ до $2^{2^{k+1}}$	От $(2^{2^k} + 1)^2$ до $(2^{2^{k+1}} + 1)^2 - 1$	2^{2^k}
От $2^{2^k} + 1$ до $2^{2^{k+1}} - 1$	От $(2^{2^k} + 1)^2$ до $2^{2^{k+1}} - 1$	2^{2^k}

Соответствующая процедура показана в листинге 11.2. Она особенно удобна при работе с малыми значениями x ($0 \leq x \leq 24$), так как при этом не требуется выполнение деления.

ЛИСТИНГ 11.2. Целочисленный квадратный корень с вычислением первого приближения путем бинарного поиска

```
int isqrt(unsigned x)
{
    int s, g0, g1;

    if (x <= 4224)
        if (x <= 24)
            if (x <= 3) return (x + 3) >> 2;
            else if (x <= 8) return 2;
            else return (x >> 4) + 3;
        else if (x <= 288)
            if (x <= 80) s = 3; else s = 4;
            else if (x <= 1088) s = 5; else s = 6;
        else if (x <= 1025*1025 - 1)
            if (x <= 257*257 - 1)
                if (x <= 129*129 - 1) s = 7; else s = 8;
                else if (x <= 513*513 - 1) s = 9; else s = 10;
            else if (x <= 4097*4097 - 1)
                if (x <= 2049*2049 - 1) s = 11; else s = 12;
            else if (x <= 16385*16385 - 1)
                if (x <= 8193*8193 - 1) s = 13; else s = 14;
            else if (x <= 32769*32769 - 1) s = 15; else s = 16;
        g0 = 1 << s; // g0 = 2**s
        // Далее все так же, как и в листинге 11.1
}
```

Время выполнения алгоритма из листинга 11.1 на компьютере с базовым набором RISC-команд в наихудшем случае составляет около $26 + (D+6)n$ тактов, где D — время выполнения команды деления в циклах, а n — количество итераций цикла `while`. Время выполнения алгоритма из листинга 11.2 составляет в наихудшем случае $27 + (D+6)n$ тактов в предположении (в обоих случаях), что команда *ветвления* выполняется за один такт. В приведенной ниже таблице показаны среднее количество итера-

ций цикла в обоих алгоритмах для равномерно распределенного в указанных диапазонах значения x .

x	Листинг 11.1	Листинг 11.2
От 0 до 9	0.80	0
От 0 до 99	1.46	0.83
От 0 до 999	1.58	1.44
От 0 до 9999	2.13	2.06
От 0 до $2^{32} - 1$	2.97	2.97

Если считать, что время деления равно 20 тактам, а x равномерно распределено от 0 до 9999, то оба алгоритма требуют около 81 такта процессорного времени.

Бинарный поиск

Поскольку алгоритм, основанный на методе Ньютона, начинается с бинарного поиска начального приближения, почему бы не воспользоваться бинарным поиском самого квадратного корня? Этот способ может начать работу с двух границ: по всей видимости, 0 и 2^{16} . Затем в качестве приближения квадратного корня берется середина между границами. Если квадрат этого значения больше аргумента x , то средняя точка становится новой верхней границей; в противном случае средняя точка становится новой нижней границей. Итерации продолжаются до тех пор, пока нижняя и верхняя границы не будут отличаться на 1. При этом нижняя граница будет представлять собой искомым результат.

Таким образом можно избежать деления, но потребуется значительное количество умножений, а именно 16, если в качестве начальных границ воспользоваться значениями 0 и 2^{16} . (Этот метод на каждой итерации вычисляет по одному биту точного значения.) В листинге 11.3 показан один из вариантов данного алгоритма, в котором используются несколько улучшенные по сравнению со значениями 0 и 2^{16} начальные границы. Кроме того, процедура в листинге 11.3 для большинства RISC-компьютеров экономит один такт на цикл, изменяя a и b так, что в программе используется сравнение $b \geq a$ вместо $b - a \geq 1$.

Листинг 11.3. Бинарный поиск целочисленного квадратного корня

```
int isqrt(unsigned x)
{
    unsigned a, b, m;    // Границы и средняя точка

    a = 1;
    b = (x >> 5) + 8;    // См. пояснения в тексте
    if (b > 65535) b = 65535;
    do {
        m = (a + b) >> 1;
        if (m*m > x) b = m - 1;
        else a = m + 1;
    } while (b >= a);
    return a - 1;
}
```

В начале каждой итерации должны выполняться предикаты $a \leq \lfloor \sqrt{x} \rfloor + 1$ и $b \geq \lfloor \sqrt{x} \rfloor$.

Начальное значение b должно быть легко вычисляемым и близким к $\lfloor \sqrt{x} \rfloor$. Вот некоторые из возможных начальных значений: x , $x+4+1$, $x+8+2$, $x+16+4$, $x+32+8$, $x+64+16$ и т.д. Выражения в начале этого списка лучше подходят для небольших значений x , выражения ближе к концу списка — для x побольше. (Значение $x+2+1$ также применимо, но не имеет смысла, так как выражение $x+4+1$ везде дает лучшую или такую же границу.)

Семь вариантов приведенного в листинге 11.3 алгоритма можно получить, заменяя a на $a+1$, b на $b-1$, используя вместо $m=(a+b)+2$ выражение $m=(a+b+1)+2$ или комбинируя перечисленные подстановки.

Время выполнения процедуры из листинга 11.3 составляет примерно $6+(M+7.5)n$ тактов, где M — время выполнения умножения в тактах, а n — количество выполняемых итераций цикла. Приведенная ниже таблица содержит средние количества выполнения циклов для равномерно распределенных в указанных интервалах значений x .

x	Среднее количество итераций цикла
От 0 до 9	3.00
От 0 до 99	3.15
От 0 до 999	4.68
От 0 до 9999	7.04
От 0 до $2^{32}-1$	16.00

Если считать, что время выполнения умножения — 5 тактов, а x равномерно распределено от 0 до 9999, то время выполнения алгоритма составит около 94 тактов. Максимальное время выполнения ($n=16$) составляет около 206 тактов.

Если компьютер оснащен командой вычисления количества ведущих нулевых битов, то начальные границы можно определить следующим образом.

$$b = (1 \ll (33 - \text{nlz}(x)) / 2) - 1;$$

$$a = (b + 3) / 2;$$

(То есть, $b = 2^{(33 - \text{nlz}(x)) / 2} - 1$.) Это очень хорошие границы для небольших x (так, при $0 \leq x \leq 15$ требуется только одна итерация), но для больших значений x дает только небольшое улучшение по сравнению с границами из листинга 11.3. При x из диапазона от 0 до 9999 среднее количество итераций примерно равно 5.45, что дает время выполнения около 74 тактов (при использовании тех же предположений, что и раньше).

Аппаратный алгоритм

Существует алгоритм, использующий при вычислении квадратного корня сдвиги и вычитания, который очень похож на алгоритм аппаратного деления из листинга 9.2 на с. 220. Будучи аппаратно встроенным в 32-битовый компьютер, этот алгоритм использу-

ет 64-битовый регистр, инициализируемый 32 нулевыми битами, за которыми следует аргумент x . При каждой итерации над 64-битовым регистром выполняется операция сдвига на два бита влево, а текущий результат y (изначально равный 0) — на один бит. Затем из левой половины 64-битового регистра вычитается значение $2y+1$. Если результат вычитания неотрицателен, он заменяет левую половину 64-битового регистра, а к значению y прибавляется 1 (для этого не требуется сумматор, так как y в этот момент заканчивается нулевым битом). Если результат вычитания отрицателен, то и 64-битовый регистр, и y остаются неизменными. Итерации выполняются 16 раз.

Этот алгоритм был описан в 1945 году [61].

Пожалуй, неожиданным результатом оказывается то, что время работы этого процесса составляет примерно половину времени, необходимого для выполнения деления $64+32 \Rightarrow 32$ с помощью упомянутого аппаратного алгоритма, так как в этом случае выполняется в два раза меньше итераций примерно той же сложности, что и в алгоритме деления.

При программном использовании данного алгоритма, вероятно, лучше всего избежать использования сдвига в двойном слове, которое требует около четырех команд сдвига. Алгоритм из листинга 11.4 [37] осуществляет это путем сдвига y и маскирования бита m справа. В среднем ему требуется выполнение 149 базовых RISC-команд. Два выражения $y \mid m$ могут быть заменены сложением $y+m$.

ЛИСТИНГ 11.4. Аппаратный алгоритм поиска целочисленного квадратного корня

```
int isqrt(unsigned x)
{
    unsigned m, y, b;

    m = 0x40000000;
    y = 0;
    while(m != 0) {           // Выполняется 16 раз
        b = y | m;
        y = y >> 1;
        if (x >= b) {
            x = x - b;
            y = y | m;
        }
        m = m >> 2;
    }
    return y;
}
```

Работа этого алгоритма похожа на школьный метод поиска квадратного корня. Ниже показан пример поиска $\lfloor \sqrt{179} \rfloor$ в соответствии с этим алгоритмом на 8-битовом компьютере.

1011 0011	x0	Изначально x = 179 (0xB3)		
- 1	b1			
<hr/>				
0111 0011	x1	0100 0000	y1	
- 101	b2	0010 0000	y2	
<hr/>				
0010 0011	x2	0011 0000	y2	
- 11 01	b3	0001 1000	y3	
<hr/>				
0010 0011	x3	0001 1000	y3	(Вычитать нельзя)
- 1 1001	b4	0000 1100	y4	
<hr/>				
0000 1010	x4	0000 1101	y4	

Результат вычислений равен 13; в регистре x остается остаток 10.

При использовании обычного трюка со *знаковым сдвигом вправо* на 31 бит можно избежать выполнения проверки $if\ x \geq b$. Можно доказать, что старший бит b всегда нулевой (на самом деле $b \leq 5 \cdot 2^{28}$), что упрощает предикат $x \geq b$ (см. раздел "Предикаты сравнения" на с. 43). В результате группу операторов if можно заменить следующими.

```
t = (int)(x | ~(x - b)) >> 31; // -1, если x >= b, иначе 0
x = x - (b & t);
y = y | (m & t);
```

Тем самым мы заменяем три такта семью в предположении, что компьютер оснащен командой *или-не*, однако это может оказаться выгодной заменой, если условный переход в данном контексте требует более пяти тактов.

Представляется, что должен быть более простой метод программного вычисления целочисленного квадратного корня, чем требующий нескольких сотен тактов процессора. Ниже приведен ряд выражений для вычисления целочисленного квадратного корня для очень малых аргументов. Эти выражения могут оказаться полезными для ускорения рассмотренных алгоритмов, когда ожидается извлечение квадратного корня из малых чисел.

Выражение	Диапазон корректности	Количество RISC-команд (расширенный набор)
x	От 0 до 1	0
$x > 0$	От 0 до 3	1
$(x+3) \div 4$	От 0 до 3	2
$x \gg \left(x \div 2 \right)$	От 0 до 3	2
$x \gg (x > 1)$	От 0 до 5	2
$(x+12) \div 8$	От 1 до 8	2
$(x+15) \div 8$	От 4 до 15	2
$(x > 0) + (x > 3)$	От 0 до 8	3
$(x > 0) + (x > 3) + (x > 8)$	От 0 до 15	5

11.2. Целочисленный кубический корень

В случае вычисления кубического корня метод Ньютона работает существенно хуже в силу более сложной итеративной формулы

$$x_{n+1} = \frac{1}{3} \left(2x_n + \frac{a}{x_n^2} \right),$$

а кроме того, как обычно, возникает проблема поиска хорошего начального приближения x_0 .

Однако существует аппаратный алгоритм, приведенный в листинге 11.5, который так же легко реализуется программно.

ЛИСТИНГ 11.5. Аппаратный алгоритм поиска целочисленного кубического корня

```
int icbrt(unsigned x)
{
    int s;
    unsigned y, b;

    y = 0;
    for (s = 30; s >= 0; s = s - 3)
    {
        y = 2*y;
        b = (3*y*(y + 1) + 1) << s;
        if (x >= b) {
            x = x - b;
            y = y + 1;
        }
    }
    return y;
}
```

Три команды *сложения* с 1 можно заменить командами *или* с 1, поскольку увеличиваемые значения четны. Однако даже при таких изменениях остается сомнительной его применимость для аппаратной реализации, в основном из-за наличия умножения $y*(y+1)$.

Этого умножения легко избежать путем определенной оптимизации. Введем еще одну беззнаковую переменную $y2$, которая будет содержать значение y в квадрате и обновляться при каждом изменении y . Непосредственно перед присвоением $y = 0$ вставим присвоение $y2 = 0$, а перед $y = 2*y$ — присвоение $y2 = 4*y2$. Присвоение значения переменной b заменим выражением $b = (3*y2+3*y+1) << s$ (и вынесем 3 за скобки), а непосредственно перед $y = y+1$ вставим $y2 = y2+2*y+1$. Полученная в результате программа не содержит умножений, за исключением умножения на малую константу, которое можно заменить *сдвигами* и *сложениями*. Эта программа содержит три *сложения* с 1, которые можно заменить командой *или*. Такая программа будет работать быстрее, если только умножение в вашем компьютере не выполняется за два такта (или быстрее).

Будьте осторожны: в [37] указывается, что этот код неработоспособен при непосредственной адаптации его для 64-битовых машин. Присвоение b в этом случае может вызвать переполнение. Этой проблемы можно избежать, убрав *сдвиг влево* на s позиций

в присвоении b , вставив после присвоения b выражение $bs = b << s$ и заменив две строки `if (x>=b) {x=x-b...` строками `if (x>=bs && b==(bs>>s)) {x=x-bs...`

11.3. Целочисленное возведение в степень

Вычисление x^n бинарным разложением n

Хорошо известен метод вычисления x^n , где n — неотрицательное целое число, использующий двоичное представление числа n . Этот метод применим к вычислению выражений вида $x \cdot x \cdot x \dots \cdot x$, где \cdot представляет собой любой ассоциативный оператор, такой, как сложение, умножение (включая умножение матриц) или конкатенация строк (с использованием записи $(ab)^3 = ababab$). В качестве примера рассмотрим вычисление $y = x^{13}$. Поскольку 13 в двоичном представлении имеет вид 1101 (т.е. 13 равно $8+4+1$), получаем

$$x^{13} = x^{8+4+1} = x^8 \cdot x^4 \cdot x^1.$$

Таким образом, x^{13} можно вычислить, как показано ниже.

$$\begin{aligned} t_1 &\leftarrow x^2 \\ t_2 &\leftarrow t_1^2 \\ t_3 &\leftarrow t_2^2 \\ y &\leftarrow t_3 \cdot t_2 \cdot x \end{aligned}$$

Для этого потребуется пять умножений, что значительно меньше двенадцати, необходимых при непосредственном перемножении значений x .

Если показатель степени — неотрицательная целая переменная, данный метод можно реализовать в виде функции, показанной в листинге 11.6.

Листинг 11.6. Вычисление x^n бинарным разложением n

```
int iexp(int x, unsigned n)
{
    int p, y;

    y = 1; // Инициализация результата
    p = x; // и переменной p
    while(1) {
        if (n & 1)
            y = p*y; // Если n нечетно, множим на p
        n = n >> 1; // Позиция очередного бита n
        if (n == 0)
            return y; // Если больше нет битов n
        p = p*p; // Степень для очередного бита n
    }
}
```

Количество умножений, которые нужно выполнить при использовании данного метода при показателе степени $n \geq 1$, равно

$$\lfloor \log_2 n \rfloor + \text{nbits}(n) - 1.$$

Это не всегда минимально необходимое число умножений. Например, при $n = 27$ метод бинарного разложения приводит к вычислению

$$x^{16} \cdot x^8 \cdot x^2 \cdot x^1,$$

что требует выполнения семи умножений. Однако если воспользоваться схемой

$$\left((x^3)^3 \right)^3,$$

то будет достаточно шести умножений. Наименьшее значение показателя степени, при котором метод бинарного разложения оказывается неоптимальным, — 15 (подсказка:

$$x^{15} = (x^3)^5).$$

Вероятно, это покажется неожиданным, но простой метод поиска оптимальной последовательности умножений для вычисления x^n для произвольного n неизвестен. Единственный известный сегодня метод включает обширный поиск. Данная проблема подробно рассматривается в [67, раздел 4.6.3].

Метод бинарного разложения имеет версию, в которой бинарное представление степени сканируется слева направо [96], что аналогично методу преобразования двоичного представления в десятичное слева направо. Инициализируем результат y значением 1 и сканируем двоичное представление показателя степени слева направо. Когда нам встречается нулевой бит, возводим y в квадрат; единичный бит приводит к возведению y в квадрат и умножению на x . Таким образом, $x^{13} = x^{1101}$ вычисляется так.

$$\left(\left((1^2 \cdot x)^2 \cdot x \right)^2 \right)^2 \cdot x$$

Этот метод требует того же количества умножений, что и метод сканирования справа налево из листинга 11.6.

2^n в Fortran

Компилятор IBM XL Fortran использует следующее определение функции.

$$\text{pow2}(n) = \begin{cases} 2^n, & 0 \leq n \leq 30 \\ -2^{31}, & n = 31 \\ 0, & n < 0 \text{ или } n \geq 32 \end{cases}$$

Здесь предполагается, что и n , и результат возведения в степень рассматриваются как знаковые целые числа. Стандарт ANSI/ISO Fortran требует, чтобы результат был нулевым при $n < 0$. Данное выше определение при $n \geq 31$ представляется разумным в том плане, что это корректный результат вычислений по модулю 2^{32} , и согласуется с тем, что дают многократные умножения.

Стандартный способ вычисления 2^n состоит в размещении числа 1 в регистре и сдвиге влево на n позиций. Этот способ не удовлетворяет определению Fortran, поскольку обычно величина сдвига рассматривается по модулю 64 или модулю 32 (на 32-разрядном компьютере), что приводит к неверному результату для больших или отрицательных значений смещения.

Если ваш компьютер оснащен командой подсчета количества ведущих нулевых битов, то вычисление $\text{pow2}(n)$ может быть выполнено следующим образом:

```
x ← nlz( $n \gg 5$ ); // x ← 32 при n из диапазона от 0 до 31; в противном случае x < 32
x ← x  $\gg 5$ ; // x ← 1 при n из диапазона от 0 до 31; в противном случае x ← 0
pow2 ← x  $\ll n$ ;
```

В данном случае операция сдвига вправо — беззнаковая, даже если n является знаковой величиной.

Если компьютер не имеет команды `nlz`, вместо нее можно использовать предикат $x = 0$ (см. раздел “Предикаты сравнения” на с. 43), заменив выражение $x \gg 5$ выражением $x \gg 31$. Возможно, лучшим методом реализации предиката $0 \leq x \leq 31$ является использование его эквивалентности предикату $x < 32$ с упрощением выражения для $x < y$ из упомянутого раздела; предикат превращается в $\neg x \& (x - 32)$. Такой метод даст нам решение из пяти команд (четырёх, если машина оснащена командой *и-не*).

```
x ←  $\neg n \& (n - 32)$ ; // x < 0 тогда и только тогда, когда  $0 \leq n \leq 31$ 
x ← x  $\gg 31$ ; // x = 1 при  $0 \leq n \leq 31$ , иначе 00
pow2 ← x  $\ll n$ ;
```

11.4. Целочисленный логарифм

Под целочисленным логарифмом подразумевается функция $\lfloor \log_b x \rfloor$, где x — положительное целое число, а b — целое число, не меньшее 2. Обычно $b = 2$ или 10; дадим таким функциям имена `ilog2` и `ilog10` соответственно. Имя `ilog` используется, когда основание логарифма не указано.

Вполне уместно расширить определение функции для $x = 0$, полагая $\text{ilog}(0) = -1$ [16]. Для такого определения имеется ряд причин.

- Функция $\text{ilog2}(x)$ оказывается при этом тесно связанной с функцией количества ведущих нулевых битов $\text{nlz}(x)$, так что если одна из этих функций реализована аппаратно, вычисление другой — задача совсем несложная.

$$\text{ilog2}(x) = 31 - \text{nlz}(x).$$

- При этом легко вычисляется значение $\lceil \log(x) \rceil$, если воспользоваться приведенной ниже формулой. Подставив $x = 1$, получаем, что $\text{ilog}(0) = -1$.

$$\lceil \log(x) \rceil = \text{ilog}(x - 1) + 1$$

- Такое определение делает справедливым для $x = 1$ (но не для $x = 0$) следующее тождество.

$$\text{ilog2}(x + 2) = \text{ilog2}(x) - 1$$

- При таком определении возможные значения функции $\text{ilog}(x)$ представляют собой небольшое компактное множество (от -1 до 31 для функции $\text{ilog2}(x)$ на 32-разрядном компьютере при беззнаковом x), что позволяет использовать их для индексации таблиц.
- Это определение естественным путем вытекает из ряда алгоритмов вычисления $\text{ilog2}(x)$ и $\text{ilog10}(x)$.

К сожалению, данное определение некорректно для определения “количества цифр числа x ”, которое равно $\text{ilog}(x) + 1$ для всех x , кроме 0 . Однако в силу большого количества преимуществ данное определение для $x = 0$ можно считать исключением из правила.

Для отрицательных x функция $\text{ilog}(x)$ не определена. Для расширения области определения данная функция рассматривается как отображающая беззнаковые значения на знаковые, а в этом случае отрицательный аргумент невозможен.

Целочисленный логарифм по основанию 2

Вычисление $\text{ilog2}(x)$, по сути, такое же, как и вычисление количества ведущих нулевых битов, которое рассматривалось в разделе “Подсчет ведущих нулевых битов” на с. 122. Все алгоритмы из этого раздела могут быть легко преобразованы для вычисления $\text{ilog2}(x)$ путем вычисления $\text{nlz}(x)$ с последующим вычитанием этого значения из 31 (в случае алгоритма из листинга 5.13 на с. 125 достаточно заменить строку `return pop(~x)` строкой `return pop(x) - 1`).

Целочисленный логарифм по основанию 10

Эта функция применяется при преобразовании числа для включения его в строку с удаленными начальными нулями. Этот процесс состоит в последовательном делении на 10 , что дает старшую десятичную цифру. Однако лучше заранее знать, из скольких цифр состоит число, с тем, чтобы избежать размещения полученной строки во временной области для подсчета ее длины с последующим перемещением в нужное место.

Для вычисления $\text{ilog10}(x)$ вполне применим метод поиска в таблице. Вполне возможно использование бинарного поиска, но таблица так мала, что, пожалуй, лучше воспользоваться последовательным поиском. Такой метод использован в приведенном в листинге 11.7 алгоритме.

ЛИСТИНГ 11.7. Целочисленный логарифм по основанию 10, метод поиска в таблице

```
int ilog10(unsigned x)
{
    int i;
    static unsigned table[11] =
        {0, 9, 99, 999, 9999,
         99999, 999999, 9999999, 99999999, 999999999,
         0xFFFFFFFF};
    };
```

```

for (i = -1; ; i++)
{
    if (x <= table[i+1]) return i;
}

```

На компьютере с базовым набором RISC-команд эта программа требует выполнения $9 + 4 \lfloor \log_{10} x \rfloor$ команд, так что максимальное количество выполняемых команд — 45, причем наиболее типичным количеством является 13 (для $10 \leq x \leq 99$).

Программа из листинга 11.7 может быть легко преобразована в версию без использования таблицы. Выполнимая часть такой программы представлена в листинге 11.8. Этот метод хорошо работает на компьютере с быстрым умножением на 10.

Листинг 11.8. Целочисленный логарифм по основанию 10, метод умножений на 10

```

p = 1;
for (i = -1; i <= 8; i++)
{
    if (x < p) return i;
    p = 10*p;
}
return i;

```

Такая программа требует выполнения примерно $10 + 6 \lfloor \log_{10} x \rfloor$ команд на компьютере с базовым RISC-набором (считая *умножение* одной командой). Для x из диапазона от 10 до 99 требуется выполнение 16 команд.

Можно применить бинарный поиск, что позволит избавиться от циклов и не использовать таблицу. Такой алгоритм может сравнивать x с 10^4 , затем с 10^2 или с 10^6 и так до тех пор, пока не будет найден показатель степени n , такой, что $10^n \leq x < 10^{n+1}$. В этом случае нам потребуется выполнение от 10 до 18 команд, четыре или пять из которых будут командами ветвления (с учетом последнего безусловного перехода).

Программа, показанная в листинге 11.9, представляет собой модификацию бинарного поиска, которая имеет максимум четыре ветвления на каждом из путей и написана в расчете на работу, в первую очередь, с небольшими значениями x . Она требует выполнения шести базовых RISC-команд для $10 \leq x \leq 99$ и от 11 до 16 команд при $x \geq 100$.

Листинг 11.9. Целочисленный логарифм по основанию 10, модифицированный бинарный поиск

```

int ilog10(unsigned x)
{
    if (x > 99)
        if (x < 1000000)
            if (x < 10000)
                return 3 + ((int)(x - 1000) >> 31);
            else
                return 5 + ((int)(x - 100000) >> 31);
        else
            if (x < 100000000)

```

```

        return 7 + ((int)(x - 10000000) >> 31);
    else
        return 9 + ((int)((x-1000000000)&-x) >> 31);
else
    if (x > 9) return 1;
    else      return ((int)(x - 1) >> 31);
}

```

Команда *сдвига* в данной программе — знаковая (именно поэтому в программе использовано приведение типа (int)). Если в вашем компьютере эта команда отсутствует, можно воспользоваться одним из описанных далее вариантов, в которых используется беззнаковый сдвиг. Далее приведены три варианта замены первого оператора return. К сожалению, первые два из них требуют наличия команды *вычитания из непосредственно заданного значения*, которая на большинстве машин отсутствует. Последний вариант использует сложение с большой константой (две команды), но это не так важно для остальных операторов return, которые в любом случае работают с большими константами. Величина этой константы равна $2^{31} - 1000$.

```

return 3 - ((x - 1000) >> 31);
return 2 + ((999 - x) >> 31);
return 2 + ((x + 2147482648) >> 31);

```

Четвертый оператор return можно заменить следующим.

```
return 8 + ((x + 1147483648) | x) >> 31;
```

Здесь большая константа равна $2^{31} - 10^9$. Такая замена позволяет избежать использования команды *и-не* и знакового сдвига.

Последняя конструкция if-else может быть заменена одной из приведенных ниже.

```

return ((int)(x - 1) >> 31) | ((unsigned)(9 - x) >> 31);
return (x > 9) + (x > 0) - 1;

```

Обе они позволяют сэкономить команду ветвления.

Если компьютер оснащен командой вычисления $\text{nlz}(x)$ или $\text{ilog2}(x)$, то для вычисления $\text{ilog10}(x)$ существуют более эффективные и интересные способы. Так, например, программа из листинга 11.10 позволяет сделать это с помощью двух поисков в таблице [16].

ЛИСТИНГ 11.10. Целочисленный логарифм по основанию 10 из логарифма по основанию 2, метод двойного поиска в таблице

```

int ilog10(unsigned x)
{
    int y;
    static unsigned char table1[33] = {
        9, 9, 9, 8, 8, 8, 7, 7, 7, 6, 6,
        6, 6, 5, 5, 5, 4, 4, 4, 3, 3, 3,
        3, 2, 2, 2, 1, 1, 1, 0, 0, 0, 0
    };
    static unsigned table2[10] = {
        1, 10, 100, 1000, 10000,
        100000, 1000000, 10000000,
        100000000, 1000000000
    };
};

```



```

y = table1[nlz(x)];
if (x < table2[y]) y = y - 1;
return y;
}

```

Таблица `table1` дает нам приближение $\text{ilog}_{10}(x)$. Это приближение в основном правильное, однако для некоторых значений x (равного 0 и лежащего в диапазонах от 8 до 9, от 64 до 99, от 512 до 9999, от 8192 до 9999 и т.д.) оно завышено на 1. Если требуется коррекция получаемого значения, то используется таблица `table2`.

В этой схеме используются 73 байта для хранения таблиц, а код состоит из шести команд на машине IBM System/370 [16] (чтобы достичь этого, значения в таблице `table1` должны быть в четыре раза больше показанных). На RISC-компьютере с наличием команды *вычисления количества ведущих нулевых битов* требуется выполнение около десяти команд. Другие рассматриваемые далее методы, по сути, являются вариантами данного.

Первый вариант состоит в том, чтобы избежать применения оператора `if`. При наличии команды *установки бита, если значение меньше беззнакового значения*, избежать ветвления можно и при использовании программы из листинга 11.10, но описываемый далее метод позволяет сделать это и на других компьютерах, не оснащенных столь экзотической командой.

Метод состоит в замене инструкции `if` вычитанием, за которым следует *сдвиг вправо* на 31 позицию, что эквивалентно вычитанию знакового бита. Сложности возникают только при больших значениях x ($x \geq 2^{31} + 10^9$); для решения этих проблем в `table2` внесен дополнительный элемент (см. листинг 11.11).

Листинг 11.11. Целочисленный логарифм по основанию 10 из логарифма по основанию 2, метод двойного поиска в таблице без условных переходов

```

int ilog10(unsigned x)
{
    int y;
    static unsigned char table1[33] = {
        10, 9, 9, 8, 8, 8, 7, 7, 7, 6, 6,
        6, 6, 5, 5, 5, 4, 4, 4, 3, 3, 3,
        3, 2, 2, 2, 1, 1, 1, 0, 0, 0, 0
    };
    static unsigned table2[11] = {
        1, 10, 100, 1000, 10000,
        100000, 1000000, 10000000,
        100000000, 1000000000, 0
    };
    y = table1[nlz(x)];
    y = y - ((x - table2[y]) >> 31);
    return y;
}

```

Эта программа требует выполнения примерно 11 команд на RISC-компьютере, оснащенном командой *вычисления количества ведущих нулевых битов*, но в остальном обладающей "базовым" набором команд. Эту программу легко преобразовать так, чтобы для $x = 0$ она возвращала не значение -1 , а значение 0 (для этого достаточно изменить последний элемент таблицы table1 на 1, т.е. заменить "0, 0, 0, 0" на "0, 0, 0, 1").

Еще один вариант состоит в замене первой таблицы вычитанием, умножением и сдвигом. Это оказывается возможным, так как $\log_{10} x$ и $\log_2 x$ связаны соотношением $\log_{10} x = \log_{10} 2 \cdot \log_2 x$, где $\log_{10} 2 \approx 0.30103\dots$. Таким образом, $\text{ilog}10(x)$ можно вычислить путем вычисления $\lfloor c \text{ilog}2(x) \rfloor$ с некоторой константой c и последующей коррекции полученного результата с использованием таблицы (подобно тому, как в листинге 11.11 для этого использовалась таблица table2).

Рассмотрим вопрос выбора константы более подробно. Итак, пусть $\log_{10} 2 = c + \epsilon$, где $c > 0$ — рациональное приближение $\log_{10} 2$, представляющее собой подходящий множитель, и $\epsilon > 0$. Тогда для $x \geq 1$

$$\begin{aligned} \text{ilog}10(x) &= \lfloor \log_{10} x \rfloor = \lfloor (c + \epsilon) \log_2 x \rfloor \\ \lfloor c \log_2 x \rfloor &\leq \text{ilog}10(x) = \lfloor c \log_2 x + \epsilon \log_2 x \rfloor \\ \lfloor c \text{ilog}2(x) \rfloor &\leq \text{ilog}10(x) \leq \lfloor c(\text{ilog}2(x) + 1) + \epsilon \log_2 x \rfloor \\ &\leq \lfloor c \text{ilog}2(x) + c + \epsilon \log_2 x \rfloor \\ &\leq \lfloor c \text{ilog}2(x) \rfloor + \lfloor c + \epsilon \log_2 x \rfloor + 1. \end{aligned}$$

Таким образом, если выбрать c так, что $c + \epsilon \log_2 x < 1$, то $\lfloor c \text{ilog}2(x) \rfloor$ аппроксимирует функцию $\text{ilog}10(x)$ с ошибкой, равной 0 или $+1$. Кроме того, если полагать, что $\text{ilog}2(0) = \text{ilog}10(0) = -1$, то $\lfloor c \text{ilog}2(0) \rfloor = \text{ilog}10(0)$, так как $0 < c \leq 1$, поэтому нет необходимости в каком-то отдельном рассмотрении этого случая (конечно, могут быть и другие определения, которые также не будут требовать каких-либо дополнительных действий, например $\text{ilog}2(0) = \text{ilog}10(0) = 0$).

Так как $\epsilon = \log_{10} 2 - c$, необходимо выбрать c таким, что

$$\begin{aligned} c + (\log_{10} 2 - c) \log_2 x &< 1 \text{ или} \\ c(\log_2 x - 1) &> (\log_{10} 2) \log_2 x - 1. \end{aligned}$$

Это условие выполняется при $x = 1$ (так как $c < 1$) и 2 . Для больших значений x требуется

$$c > \frac{(\log_{10} 2) \log_2 x - 1}{\log_2 x - 1}.$$

Наиболее строгое условие накладывается на c при больших x . На 32-разрядном компьютере $x < 2^{32}$, так что c должно удовлетворять условию

$$c > \frac{0.30103 \cdot 32 - 1}{32 - 1} \approx 0.27848.$$

Так как $\epsilon > 0$, $c < 0.30103$, и из соображений удобства можно выбрать значение $c = 9/32 = 0.28125$. Эксперименты показывают, что более грубые значения $5/16$ и $1/4$ действительно не подходят для нашего приближения c .

Это приводит нас к схеме, показанной в листинге 11.12, в которой используется заниженная оценка, корректируемая прибавлением 1. Она требует выполнения примерно 11 RISC-команд, среди которых команда *вычисления количества ведущих нулевых битов* (умножение считается одной командой).

Листинг 11.12. Целочисленный логарифм по основанию 10 из логарифма по основанию 2, метод поиска в таблице

```
static unsigned table2[10] = {
    0, 9, 99, 999, 9999,
    99999, 999999, 9999999,
    99999999, 999999999
};

y = (9 * (31 - nlz(x))) >> 5;
if (x > table2[y+1]) y = y + 1;
return y;
```

В этой программе можно избежать условных переходов, но при этом вновь возникнут сложности при больших значениях $x > 2^{31} + 10^9$, справиться с которыми можно различными способами. Один из способов состоит в использовании другого множителя, а именно $19/64$, и несколько измененной таблицы. Соответствующая программа показана в листинге 11.13. Она требует выполнения примерно 11 RISC-команд, среди которых команда *вычисления количества ведущих нулевых битов* (умножение считается одной командой).

Листинг 11.13. Целочисленный логарифм по основанию 10 из логарифма по основанию 2, метод поиска в таблице, без команд ветвления

```
int ilog10(unsigned x)
{
    int y;
    static unsigned table2[11] = {
        0, 9, 99, 999, 9999,
        99999, 999999, 9999999,
        99999999, 999999999,
        0xFFFFFFFF
    };
    y = (19 * (31 - nlz(x))) >> 6;
    y = y + ((table2[y+1] - x) >> 31);
    return y;
}
```

Другой способ состоит в использовании команды *или* с *x* и разностью в качестве операндов для обеспечения установки знакового бита при $x \geq 2^{31}$. Таким образом, при использовании этого способа вторая выполняемая строка в листинге 11.12 должна быть заменена следующей.

```
y = y + (((table2[y+1] - x) | x) >> 31);
```

Этот способ предпочтительнее, если умножение на 19 оказывается существенно сложнее умножения на 9 (реализуемые в виде последовательностей *сдвигов* и *сложений*).

На 64-разрядном компьютере следует выбрать *c*, удовлетворяющее такому условию.

$$c > \frac{0.30103 \cdot 64 - 1}{64 - 1} \approx 0.28993$$

Таким значением может быть $19/64 = 0.296875$, и, как показывают эксперименты, более грубого приближения не существует. В результате будет получена следующая программа.

```
unsigned table2[20] = {0, 9, 99, 999, 9999, ...,
                      99999999999999999999};
y = ((19 * (63 - nlz(x)) >> 6);
y = y + ((table2[y+1] - x) >> 63);
return y;
```

Упражнения

1. Является ли корректным корень четвертой степени целого числа *x*, вычисляемый как целочисленный квадратный корень целочисленного квадратного корня *x*? Иными словами, справедливо ли следующее соотношение?

$$\left\lfloor \sqrt{\left\lfloor \sqrt{x} \right\rfloor} \right\rfloor = \left\lfloor \sqrt[4]{x} \right\rfloor$$

2. Запишите код 64-битовой версии программы извлечения кубического корня, упомянутой в конце раздела "Целочисленный кубический корень". Воспользуйтесь типом данных `long long` языка программирования C. Не видите ли вы альтернативный метод обработки переполнения *b*, который мог бы привести к более быстрой программе?
3. Сколько умножений требуется для вычисления x^{2^3} (по модулю 2^W , где *W* — размер слова компьютера)?
4. Опишите простыми словами функции (а) $2^{\log_2(x)}$ и (б) $2^{\log_2(x-1)+1}$, где *x* — целое число, большее 0.

СИСТЕМЫ СЧИСЛЕНИЯ С НЕОБЫЧНЫМИ ОСНОВАНИЯМИ

В этой главе рассматривается несколько необычных позиционных систем счисления. В основном такие системы представляют собой не более чем интересный курьез, не имеющий практического применения. Наше рассмотрение ограничивается целыми числами, однако его легко распространить и на цифры после точки, что обычно (хотя и не всегда) обозначает нецелые числа.

12.1. Основание -2

Использование системы счисления по основанию -2 дает возможность выражать как положительные, так и отрицательные числа без явного указания их знака или, например, указания отрицательного веса старшего значащего бита [68]. При такой записи используются цифры 0 и 1, как и в случае системы счисления по основанию 2. Таким образом, значение, представленное строкой из нулей и единиц, трактуется как

$$(a_n \dots a_3 a_2 a_1 a_0) = a_n (-2)^n + \dots + a_3 (-2)^3 + a_2 (-2)^2 + a_1 (-2)^1 + a_0.$$

Из этого определения видно, что процедура поиска представления числа в системе счисления по основанию -2 должна выполнять последовательное деление на -2 , записывая получаемые остатки. Используемое деление должно быть таким, которое всегда дает в остатке 0 или 1 (используемые в записи чисел цифры), т.е. это должно быть модульное деление. В качестве примера покажем, как найти представление числа -3 по основанию -2 .

$$\begin{aligned} \frac{-3}{-2} &= 2 \text{ rem } 1 \\ \frac{2}{-2} &= -1 \text{ rem } 0 \\ \frac{-1}{-2} &= 1 \text{ rem } 1 \\ \frac{1}{-2} &= 0 \text{ rem } 1 \end{aligned}$$

Поскольку мы достигли нулевого частного, процесс на этом завершается (при его продолжении все остальные частные и остатки будут равны 0). Таким образом, считывая значения остатков снизу вверх, получаем, что число -3 , записанное в системе счисления с основанием -2 , равно 1101.

В табл. 12.1 слева показаны десятичные числа, соответствующие каждой из битовых последовательностей от 0000 до 1111 в системе счисления по основанию -2 , а справа — представление в этой системе счисления десятичных чисел в диапазоне от -15 до $+15$.

ТАБЛИЦА 12.1. ПРЕОБРАЗОВАНИЕ ДЕСЯТИЧНЫХ ЧИСЕЛ В ЧИСЛА ПО ОСНОВАНИЮ -2

n (основание -2)	n (десятичное)	n (десятичное)	n (основание -2)	$-n$ (основание -2)
0	0	0	0	0
1	1	1	1	11
10	-2	2	110	10
11	-1	3	111	1101
100	4	4	100	1100
101	5	5	101	1111
110	2	6	11010	1110
111	3	7	11011	1001
1000	-8	8	11000	1000
1001	-7	9	11001	1011
1010	-10	10	11110	1010
1011	-9	11	11111	110101
1100	-4	12	11100	110100
1101	-3	13	11101	110111
1110	-6	14	10010	110110
1111	-5	15	10011	110001

Не так очевидно, что 2^n возможных битовых строк длиной n однозначно представляют все целые числа из некоторого диапазона, но это свойство может быть доказано по индукции. Индуктивная гипотеза в данном случае состоит в том, что n -битовые слова представляют все целые числа из диапазона

$$\text{от } -(2^{n-1} - 2)/3 \text{ до } (2^n - 1)/3 \text{ для четных } n, \quad (1,а)$$

$$\text{от } -(2^n - 2)/3 \text{ до } (2^{n+1} - 1)/3 \text{ для нечетных } n. \quad (1,б)$$

Предположим сначала, что n четно. Для $n = 2$ представляемыми числами являются 10, 11, 00 и 01 по основанию -2, т.е. десятичные числа -2, -1, 0, 1. Это согласуется с формулой (1,а), причем каждое из чисел диапазона представлено только один раз.

Слово из $n+1$ бит при ведущем бите 0 может представлять все числа, задаваемые формулой (1,а). Кроме того, если ведущий бит равен 1, то это слово может представлять все эти же целые числа, смещенные на $(-2)^n = 2^n$. Таким образом, новый диапазон представляет собой

$$\text{от } 2^n - (2^{n-1} - 2)/3 \text{ до } 2^n + (2^n - 1)/3$$

или

$$\text{от } (2^n - 1)/3 + 1 \text{ до } (2^{n+1} - 1)/3.$$

Этот диапазон является смежным с диапазоном (1,а), так что слово размером $n+1$ бит однозначно представляет все целые числа из диапазона

от $-(2^{n+1}-2)/3$ до $(2^{n+2}-1)/3$.

Это согласуется с (1,6), если заменить n на $n+1$.

Доказательство того, что (1,а) следует из (1,б) при нечетном n и что все целые числа диапазона представлены однозначно, проводится аналогично.

В случае сложения и вычитания используются обычные правила, такие как $0+1=1$ и $1-1=0$. Поскольку 2 записывается как 110, а -1 — как 11, применимы дополнительные правила, которых наряду с обычными вполне достаточно для проведения арифметических вычислений.

$$\begin{aligned} 1+1 &= 110 \\ 11+1 &= 0 \\ 1+1+1 &= 111 \\ 0-1 &= 11 \\ 11-1 &= 10 \end{aligned}$$

При сложении и вычитании иногда имеется два бита переноса. Биты переноса *прибавляются* к столбцу даже при вычитании. Их удобно располагать над следующим битом слева и использовать (по возможности) упрощение $11+1=0$. Если 11 переносится в столбец, который содержит два 0, записываем в нем 1 и переносим 1 дальше. Приведем конкретные примеры сложения и вычитания.

Сложение						Вычитание							
11	1	11	11	11	19	1	11	1	1	1	21		
+ 1	1	0	1	1	1	- 1	0	1	1	1	0		
0	1	1	0	0	0	1	0	0	1	1	1		
						8							59

Возможны только следующие значения битов переносов: 0, 1 и 11. Переполнение происходит тогда, когда имеется перенос (1 или 11) из старшей позиции. Это замечание справедливо как для сложения, так и для вычитания.

Поскольку возможны три варианта переноса, сумматор по основанию счисления -2 существенно сложнее сумматора в дополнительном коде.

Существует два способа изменить знак числа. Его можно прибавить к самому себе, сдвинутому влево на одну позицию (т.е. умножить на -1), либо вычесть из 0. В этом случае нет такого простого и удобного правила, как “дополнить и прибавить 1”, присущего арифметике на основе дополнительного кода. При работе с дополнительным кодом это правило используется для создания вычитающего модуля на основе сумматора (для вычисления $A-B$ формируется сумма $A+\bar{B}+1$).

В случае работы с числами в системе счисления по основанию -2 построить такое простое устройство невозможно, но почти столь же простой метод заключается в дополнении уменьшаемого (т.е. инверсии каждого его бита), сложении полученного значения с вычитаемым и последующем дополнении суммы [77]. Вот пример, демонстрирующий вычитание 13 из 6 в соответствии с этой схемой на 8-разрядной машине.


```

00011010    6
00011101    13
11100101    6 дополненное
-----
11110110    (Дополненное 6) + 13
00001001    Дополнение суммы (-7)

```

Этот метод использует соотношение

$$A - B = I - ((I - A) + B)$$

в арифметике по основанию -2 , где I — слово, в котором все биты равны 1.

Умножение в системе счисления по основанию -2 достаточно простое. При этом используются простые правила: $1 \times 1 = 1$ и умножение на 0 дает 0; сложение в столбик выполняется с применением правил сложения чисел в системе счисления по основанию -2 .

Деление, однако, оказывается существенно сложнее. Это действительно очень интересная и сложная задача — разработка реального аппаратного алгоритма деления, т.е. алгоритма, основанного на повторяющихся вычитаниях и сдвигах. В листинге 12.1 показан алгоритм, предназначенный для работы на 8-битовом компьютере. Он реализует модульное деление (когда остаток от деления неотрицателен).

ЛИСТИНГ 12.1. Деление в системе счисления по основанию -2

```

int divbm2(int n, int d)    // q = n/d по основанию -2
{
    int r, dw, c, q, i;

    r = n;                  // Инициализация остатка
    dw = (-128)*d;         // Позиция d
    c = (-43)*d;           // Инициализация компаранда
    if (d > 0) c = c + d;
    q = 0;                  // Инициализация частного
    for (i = 7; i >= 0; i--) {
        if (d > 0 ^ (i&1) == 0 ^ r >= c) {
            q = q | (1 << i); // Установка бита частного
            r = r - dw;       // Вычитание сдвинутого d
        }
        dw = dw / (-2);     // Позиция d
        if (d > 0) c = c - 2*d; // Установка компаранда
        else c = c + d;     // для следующей итерации
        c = c / (-2);
    }
    return q;              // Возврат частного
                           // по основанию -2
                           // Остаток r, 0 <= r < |d|
}

```

Хотя эта программа написана на C и протестирована на машине с дополнительным кодом, это не играет никакой роли — данный код следует рассматривать абстрактно. Входные величины n и d и все внутренние переменные, за исключением q , являются просто числами без какого-либо конкретного представления. Выходная величина q является строкой битов, интерпретируемой как число в системе счисления по основанию -2 .

Здесь требуется небольшое пояснение. Если бы входные величины были числами в системе счисления по основанию -2 , то алгоритм было бы очень трудно выразить в выполнимом виде. Например, проверка "if ($d > 0$)" должна была бы установить, находится ли старший значащий бит числа d в четной позиции. Сложение в выражении " $c = c + d$ " также должно было быть реализовано по правилам сложения в системе счисления по основанию -2 . Такой код был бы слишком сложен для чтения. Поэтому, рассматривая данный алгоритм, вы должны воспринимать n и d как числа без конкретного представления. Приведенный код просто показывает, какие именно арифметические операции должны быть выполнены, независимо от используемого способа кодирования. Если бы числа были представлены в системе счисления по основанию -2 , как при аппаратной реализации данного алгоритма, умножение на -128 представляло бы собой сдвиг влево на семь позиций, а деление на -2 — сдвиг вправо на одну позицию.

В качестве примеров приведем вычисляемые этим кодом значения.

$$\begin{aligned} \text{divbm2}(6, 2) &= 7 && \text{(шесть, деленное на 2, равно } 111_{-2}\text{)} \\ \text{divbm2}(-4, 3) &= 2 && \text{(минус четыре, деленное на 3, равно } 10_{-2}\text{)} \\ \text{divbm2}(-4, -3) &= 6 && \text{(минус четыре, деленное на минус три, равно } 110_{-2}\text{)} \end{aligned}$$

Шаг $q = q \ll (1 < i)$; представляет собой просто установку i -го бита числа q . Следующая строка — $r = r - dw$; — представляет уменьшение остатка на сдвинутое влево значение делителя d .

Этот алгоритм трудно описать детально, но постараемся представить главную идею.

Рассмотрим определение значения первого бита частного, бита 7 числа q . При основании системы счисления, равном -2 , 8-битовое число, старший бит которого равен 1, находится в диапазоне от -170 до -43 . Таким образом, игнорируя возможность переполнения, первый (старший) бит частного будет равен 1 тогда и только тогда, когда частное алгебраически не превышает -43 .

Поскольку $n = qd + r$ и для положительного делителя справедливо соотношение $r \leq d - 1$, для положительного делителя первый бит частного будет равен 1 тогда и только тогда, когда $n \leq -43d + (d - 1)$ или $n < -43d + d$. Для отрицательного делителя первый бит частного будет равен 1 тогда и только тогда, когда $n \geq -43d$ (при модульном делении $r \geq 0$).

Таким образом, первый бит частного равен 1 тогда и только тогда, когда

$$(d > 0 \& \neg(n \geq -43d + d)) \vee (d < 0 \& n \geq -43d).$$

Игнорируя возможное нулевое значение d , можем записать это выражение как

$$d > 0 \oplus n \geq c,$$

где $c = -43d + d$ при положительном d и $c = -43d$ при d отрицательном.

Так выглядит логика определения бита частного для нечетной позиции. Для четной позиции логика обратная. Из-за этого условие в операторе if включает проверку $(i \& 1) == 0$ (напомним, что символ \wedge в языке C означает операцию *исключающее или*).

При каждой итерации с устанавливается равным наименьшему (наиболее близкому к 0) целому числу с единичным битом в позиции i после деления на d . Если текущий остаток r превышает это значение, бит i числа q устанавливается равным 1, а r корректируется путем вычитания числа, представляющего собой 1 в этой позиции, умноженного на делитель d . Реальное умножение при этом не требуется — d просто соответствующим образом позиционируется, после чего выполняется вычитание.

Алгоритм не назовешь элегантным. Его очень трудно реализовать из-за ряда сложений, вычитаний, сравнений и даже умножения (на константу), которые должны быть выполнены в начале. Надеяться на существование “однородного” алгоритма, т.е. алгоритма, в котором действия одинаковы вне зависимости от знаков аргументов, для системы счисления по основанию -2 , по-видимому, не приходится. Причина в том, что деление представляет собой, по сути, неоднородный процесс. Рассмотрим простейший алгоритм на основе *сдвигов* и *вычитаний*. Такой алгоритм может вовсе не использовать сдвиги и для положительных аргументов использовать только итеративное вычитание делителя из делимого с подсчетом количества вычитаний до тех пор, пока остаток не окажется меньше делителя. Однако, если делимое отрицательно (а делитель положителен), процесс будет состоять в итеративном прибавлении делителя к делимому, пока остаток не станет неотрицательным (частное при этом будет представлять собой количество сложений с обратным знаком). Процесс будет иным, если делитель отрицателен.

Несмотря на это деление является однородным процессом для представления чисел в прямом коде со знаком. При таком представлении значения чисел положительны, так что алгоритм может просто осуществлять вычитания до тех пор, пока остаток не станет отрицательным, а затем установить бит знака частного равным *исключающему или* знаковых битов делимого и делителя, а знаковый бит остатка — равным знаку делимого (что дает нам обычное деление с отсечением).

Алгоритм, приведенный выше, можно сделать в определенном смысле более однородным, дополняя делитель, если он отрицателен, а затем выполняя упрощенные для случая $d > 0$ шаги. В конце алгоритма выполняется окончательная коррекция. Для модульного деления коррекция изменяет знак частного и оставляет неизменным остаток. Такое изменение алгоритма выносит некоторые проверки за пределы цикла, но в целом особой красоты алгоритму не прибавляет.

Интересно сравнить распространенные представления чисел и числа в системе счисления по основанию -2 в аспекте однородности при выполнении четырех арифметических операций. Точного определения понятия “однородность” нет, но можно понимать под этим наличие или отсутствие операций, выполняемых в зависимости от знака аргументов. Мы считаем операцию однородной, если знаковый бит результата равен *исключающему или* от знаковых битов аргументов операции. В табл. 12.2 показано, какие из операций и при работе с какими представлениями чисел рассматривают свои аргументы однородно.

ТАБЛИЦА 12.2. ОДНОРОДНОСТЬ ОПЕРАЦИЙ ПРИ ИСПОЛЬЗОВАНИИ РАЗНЫХ ПРЕДСТАВЛЕНИЙ ЧИСЕЛ

Операция	Знаковые величины	Дополнение до 1	Дополнение до 2	По основанию -2
Сложение	Нет	Да	Да	Да
Вычитание	Нет	Да	Да	Да
Умножение	Да	Нет	Нет	Да
Деление	Да	Нет	Нет	Нет

Сложение и вычитание чисел в дополнительном к 1 представлении выполняются однородно с помощью приема “с переносом в конец слова”. В случае сложения все биты, включая знаковый, суммируются по обычным правилам двоичной арифметики, и перенос из крайнего слева (знакового) бита добавляется к младшему биту результата. Этот процесс всегда корректно завершается (т.е. такое добавление переноса не может привести к новому переносу в позиции знакового бита).

В случае умножения чисел в дополнительном к 2 представлении запись “Да” справедлива только в случае, если нас интересует только правая половина двойного слова произведения.

Завершим это рассмотрение системы счисления по основанию -2 некоторыми наблюдениями, касающимися выполнения преобразований между числами в этой системе счисления и обычными двоичными числами.

Для того чтобы преобразовать число из системы счисления по основанию -2 в двоичное число, формируем слово, состоящее только из битов с положительными весами, и вычитаем из него слово, состоящее из битов с отрицательными весами, используя при этом правила вычитания для двоичной арифметики. Другой метод, который может показаться немного проще, состоит в выделении битов, находящихся в позициях с отрицательными весами, сдвиге их на одну позицию влево и вычитании выделенного числа из исходного с использованием обычных правил вычитания двоичной арифметики.

Для преобразования двоичного числа в число в системе счисления по основанию -2 нужно выделить биты из нечетных позиций (позиций с весом 2^n , где n нечетно), сдвинуть их на одну позицию влево и сложить два числа с использованием правил сложения чисел в системе счисления по основанию -2. Приведем два примера преобразования чисел.

Двоичное из системы счисления по основанию -2	В систему счисления по основанию -2 из двоичного
110111 (-13)	110111 (55)
- 1 0 1 (двоичное вычитание)	+ 1 0 1 (сложение по основанию -2)
...111110011 (-13)	1001011 (55)

В компьютере, с его фиксированным размером слова, эти преобразования работают и для отрицательных чисел, если просто игнорировать перенос из старшей позиции. Для иллюстрации этого можно рассмотреть приведенный выше справа пример как преобразование двоичного числа -9 (размер слова — 6 бит) в систему счисления по основанию -2.

Приведенный выше алгоритм преобразования в число в системе счисления по основанию -2 не поддается легкой программной реализации на двоичном компьютере, поскольку требует выполнения сложения в соответствии с правилами сложения чисел в системе счисления по основанию -2 . Шреппель (Schroepfel) [44, item 128] преодолел это препятствие, разработав более ясный и практичный метод выполнения преобразований в обоих направлениях. Для преобразования в двоичное число его метод выглядит следующим образом.

$$B \leftarrow (N \oplus 0b10\dots1010) - 0b10\dots1010$$

Чтобы убедиться, что этот метод работает, рассмотрим число, которое в системе счисления по основанию -2 имеет вид $abcd$. Тогда, если (ошибочно) рассматривать его как обычное двоичное число, его значение равно $8a + 4b + 2c + d$. После выполнения операции *исключающее или* это число, рассматриваемое как двоичное, равно $8(1-a) + 4b + 2(1-c) + d$. После (двоичного) вычитания $8+2$ получаем значение $-8a + 4b - 2c + d$, которое представляет собой значение числа $abcd$ в системе счисления по основанию -2 .

Формула Шреппеля может быть легко разрешена для выполнения обратного преобразования, давая нам метод преобразования в обратном направлении, требующий выполнения трех команд. Ниже приведены формулы для преобразования в двоичное число на 32-битовой машине.

$$\begin{aligned} B &\leftarrow (N \& 0x55555555) - (N \& -0x55555555) \\ B &\leftarrow N - ((N \& 0xAAAAAAAA) \ll 1) \\ B &\leftarrow (N \oplus 0xAAAAAAAA) - 0xAAAAAAAA \end{aligned}$$

Для преобразования числа в систему счисления по основанию -2 используется следующая формула.

$$N \leftarrow (B + 0xAAAAAAAA) \oplus 0xAAAAAAAA$$

12.2. Основание $-1+i$

Используя в качестве основания системы счисления $-1+i$, где i обозначает $\sqrt{-1}$, все комплексные целые числа (т.е. комплексные числа, действительная и мнимая части которых целые) можно выразить в виде одного "числа" без явного использования знака или других вспомогательных средств. Приятной неожиданностью оказывается то, что это можно осуществить только с помощью цифр 0 и 1, причем все целые числа имеют при этом однозначное представление. Не будем здесь доказывать этот факт, как и многие другие, только опишем вкратце свойства данной системы счисления.

Не таким уж тривиальным оказывается выяснение того, как следует записать целое число 2 в этой системе счисления¹. Однако эта задача вполне разрешима путем последовательного деления 2 на основание системы счисления и записи остатков. Но что такое

¹ Попробуйте сделать это самостоятельно, отложив на время книгу.

остаток от деления в данном контексте? Мы хотим, чтобы остаток после деления на $-1+i$ был равен 0 или 1, если это возможно (так, чтобы цифрами были только 0 и 1). Для того чтобы увидеть, что это всегда возможно, предположим, что выполняется деление произвольного комплексного целого числа $a+bi$ на $-1+i$. Требуется найти q и r , такие, что q — комплексное целое число, r равно 0 или 1 и

$$a+bi = (q_r + q_i i)(-1+i) + r,$$

где q_r и q_i означают действительную и мнимую части q соответственно. Приравнивая действительные и мнимые части и решая одновременно полученные уравнения, получим

$$q_r = \frac{b-a+r}{2},$$

$$q_i = \frac{-a-b+r}{2}.$$

Ясно, что если a и b оба четны или оба нечетны, то выбор $r=0$ дает нам комплексное целое число q . Если же одно из чисел a и b четно, а второе — нет, то комплексное целое число q можно получить при выборе $r=1$.

Таким образом, целое число 2 может быть приведено к основанию $-1+i$ так, как показано ниже.

Поскольку действительная и мнимая части целого числа 2 четны, просто выполняем деление, зная, что остаток равен 0.

$$\frac{2}{-1+i} = \frac{2(-1-i)}{(-1+i)(-1-i)} = -1-i \text{ rem } 0$$

Действительная и мнимая части числа $-1-i$ нечетны, так что остаток опять равен 0.

$$\frac{-1-i}{-1+i} = \frac{(-1-i)(-1-i)}{(-1+i)(-1-i)} = i \text{ rem } 0$$

Поскольку теперь действительная и мнимая части i являются соответственно четной и нечетной, остаток будет равен 1. Проще всего учесть его, вычитая 1 из делимого.

$$\frac{i-1}{-1+i} = 1 \text{ (остаток 1)}$$

Теперь действительная и мнимая части 1 являются соответственно нечетной и четной, так что остаток будет равен 1. Вычитая его из делимого, получаем следующее.

$$\frac{1-1}{-1+i} = 0 \text{ (остаток 1)}$$

Таким образом получено нулевое частное, и процесс на этом завершается. Считывая полученные при делении остатки в обратном направлении, получаем, что представление 2 в системе счисления по основанию $-1+i$ имеет вид 1100.

В табл. 12.3 показаны все битовые строки от 0000 до 1111 и их интерпретация в системе счисления по основанию $-1+i$, а также представление в этой системе десятичных чисел от -15 до $+15$.

ТАБЛИЦА 12.3. ПРЕОБРАЗОВАНИЯ МЕЖДУ ДЕСЯТИЧНОЙ СИСТЕМОЙ СЧИСЛЕНИЯ И ПО ОСНОВАНИЮ $-1+i$

n (основание $-1+i$)	n (десятичное)	n (десятичное)	n (основание $-1+i$)	$-n$ (основание $-1+i$)
0	0	0	0	0
1	1	1	1	11101
10	$-1+i$	2	1100	11100
11	i	3	1101	10001
100	$-2i$	4	111010000	10000
101	$1-2i$	5	111010001	11001101
110	$-1-i$	6	111011100	11001100
111	$-i$	7	111011101	11000001
1000	$2+2i$	8	111000000	11000000
1001	$3+2i$	9	111000001	11011101
1010	$1+3i$	10	111001100	11011100
1011	$2+3i$	11	111001101	11010001
1100	2	12	100010000	11010000
1101	3	13	100010001	1110100001101
1110	$1+i$	14	100011100	1110100001100
1111	$2+i$	15	100011101	1110100000001

Правила сложения в системе счисления по основанию $-1+i$ (кроме тривиальных правил сложения с участием нулевых битов) выглядят следующим образом.

$$\begin{aligned}
 1 + 1 &= 1100 \\
 1 + 1 + 1 &= 1101 \\
 1 + 1 + 1 + 1 &= 111010000 \\
 1 + 1 + 1 + 1 + 1 &= 111010001 \\
 1 + 1 + 1 + 1 + 1 + 1 &= 111011100 \\
 1 + 1 + 1 + 1 + 1 + 1 + 1 &= 111011101 \\
 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 &= 111000000
 \end{aligned}$$

При сложении двух чисел наибольшее количество переносов, которое может возникнуть в одном столбце, — шесть, так что максимальная сумма, которая может быть получена в одном столбце, — 8 (111000000). Это делает построение сумматора крайне сложной задачей. Если бы нашелся кто-то разрабатывающий компьютер с комплексной арифметикой, то, вне всяких сомнений, гораздо лучше было бы хранить действительную и мнимую части отдельно², причем каждая часть имела бы более разумное представление, например в дополнительном к 2 коде.

² Этот способ был использован в 1940 году в Bell Labs в Complex Number Calculator Джорджа Стибиза (George Stibitz) [60].

12.3. Другие системы счисления

Система счисления по основанию $-1-i$ имеет, по сути, все те же свойства, что и система счисления по основанию $-1+i$, рассмотренная выше. Если некоторая битовая строка представляет число $a+bi$ в одной из этих систем, то в другой та же битовая строка представляет число $a-bi$.

Основания систем счисления $1+i$ и $1-i$ также могут представлять все комплексные целые числа с использованием цифр 1 и 0. Эти две системы счисления имеют такое же соотношение между представлением чисел в них, как и системы по основанию $-1\pm i$. В системах счисления по основанию $1\pm i$ представление некоторых целых чисел имеет бесконечную строку единиц в левой части, аналогично представлению отрицательных чисел в дополнительном к 2 коде. Это естественным образом вытекает из использования однородных правил для сложения и вычитания, как и в случае дополнительного кода. Одним из таких чисел является число 2, которое (в любой из этих систем счисления) имеет вид $\dots 11101100$. Таким образом, сложение в этих системах счисления имеет очень сложное правило: $1+1 = \dots 11101100$.

Группируя в пары биты в представлении целого числа в системе счисления по основанию -2 , можно получить представление положительных и отрицательных чисел в системе счисления по основанию 4 с использованием цифр -2 , -1 , 0 и 1 .

$$-14_{10} = 110110_{-2} = (-1)(1)(-2)_4 = -1 \cdot 4^2 + 1 \cdot 4^1 - 2 \cdot 4^0$$

Аналогично группированием пар битов в представлении числа в системе счисления по основанию $-1+i$ будет получено представление комплексного числа в системе счисления по основанию $-2i$ с использованием цифр 0 , 1 , $-1+i$ и i . Однако оно слишком сложное, чтобы представлять интерес.

Аналогична и "мнимочетверичная" система счисления [67]. Она представляет комплексные числа с использованием в качестве основания системы счисления величины $2i$ и цифр 0 , 1 , 2 и 3 (без знака). Для представления некоторых целых чисел, а именно с нечетным мнимым компонентом, при этом приходится использовать цифры справа от десятичной точки. Например, i в системе счисления по основанию $2i$ имеет вид 10.2 .

12.4. Какое основание наиболее эффективно

Предположим, что вы разрабатываете компьютер и хотите решить, какое основание системы счисления следует использовать для представления целых чисел. У вас есть различные схемные решения, которые позволяют использовать регистры с двумя состояниями (бинарные), тремя, четырьмя и т.д. Какое же решение вам следует принять?

Будем считать, что стоимость схемы с b состояниями пропорциональна b , так что схема с тремя состояниями на 50% дороже бинарной, а с четырьмя — в два раза дороже бинарной.

Предположим, что вы хотите хранить в регистрах целые числа от 0 до некоторого максимального значения M . Представление целых чисел от 0 до M в системе счисления по основанию b требует $\lceil \log_b (M+1) \rceil$ цифр (например, для представления всех целых чисел от 0 до 999999 в десятичной системе счисления требуется $\log_{10} (1000000) = 6$ цифр).

Таким образом, следует ожидать, что стоимость регистра равна произведению количества требующихся цифр на стоимость представления каждой цифры.

$$c = k \log_b (M + 1) \cdot b$$

Здесь c — стоимость регистра, а k — константа. Наша задача — найти значение b , которое минимизирует стоимость для данного значения M .

Минимум этой функции находится в той точке, где значение производной $dc/db = 0$. Таким образом, получаем следующее уравнение.

$$\frac{d}{db} (kb \log_b (M + 1)) = \frac{d}{db} \left(kb \frac{\ln(M + 1)}{\ln b} \right) = k \ln(M + 1) \frac{\ln b - 1}{(\ln b)^2} = 0.$$

Производная равна нулю, когда $\ln b = 1$, т.е. $b = e$.

Не очень удовлетворительный результат. Так как $e \approx 2.71828$, наиболее эффективными основаниями систем счисления являются 2 и 3. Какое из них более эффективно? Отношение стоимости регистра при использовании основания 2 к стоимости регистра при использовании основания системы счисления, равного 3, составляет

$$\frac{c(2)}{c(3)} = \frac{k \cdot 2 \log_2 (M + 1)}{k \cdot 2 \log_3 (M + 1)} = \frac{2 \ln(M + 1) / (\ln 2)}{3 \ln(M + 1) / (\ln 3)} = \frac{2 \ln 3}{3 \ln 2} \approx 1.056.$$

Таким образом, использование системы счисления по основанию 2 несколько дороже использования системы счисления по основанию 3, но на очень небольшую величину.

Аналогичный анализ приводит к выводу, что использование системы счисления по основанию 2 больше использования системы счисления по основанию e примерно в 1.062 раза.

Упражнения

1. Формула Шреппеля для преобразования числа из системы счисления с основанием -2 в бинарное включает дуальное использование константы 0x55555555. Сможете ли вы его найти?
2. Покажите, как добавить 1 в системе счисления с основанием -2 , используя арифметические и логические операции бинарного компьютера. Например, 0b111 \Rightarrow 0b100.
3. Покажите, как округлить число в системе счисления с основанием -2 в отрицательном направлении к кратному 16, используя арифметические и логические операции бинарного компьютера. Например, 0b10 \Rightarrow 0b110000.
4. Напишите программу на любом языке программирования для приведения целого числа в системе счисления с основанием $-1+i$ к виду $a+bi$, где a и b — действительные числа. Например, если вы передадите программе целое число 33, или 0x21, она должна вывести что-то вроде $5-4i$.
5. Как изменить знак числа в системе счисления с основанием $-1+i$? Выделить действительную часть? Выделить мнимую часть? Получить комплексно сопряженное число? (Комплексно сопряженным к $a+bi$ является число $a-bi$.)

ГЛАВА 13

КОД ГРЕЯ

13.1. Код Грея

Возможно ли циклически перебрать все 2^n комбинаций из n бит, изменяя при каждой итерации только один бит? Ответ на этот вопрос — “Да”; и именно это свойство определяет коды Грея (Gray codes). Итак, код Грея — это такой способ кодирования целых чисел, при котором закодированное число и число, следующее за ним, отличаются только одним битом. Эта концепция может быть обобщена для любого основания системы счисления, например для десятичных чисел, но в этой главе рассматриваются только бинарные коды Грея.

Хотя существует множество различных вариантов кодов Грея, рассмотрим только один: “двоичный отраженный (рефлексный) код Грея”. Именно этот код обычно имеется в виду, когда говорят о неконкретном “коде Грея”. Покажем (как обычно, без доказательств), как выполняются основные операции при таком представлении целых чисел, и расскажем о некоторых интересных свойствах кода Грея.

Отраженный двоичный код Грея строится следующим образом. Начинаем со строк 0 и 1, которые представляют соответственно целые числа 0 и 1.

0
1

Возьмем отражение этих строк относительно горизонтальной оси после приведенного списка и поместим слева от новых записей списка 1, а слева от уже имевшихся — 0.

00
01
11
10

Таким образом получен отраженный код Грея для $n = 2$. Чтобы получить код для $n = 3$, повторим описанную процедуру.

000
001
011
010
110
111
101
100

При таком способе построения легко увидеть по индукции по n , что, во-первых, каждая из 2^n комбинаций битов появляется в списке, причем только один раз; во-вторых, при переходе от одного элемента списка к рядом стоящему изменяется только один бит; в-третьих, только один бит изменяется при переходе от последнего элемента списка к первому. Коды Грея, обладающие последним свойством, называются циклическими, и отраженный код Грея обязательно является таковым.

При $n > 2$ существуют нециклические коды Грея, которые состоят из всех 2^n значений, взятых по одному разу. Примером такого кода может служить следующий: 000 001 011 010 110 100 101 111.

В схеме на рис. 13.1 показаны четырехбитовые числа, закодированные в обычной двоичной системе счисления и кодами Грея. Показаны также формулы для преобразования одного представления в другое на побитовом уровне, применимом в аппаратной реализации.

Бинарный код	Код Грея		
<i>abcd</i>	<i>efgh</i>		
0000	0000		
0001	0001	Код Грея из бинарного кода	Бинарный код из кода Грея
0010	0011	$e = a$	$a = e$
0011	0010	$f = a \oplus b$	$b = e \oplus f$
0100	0110	$g = b \oplus c$	$c = e \oplus f \oplus g$
0101	0111	$h = c \oplus d$	$d = e \oplus f \oplus g \oplus h$
0110	0101		
0111	0100		
1000	1100		
1001	1101		
1010	1111		
1011	1110		
1100	1010		
1101	1011		
1110	1001		
1111	1000		

Рис. 13.1. Четырехбитовый код Грея и формулы преобразования

Заметим, что циклический код Грея из n бит остается таковым как после любого циклического сдвига кода по вертикали, так и при любом изменении порядка столбцов. Любая из этих операций дает новый циклический код Грея, отличающийся от остальных, полученных таким же методом. Таким образом, имеется как минимум $2^n \cdot n!$ циклических бинарных кодов Грея из n бит (на самом деле их количество превышает данную оценку для $n \geq 3$).

Код Грея и двоичное представление обладают следующим дуальным отношением, очевидным из приведенных на рис. 13.1 формул.

- Бит i числа в коде Грея представляет собой четность бита i и бита слева от него в соответствующем двоичном числе (если слева от бита i нет битов, используем значение 0).

- Бит i двоичного числа представляет собой четность всех битов в позиции i и слева от нее в соответствующем числе кода Грея.

Преобразование двоичного числа в код Грея может быть выполнено с помощью всего двух команд.

$$G \leftarrow B \oplus (B \gg 1)$$

Преобразование кода Грея в двоичное число существенно сложнее.

$$B \leftarrow \bigoplus_{i=0}^{n-1} G \gg i.$$

Вы уже встречались с этой формулой в разделе 5.2, "Четность", на с. 119. Как упоминалось в этом разделе, вычисления по данной формуле можно выполнять так, как показано ниже для $n = 32$.

```

B = G ^ (G >> 1);
B = B ^ (B >> 2);
B = B ^ (B >> 4);
B = B ^ (B >> 8);
B = B ^ (B >> 16);

```

Таким образом, в общем случае требуется выполнение $2 \cdot \lceil \log_2 n \rceil$ команд.

Поскольку преобразование двоичных чисел в код Грея выполняется очень просто, генерация последовательных чисел кода Грея представляет собой тривиальную задачу.

```

for (i = 0; i < n; i++) {
    G = i ^ (i >> 1);
    output G;
}

```

13.2. Увеличение чисел кода Грея

Логика увеличения четырехбитового двоичного числа $abcd$ может быть выражена следующим образом с использованием обозначений булевой алгебры.

$$\begin{aligned}
 d' &= \bar{d} \\
 c' &= c \oplus d \\
 b' &= b \oplus cd \\
 a' &= a \oplus bcd
 \end{aligned}$$

Таким образом, один способ построения аппаратного счетчика кода Грея состоит в создании бинарного счетчика с использованием описанной логики и преобразовании его вывода a' , b' , c' и d' в код Грея с помощью операции *исключающего или* над соседними битами, как показано на рис. 13.1.

Приведем еще один способ, который может оказаться несколько эффективнее.

$$\begin{aligned}
 p &= e \oplus f \oplus g \oplus h \\
 h' &= h \oplus \bar{p} \\
 g' &= g \oplus hp \\
 f' &= f \oplus g\bar{h}p \\
 e' &= e \oplus f\bar{g}\bar{h}p
 \end{aligned}$$

Таким образом, в общем случае

$$G'_n = G_n \oplus (G_{n-1}\bar{G}_{n-2}\dots\bar{G}_0p), \quad n \geq 2.$$

Поскольку четность p чередуется между 0 и 1, схема счетчика может поддерживать p как отдельный однобитовый регистр и просто инвертировать его при каждой итерации.

Программно наилучший способ поиска последующего элемента G' для целого числа G в коде Грея, вероятно, состоит в конвертации G в двоичное число, его увеличении и обратном преобразовании в код Грея. Еще один интересный и почти такой же по эффективности путь — это определение того, какой бит должен быть изменен в числе G . Вот как выглядит последовательность этих битов, представленная в виде слов, которые с помощью операции *исключающего или* с числом G дают очередное число кода.

1 2 1 4 1 2 1 8 1 2 1 4 1 2 1 16

Внимательный читатель распознает в этой последовательности маску, указывающую положение крайнего слева бита, который изменяется при увеличении целого числа 0, 1, 2, 3, ..., соответствующего позиции в приведенной последовательности. Таким образом, при увеличении числа G из кода Грея позиция инвертируемого бита определяется крайним слева битом, который изменяется при прибавлении 1 к двоичному числу, соответствующему G .

Рассмотренные способы приводят нас к алгоритмам увеличения чисел кода Грея, показанным в листинге 13.1 (оба они сначала преобразуют число G в двоичное, что обозначено как $\text{index}(G)$).

Листинг 13.1. Увеличение числа кода Грея

```

B = index(G);           B = index(G);
B = B + 1;             M = ~B & (B + 1);
Gp = B ^ (B >> 1);    Gp = G ^ M;

```

Метод увеличения чисел кода Грея “с карандашом в руке” выглядит следующим образом.

Начиная справа, найти первый бит, для которого четность данного бита и всех битов слева положительна. Инвертировать бит в данной позиции.

Можно воспользоваться другим, эквивалентным, правилом.

Пусть p — четность слова G . Если p четно, инвертировать младший бит. Если же p нечетно, инвертировать бит слева от крайнего единичного бита справа.

Последнее правило непосредственно выражено приведенными выше булевыми выражениями.

13.3. Отрицательно-двоичный код Грея

Если вы запишете целые числа в системе счисления по основанию -2 и преобразуете их с помощью *сдвига и исключенияющего или* так же, как обычное двоичное число преобразуется в число кода Грея, то полученные числа также дадут код Грея. Трехбитовый код Грея имеет индексы в диапазоне трехбитовых чисел в системе счисления по основанию -2 , а именно от -2 до 5 ; четырехбитовый — от -10 до 5 . Полученный таким образом код Грея не является отраженным, но очень близок к нему. Такой четырехбитовый код Грея можно построить, начиная с 0 и 1 , отражая их относительно *верха* списка, затем относительно *низа* списка и так далее, чередуя отражения.

Для преобразования чисел такого кода Грея обратно в систему счисления по основанию -2 правила, разумеется, те же, что и для преобразования обычного кода Грея в двоичные числа (так как эти операции обратны одна другой, не имеет значения, как именно трактуются битовые строки, участвующие в операциях).

13.4. Краткая история и применение

Коды Грея получили свое название по имени Франка Грея (Frank Gray), физика из Bell Telephone Laboratories, который в 1930-х годах изобрел метод, в настоящее время используемый для передачи цветного телевизионного сигнала, совместимого с существующими методами передачи и получения черно-белого сигнала (т.е. при получении цветного сигнала черно-белым приемником изображение выводится оттенками серого цвета).

Мартин Гарднер (Martin Gardner) [32] рассматривал применение кодов Грея для решения головоломок с китайскими кольцами и ханойскими башнями, а также для построения гамильтоновых путей в графах, представляющих гиперкубы. Он также показал, каким образом можно преобразовывать числа из десятичного представления в десятичные коды Грея.

Коды Грея используются в датчиках положения. Представьте себе полосу материала, которая разделена на проводящие и непроводящие области, соответствующие нулям и единицам целых чисел кода Грея. Каждая такая полоса контактирует с проводящей щеткой. Если щетка располагается над разделяющей линией между двумя областями так, что получается неоднозначное считывание позиции, то не имеет никакого значения, как именно будет разрешена данная неоднозначность, поскольку неоднозначным может быть положение только одной щетки (и как бы не определялось ее положение, это будут две соседние области).

Полоса материала может быть серией концентрических круговых дорожек, как показано на рис. 13.2 (четыре точки указывают четыре считывающие щетки); в этом случае будет получен датчик углового положения. В данном случае, очевидно, требуется применение циклического кода Грея.

Можно создать циклические коды Грея для датчика поворота, которому требуется только одно кольцо из проводящих и непроводящих областей, хотя и ценой снижения точности для заданного количества щеток. В этом случае щетки располагаются вдоль кольца, а не вдоль радиуса. Такие коды Грея называются *кодами Грея для одной дорожки* (single track Gray codes — STGC).

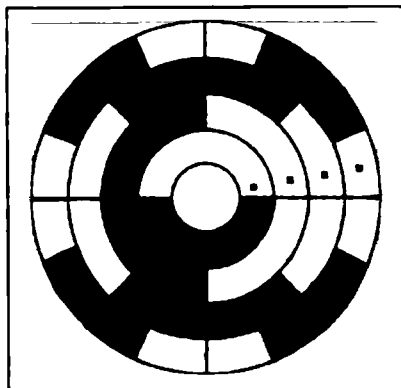


Рис. 13.2. Датчик углового положения

Идея заключается в поиске кода, для которого, если записать его как на рис. 13.1, каждый столбец будет представлять собой циклический сдвиг первого столбца (такой код будет циклическим, что важно в случае поворотного устройства). Отраженный код Грея при $n = 2$ тривиально является STGC. Коды Грея для одной дорожки для n от 2 до 4 показаны ниже.

$n = 2$	$n = 3$	$n = 4$
00	000	0000
01	001	0001
11	011	0011
10	111	0111
	110	1111
	100	1110
		1100
		1000

Такие коды позволяют создавать более компактные датчики углового положения. На рис. 13.3 показано соответствующее устройство для $n = 3$.

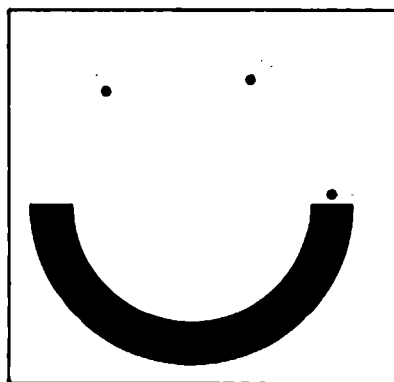


Рис. 13.3. Датчик углового положения с одной дорожкой

Все эти коды очень похожи и представляют собой весьма интересный шаблон. Следуя этому шаблону, однокоричечный код для $n = 5$ состоит из 10 слов и обеспечивает точность определения углового положения, равную 36° . Однако можно добиться большего. Код на рис. 13.4 для $n = 5$ состоит из 30 слов и обеспечивает точность определения углового положения, равную 12° . Это близко к оптимальному значению (32 кодовых слова).

10000	01000	00100	00010	00001
10100	01010	00101	10010	01001
11100	01110	00111	10011	11001
11110	01111	10111	11011	11101
11010	01101	10110	01011	10101
11000	01100	00110	00011	10001

Рис. 13.4. Датчик углового положения с одной дорожкой

Все однокоричечные коды, приведенные в данном разделе, являются наилучшими возможными в том смысле, что для n от 2 до 5 наибольшее количество возможных кодовых слов равно 4, 6, 8 и 30.

Для случая $n = 9$ имеется однокоричечный код ровно из 360 слов (это наименьшее возможное для этого значение n , так как ни один код при $n = 8$ не может состоять более чем из 256 слов) [51].

Упражнения

- Покажите, что если целое число x четно, то $G(x)$ (отраженный бинарный код Грея для x) имеет четное число единичных битов, а если x нечетно, то число единичных битов в $G(x)$ нечетно.
- Сбалансированным кодом Грея называется циклический код Грея, в котором количество изменяющихся битов при одном проходе по коду одинаково во всех столбцах.
 - Покажите, что однокоричечный код всегда сбалансирован.
 - Сможете ли вы найти сбалансированный код Грея для $n = 3$, состоящий из восьми кодовых слов?
- Разработайте циклический код Грея для кодирования целых чисел от 0 до 9.
- [71] Пусть имеется разложение числа на простые множители. Покажите, как перечислить все делители этого числа так, чтобы каждый делитель в списке получался из предыдущего путем одного умножения или деления на простое число.

ЦИКЛИЧЕСКИЙ ИЗБЫТОЧНЫЙ КОД

14.1. Введение

Циклический избыточный код (cyclic redundancy check — CRC) представляет собой метод обнаружения ошибок в цифровых данных, но не внесения исправлений при обнаружении таковых. В методе CRC к передаваемому сообщению добавляется некоторое количество битов проверки, часто именуемых контрольной суммой (checksum) или хеш-кодом. Получатель может определить, согласуются ли проверочные биты с полученными данными, и тем самым с некоторой степенью достоверности выяснить, не произошла ли ошибка в процессе передачи данных. Если ошибка имеется, то получатель шлет отправителю сигнал отсутствия подтверждения приема (negative acknowledgement — NAK), запрашивая тем самым повторную передачу сообщения.

Этот метод иногда применяется и в устройствах хранения данных, таких как диски. В этой ситуации у каждого блока диска имеются контрольные биты, и аппаратное обеспечение может автоматически выполнить повторное чтение блока в случае обнаружения ошибки (или сообщить об ошибке программному обеспечению).

В представленном далее материале использованы термины “отправитель” и “получатель” некоторого “сообщения”, но должно быть понятно, что материал применим, например, и к записи информации на диск, и к чтению с него.

В разделе 14.2 описывается теория, лежащая в основе методологии CRC. В разделе 14.3 показано, как эта теория воплощается в практические схемы, и приводится программная реализация популярного метода, известного как CRC-32.

Основы

Имеется ряд методов генерации проверочных битов, которые могут быть добавлены в сообщение. Пожалуй, простейшим является добавление единичного бита (именуемого битом четности (parity bit)), который делает общее количество единичных битов *кодového вектора* (code vector), т.е. сообщения с добавленным битом четности, четным (или нечетным). Если при передаче сообщения изменяется единственный бит, это приведет к изменению четности на противоположное значение. Отправитель генерирует бит четности, просто суммируя биты сообщения по модулю 2, т.е. выполняя над всеми ними операцию *исключающего или*. Затем к сообщению добавляется бит четности (или дополнительный к нему). Получатель может проверить сообщение, выполняя суммирование битов сообщения по модулю 2 и проверяя согласованность полученного значения с битом четности. Или, что то же самое, получатель может выполнить суммирование всех битов (сообщения и бита четности) и проверить, получился ли нулевой результат (при использовании четной четности).

Зачастую говорят, что эта простая методика обнаруживает ошибку в одном бите. На самом деле она обнаруживает ошибку в любом нечетном количестве битов (включая бит четности), но малоприятно знать, что вы в состоянии обнаружить ошибку в трех битах и не в состоянии — в двух.

Аппаратное обеспечение генерации и проверки единственного бита четности при отправке и получении последовательности битов очень простое. Оно состоит из одной схемы *исключающего или* и небольшой схемы управления. В случае параллельной передачи битов можно использовать дерево схем *исключающего или*, как показано на рис. 14.1. Эффективный способ программного вычисления бита четности приведен в разделе 5.2 на с. 114.

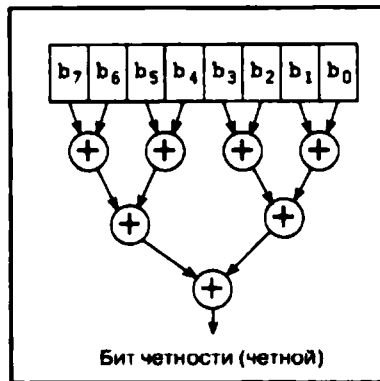


Рис. 14.1. Дерево схем *исключающего или*

Другие методы вычисления контрольных сумм состоят в вычислении *исключающего или* всех байтов в сообщении или в вычислении суммы всех байтов с циклическим переносом. В последнем случае перенос каждой 8-битовой суммы прибавляется к младшему биту накопителя. Считается, что этот метод лучше отлавливает ошибки, чем простое *исключающее или* или суммирование байтов с отбрасыванием переноса.

Еще лучшим методом с точки зрения обнаружения ошибок, и при этом легко реализуемым аппаратно, является циклический избыточный код. Это еще один способ вычисления контрольной суммы, обычно длиной 8, 16 или 32 бита, которая добавляется к сообщению. Мы кратко рассмотрим теорию, покажем, как она реализуется в аппаратном обеспечении, и приведем программную реализацию, обычно применяемую для вычисления 32-битовой контрольной суммы CRC.

Следует упомянуть о наличии гораздо более сложных способов вычисления контрольных сумм, или хеш-кодов, данных. Примерами могут служить такие хеш-функции, как MD5 и SHA-1, дающие коды длиной 128 и 160 бит соответственно. Эти методы используются в основном в криптографических приложениях и существенно более сложны в реализации — как аппаратной, так и программной, — чем описываемая здесь методика CRC. Тем не менее SHA-1 используется, например, в ряде систем управления версиями (Git и другие) просто для проверки целостности данных.

14.1. Теория

Технология CRC основана на полиномиальной арифметике, в частности на вычислении остатка при делении полиномов в $GF(2)$ (поле Галуа с двумя элементами). Это в определенной степени схоже с тем, как если бы сообщение рассматривалось как очень большое бинарное число, и вычислялся остаток от его деления на большое простое число, такое как $2^{32} - 5$. Интуитивно следует ожидать получения таким способом надежной контрольной суммы.

Полином в $GF(2)$ представляет собой полином от одной переменной x с коэффициентами 0 и 1. Сложение и вычитание выполняются по модулю 2, т.е. как если бы это были операции *исключающего или*. Например, сумма полиномов

$$\begin{aligned} &x^3 + x + 1 \quad \text{и} \\ &x^4 + x^3 + x^2 + x \end{aligned}$$

равна, как и их разность, $x^4 + x^2 + 1$. Обычно при записи таких полиномов знак “минус” не используется, поскольку коэффициент -1 эквивалентен коэффициенту 1.

Умножение таких полиномов выполняется просто и непосредственно. Произведение одного коэффициента на другой эквивалентно их комбинации с помощью оператора логического *и*, а промежуточные произведения суммируются с использованием *исключающего или*. При вычислении контрольных сумм CRC произведение полиномов не применяется.

Деление полиномов над $GF(2)$ можно выполнять почти так же, как длинное деление полиномов с целыми коэффициентами. Вот пример такого деления.

$$\begin{array}{r} x^7 + x^6 + x^5 + \\ x^7 + \\ \hline + x^6 + \\ + x^6 + + x^3 \\ \hline + x^3 + x^2 + x \\ + x^3 + + 1 \\ \hline + x^2 + 1 \end{array}$$

Вы можете самостоятельно убедиться в том, что частное $x^4 + x^3 + 1$, будучи умноженным на делитель $x^3 + x + 1$, после прибавления остатка $x^2 + 1$ даст делимое.

Метод CRC рассматривает сообщение как полином в $GF(2)$. Например, сообщение 11001001 при передаче слева направо (110...) рассматривается как представление полинома $x^7 + x^6 + x^3 + 1$. Отправитель и получатель согласовывают между собой некоторый фиксированный полином, именуемый *генератором*. Например, для 16-битового CRC комитет ССИТ (Le Comité Consultatif International Télégraphique et Téléphonique)¹ выбрал

¹ В настоящее время переименован в ITU-TSS (International Telecommunications Union — Telecommunications Standards Sector).

генератор $x^{16} + x^{12} + x^5 + 1$, широко используемый в настоящее время для вычисления 16-битовой контрольной суммы CRC. Для вычисления r -битовой контрольной суммы CRC генератор должен иметь степень r . Отправитель добавляет r нулевых битов к m -битовому сообщению и делит получившийся полином степени $m+r-1$ на генератор. Эта процедура дает нам в остатке полином степени $r-1$ (или меньшей). Полином-остаток имеет r коэффициентов, которые и являются контрольной суммой. Полином-частное игнорируется. Передаваемые данные (кодированный вектор) представляют собой исходное m -битовое сообщение, за которым следует r -битовая контрольная сумма.

Имеется два способа проверки корректности переданного сообщения получателем. Первый заключается в вычислении контрольной суммы первых m бит полученных данных и сравнении ее с последними полученными r битами. Второй (распространенный на практике) вариант состоит в том, чтобы поделить все $m+r$ бит на полином-генератор и проверить, получился ли в результате нулевой r -битовый остаток. Рассмотрим, почему остаток при этом должен быть нулевым. Пусть M — полиномиальное представление сообщения, а R — полиномиальное представление остатка, вычисленное отправителем. Тогда переданные данные соответствуют полиному $Mx^r - R$ (или, что то же самое, $Mx^r + R$). По способу вычисления R мы знаем, что $Mx^r = QG + R$, где G — полином-генератор, а Q — (игнорируемое) частное. Следовательно, передаваемые данные $Mx^r - R$ есть не что иное, как QG , что, очевидно, кратно G . Если получатель создан так же, как и отправитель, то он добавляет к полученным данным r нулевых битов, а затем вычисляет остаток R . Полученные данные после добавления r нулевых битов остаются кратными G , так что вычисленный остаток будет нулевым.

Так выглядит основная идея, но в действительности процесс несколько видоизменен, чтобы исправить имеющиеся недостатки. Например, описанный метод нечувствителен к наличию ведущих и завершающих нулевых битов в передаваемых данных. В частности, если при передаче будут обнулены все полученные данные, включая контрольную сумму, эти данные успешно пройдут проверку.

Выбор “хорошего” генератора в определенной мере представляет собой искусство и находится за пределами тематики данной книги. Упомянем только два соображения: в случае r -битовой контрольной суммы генератор G должен быть полиномом степени r (поскольку иначе первый бит контрольной суммы всегда будет нулевым, так что он просто окажется потерянным), а последний коэффициент генератора должен быть равен единице (т.е. генератор G не должен делиться на x), поскольку в противном случае последний бит контрольной суммы всегда будет нулевым (так как $Mx^r = QG + R$, если G делится на x , то и R делится на x). В работах [91] и [105] доказаны следующие факты о генераторах.

- Если G состоит из двух или большего числа членов, обнаруживаются все единичные ошибки.
- Если G не делится на x (т.е. если последний член — 1) и e — наименьшее положительное целое число, такое, что G делит $x^e + 1$, то обнаруживаются все двойные ошибки в кадре из e бит. Особенно хорошим полиномом в этом отношении является полином $x^{15} + x^{14} + 1$, для которого $e = 32\,767$.

- Если $x+1$ является множителем G , то обнаруживаются все ошибки из нечетного количества битов.
- r -битовая контрольная сумма CRC обнаруживает все пакетные ошибки длиной $\leq r$ (пакетная ошибка длиной r представляет собой строку из r бит, в которой первый и последний биты ошибочны, а промежуточные $r-2$ бит могут быть как ошибочными, так и корректными).

Генерирующий полином $x+1$ создает контрольную сумму единичной длины, которая вычисляет четную четность сообщения (*подсказка для доказательства*: чему равен остаток от деления x^k на $x+1$ для произвольного $k \geq 0$?).

Интересно заметить, что если код любого типа в состоянии обнаружить все двойные и одинарные ошибки, то он в принципе в состоянии исправлять одинарные ошибки. Чтобы убедиться в этом, предположим, что получены данные с ошибкой в одном бите. Представим, что мы по одному меняем биты на противоположные. Во всех случаях, кроме одного, мы получим двойную ошибку, которая обнаруживается рассматриваемым кодом. Но при изменении неверного бита получится корректное сообщение, что и будет распознано нашим кодом. Несмотря на это метод CRC не предназначен для исправления одиночных ошибок. Вместо этого при обнаружении ошибки отправитель все сообщение отправляет заново.

14.3. Практика

В табл. 14.1 приведены генерирующие полиномы для некоторых распространенных стандартов CRC. В столбце "Значение" приведено шестнадцатеричное представление генератора; старший бит в нем опущен — он всегда равен 1.

Таблица 14.1. Полиномы для некоторых распространенных стандартов CRC

Обычное название	r	Генератор	
		Полином	Значение
CRC-12	12	$x^{12} + x^{11} + x^3 + x^2 + x + 1$	80F
CRC-16	16	$x^{16} + x^{15} + x^2 + 1$	8005
CRC-CCITT	16	$x^{16} + x^{12} + x^5 + 1$	1021
CRC-32	32	$x^{32} + x^{26} + x^{23} + x^{22} + x^{16} + x^{12} + x^{11} + x^{10} + x^8 + x^7 + x^5 + x^4 + x^2 + x + 1$	04C11DB7

Стандарты CRC отличаются не только выбором генерирующих полиномов. Большинство из них предполагает, что сообщению предшествуют некоторые ненулевые биты, но у остальных такая инициализация отсутствует. Большинство стандартов при передаче требует, чтобы первым был передан младший бит, но некоторые начинают передачу байта со старшего бита. Большинство стандартов первым добавляют младший байт

контрольной суммы, но и здесь некоторые стандарты сначала добавляют старший байт. Некоторые стандарты выполняют также дополнение контрольной суммы.

CRC-12 используется для передачи 6-битового символического потока, остальные работают с 8-битовыми символами или 8-битовыми байтами произвольных данных. CRC-16 используется в коммуникационном стандарте IBM BISYNCH. Полином CRC-CCITT, известный также как ITU-TSS, применяется в таких коммуникационных протоколах, как XMODEM, X.25, IBM SDLC и ISO HDLC [105]. CRC-32 известен также как AUTODIN-II и ITU-TSS (ITU-TSS определяет как 16-, так и 32-битовый полиномы). Он применяется в PKZip, Ethernet, AAL5 (ATM Adaptation Layer 5), FDDI (Fiber Distributed Data Interface), стандарте IEEE-802 LAN/MAN и в некоторых приложениях DOD. В этой главе приводится программная реализация именно этого метода.

Одночлен $x+1$ является делителем для первых трех полиномов из табл. 14.1, но не для CRC-32.

Для того чтобы обнаруживать ошибку, связанную с некорректной вставкой или удалением ведущих нулевых битов, некоторые протоколы добавляют к сообщению один или несколько ненулевых битов. Сами эти биты не передаются вместе с сообщением; они просто используются для инициализации ключевого регистра (описан ниже), применяемого при вычислении CRC. Похоже, универсальным является значение, состоящее из r единичных битов. Получатель инициализирует свой регистр тем же способом.

Проблема завершающих нулей несколько сложнее. Похоже, никаких проблем не должно было бы возникнуть, если бы получатель работал путем сравнения остатков, основанных только на битах сообщения. Но, пожалуй, для получателя проще вычислять остаток для всех полученных битов (сообщения и контрольной суммы) плюс r добавленных нулевых битов. Остаток при этом должен быть нулевым. Но в случае нулевого остатка при добавлении в конец сообщения или удаления из его конца нулевых битов остаток по-прежнему будет нулевым, так что эта ошибка пройдет незамеченной.

Обычное решение этой проблемы заключается в том, что отправитель выполняет дополнение контрольной суммы перед тем, как добавить ее к сообщению. Поскольку это делает остаток, вычисленный получателем, ненулевым (как правило), этот остаток изменится при вставке или удалении завершающих нулевых битов. Как в этом случае получатель распознает, что сообщение передано без ошибок?

Используя обозначение "mod" для остатка, можем записать

$$(Mx^r + R) \bmod G = 0.$$

Обозначая "дополнение" полинома R как \bar{R} , имеем

$$\begin{aligned} (Mx^r + \bar{R}) \bmod G &= (Mx^r + (x^{r-1} + x^{r-2} + \dots + 1 - R)) \bmod G \\ &= ((Mx^r + R) + x^{r-1} + x^{r-2} + \dots + 1) \bmod G \\ &= (x^{r-1} + x^{r-2} + \dots + 1) \bmod G. \end{aligned}$$

Таким образом, контрольная сумма, вычисленная получателем для сообщения, переданного без ошибок, должна быть равна

$$(x^{r-1} + x^{r-2} + \dots + 1) \bmod G.$$

Это значение является константой для заданного G . Для CRC-32 этот полином, именуемый *вычетом* (residue), равен

$$x^{31} + x^{30} + x^{26} + x^{25} + x^{24} + x^{18} + x^{15} + x^{14} + \\ x^{12} + x^{11} + x^{10} + x^8 + x^6 + x^5 + x^4 + x^3 + x + 1,$$

или, в шестнадцатеричной записи, C704DD7B [11].

Аппаратное обеспечение

Для разработки электронной схемы для вычисления контрольной суммы CRC мы сведем процесс деления полиномов к основным составляющим.

Процесс использует регистр сдвига, который мы обозначим как CRC. Он имеет длину r (степень G) бит, а не $r + 1$, как можно было бы ожидать. При выполнении вычитания (*исключающего или*) нет необходимости в представлении старшего бита, поскольку старшие биты G и вычитаемого из него значения оба равны единице. Неформально процесс деления можно описать следующим образом.

Инициализируем регистр CRC нулевыми битами.

Получаем первый/очередной бит m .

Если старший бит CRC равен 1,

сдвигаем CRC и m влево на одну позицию и выполняем операцию *исключающего или* над получившимся в результате сдвига значением и r младшими битами G .

В противном случае

просто сдвигаем CRC и m влево на одну позицию.

Если в сообщении есть необработанные биты, повторяем цикл.

Может показаться, что сначала должно выполняться вычитание, а затем сдвиг. Это было бы верно, если бы регистр CRC хранил генератор целиком, все $r + 1$ бит. Но поскольку у нас в регистре CRC хранятся только r младших битов G , сначала выполняется сдвиг — для корректного выравнивания значений.

Ниже показано содержимое регистра CRC для генератора $G = x^3 + x + 1$ и сообщения $M = x^7 + x^6 + x^5 + x^2 + x$. В бинарном выражении $G = 1011$ и $M = 11100110$.

- | | |
|-----|--|
| 000 | Начальное содержимое CRC. Старший бит 0, так что добавляем первый бит сообщения. |
| 001 | Старший бит равен 0, так что просто добавляем второй бит сообщения, что дает: |
| 011 | Старший бит опять равен 0, так что просто добавляем третий бит сообщения, что дает: |
| 111 | Старший бит равен 1, так что выполняем сдвиг и <i>исключающее или</i> с 011, что дает: |

- 101 Старший бит равен 1, так что выполняем сдвиг и исключающее или с 011, что дает:
 001 Старший бит равен 0, так что просто добавляем шестой бит сообщения, что дает:
 011 Старший бит равен 0, так что просто добавляем седьмой бит сообщения, что дает:
 111 Старший бит равен 1, так что выполняем сдвиг и исключающее или с 011, что дает:
 101 Больше битов в сообщении нет, так что нами получен остаток.

Эти шаги можно реализовать с помощью (упрощенной) схемы, показанной на рис. 14.2, которая известна как *сдвиговый регистр с обратной связью* (feedback shift register). Квадратиками на схеме показаны три бита регистра CRC. Когда сообщение поступает на вход, если старший бит (квадратик x^2) равен 0, то одновременно бит сообщения сдвигается в квадрат x^0 , бит x^0 сдвигается в x^1 , бит из x^1 сдвигается в x^2 , а бит из x^2 отбрасывается. Если же старший бит регистра CRC равен 1, то эта единица поступает на нижний вход обеих схем *исключающего или*. При поступлении очередного бита сообщения выполняются те же сдвиги, но над тремя битами в регистре CRC выполняется операция *исключающего или* с бинарным значением 011. По окончании обработки всех битов сообщения регистр CRC хранит значение $M \bmod G$.

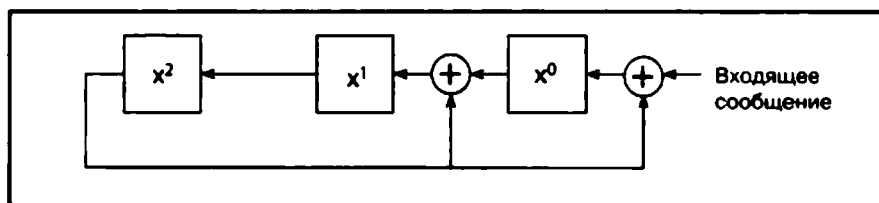
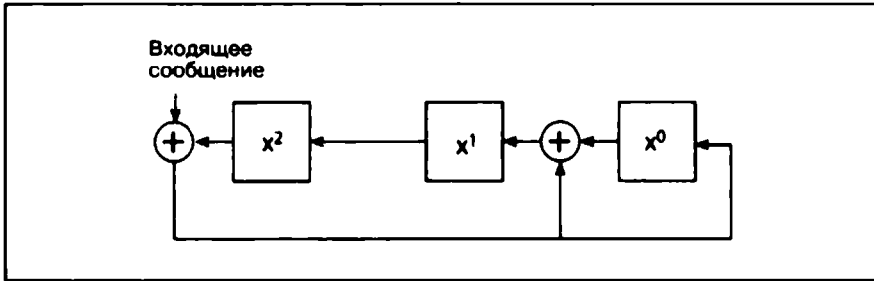


Рис. 14.2. Схема полиномиального деления для $G = x^3 + x + 1$

Если применить схему на рис. 14.2 для вычисления CRC, то по окончании обработки сообщения в нее надо послать r (в нашем случае 3) нулевых бита. Тогда регистр CRC будет содержать искомую контрольную сумму, $Mx^r \bmod G$. Однако имеется способ избежать этого шага простой перестановкой схемы.

Вместо подачи сообщения с правого конца схемы его можно подавать с левого конца, так, как показано на рис. 14.3. Влияние этой перестановки эквивалентно предварительному умножению сообщения M на x^r . Но предумножение и постумножение для полиномов являются одним и тем же. Следовательно, при поступлении каждого бита сообщения содержимое регистра CRC представляет собой остаток для обработанной части сообщения, как если бы к этой части были добавлены r нулевых битов.

Рис. 14.3. Схема CRC для $G = x^3 + x + 1$

На рис. 14.4 показана схема для полинома CRC-32.

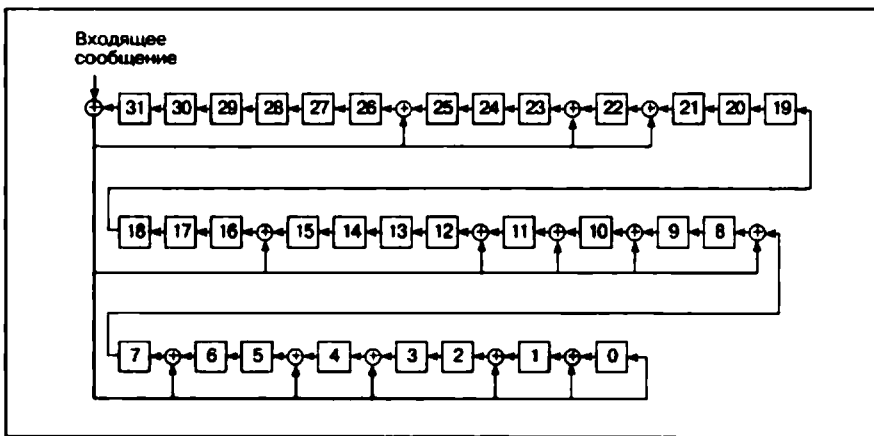


Рис. 14.4. Схема CRC для CRC-32

Программная реализация

В листинге 14.1 показана базовая программная реализация CRC-32. Протокол CRC-32 инициализирует все биты регистра CRC единицами, пересылает каждый байт начиная с его младшего бита и выполняя по окончании дополнение контрольной суммы. Мы считаем, что сообщение состоит из целого числа байтов.

Листинг 14.1. Базовый алгоритм CRC-32

```
unsigned int crc32(unsigned char *message) {
    int i, j;
    unsigned int byte, crc;

    i = 0;
    crc = 0xFFFFFFFF;
    while (message[i] != 0)
    {
        byte = message[i];           // Очередной байт
        byte = reverse(byte);       // 32-битовое обращение
        for (j = 0; j <= 7; j++)   // Восемь раз
        {
```

```

    if ((int)(crc ^ byte) < 0)
        crc = (crc << 1) ^ 0x04C11DB7;
    else crc = crc << 1;
    byte = byte << 1;      // Следующий бит байта
}
i = i + 1;
}
return reverse(~crc);
}

```

Следуя, насколько это возможно, рис. 14.4, программа использует левые сдвиги. Для этого требуется реверс каждого байта сообщения и позиционирование его с левой стороны 32-битового регистра, который в программе обозначен как `byte`. Можно воспользоваться программой реверса на уровне слова, приведенной в листинге 7.1 на с. 155 (хотя она и не очень эффективна, так как нам нужно выполнить реверс только восьми битов).

Код в листинге 14.1 приведен только в иллюстративных целях. Его можно существенно улучшить, сохранив при этом его побитовый характер. Начнем с того, что заметим, что восемь битов реверсного значения `byte` используются в инструкции `if` внутреннего цикла, после чего отбрасываются. Кроме того, старшие восемь битов `crc` во внутреннем цикле не изменяются (не считая сдвига). Таким образом, можно установить `crc = crc ^ byte` перед внутренним циклом, упрощая тем самым инструкцию `if` и опуская левый сдвиг переменной `byte` в конце цикла.

Двух реверсов можно избежать, применяя сдвиг вправо вместо сдвига влево. Это изменение требует реверса шестнадцатеричной константы, представляющей полином CRC-32, и проверки младшего бита `crc`. Наконец проверку `if` можно заменить простой логикой для уменьшения количества ветвлений. Результат показан в листинге 14.2.

ЛИСТИНГ 14.2. Усовершенствованный побитовый алгоритм CRC-32

```

unsigned int crc32(unsigned char *message) {
    int i, j;
    unsigned int byte, crc, mask;
    i = 0;
    crc = 0xFFFFFFFF;
    while (message[i] != 0)
    {
        byte = message[i];      // Очередной байт
        crc = crc ^ byte;
        for (j = 7; j >= 0; j--) // Восемь раз
        {
            mask = -(crc & 1);
            crc = (crc >> 1) ^ (0xEDB88320 & mask);
        }
        i = i + 1;
    }
    return ~crc;
}

```

Полное разворачивание внутреннего цикла неразумно. Если поступить таким образом, программа в листинге 14.2 будет выполнять около 46 команд на байт входного сообщения, включая загрузки и ветвления. (Мы считаем, что компилятор обычно генери-

рует две загрузки `message[i]` и преобразует цикл `while` так, что имеется только одно ветвление в конце цикла.)

Очередная версия программной реализации использует выборку из таблицы. Это обычный способ вычисления CRC-32. Хотя приведенные выше программы работают по одному биту за раз, метод выборки из таблицы (как он обычно реализуется) обрабатывает сообщение побайтово. При этом используется таблица из 256 константных слов.

Внутренний цикл в листинге 14.2 сдвигает регистр `crc` вправо восемь раз, в то время как операция *исключающего или* с константой применяется только тогда, когда младший бит `crc` равен 1. Эти шаги можно заменить одним правым сдвигом на восемь позиций, за которым следует одна операция *исключающего или* с маской, зависящей от того, как располагаются единичные биты в крайних справа восьми битах регистра `crc`.

Оказывается, что вычисления, необходимые для заполнения таблицы, те же, что и для вычисления CRC одного байта. Соответствующий код приведен в листинге 14.3. Чтобы программа была самодостаточна, она включает заполнение таблицы при первом вызове, но на практике эти шаги обычно входят в отдельную функцию, чтобы сделать вычисление CRC как можно более простым. Как вариант таблица может быть определена как длинная последовательность данных для инициализации массива. При компиляции компилятором GCC на базовом RISC для обработки одного входного байта выполняется 13 команд, в которые входят две загрузки и одно ветвление.

ЛИСТИНГ 14.3. Алгоритм расчета CRC с выборкой из таблицы

```
unsigned int crc32(unsigned char* message)
{
    int i, j;
    unsigned int byte, crc, mask;
    static unsigned int table[256];

    /* Заполнение таблицы при необходимости */
    if (table[1] == 0)
    {
        for (byte = 0; byte <= 255; byte++)
        {
            crc = byte;
            for (j = 7; j >= 0; j--) // Восемь раз
            {
                mask = -(crc & 1);
                crc = (crc >> 1) ^ (0xEDB88320 & mask);
            }
            table[byte] = crc;
        }
    }

    /* Когда таблица заполнена, вычисляем CRC. */

    i = 0;
    crc = 0xFFFFFFFF;
    while ((byte = message[i]) != 0)
    {
        crc = (crc >> 8) ^ table[(crc ^ byte) & 0xFF];
        i = i + 1;
    }
}
```

```
    }  
    return ~crc;  
}
```

Более быстрые версии этих программ строятся с применением стандартных методик оптимизации, но какие-то особо выдающиеся результаты автору не известны. Можно развернуть цикл и спланировать загрузки более аккуратно, чем это автоматически сделает компилятор. Можно загружать строку сообщения по полслова или по целому слову за один раз (с соответствующим учетом вопросов выравнивания), чтобы уменьшить количество загрузок сообщения и операций *исключающего* или регистра `crc` с сообщением (см. упр. 1). Метод выборки из таблицы позволяет обрабатывать сообщение по два байта за раз ценой таблицы размером 65536 слов. Это может сделать программу более быстрой или более медленной — в зависимости от размера кеша данных и дополнительных расходов при отсутствии данных в кеше.

Упражнения

1. Покажите, что если генератор G содержит два или большее количество членов, обнаруживаются все одиночные ошибки.
2. Используя листинг 14.3, покажите, как закодировать основной цикл так, чтобы загружать данные сообщения по одному слову за раз. Для простоты считайте, что сообщение выровнено по границе слова и содержит целое число слов, после чего нулевой байт указывает конец сообщения.

ГЛАВА 15

КОДЫ С КОРРЕКЦИЕЙ ОШИБОК

15.1. Введение

Этот раздел представляет собой краткое введение в теорию и практику кодов с коррекцией ошибок (error-correcting codes — ECC). Мы ограничимся бинарными блочными кодами с прямой коррекцией ошибок (forward error-correcting — FEC). Это название означает, что алфавит символов состоит только из двух символов, обозначаемых как 0 и 1, которые получатель может корректировать в процессе получения без повторной посылки сообщения или запроса дополнительной информации от отправителя, и что передача осуществляется блоками фиксированной длины, именуемыми *кодowymi словами* (code word).

В разделе 15.2. описан код, независимо изобретенный Р.У. Хэммингом (R.W. Hamming) и М. Голеем (M.J.E. Golay) до 1950 года [45]. Это код с коррекцией одиночной ошибки (single error-correcting — SEC), а его простое расширение, также изобретенное Хэммингом, одновременно с коррекцией одной ошибки позволяет обнаруживать двойные ошибки (double error-detecting — SEC-DED).

В разделе 15.4 мы вернемся к вопросу о других возможных решениях в области прямой коррекции ошибок. Ограничиваясь бинарными блочными кодами с прямой коррекцией ошибок, мы зададимся вопросом о том, сколько различных кодовых слов можно закодировать в случае заданных размера блока (или *длины кода*) и уровня обнаружения и коррекции ошибок.

Раздел 15.2 предназначен для читателей, которые интересуются тем, как коды ECC работают в памяти компьютера, в то время как раздел 15.4 — для тех, кого в большей степени интересуют математические основы данных кодов и привлекают нерешенные математические задачи.

Следует заметить, что за последние 50 лет коды с коррекцией ошибок выросли в отдельную большую тему многочисленных исследований. Опубликовано множество работ, посвященных данной тематике (можно упомянуть такие, как [50, 78, 87 и 98]). Здесь мы только вкратце коснемся данной области, познакомив читателя только с двумя важными ее разделами и используемой в ней терминологии. Хотя многие разделы теории кодов с коррекцией ошибок зависят от линейной алгебры (и это действительно отличное практическое применение для абстрактной теории), мы старались избегать ее, чтобы не усложнять понимание материала для тех, кто с ней не знаком.

В данной главе используются следующие обозначения (определения терминов будут даны в последующих разделах).

m	Количество битов “информации” или “сообщения”
k	Количество битов контроля четности (проверочных битов)
n	Длина кода, $n = m + k$
u	Битовый вектор информации, u_0, u_1, \dots, u_{m-1}

p Битовый вектор четности, p_0, p_1, \dots, p_{t-1}
 s Вектор синдрома, s_0, s_1, \dots, s_{t-1}

15.2. Код Хэмминга

Разработка Хэмминга [45] представляет собой очень простое построение кода, который позволяет корректировать единичные ошибки. В предположении, что передаваемые данные состоят из некоторого количества *информационных битов* m , он добавил к ним ряд *проверочных битов* p , таких, что если полученный блок содержал не более одного ошибочного бита, то p идентифицирует этот ошибочный бит (который может быть одним из проверочных битов). Конкретно в коде Хэмминга в случае отсутствия ошибок целочисленное значение p равно нулю, а в случае однократной ошибки оно представляет собой позицию ошибочного бита (при этом первый бит имеет номер 1). Пусть количество использованных информационных битов — m , а количество проверочных битов — k . Поскольку k проверочных битов должны проверять как себя, так и информационные биты, значение p , рассматриваемое как целое число, должно находиться в диапазоне от 0 до $m+k$, т.е. принимать $m+k+1$ различных значений. Поскольку k бит позволяют различать 2^k случаев, должно выполняться соотношение

$$2^k \geq m+k+1. \quad (1)$$

Это соотношение называется правилом Хэмминга. Оно применимо к любому бинарному блочному коду с прямой коррекцией ошибок, в котором должны проверяться все передаваемые биты. Проверочные биты располагаются среди информационных битов так, как описано ниже.

Поскольку p индексирует ошибочный бит (если таковой имеется), младший бит p должен быть равен 1, если ошибочный бит находится в нечетной позиции, и 0, если ошибочный бит находится в четной позиции или если ошибок нет. Простой способ достичь этого — положить младший бит p (бит p_0) равным четности нечетных позиций блока и разместить p_0 в нечетной позиции. Получатель сообщения проверяет четность нечетных позиций, включая бит p_0 . Если результат равен 1, ошибка находится в нечетной позиции, а если равен 0, то либо ошибки нет, либо она находится в четной позиции. Описанное свойство удовлетворяет условию, согласно которому p должно быть индексом ошибочного бита или нулем, если ошибок нет.

Аналогично положим следующий по старшинству бит p (бит p_1) равным четной четности позиций 2, 3, 6, 7, 10, 11, ... (бинарные значения 10, 11, 110, 111, 1010, 1011, ...) и поместим p_1 в одну из этих позиций. Бинарные представления номеров всех этих позиций имеют 1 во второй позиции (считая младший бит находящимся в первой позиции). Получатель сообщения проверяет четность этих позиций (включая позицию p_1). Если результат равен 1, ошибка произошла в одной из упомянутых позиций, а если результат равен 0, то либо ошибки нет, либо она находится в некоторой иной позиции.

Продолжая подобным образом, делаем третий проверочный бит p_2 равным четной четности позиций, бинарные представления номеров которых имеют 1 в третьей позиции, а именно — позиций 4, 5, 6, 7, 12, 13, 14, 15, 20, ..., и помещаем p_2 в одну из этих позиций.

Размещение проверочных битов в позициях, соответствующих степени 2 (1, 2, 4, 8, ...) обладает тем преимуществом, что они являются независимыми, т.е. отправитель может вычислить p_0 независимо от p_1, p_2, \dots и, в общем случае, вычислять каждый проверочный бит независимо от других.

В качестве примера разработаем код, корректирующий одиночную ошибку, для $m = 4$. Решение (1) дает $k = 3$ (при этом выполняется равенство). Это означает, что используются все 2^k возможных значений для k проверочных битов, что дает особенно эффективный код. Код, обладающий таким свойством, называется *идеальным* (perfect).¹

Этот код называется кодом Хэмминга (7,4), что означает, что длина кода равна 7, а количество информационных битов — 4. Позиции проверочных p_i и информационных m_i битов показаны ниже.

p_0	p_1	m_3	p_2	m_2	m_1	m_0
1	2	3	4	5	6	7

В табл. 15.1 показан код полностью. В 16 строках показаны 16 возможных значений информационных битов и значения проверочных битов, вычисленных методом Хэмминга.

ТАБЛИЦА 15.1. Код Хэмминга (7,4)

Информация	1	2	3	4	5	6	7
	p_0	p_1	m_3	p_2	m_2	m_1	m_0
0	0	0	0	0	0	0	0
1	1	1	0	1	0	0	1
2	0	1	0	1	0	1	0
3	1	0	0	0	0	1	1
4	1	0	0	1	1	0	0
5	0	1	0	0	1	0	1
6	1	1	0	0	1	1	0
7	0	0	0	1	1	1	1
8	1	1	1	0	0	0	0
9	0	0	1	1	0	0	1

¹ Идеальный код существует для $m = 2^k - k - 1$ при целых k , т.е. для $m = 1, 4, 11, 26, 57, 120, \dots$

Окончание табл. 15.1

Информация	1	2	3	4	5	6	7
	p_0	p_1	m_1	p_2	m_2	m_1	m_0
10	1	0	1	1	0	1	0
11	0	1	1	0	0	1	1
12	0	1	1	1	1	0	0
13	1	0	1	0	1	0	1
14	0	0	1	0	1	1	0
15	1	1	1	1	1	1	1

Чтобы проиллюстрировать, как получатель корректирует ошибку в одном бите, предположим, что получено кодовое слово 1001110. Это строка 4 табл. 15.1 с измененным битом 6. Получатель вычисляет *исключающее или* битов в нечетных позициях и получает 0. Затем он вычисляет *исключающее или* битов 2, 3, 6 и 7 и получает 1. И наконец он вычисляет *исключающее или* битов 4, 5, 6 и 7 и получает 1. Таким образом, указатель ошибки, который называется *синдромом* (syndrome), имеет бинарное представление 110, или значение 6. Следовательно, для коррекции блока получатель должен инвертировать бит в позиции 6.

Код SEC-DED

Для многих приложений кода с коррекцией единичной ошибки недостаточно, поскольку он принимает все полученные блоки. Код SEC-DED выглядит более безопасным, и именно этот уровень коррекции и обнаружения наиболее часто применяется в компьютерных запоминающих устройствах.

Код Хэмминга можно превратить в код SEC-DED путем добавления одного бита — бита четности (предположим, четной) для всего кодового слова SEC. Этот код называется *расширенным кодом Хэмминга* [50, 87]. То, что при этом получается код SEC-DED, не очевидно. Чтобы убедиться в этом, рассмотрим табл. 15.2. Будем считать априори, что произошла одна из ошибок 0, 1 и 2. Как указано в табл. 15.2, если ошибок нет, то общая четность (четность всего n -битового кодового слова) будет четной, а синдром $(n-1)$ -битовой SEC-части блока будет нулевым. В случае одной ошибки общая четность полученного блока будет нечетной. Если ошибка произошла в бите общей четности, синдром будет нулевым. Если же ошибка произошла в каком-то другом месте, синдром будет ненулевым и будет указывать ошибочный бит. В случае двух ошибок общая четность останется четной. Если одна из двух ошибок приходится на бит общей четности, то вторая приходится на SEC-часть блока. В этом случае синдром будет ненулевым (и будет указывать на ошибочный бит в SEC-части блока). Если же обе ошибки находятся в SEC-части блока, то синдром также будет ненулевым, хотя пояснить, почему это так, немного сложнее.

ТАБЛИЦА 15.2. ПОЛУЧЕНИЕ SEC-DED КОДА ДОБАВЛЕНИЕМ БИТА ЧЕТНОСТИ

Возможные варианты			Заключение получателя
Количество ошибок	Общая четность	Синдром	
0	Четная	0	Ошибка нет
1	Нечетная	0	Ошибка в бите общей четности
		$\neq 0$	Синдром указывает ошибочный бит
2	Четная	$\neq 0$	Двойная ошибка (коррекции не подлежит)

Причина заключается в том, что должен иметься бит, который проверяет одну из двух битовых позиций, но не другую. Четность этого проверочного бита и битов, которые он проверяет, должна быть нечетна, чтобы давать ненулевой синдром. Но почему должен существовать такой бит, который проверяет один ошибочный бит, но не оба? Чтобы разобраться в этом вопросе, сначала предположим, что один ошибочный бит находится в четной позиции, а второй — в нечетной. Тогда, поскольку один из проверочных битов (p_0) проверяет все четные позиции и не проверяет ни одну из нечетных, четность битов в нечетных позициях будет нечетной, что даст ненулевой синдром. Рассмотрим более общий случай, предположив, что ошибочные биты находятся в позициях i и j ($i \neq j$). Тогда, поскольку бинарные представления i и j должны отличаться в некоторой битовой позиции, один из них будет иметь в этой позиции 0, а другой — 1. Проверочный бит, соответствующий этой позиции в бинарных целых числах, проверяет биты в тех позициях кодового слова, которые в данном месте номера позиции имеют 1, но не тех, в номерах которых в данном месте находится нулевой бит. Биты, охваченные такой проверкой, будут иметь нечетную четность, так что синдром окажется ненулевым. В качестве примера предположим, что ошибочные биты находятся в позициях 3 и 7. Бинарные представления этих позиций — 0...0011 и 0...0111. Эти числа отличаются в третьей позиции справа, где число 7 имеет единичный бит, а число 3 — нулевой. Следовательно, биты, проверяемые третьим проверочным битом (это биты 4, 5, 6, 7, 12, 13, 14, 15, ...), будут иметь нечетную четность.

Таким образом, как указано в табл. 15.2, общая четность и синдром вместе единственным образом определяют, произошла ли одна или две ошибки или передача прошла без ошибок. В случае двух ошибок получатель не может определить, произошла ли в SEC-части сообщения только одна из ошибок (и в таком случае ее можно корректировать) или в этой части находятся обе ошибки (в таком случае попытка коррекции приведет к неверному результату).

Бит общей четности может проверять четность только в четных позициях, так как общая четность может быть легко вычислена из этой четности и четности в нечетных позициях (которая представлена младшим проверочным битом). В общем случае бит общей четности может проверять четность дополненного множества битов, проверяе-

мых любым из битов четности SEC. Данное наблюдение может позволить сэкономить несколько логических схем при аппаратной реализации.

Должно быть очевидно, что SEC-код Хэмминга имеет минимальную избыточность, т.е. для заданного количества информационных битов этот код добавляет минимальное количество проверочных битов, обеспечивающих коррекцию одной ошибки. Это свойство выполняется, поскольку по самому принципу построения кода добавляется количество проверочных битов, минимально необходимое для того, чтобы, будучи интерпретированными как целое число, они могли индексировать любой бит кода плюс одно состояние, обозначающее отсутствие ошибок. Другими словами, код удовлетворяет неравенству (1). Хэмминг показал, что SEC-DED код, построенный из SEC-кода путем добавления бита общей четности, также обладает минимальной избыточностью. Его доказательство состоит в предположении, что имеется SEC-DED код с меньшим числом проверочных битов, после чего из этого факта он получает противоречие с исходным предположением, что начальный SEC-код обладает минимальной избыточностью.

Минимально необходимое количество проверочных битов

В среднем столбце табл. 15.3 показаны минимальные решения неравенства (1) для диапазонов значений m . В правом столбце просто добавлен один бит, необходимый для получения SEC-DED кода. Из этой таблицы видно, например, что для обеспечения кода с коррекцией ошибок уровня SEC-DED для слова памяти, содержащего 64 информационных бита, требуется восемь проверочных битов, что дает в сумме слово памяти размером 72 бита.

ТАБЛИЦА 15.3. ДОПОЛНИТЕЛЬНЫЕ БИТЫ ДЛЯ КОРРЕКЦИИ/ОБНАРУЖЕНИЯ ОШИБКИ

Количество информационных битов m	k для SEC	k для SEC-DED
1	2	3
От 2 до 4	3	4
От 5 до 11	4	5
От 12 до 26	5	6
От 27 до 57	6	7
От 58 до 120	7	8
От 121 до 247	8	9
От 238 до 502	9	10

Заключительные замечания

В более математически ориентированной литературе по кодам с коррекцией ошибок термин “код Хэмминга” относится только к описанным выше идеальным кодам, т.е. к кодам с $(n, m) = (3, 1), (7, 4), (15, 11), (31, 26)$ и т.д. Аналогично расширенными кодами Хэмминга называют описанные выше идеальные SEC-DED коды. Специалисты по вычис-

лительным архитектурам и инженеры часто используют эти термины по отношению к любым описанным выше кодам и некоторым их вариантам. Термин "расширенный" зачастую просто подразумевается.

Первой ЭВМ IBM с использованием кодов Хэмминга стала машина IBM Stretch (model 7030), созданная в 1961 году [78]. Она использовала (72,64) SEC-DED код (не являющийся идеальным). Следующая машина, известная как Harvest (model 7950) и созданная в 1962 году, была оснащена 22-дорожечными ленточными накопителями, применявшими (22,16) SEC-DED код. Коды с коррекцией ошибок на современных машинах обычно кодами Хэмминга не являются; здесь применяются коды, разработанные для получения некоторых логических или схемотехнических свойств, таких как минимизация деревьев проверки четности или достижение их одинаковой длины. Такие коды для определения места ошибки вместо простого метода Хэмминга используют аппаратную выборку из таблицы.

На момент написания этой книги в 2012 году большинство ноутбуков не имели средств проверки ошибок в системах памяти. Настольные компьютеры либо не имеют их вовсе, либо используют только простую проверку четности. Компьютеры серверного класса обычно оснащены коррекцией ошибок на уровне SEC-DED.

В ранних твердотельных компьютерах, оснащенных ECC-памятью, память обычно представляет собой восемь проверочных битов и 64 информационных бита. Модуль памяти (группа микросхем) обычно строится из 8-битовых схем. При обращении к слову (72 бита, включая проверочные) происходит выборка восьми битов из каждой из девяти микросхем. Каждая схема спроектирована так, чтобы восемь битов, к которым происходит обращение при выборке одного слова, были физически удалены друг от друга. Таким образом, обращение к слову приводит к обращению к 72 физически разнесенным в пространстве битам. При таком подходе, если несколько расположенных рядом битов будут изменены, например, пролетевшей альфа-частицей или какой-то частицей из космических лучей, у нескольких слов окажется всего лишь по одной ошибке, которая может быть исправлена. Некоторые большие модули памяти используют технологию, известную как Chipkill. Она позволяет компьютеру продолжать работу даже при выходе из строя микросхемы целиком, например из-за сбоя ее питания.

Такая технология "разнесения" может использоваться в приложениях связи для исправления пакетных ошибок путем разнесения битов сообщения во времени.

Современная организация памяти с коррекцией ошибок зачастую более сложна, чем наличие 8 проверочных битов для 64 информационных. Современная серверная память может иметь 16 или 32 информационных байта (128 или 256 бит), проверяемых как единое ECC-слово. Каждая микросхема DRAM может хранить два, три или четыре бита в соседних позициях. Соответственно, ECC работает с алфавитом из 4, 8 или 16 символов (этот материал выходит за рамки нашего рассмотрения). Поскольку обычно микросхемы DRAM имеют 8- или 16-битовую конфигурацию, модуль памяти зачастую предоставляет более чем достаточное количество бит для функционирования ECC. Лишние биты могут использоваться для иных функций, таких как один или два бита четности на адрес памяти. Это позволяет памяти убедиться, что полученный адрес (вероятно) является адресом, сгенерированным процессором.

В современных серверных машинах коды с коррекцией ошибок могут использоваться на разных уровнях, а также в основной памяти. Они могут применяться и в других областях, например при работе с шинами.

15.3. Программная реализация SEC-DED для 32 информационных битов

В этом разделе рассматривается код, для которого можно реализовать эффективное программное обеспечение на машине с базовым набором команд RISC. Он выполняет коррекцию одной ошибки и обнаружение двух ошибок в блоке из 32 информационных битов. Используемый метод базируется на разработках Хэмминга.

Мы следуем Хэммингу в использовании проверочных битов таким образом, что получатель в состоянии легко определить, произошла ли пересылка без ошибок, с одной или с двумя ошибками, и в случае одной ошибки легко исправить ее. Мы также следуем идее Хэмминга о добавлении одного бита общей четности для превращения SEC-кода в SEC-DED код и выбираем проверочные биты таким образом, чтобы проверочный бит и биты, которые он проверяет, обладали четной четностью. Всего нам требуется семь проверочных битов (см. табл. 15.3).

Рассмотрим сначала свойство SEC, без DED. Для обеспечения SEC достаточно шести битов. При программной реализации основная сложность при применении метода Хэмминга заключается в том, что следует перемешать шесть проверочных битов с 32 информационными, получая в результате 38-битовое значение. Для отправителя очень неудобно распределять информационные биты среди 38 бит и вычислять проверочные биты в позициях, указанных методом Хэмминга. Получатель также столкнется с аналогичными трудностями. Проверочные биты можно было бы перенести в отдельное слово или регистр, а 32 информационных бита хранить в другом слове или регистре. Но при этом мы получим нерегулярный спектр позиций, контролируемых каждым проверочным битом. В представленной схеме эти диапазоны сохраняют большую часть регулярности, которую они проявляют в схеме Хэмминга (игнорирующей границы слов), а регулярность приводит к упрощению вычислений.

Позиции, проверяемые каждым проверочным битом, показаны в табл. 15.4. В этой таблице используется прямой порядок битов с нумерацией, начинающейся с позиции 0 младшего бита (в отличие от нумерации Хэмминга).

ТАБЛИЦА 15.4. Позиции, КОНТРОЛИРУЕМЫЕ ПРОВЕРОЧНЫМИ БИТАМИ

Проверочный бит	Контролируемые позиции
P_0	0, 1, 3, 5, 7, 9, 11, ..., 29, 31
P_1	0, 2-3, 6-7, 10-11, ..., 30-31
P_2	0, 4-7, 12-15, 20-23, 28-31
P_3	0, 8-15, 24-31
P_4	0, 16-31
P_5	1-31

Заметим, что каждая из 32 позиций битов информационного слова проверяется как минимум двумя проверочными битами. Например, позиция 6 проверяется битами p_1 и p_2 (а также p_5). Таким образом, если два информационных слова отличаются в одной позиции, кодовые слова отличаются минимум в трех позициях (поврежденный информационный бит и два или более проверочных битов), так что кодовые слова отстоят одно от другого на расстояние не менее 3 (см. раздел "Расстояние Хэмминга" на с. 370). Кроме того, если два информационных слова отличаются в двух битовых позициях, то как минимум один из битов $p_0 - p_5$ проверяет одну из позиций, но не проверяет другую, так что эти слова, опять же, располагаются на расстоянии не менее 3 одно от другого. Следовательно, предложенная схема представляет код с минимальным расстоянием 3 (SEC-код).

Предположим, что кодовое слово передается получателю. Пусть m обозначает полученные информационные биты, p обозначает полученные проверочные биты, а s (синдром) обозначает результат операции *исключающего или* над p и проверочными битами, вычисленными получателем на основании m . Тогда изучение табл. 15.4 показывает, что s в случае одной ошибки (или отсутствия таковой) в кодовом слове будет установлено так, как показано в табл. 15.5.

ТАБЛИЦА 15.5. Синдром для отсутствия ошибок или одной ошибки

Ошибка в бите	Синдром $s_5 \dots s_0$
Ошибок нет	000000
m_0	011111
m_1	100001
m_2	100010
m_3	100011
m_4	100100
...	...
m_{10}	111110
m_{11}	111111
p_0	000001
p_1	000010
p_2	000100
p_3	001000
p_4	010000
p_5	100000

В качестве примера предположим, что при передаче искажен бит u_4 . В табл. 15.4 показано, что бит u_4 контролируется проверочными битами p_2 и p_1 . Следовательно, проверочные биты, вычисленные отправителем и получателем, будут отличаться битами p_2 и p_1 . В этом сценарии проверочные биты приходят к получателю в неизменном виде, так что в синдроме установлены биты 2 и 5, т.е. он имеет значение 100100.

Если в процессе передачи искажен один из проверочных битов (а среди информационных битов ошибок нет), то полученные при передаче и вычисленные проверочные биты будут отличаться в одном поврежденном проверочном бите, как показано в шести последних строках табл. 15.5.

Синдромы, приведенные в табл. 15.5, различны для всех 39 возможных вариантов отсутствия ошибки и одиночных ошибок в любом месте кодового слова. Следовательно, синдром определяет, произошла ли ошибка, и если произошла, то в какой позиции кодового слова это случилось. Кроме того, в случае одиночной ошибки достаточно легко вычислить (не прибегая к выборке из таблицы) ошибочный бит и исправить его. Вот как это делается.

Если $s = 0$, ошибки нет.

Если $s = 011111$, неверен бит u_6 .

Если $s = 1xxxxx$, где $xxxxx$ имеет ненулевое значение, ошибка в u находится в позиции $xxxxx$.

В противном случае установлен единственный бит s , ошибка произошла в проверочном бите, и корректные проверочные биты получаются путем применения *исключающего или* к синдрому и полученным при передаче проверочным битам (или путем вычисления проверочных битов заново).

В предположении, что корректировать ошибку в проверочных битах смысла не имеет, описанную логику можно реализовать так, как показано ниже, где b — номер корректируемого бита.

```

if (s & (s - 1)) = 0 then ...           // Коррекция не нужна
else do
  if s = 0b011111 then b ← 0
  else b ← s & 0b011111
  u ← u ⊕ (1 ≪ b)                       // Дополнение бита b в u
end

```

Имеется прием, позволяющий заменить вторую конструкцию if-then-else, показанную выше, инструкцией присваивания.

Для распознавания двойной ошибки вычисляется бит общей четности (четности $u_{31,0}$ и $p_{5,0}$) и помещается для передачи в позицию 6 проверочных битов p . Двойная ошибка определяется тем, что в этом случае бит общей четности корректен, но значение синдрома $s_{5,0}$ при этом ненулевое. Причины, по которым в случае двойной ошибки синдром

имеет ненулевое значение, те же, что и в случае расширенного кода Хэмминга, приведенные на с. 361.

Программная реализация этого кода показана в листингах 15.1 и 15.2. Мы рассматриваем простой случай отправителя и получателя, и получатель при этом не должен выполнять коррекцию ошибки в проверочных битах или бите общей четности.

ЛИСТИНГ 15.1. Вычисление проверочных битов

```

unsigned int checkbits(unsigned int u)
{
    /* Вычисление шести проверочных битов четности для
    информационных битов из заданного 32-битового слова u.
    Проверочными являются биты p[5:0]. При отправке (другим
    процессом) в начало p добавляется бит общей четности.
    Проверочные биты контролируют следующие биты слова u:
    p[0] 0, 1, 3, 5, ..., 31 (0 и нечетные позиции).
    p[1] 0, 2-3, 6-7, ..., 30-31 (0 и позиции xxx1x).
    p[2] 0, 4-7, 12-15, 20-23, 28-31 (0 и позиции xx1xx).
    p[3] 0, 8-15, 24-31 (0 и позиции x1xxx).
    p[4] 0, 16-31 (0 и позиции 1xxxx).
    p[5] 1-31 */

    unsigned int p0, p1, p2, p3, p4, p5, p6, p;
    unsigned int t1, t2, t3;

    // Сначала вычисляем p[5:0], игнорируя u[0].
    p0 = u ^ (u >> 2);
    p0 = p0 ^ (p0 >> 4);
    p0 = p0 ^ (p0 >> 8);
    p0 = p0 ^ (p0 >> 16); // p0 в позиции 1.

    t1 = u ^ (u >> 1);
    p1 = t1 ^ (t1 >> 4);
    p1 = p1 ^ (p1 >> 8);
    p1 = p1 ^ (p1 >> 16); // p1 в позиции 2.

    t2 = t1 ^ (t1 >> 2);
    p2 = t2 ^ (t2 >> 8);
    p2 = p2 ^ (p2 >> 16); // p2 в позиции 4.

    t3 = t2 ^ (t2 >> 4);
    p3 = t3 ^ (t3 >> 16); // p3 в позиции 8.

    p4 = t3 ^ (t3 >> 8); // p4 в позиции 16.

    p5 = p4 ^ (p4 >> 16); // p5 в позиции 0.

    p = ((p0>>1) & 1) | ((p1>>1) & 2) | ((p2>>2) & 4) |
        ((p3>>5) & 8) | ((p4>>12) & 16) | ((p5 & 1) << 5);

    p = p ^ (-(u & 1) & 0x3F); // Учитываем u[0].
    return p;
}

```


Листинг 15.2. Действия получателя

```

int correct(unsigned int pr, unsigned int* ur)
{
    /* Эта функция исследует семь полученных проверочных
    битов и 32 информационных бита (pr и ur) и определяет
    количество имеющихся ошибок (предположительно их может
    быть 0, 1 или 2). Функция возвращает 0, 1 или 2 (что
    означает отсутствие ошибок, одну или две ошибки
    соответственно. В случае одной ошибки функция исправляет
    полученное информационное слово (ur). */

    unsigned int po, p, syn, b;

    po = parity(pr ^ *ur); // Общая четность полученных данных

    p = checkbits(*ur);    // Вычисление проверочных битов
                          // для полученной информации
    syn = p ^ (pr & 0x3F); // Синдром (исключая бит общей
                          // четности)

    if (po == 0)
    {
        if (syn == 0)
        {
            return 0;      // Ошибок нет, возврат 0
        }
        else
        {
            return 2;      // Две ошибки, возврат 2
        }
    }

    // Одна ошибка
    if (((syn - 1) & syn) == 0) // Если в syn установлено
    {                          // нуль или один бит, то
        return 1;              // ошибка в проверочных битах
    }                          // или бите общей четности
                              // и коррекция не нужна

    // Ошибка одна, и биты syn 5:0 указывают ее положение в ur
    b = syn - 31 - (syn >> 5); // Отображение syn в диапазон
                              // от 0 до 31
    // if (syn == 0x1f) b = 0; // (Эти две строки эквивалентны
    // else b = syn & 0x1f;    // одной строке выше)
    *ur = *ur ^ (1 << b);    // Коррекция бита
    return 1;
}

```

Для вычисления проверочных битов функция `checkbits` сначала игнорирует информационный бит m_i и, исключая строку i при вычислении проверочного бита k_i , $0 \leq i \leq 4$, вычисляет

$$0. \left(x \leftarrow u \oplus \left(u \gg 1 \right) \right)$$

$$1. \left(x \leftarrow x \oplus \left(x \gg 2 \right) \right)$$

$$2. \left(x \leftarrow x \oplus \left(x \gg 4 \right) \right)$$

$$3. \left(x \leftarrow x \oplus \left(x \gg 8 \right) \right)$$

$$4. \left(x \leftarrow x \oplus \left(x \gg 16 \right) \right)$$

Эти вычисления помещают p_i в различные позиции слова x , как показано в листинге 15.1. При вычислении p_i выполняются все показанные выше присваивания. Здесь срабатывает регулярность шаблона информационных битов, контролируемых каждым проверочным битом. Она позволяет использовать существенную унификацию кода, уменьшая количество таких присваиваний с $4 \times 5 + 5 = 25$ до 15, как показано в листинге 15.1.

Кстати, если компьютер оснащен командой для вычисления четности слова или имеет команду для *вычисления степени заполнения* (которая помещает четность слова в младший бит целевого регистра), регулярный шаблон не является необходимым. На такой машине проверочные биты можно вычислить следующим образом.

```
p0 = pop(u ^ 0xAAAAAAAA) & 1;
p1 = pop(u & 0xCCCCCDD) & 1;
```

И так далее.

После упаковки шести проверочных битов в одно значение p функция `checkbits` учитывает значение информационного бита m_0 , выполняя дополнение всех шести битов, если $m_0 = 1$ (см. табл. 15.4; бит p_i должен быть дополнен, поскольку бит m_0 до этого момента ошибочно включен в вычисление бита p_i).

15.4. Общее рассмотрение задачи коррекции ошибок

В этом разделе продолжается рассмотрение бинарных блочных кодов с прямой коррекцией ошибок, но несколько более обобщенное по сравнению с описанием кодов в разделе 15.2. Мы отбросим предположение о том, что блок состоит из множества информационных битов и другого множества проверочных битов, а также следствие, заключающееся в том, что количество кодовых слов должно быть равно степени 2. Мы также рассмотрим уровни коррекции и обнаружения ошибок, превышающие SEC и SEC-DED. Предположим, например, что нам нужен код, исправляющий две ошибки, для бинарного представления десятичных цифр. Если код имеет 16 кодовых слов (десять из них используются для представления десятичных цифр и шесть — не используются), то длина кодового слова должна быть не менее 11 бит. Но если код использует всего лишь

десять кодовых слов, то длина кодового слова может быть равна 10 битам (см. в табл. 15.8 на с. 377 столбец для $d = 5$; все пояснения приведены ниже).

Код представляет собой простое множество *кодовых слов*, и для наших целей кодовые слова являются бинарными строками одной и той же длины, которая, как упоминалось ранее, называется *длиной кода*. Количество кодовых слов в множестве называется *размером кода*. Мы не интерпретируем кодовые слова; они могут представлять буквы и цифры или, например, пиксели изображения.

В качестве тривиального примера можно привести код, который состоит из бинарных чисел от 0 до 7, с повторением каждого бита три раза.

{000000000, 000000111, 000111000, 000111111, 111000000, ..., 111111111}

Еще одним примером может служить код *два-из-пяти*, в котором каждое кодовое слово имеет ровно два единичных бита.

{00011, 00101, 00110, 01001, 01010, 10000, 10001, 10010, 10100, 11000}

Размер этого кода равен 10; таким образом, он подходит для представления десятичных цифр. Заметим, что если кодовое слово 00110 рассматривать как представление десятичной цифры 0, то остальные значения можно декодировать в цифры от 1 до 9, если присвоить битам веса 6, 3, 2, 1 и 0 в порядке слева направо.

Кодовая скорость (code rate) представляет собой меру эффективности кода. Для кода наподобие кода Хэмминга ее можно определить как количество информационных битов, деленное на длину кода. Для рассмотренного ранее кода Хэмминга она равна $4/7 \approx 0.57$. В общем случае кодовую скорость можно определить как логарифм по основанию 2 размера кода, деленный на длину кода. Приведенные выше примеры кодов имеют скорости $\log_2(8)/9 \approx 0.33$ и $\log_2(10)/5 \approx 0.66$ соответственно.

Расстояние Хэмминга

Центральной концепцией теории кодов с коррекцией ошибок является *расстояние Хэмминга*. Расстояние Хэмминга между двумя словами (равной длины) представляет собой количество битов, которыми эти слова отличаются. Иначе говоря, это степень заполнения результата операции *исключающего или* над двумя словами. Эта величина с полным правом может называться расстоянием, поскольку удовлетворяет определению функции расстояния из линейной алгебры.

$$d(x, y) = d(y, x)$$

$$d(x, y) \geq 0$$

$$d(x, y) = 0 \text{ тогда и только тогда, когда } x = y$$

$$d(x, y) + d(y, z) \geq d(x, z) \text{ (неравенство треугольника)}$$

Здесь $d(x, y)$ обозначает расстояние Хэмминга между словами x и y , которое для краткости мы будем называть просто расстоянием между x и y .

Предположим, что код имеет минимальное расстояние, равное 1, т.е. в множестве кодовых слов имеется два слова, x и y , отличающиеся в одной битовой позиции. Понятно, что при передаче x и изменении значения бита, отличающего его от слова y , в результате ошибки при передаче получатель не сумеет отличить полученное с ошибкой слово x от слова y , переданного без ошибки. Следовательно, в общем случае в таком коде невозможно обнаружить даже единичную ошибку.

Предположим теперь, что код имеет минимальное расстояние, равное 2. Тогда, если при передаче искажается только один бит, получается неверное кодовое слово, так что получатель способен (в принципе) обнаружить ошибку. Если же при передаче меняются два бита, одно корректное кодовое слово может превратиться в другое корректное кодовое слово. Таким образом, двойная ошибка обнаружена быть не может. Кроме того, одинарная ошибка не может быть *исправлена*. Это связано с тем, что если получателю поступает слово с одной ошибкой, то могут иметься два слова, каждое из которых при изменении одного бита даст полученное некорректное слово, — так что нельзя определить, ошибка в каком из этих слов привела к полученному ошибочному слову.

К этой категории относится код, полученный добавлением к слову одного бита четности. Ниже показан случай трех информационных битов ($m = 3$). Крайний справа бит — бит четности, выбираемый так, чтобы четность всех четырех битов была четной. Читатель может самостоятельно убедиться, что минимальное расстояние между кодовыми словами равно 2.

```
0000
0011
0101
0110
1001
1010
1100
1111
```

В действительности добавление бита четности позволяет определять любое нечетное количество ошибок, но когда мы говорим о том, что код позволяет обнаруживать k ошибок, мы имеем в виду *все* ошибки до k бит включительно.

Теперь рассмотрим случай, когда минимальное расстояние между кодовыми словами равно 3. Если при передаче поменяют значение один или два бита, в результате будет получено неверное кодовое слово. Если будет изменен только один бит, получатель может (по крайней мере, в принципе) попытаться изменить каждый бит слова по очереди, и только в одном случае будет получено корректное кодовое слово. Следовательно, такой код позволяет получателю обнаруживать и корректировать одинарную ошибку. Двойная ошибка может иметь тот же вид, что и единичная ошибка в другом кодовом слове, так что получатель не в состоянии обнаруживать двойные ошибки.

Аналогично несложно доказать, что если минимальное расстояние кода равно 4, то получатель в состоянии корректировать все одинарные ошибки и обнаруживать все двойные ошибки (т.е. мы получаем SEC-DED код). Как упоминалось выше, этот уровень часто применяется в компьютерной памяти.

В табл. 15.6 приведены итоговые результаты по способности к обнаружению и коррекции ошибок блочных кодов в зависимости от их минимального расстояния.

ТАБЛИЦА 15.6. КОЛИЧЕСТВО ИСПРАВЛЯЕМЫХ/ОБНАРУЖИВАЕМЫХ ОШИБОК

Минимальное расстояние	Коррекция	Обнаружение
1	0	0
2	0	1
3	1	1
4	1	2
5	2	2
6	2	3
7	3	3
8	3	4
d	$\lfloor (d-1)/2 \rfloor$	$\lfloor d/2 \rfloor$

Возможностью коррекции ошибок можно пожертвовать в пользу возможности обнаружения ошибок. Например, если минимальное расстояние кода равно 3, то его избыточность может использоваться для обнаружения одинарных и двойных ошибок ценой отказа от коррекции ошибок. Если минимальное расстояние равно 5, код может использоваться как для коррекции одинарных и обнаружения тройных ошибок, так и для обнаружения четверных ошибок при отсутствии возможности коррекции. Все, что вычитается из столбца “Коррекция” в табл. 15.6, можно добавить в столбец “Обнаружение”.

Основная задача теории кодирования

До этого момента мы задавались вопросом “Сколько проверочных битов требуется для заданного количества информационных битов m , чтобы получить код с минимальным расстоянием d ?” С целью обобщения “развернем” вопрос и сформулируем его как “Сколько кодовых слов может быть у кода с длиной n и минимальным расстоянием d ?” Таким образом, количество кодовых слов не обязано быть целой степенью числа 2.

Следуя [98] и другим работам, обозначим через $A(n, d)$ наибольший возможный размер (бинарного) кода с длиной n и минимальным расстоянием d . Остаток этого раздела посвящен изучению известных фактов о данной функции. Определение ее значений носит название *основной задачи теории кодирования* [50, 98]. Везде в этом разделе полагается, что $n \geq d \geq 1$.

Практически тривиальным результатом является

$$A(n, 1) = 2^n, \quad (2)$$

поскольку имеется 2^n различных слов длиной n .

Для минимального расстояния 2 из примера с одним битом четности нам известно, что $A(n, 2) \geq 2^{n-1}$. Но $A(n, 2)$ не может превышать 2^{n-1} по следующей причине. Предположим, что имеется код длиной n с минимальным расстоянием 2, количество кодовых слов которого превышает 2^{n-1} . Удалим любой один столбец из кодовых слов. (Мы предполагаем, что кодовые слова организованы в матрицу наподобие того, как это сделано в табл. 15.1 на с. 359.) Так мы получим код длиной $n-1$ с минимальным расстоянием не менее 1 (удаление столбца может уменьшить минимальное расстояние не более чем на 1), размер которого превышает 2^{n-1} . Таким образом, у такого кода $A(n-1, 1) > 2^{n-1}$, что противоречит уравнению (2). Следовательно,

$$A(n, 2) = 2^{n-1}.$$

Пока что все было несложно. Но что можно сказать об $A(n, 3)$? Это задача нерешенная, в том смысле, что неизвестна ни формула, ни простой способ вычисления этой величины. Конечно, известны многие конкретные значения $A(n, 3)$, известны некоторые границы, но в большинстве случаев точное значение неизвестно.

Когда в соотношении (1) выполняется равенство, оно представляет решение данной задачи для $d = 3$. Положив $n = m + k$, можно переписать (1) как

$$2^m \leq \frac{2^n}{n+1}. \quad (3)$$

Здесь m представляет собой количество информационных битов, так что 2^m является максимальным количеством кодовых слов. Следовательно, мы имеем

$$A(n, 3) \leq \frac{2^n}{n+1},$$

причем равенство выполняется, когда $2^n/(n+1)$ является целым числом (согласно построению Хэмминга).

Для $n = 7$ это дает $A(7, 3) = 16$, результат, который уже известен из раздела 15.2. Для $n = 3$ это дает $A(3, 3) \leq 2$, и предельного значения 2 можно достичь с кодовыми словами 000 и 111. Для $n = 4$ мы получаем $A(4, 3) \leq 3.2$, и, немного поэкспериментировав, можно увидеть, что невозможно получить три кодовых слова длиной 4 с $d = 3$. Таким образом, когда в (3) не выполняется равенство, это соотношение дает нам просто верхнюю границу (вполне возможно — недостижимую) для максимального количества кодовых слов.

Для $n \geq 2$ имеется интересное соотношение

$$A(n, d) \leq 2A(n-1, d). \quad (4)$$

Таким образом, добавление 1 к длине кода не более чем удваивает количество кодовых слов, возможных при том же минимальном расстоянии d . Чтобы увидеть это, предположим, что есть код длиной n с расстоянием d и размером $A(n, d)$. Выберем произвольный столбец кода. В этой позиции либо не менее половины битов имеют значение 0, либо не

менее половины битов имеют значение 1. Выберем из этих двух подмножеств то, в котором имеется не менее $A(n, d)/2$ кодовых слов, образуем новый код, состоящий из этого подмножества, и удалим выбранный столбец (в котором все биты равны либо 0, либо 1). Получающееся в результате множество кодовых слов имеет длину $n-1$, то же минимальное расстояние d и как минимум $A(n, d)/2$ кодовых слов. Таким образом, $A(n-1, d) \geq A(n, d)/2$, откуда следует неравенство (4).

В случае четного d имеется полезное соотношение

$$A(n, d) = A(n-1, d-1). \quad (5)$$

Чтобы убедиться в этом, предположим, что у нас есть код C длиной n с минимальным расстоянием d , где d нечетно. образуем новый код, добавляя к каждому слову C бит четности, делая общую четность каждого слова, скажем, четной. Новый код имеет длину $n+1$, то же количество кодовых слов, что и код C , и минимальное расстояние $d+1$. Если два слова C находятся на расстоянии x , где x нечетно, то одно слово должно иметь четную четность, а другое — нечетную. Таким образом, мы добавляем 0 в первом случае и 1 — во втором, что увеличивает расстояние между словами до $x+1$. Если x четно, мы добавляем к обоим словам 0, что не изменяет расстояния между ними. Поскольку d нечетно, все пары слов, имеющие расстояние d между собой, оказываются на расстоянии $d+1$ одно от другого. Расстояние между двумя словами, большее, чем d , либо не изменяется, либо увеличивается. Следовательно, новый код имеет минимальное расстояние $d+1$. Это показывает, что если d нечетно, то $A(n+1, d+1) \geq A(n, d)$, или, что то же самое, $A(n, d) \geq A(n-1, d-1)$ для четного $d \geq 2$.

Теперь предположим, что у нас есть код длиной n с минимальным расстоянием $d \geq 2$ (d может быть четным или нечетным). образуем новый код путем удаления любого одного столбца. Новый код имеет длину $n-1$, минимальное расстояние не менее $d-1$ и тот же размер, что и исходный код (все кодовые слова нового кода различны, поскольку новый код имеет минимальное расстояние, не меньшее 1). Следовательно, $A(n-1, d-1) \geq A(n, d)$. Таким образом, уравнение (5) доказано.

Сферы

Верхняя и нижняя границы $A(n, d)$ для любого $d \geq 1$ можно вывести путем рассмотрения n -мерных сфер. Будем рассматривать заданное кодовое слово как центр "сферы" радиусом r , состоящей из всех слов на расстоянии Хэмминга не более r от центра.

Сколько точек (слов) находится в сфере радиуса r ? Сначала рассмотрим, сколько точек находится в оболочке на расстоянии, равном ровно r , от центрального слова. Это количество определяется числом способов выбора r различных элементов из n без учета порядка выбора. Представим r выбранных битов как дополненных для образования слова на расстоянии, в точности равном r , от центральной точки. Это число сочетаний, часто записываемое как $\binom{n}{r}$, можно вычислить по формуле²

² Это число называют также биномиальным коэффициентом, поскольку $\binom{n}{r}$ является коэффициентом при члене $x^r y^{n-r}$ в разложении бинома $(x+y)^n$.

$$\binom{n}{r} = \frac{n!}{r!(n-r)!}.$$

Таким образом, $\binom{n}{0} = 1$, $\binom{n}{1} = n$, $\binom{n}{2} = n(n-1)/2$, $\binom{n}{3} = n(n-1)(n-2)/6$ и т.д.

Общее количество точек в сфере с радиусом r равно сумме точек в оболочках радиусом от 0 до r :

$$\sum_{i=0}^r \binom{n}{i}.$$

Похоже, простой формулы для этой суммы не существует [66].

Отсюда легко получить границы для $A(n, d)$. Сначала предположим, что у нас есть код длиной n с минимальным расстоянием d , и он состоит из M кодовых слов. Окружим каждое кодовое слово сферой одного и того же максимального радиуса, так, что никакие две сферы не имеют общих точек. Этот радиус равен $(d-1)/2$, если d нечетно, и $(d-2)/2$, если d четно (рис. 15.1). Поскольку каждая точка находится не более чем в одной сфере, общее количество точек в M сферах не должно превышать общего количества точек в пространстве, т.е.

$$M \sum_{i=0}^{\lfloor (d-1)/2 \rfloor} \binom{n}{i} \leq 2^n.$$

Это неравенство справедливо для любого M , следовательно, и для $M = A(n, d)$, так что

$$A(n, d) \leq \frac{2^n}{\sum_{i=0}^{\lfloor (d-1)/2 \rfloor} \binom{n}{i}}.$$

Это неравенство носит название *граница упаковки сфер* или *граница Хэмминга*.

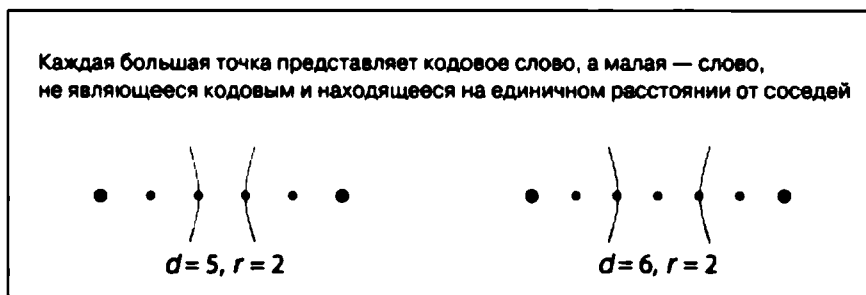


Рис. 15.1. Максимальный радиус, обеспечивающий коррекцию точек сферы

Идея со сферами легко дает и нижнюю границу $A(n, d)$. Вновь предположим, что у нас есть код длиной n с минимальным расстоянием d , который имеет максимально возможное количество кодовых слов, т.е. $A(n, d)$. Окружим каждое слово сферой с ра-

диусом $d-1$. Тогда эти сферы должны покрывать *все* 2^n точек пространства (возможно, с перекрытием). Если же это не так, то должна иметься точка на расстоянии d или более от всех кодовых слов, что невозможно, поскольку такая точка должна быть кодовым словом. Так мы получаем слабую форму границы Гильберта–Варшамова.

$$A(n, d) \sum_{i=0}^{d-1} \binom{n}{i} \geq 2^n$$

Имеется сильная форма границы Гильберта–Варшамова, применимая к *линейным* кодам. Ее вывод основан на методах линейной алгебры, которые, будучи важными при рассмотрении линейных кодов, не могут быть рассмотрены в кратком введении, посвященном кодам с коррекцией ошибок. Достаточно сказать, что линейный код представляет собой код, в котором сумма (*исключающее или*) любых двух кодовых слов также является кодовым словом. Код Хэмминга в табл. 15.1 является линейным кодом. Поскольку граница Гильберта–Варшамова представляет собой нижнюю границу для линейных кодов, она также является нижней границей для рассматриваемых здесь неограниченных кодов. При больших n это наилучшая известная нижняя граница как для линейных, так и для неограниченных кодов.

Строгая граница Гильберта–Варшамова утверждает, что $A(n, d) \geq 2^m$, где m — наибольшее целое число, такое, что

$$2^m < \frac{2^n}{\sum_{i=0}^{d-1} \binom{n-1}{i}}.$$

Иначе говоря, значение в правой части этого неравенства округляется в сторону уменьшения к строго меньшей целочисленной степени 2. “Строгость” важна для таких случаев, как $(n, d) = (8, 3)$, $(16, 3)$ и (вырожденный случай) $(6, 7)$.

Комбинируя полученные результаты, можем записать:

$$\text{GP2LT} \left(\frac{2^n}{\sum_{i=0}^{d-2} \binom{n-1}{i}} \right) \leq A(n, d) \leq \frac{2^n}{\sum_{i=0}^{\lfloor (d-1)/2 \rfloor} \binom{n}{i}}, \quad (6)$$

где GP2LT обозначает функцию, представляющую собой наибольшую целую степень 2, (строго) меньшую, чем аргумент функции.

В табл. 15.7 приведены значения этих границ для некоторых малых значений n и d . Единственное число в записи означает, что нижняя и верхняя границы в (6) совпадают.

ТАБЛИЦА 15.7. Границы Гильберта–Варшамова и Хэмминга для $A(n, d)$

n	$d = 4$	$d = 6$	$d = 8$	$d = 10$	$d = 12$	$d = 14$	$d = 16$	n
6	4-5	2	-	-	-	-	-	5
7	8-9	2	-	-	-	-	-	6
10	32-51	4-11	2-3	2	-	-	-	9
13	256-315	16-51	2-13	2-5	2	-	-	12
16	2048	64-270	8-56	2-16	2-6	2-3	2	15
19	8192-13797	256-1524	16-265	4-64	2-20	2-8	2-4	18
22	65536-95325	1024-9039	64-1342	8-277	4-75	2-25	2-10	21
25	2^{19} -671088	4096-55738	256-7216	32-1295	8-302	2-88	2-31	24
28	2^{22} -4793490	32768-354136	1024-40622	128-6436	16-1321	4-337	2-104	27
	$d = 3$	$d = 5$	$d = 7$	$d = 9$	$d = 11$	$d = 13$	$d = 15$	

Если d четно, границы можно вычислить непосредственно из (6) или, воспользовавшись (5), из (6), заменяя d на $d-1$, а n на $n-1$ в двух выражениях для границ. Оказывается, что последний способ всегда дает более точные границы. Поэтому записи в табл. 15.7 вычислялись только для нечетных d . Чтобы получить значения для четных d , воспользуйтесь значениями d в верхней строке и значениями n в левом столбце.

Как видите, границы (6) довольно нечеткие, в особенности для больших d . Отношение верхней границы к нижней расходится к бесконечности с ростом n . Нижняя граница особенно слаба. Написано более тысячи работ, в которых предлагаются методы улучшения этих границ, и достигнутые результаты показаны в табл. 15.8 [2, 13; там, где результаты этих работ различны, приводится более строгое значение].

Случаи $(n, d) = (7, 3)$, $(15, 3)$ и $(23, 7)$ являются *идеальными* кодами, что означает, что они достигают верхней границы, указанной в (6). Это определение представляет собой обобщение, данное на с. 359. Коды, у которых n нечетно и $n = d$, также являются идеальными (см. упр. 8).

Мы завершим данную главу, указав, что идея минимального расстояния для всего кода, которая приводит к обнаружению p -кратных и коррекции q -кратных ошибок для некоторых p и q , не является единственным критерием "мощности" бинарного блочного кода с прямой коррекцией ошибок. Например, была проделана работа по кодам, направленным на исправление пакетных ошибок. В работе [30] продемонстрирован код $(16, 11)$ (и некоторые другие), который может исправить любые однократные ошибки и любые ошибки в двух последовательных битах, и является идеальным в некотором смысле, который здесь не рассматривается. Однако исправлять двойные ошибки в общем случае этот код не в состоянии. Расширенный код Хемминга $(16, 11)$ является идеальным SEC-DED кодом. Таким образом, его код обеспечивает *обнаружение* двойных ошибок в об-

шем случае вместо *исправления* двойных ошибок в последовательных битах. Это достаточно важно и интересно, потому что во многих приложениях ошибки происходят в пределах коротких пакетов.

ТАБЛИЦА 15.8. НАИЛУЧШИЕ ИЗВЕСТНЫЕ ГРАНИЦЫ ДЛЯ $A(n, d)$

n	$d = 4$	$d = 6$	$d = 8$	$d = 10$	$d = 12$	$d = 14$	$d = 16$	n
6	4	2	-	-	-	-	-	5
7	8	2	-	-	-	-	-	6
8	16	2	2	-	-	-	-	7
9	20	4	2	-	-	-	-	8
10	40	6	2	2	-	-	-	9
11	72	12	2	2	-	-	-	10
12	144	24	4	2	2	-	-	11
13	256	32	4	2	2	-	-	12
14	512	64	8	2	2	2	-	13
15	1024	128	16	4	2	2	-	14
16	2048	256	32	4	2	2	2	15
17	2720-3276	256-340	36	6	2	2	2	16
18	5312-6552	512-673	64-72	10	4	2	2	17
19	10496-13104	1024-1237	128-135	20	4	2	2	18
20	20480-26168	2048-2279	256	40	6	2	2	19
21	36864-43688	2560-4096	512	42-47	8	4	2	20
22	73728-87376	4096-6941	1024	64-84	12	4	2	21
23	147456-173015	8192-13674	2048	80-150	24	4	2	22
24	294912-344308	16384-24106	4096	128-268	48	6	4	23
25	2^{19} -599184	16384-47538	4096-5421	192-466	52-55	8	4	24
26	2^{20} -1198368	32768-84260	4104-9672	384-836	64-96	14	4	25
27	2^{21} -2396736	65536-157285	8192-17768	512-1585	128-169	28	6	26
28	2^{22} -4792950	131072-291269	16384-32151	1024-3170	178-288	56	8	27
	$d = 3$	$d = 5$	$d = 7$	$d = 9$	$d = 11$	$d = 13$	$d = 15$	

Упражнения

1. Приведите код Хэмминга для $m = 3$ (постройте таблицу, подобную табл. 15.1).
2. В некоторых приложениях SEC-кодов не требуется выполнять коррекцию проверочных битов. Следовательно, k проверочных битов нужны только для проверки информационных битов, но не самих проверочных битов. Для m информационных битов k должно быть достаточно велико, чтобы получатель мог различать $m + 1$ случай: в каком из m бит ошибка или когда ошибок нет вовсе. Таким образом, требуемое количество проверочных битов задается соотношением $2^k \geq m + 1$. Это более слабое ограничение, накладываемое на k , чем правило Хэмминга, так что должно быть возможным построить для некоторых значений m такой SEC-код, у которого проверочных битов будет меньше, чем требуется правилом Хэмминга. Как вариант можно иметь одно значение, которое указывает, что ошибка произошла в проверочных битах, но не указывает, где именно. В таком случае получается правило $2^k \geq m + 2$.

Где в приведенных рассуждениях скрывается ошибка?

3. Как для заданного m найти наименьшее k , удовлетворяющее неравенству (1)?
4. Покажите, что функция расстояния Хэмминга для любого бинарного блочного кода удовлетворяет неравенству треугольника: если x и y являются кодовыми векторами и $d(x, y)$ обозначает расстояние Хэмминга между ними, то

$$d(x, z) \leq d(x, y) + d(y, z).$$

5. Докажите неравенство $A(2n, 2d) \geq A(n, d)$.
6. Докажите "границу синглтона" $A(n, d) \leq 2^{n-d+1}$.
7. Покажите, что понятие идеального кода, как достижение равенства в правой части неравенства (6), является обобщением правила Хэмминга.
8. Чему равно значение $A(n, d)$ при $n = d$? Покажите, что для нечетных n эти коды являются идеальными.
9. Покажите, что если n кратно 3, а $d = 2n/3$, то $A(n, d) = 4$.
10. Покажите, что если $d > 2n/3$, то $A(n, d) = 2$.
11. Схема двумерной проверки четности для 64 информационных битов размещает информационные биты $m_0 \dots m_{63}$ в массиве 8×8 , и к каждой строке и каждому столбцу добавляет бит четности, как показано ниже.

$$\begin{bmatrix} M_0 & M_1 & \dots & M_6 & M_7 & F_0 \\ M_8 & M_9 & \dots & M_{14} & M_{15} & F_1 \\ & & \dots & & & \\ M_{48} & M_{49} & \dots & M_{54} & M_{55} & F_6 \\ M_{56} & M_{57} & \dots & M_{62} & M_{63} & F_7 \\ c_0 & c_1 & \dots & c_6 & c_7 & F_8 \end{bmatrix}$$

Биты r_i представляют собой биты проверки четности в строках, а c_i — биты проверки четности в столбцах. “Угловой” проверочный бит может быть либо битом четности строки, либо битом четности столбца, но не обоих одновременно; на рисунке он показан как бит четности нижней строки (проверочных битов с c_0 по c_7).

Прокомментируйте эту схему. В частности, является ли она SEC-DED кодом? Существенно ли изменятся ее возможности по обнаружению и коррекции ошибок, если будет опущен угловой бит r_8 ? Имеется ли простое соотношение между значениями углового бита, когда он представляет собой сумму строки и сумму столбца?

ГЛАВА 16

КРИВАЯ ГИЛЬБЕРТА

В 1890 году Джузеппе Пеано (Giuseppe Peano) открыл плоскую кривую¹ с удивительным свойством “заполнения пространства”. Такая кривая заполняла единичный квадрат и проходила через каждую его точку (x, y) по меньшей мере один раз.

Кривая Пеано основана на разделении каждой стороны единичного квадрата на три равные части, которые делят его на девять меньших квадратов. Кривая проходит эти девять квадратов в определенном порядке. Затем каждый из девяти малых квадратов аналогично делится на девять частей, и кривая модифицируется таким образом, чтобы обойти все части в определенном порядке. Эта кривая может быть описана с использованием дробных чисел, записанных в системе счисления по основанию 3; Пеано описал ее впервые именно так.

В 1891 году Давид Гильберт (David Hilbert) [49] открыл вариант кривой Пеано, основанной на делении каждой стороны единичного квадрата на две равные части, что делит квадрат на четыре меньшие части. Затем каждый из четырех получившихся квадратов, в свою очередь, аналогично делится на четыре меньших квадрата и т.д. На каждой стадии такого деления Гильберт строил кривую, которая обходила все имеющиеся квадраты. Кривая Гильберта (которую иногда называют кривой Пеано–Гильберта) представляет собой предельную кривую, полученную в результате такого построения. Ее можно описать с помощью дробей, записанных в системе счисления по основанию 2.

На рис. 16.1 показаны первые три шага последовательности, приводящей к получению заполняющей пространство кривой Гильберта, в том виде, в котором они были показаны в его статье в 1891 году.

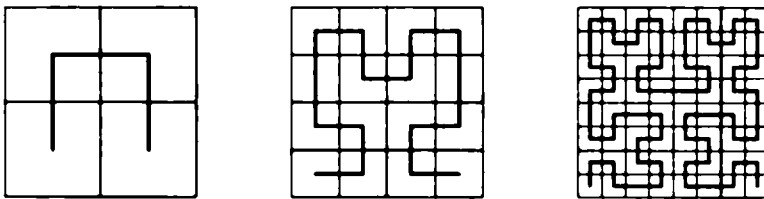


Рис. 16.1. Первые три кривые в последовательности, определяющей кривую Гильберта

Здесь мы поступим немного иначе. Используем термин “кривая Гильберта” для любой кривой из последовательности построения заполняющей пространство кривой Гильберта, и под кривой Гильберта n -го порядка будет подразумеваться n -я кривая последовательности (на рис. 16.1 показаны кривые первого, второго и третьего порядков). Каждая кривая порядка n масштабируется множителем 2^n с тем, чтобы координаты углов кривой представляли собой целые числа. Таким образом, наша кривая Гильберта порядка n имеет углы с координатами

¹ Напомним: кривая представляет собой непрерывное отображение одномерного пространства на n -мерное.

натами от 0 до $2^n - 1$ по осям x и y . Примем положительным направление вдоль кривой от точки $(x, y) = (0, 0)$ к $(2^n - 1, 0)$. На рис. 16.2 показаны рассматриваемые нами масштабированные кривые Гильберта от первого до шестого порядка.

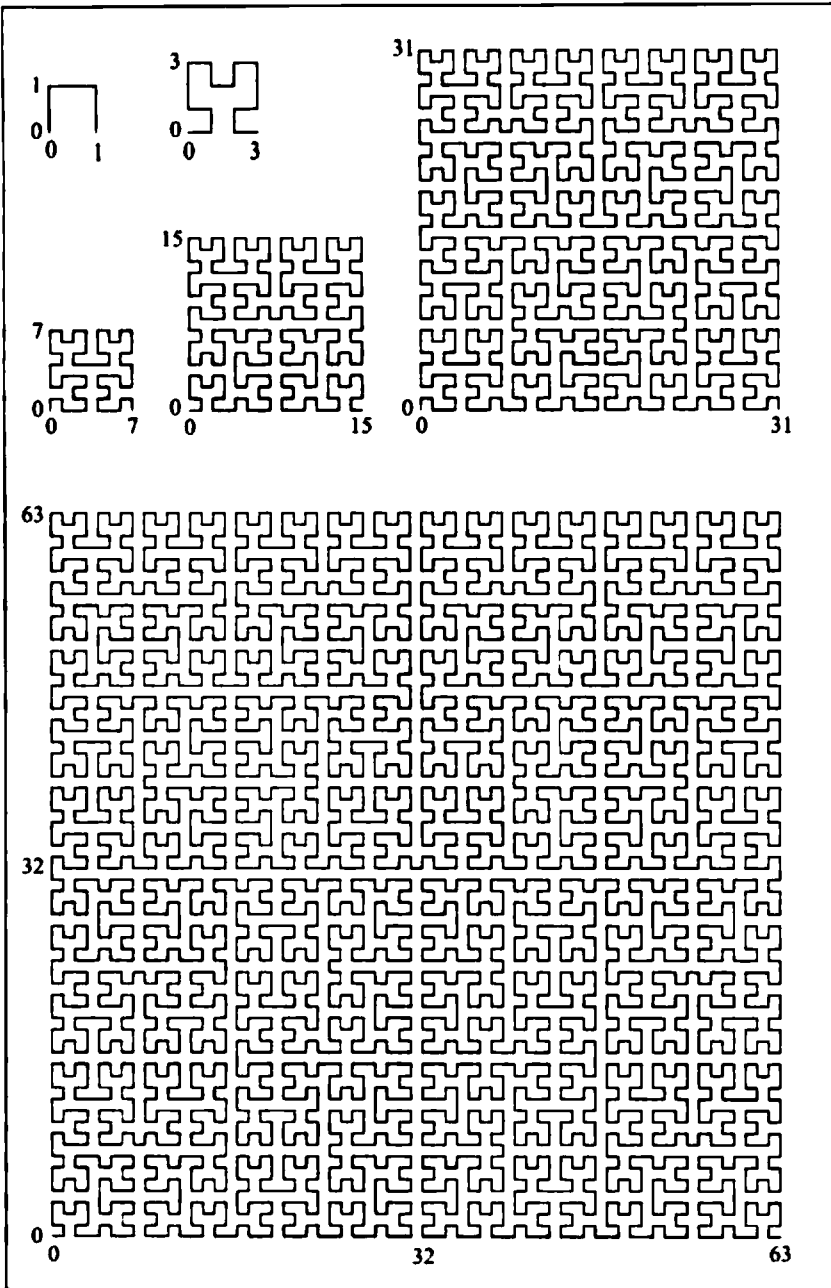


Рис. 16.2. Кривые Гильберта порядка 1–6

16.1. Рекурсивный алгоритм построения кривой Гильберта

Для того чтобы понять, каким образом строится кривая Гильберта, рассмотрим внимательно кривые на рис. 16.2. Кривая порядка 1 состоит из отрезков, направленных вверх, вправо и вниз. Кривая порядка 2 следует тому же общему шаблону. Сначала изображается U-образная кривая, идущая вверх, затем от нее вверх идет отрезок, который соединяется с очередной U-образной кривой, направленной вправо. От нее идет отрезок вправо, к такой же U-образной кривой, от которой отрезок вниз ведет нас к последней U-образной кривой, ведущей вниз. В результате из перевернутой U-образной кривой первого порядка будет получена Y-образная кривая второго порядка.

Кривую Гильберта любого порядка можно рассматривать как серию U-образных кривых разной ориентации, за каждой из которых, за исключением последней, следует отрезок в определенном направлении. При преобразовании кривой Гильберта некоторого порядка в кривую следующего порядка каждая U-образная кривая преобразуется в Y-образную кривую с той же общей ориентацией, а каждый соединяющий отрезок — в отрезок в том же направлении.

Преобразование кривой Гильберта первого порядка (кривая U с общим направлением вправо и вращательной ориентацией по часовой стрелке) в кривую второго порядка происходит следующим образом.

1. Рисуем U по направлению вверх против часовой стрелки.
2. Рисуем отрезок, направленный вверх.
3. Рисуем U по направлению вправо по часовой стрелке.
4. Рисуем отрезок, направленный вправо.
5. Рисуем U по направлению вправо по часовой стрелке.
6. Рисуем отрезок, направленный вниз.
7. Рисуем U по направлению вниз против часовой стрелки.

Рассматривая кривые разных порядков, можно видеть, что все U, ориентированные так же, как и в кривой Гильберта первого порядка, трансформируются таким же образом. Сходный набор правил трансформации может быть создан для каждой U-образной кривой с другой ориентацией. Все эти правила легко укладываются в рекурсивную программу построения кривой Гильберта, показанную в листинге 16.1 [107]. В приведенной программе ориентация U характеризуется двумя целыми числами, которые определяют направление и вращательное направления следующим образом.

dir = 0 :	вправо	rot = +1 :	по часовой стрелке
dir = 1 :	вверх	rot = -1 :	против часовой стрелки
dir = 2 :	влево		
dir = 3 :	вниз		

ЛИСТИНГ 16.1. Генерация кривой Гильберта

```

void step(int);

void hilbert(int dir, int rot, int order)
{
    if (order == 0) return;

    dir = dir + rot;
    hilbert(dir, -rot, order - 1);
    step(dir);
    dir = dir - rot;
    hilbert(dir, rot, order - 1);
    step(dir);
    hilbert(dir, rot, order - 1);
    dir = dir - rot;
    step(dir);
    hilbert(dir, -rot, order - 1);
}

```

В действительности переменная `dir` может принимать и другие значения, вне диапазона 0–3, но в этом случае необходимо брать значение `dir` по модулю 4.

В листинге 16.2 приведены программа-оболочка и функция `step`, которая используется в функции `hilbert`. Эта программа получает в качестве аргумента порядок генерируемой кривой Гильберта и выводит список отрезков, указывая для каждого направление перемещения и длину кривой до конца отрезка, а также координаты конца текущего отрезка. Например, для кривой Гильберта второго порядка вывод программы оказывается следующим.

```

0 0000 00 00
0 0001 01 00
1 0010 01 01
2 0011 00 01
1 0100 00 10
1 0101 00 11
0 0110 01 11
-1 0111 01 10
0 1000 10 10
1 1001 10 11
0 1010 11 11
-1 1011 11 10
-1 1100 11 01
-2 1101 10 01
-1 1110 10 00
0 1111 11 00

```

ЛИСТИНГ 16.2. Программа-оболочка для генератора кривой Гильберта

```

#include <stdio.h>
#include <stdlib.h>

int x = -1, y = 0; // Глобальные переменные
int i = 0;        // Расстояние вдоль кривой
int blen;        // Длина вывода

```

```

void hilbert(int dir, int rot, int order);

void binary(unsigned k, int len, char* s)
{
    /* Преобразование беззнакового целого k в строку
       символов. Результат хранится в строке s длиной len.
    */
    int i;
    s[len] = 0;
    for (i = len - 1; i >= 0; i--)
    {
        if (k & 1) s[i] = '1';
        else      s[i] = '0';
        k = k >> 1;
    }
}

void step(int dir)
{
    char ii[33], xx[17], yy[17];

    switch (dir & 3)
    {
        case 0:  x = x + 1;  break;
        case 1:  y = y + 1;  break;
        case 2:  x = x - 1;  break;
        case 3:  y = y - 1;  break;
    }

    binary(i, 2*blen, ii);
    binary(x, blen, xx);
    binary(y, blen, yy);
    printf("%5d  %s  %s  %s\n", dir, ii, xx, yy);
    i = i + 1;      // Увеличение расстояния
}

int main(int argc, char* argv[])
{
    int order;
    order = atoi(argv[1]);
    blen = order;
    step(0);        // Вывод начальной точки
    hilbert(0, 1, order);
    return 0;
}

```

16.2. Преобразование расстояния вдоль кривой Гильберта в координаты





Для того чтобы найти координаты (x, y) точки, определяемой расстоянием s вдоль кривой Гильберта порядка n , заметим, что старшие два бита $2n$ -битового числа s определяют, в каком квадранте находится точка. Это происходит потому, что кривая Гильберта любого порядка следует общему шаблону кривой первого порядка. Если два старших бита s равны 00, точка находится где-то в нижнем левом квадранте; 01 указывает на



верхний левый квадрант, 10 — на верхний правый квадрант и 11 — на нижний правый квадрант. Таким образом, два старших бита z определяют старшие биты n -битовых чисел x и y , как показано ниже.

Два старших бита z	Старшие биты (x, y)
00	(0, 0)
01	(0, 1)
10	(1, 1)
11	(1, 0)

В любой кривой Гильберта встречаются только четыре из восьми возможных U-образных кривых. Они графически показаны в табл. 16.1, где приведено также отображение двух битов z на единичные биты x и y .

ТАБЛИЦА 16.1. ЧЕТЫРЕ ВОЗМОЖНЫХ ОТОБРАЖЕНИЯ

A	B	C	D
			
00 → (0, 0)	00 → (0, 0)	00 → (1, 1)	00 → (1, 1)
01 → (0, 1)	01 → (1, 0)	01 → (1, 0)	01 → (0, 1)
10 → (1, 1)	10 → (1, 1)	10 → (0, 0)	10 → (0, 0)
11 → (1, 0)	11 → (0, 1)	11 → (0, 1)	11 → (1, 0)

Заметим, что на рис. 16.2 во всех случаях U-образная кривая, представленная отображением A () становится на следующем уровне детализации U-образной кривой, представленной отображением B, A, A или D (в зависимости от расстояния от начала исходного отображения A — 0, 1, 2 или 3). Аналогично U-образная кривая, представленная отображением B () на следующем уровне детализации становится представленной отображением A, B, B или C (в зависимости от расстояния вдоль исходной кривой — 0, 1, 2 или 3). Результатом таких наблюдений является табл. 16.2, в которой представлены переходы между состояниями, соответствующими отображениям, показанным в табл. 16.1.

Для использования этой таблицы начнем с состояния A. Целое z следует дополнить необходимым количеством нулевых битов слева с тем, чтобы его длина стала равной $2n$ бит, где n — порядок кривой Гильберта. После этого сканируем попарно биты z слева направо. Первая строка табл. 16.2 означает, что если текущим состоянием является A и сканированная пара битов z равна 00, то мы получаем на выходе (0, 0) и переходим к состоянию B, после чего считываем очередные два бита слова z . Аналогично вторая строка табл. 16.2 означает, что если текущим состоянием является A, а сканированная пара битов z равна 01, то на выходе мы получаем (0, 1) и остаемся в состоянии A.

ТАБЛИЦА 16.2. ТАБЛИЦА ПЕРЕХОДОВ СОСТОЯНИЙ ДЛЯ ВЫЧИСЛЕНИЯ (x, y) ИЗ s

Если текущее состояние	и следующие (справа) биты s	то добавляем к (x, y)	и переходим к состоянию
A	00	(0, 0)	B
A	01	(0, 1)	A
A	10	(1, 1)	A
A	11	(1, 0)	D
B	00	(0, 0)	A
B	01	(1, 0)	B
B	10	(1, 1)	B
B	11	(0, 1)	C
C	00	(1, 1)	D
C	01	(1, 0)	C
C	10	(0, 0)	C
C	11	(0, 1)	B
D	00	(1, 1)	C
D	01	(0, 1)	D
D	10	(0, 0)	D
D	11	(1, 0)	A

Получаемые на выходе биты накапливаются в порядке слева направо. Когда s оказывается сканированным до конца, мы получаем на выходе n -битовые величины x и y . Предположим в качестве примера, что $n = 3$ и $s = 110100$. Поскольку процесс начинается в состоянии A и первые сканированные биты равны 11, мы получаем на выходе (1,0) и переходим в состояние D (четвертая строка табл. 16.2). Затем, находясь в состоянии D и получая очередную пару битов s , равную 01, в соответствии со строкой 14 табл. 16.2 мы получаем на выходе (0,1) и остаемся в состоянии D. И наконец последняя пара битов 00 дает нам на выходе (1,1) и переводит нас в состояние C (хотя, поскольку строка s сканирована до конца, это состояние не играет никакой роли).

Итак, в результате сканирования и переходов получена пара (101,011), т.е. $x = 5$ и $y = 3$.

Реализующая описанный алгоритм программа на языке программирования C приведена в листинге 16.3. В этой программе текущее состояние представлено целым числом из диапазона 0–3, соответствующего состояниям A–D. В присвоении переменной `row` текущее состояние объединяется с очередными двумя битами `s`, давая целое число от 0 до 15, которое представляет собой номер строки в табл. 16.2. Переменная `row` используется при обращении к целым числам, которые представляют собой два крайних справа столбца табл. 16.2, т.е., по сути, осуществляется выборка из таблицы в регистре. В используемых шестнадцатеричных значениях биты слева направо соответствуют значениям табл. 16.2, просматриваемым снизу вверх.

Листинг 16.3. Программа для вычисления x и y по значению s

```
void hil_xy_from_s(unsigned s, int n,
                  unsigned* xp, unsigned* yp)
{
    int i;
    unsigned state, x, y, row;

    state = 0; // Инициализация
    x = y = 0;

    for (i = 2*n - 2; i >= 0; i -= 2) // Выполнить n раз
    {
        row = 4*state | (s >> i) & 3; // Строка в таблице
        x = (x << 1) | (0x936C >> row) & 1;
        y = (y << 1) | (0x39C6 >> row) & 1;
        state =
            (0x3E6B94C1 >> 2*row) & 3; // Новое состояние
    }

    *xp = x; // Возврат
    *yp = y; // результатов
}
```

В [82] приводится несколько иной алгоритм. В отличие от алгоритма из листинга 16.3, он сканирует биты s справа налево. Алгоритм базируется на том, что можно отобразить младшие значащие биты s на (x, y) , основываясь на кривой Гильберта первого порядка, после чего переходить к тестированию следующей пары битов s . Если они равны 00, то только что вычисленные значения (x, y) должны поменяться местами (что соответствует отражению относительно прямой $x = y$ (см. рис. 16.1, кривые первого и второго порядков). Если эти биты равны 01 или 10, то значения x и y остаются неизменными. И наконец, если значение пары битов равно 11, то значения x и y обмениваются и к ним применяется операция дополнения. Те же правила применяются и при дальнейшем сканировании битовой строки s . Этот алгоритм схематично представлен в табл. 16.3, а реализующий его код — в листинге 16.4. Интересно, что сначала можно добавлять биты к x и y , а уже затем выполнять операции обмена и дополнения над полученными значениями, включающими только что добавленные биты; результат окажется тем же, что и при добавлении битов после обмена и дополнения.

ТАБЛИЦА 16.3. МЕТОД ЛАМА И ШАПИРО ВЫЧИСЛЕНИЯ (x, y) ПО ЗНАЧЕНИЮ s

Если следующие два бита s слева равны	то	и добавляем спереди к (x, y)
00	Обмениваем x и y	$(0, 0)$
01	Оставляем x и y без изменений	$(0, 1)$
10	Оставляем x и y без изменений	$(1, 1)$
11	Обмениваем x и y и выполняем операцию дополнения	$(1, 0)$

ЛИСТИНГ 16.4. Метод Лама и Шапиро вычисления (x, y) по значению s

```

void hil_xy_from_s(unsigned s, int n,
                  unsigned* xp, unsigned* yp)
{
    int i, sa, sb;
    unsigned x, y, temp;

    for (i = 0; i < 2*n; i += 2)
    {
        sa = (s >> (i+1)) & 1; // Бит i+1 строки s
        sb = (s >> i) & 1;     // Бит i строки s

        if ((sa ^ sb) == 0)    // Если sa, sb = 00 или 11,
        {
            temp = x;         // обмениваем x и y
            x = y ^ (-sa);    // и, если sa = 1,
            y = temp ^ (-sa); // дополняем их
        }

        x = (x >> 1) |
            (sa << 31);       // Добавляем sa к x и
        y = (y >> 1) |
            ((sa ^ sb) << 31); // (sa^sb) к y (спереди)
    }

    *xp = x >> (32 - n);     // Выравниваем x и y вправо
    *yp = y >> (32 - n);     // и возвращаем их
}

```

Переменные x и y в листинге 16.4 не инициализированы, что в некоторых компиляторах может привести к сообщению об ошибке. Однако приведенный код корректно функционирует независимо от того, какие значения имели переменные x и y изначально.

Условного перехода в цикле в листинге 16.4 можно избежать, если воспользоваться рассматривавшимся в разделе 2.20 на с. 68 приемом “трех *исключающих или*”. Блок `if` при этом можно заменить следующим кодом (здесь `swap` и `str1` — беззнаковые целые переменные).

```

swap = (sa ^ sb) - 1;
// Если требуется обмен, -1; в противном случае 0

cmpl = -(sa & sb);
// Если требуется дополнение, -1; в противном случае 0

x = x ^ y;
y = y ^ (x & swap) ^ cmpl;
x = x ^ y;
    
```

Однако при этом требуется выполнение девяти команд, в то время как при использовании условного перехода достаточно только двух или шести, в зависимости от выполнения условия, так что, возможно, отказ от условного перехода в данном случае — не лучший выбор.

Идея "обмена и дополнения" из [82] подсказывает логическую схему для построения кривой Гильберта. Идея, лежащая в основе описываемой далее схемы, состоит в том, чтобы, проходя вдоль кривой n -го порядка, вы отображали пары битов s в координаты (x, y) в соответствии с отображением A из табл. 16.1. При прохождении различных областей координаты, полученные в результате отображения, могут меняться местами и дополняться либо над ними могут выполняться обе эти операции. Схема, показанная на рис. 16.3, отслеживает необходимость выполнения данных операций на каждом шаге, использует соответствующее отображение двух битов s на (x, y) и генерирует сигналы обмена и дополнения для следующего шага.

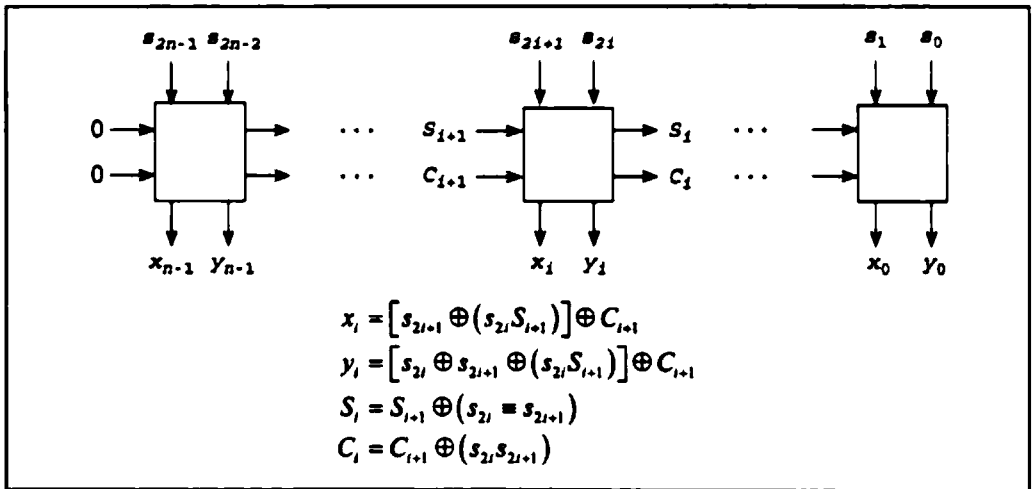


Рис. 16.3. Логическая схема для вычисления (x, y) из значения s

Предположим, что имеется регистр, содержащий длину пути s вдоль кривой Гильберта и схемы для ее увеличения. Тогда для поиска следующей точки кривой Гильберта следует сначала увеличить значение s , а затем преобразовать его так, как показано в табл. 16.4. Преобразования s в таблице выполняются слева направо, что создает определенные трудности, так как увеличение s — процесс, затрагивающий биты в порядке

справа налево. Таким образом, время, необходимое для генерации очередной точки кривой Гильберта n -го порядка, пропорционально $2n$ (для увеличения s) плюс n (для преобразования s в пару координат (x, y)).

ТАБЛИЦА 16.4. Логика вычислений (x, y) по значению s

Если очередная (справа) пара битов s равна	то добавляем к (x, y)	и устанавливаем
00	$(0, 0)^*$	$\text{swap} = \overline{\text{swap}}$
01	$(0, 1)^*$	Без изменений
10	$(1, 1)^*$	Без изменений
11	$(1, 0)^*$	$\text{swap} = \overline{\text{swap}}, \text{cmp1} = \overline{\text{cmp1}}$

* Возможно, обмененные или дополненные.

На рис. 16.3 эти вычисления показаны в виде логической схемы (S обозначает сигнал обмена, а C — сигнал дополнения).

В схеме на рис. 16.3 предлагается другой путь вычисления (x, y) по значению s . Обратите внимание на то, как сигналы обмена и дополнения распространяются слева направо через n стадий. Таким образом, здесь можно применить метод параллельного префикса для быстрого (не за $n-1$, а за $\log_2 n$ шагов) распространения информации об обмене и дополнении по всем стадиям, а затем использовать некоторые побитовые операции для вычисления x и y с использованием приведенных на рис. 16.3 уравнений. Значения x и y перемешаны, и их биты находятся в четных и нечетных позициях слова, так что их следует извлечь оттуда с использованием информации из раздела 7.2 на с. 166. Это может показаться несколько сложным и окупающимся только для больших значений n , однако давайте посмотрим, как все происходит на самом деле.

Процедура, реализующая описанную операцию, приведена в листинге 16.5 [37]. Эта процедура оперирует величинами, представляющими собой полные слова, так что сначала входное значение s дополняется слева битами '01' (эта битовая комбинация не влияет на количество обменов и дополнений). Затем вычисляется величина cs (complement-swap — дополнение-обмен). Это слово имеет вид $cscs\dots cs$, где каждый отдельный бит c , будучи равным 1, означает, что соответствующая пара битов должна быть дополнена, а s — что соответствующая пара битов должна быть обменена в соответствии с табл. 16.3. Другими словами, каждая пара битов s отображается в пару cs следующим образом.

s_{2i-1}	s_{2i}	cs
0	0	01
0	1	00
1	0	00
1	1	11

Это и есть та величина, к которой мы хотим применить операцию параллельного префикса, а именно PP-XOR. Операция применяется слева направо, поскольку последовательные единичные биты, означающие дополнение или обмен, имеют те же логические свойства, что и операция *исключающего или*: два последовательных единичных бита взаимно исключают друг друга.

Листинг 16.5. Метод параллельного префикса для вычисления (x, y) по значению s

```
void hil_xy_from_s(unsigned s, int n,
                  unsigned *xp, unsigned *yp)
{
    unsigned comp, swap, cs, t, sr;

    s = s | (0x55555555 << 2*n); // Заполняем s слева парами
    sr = (s >> 1) & 0x55555555; // 01 (без изменений)
    cs = ((s & 0x55555555) + sr) // Вычисляем пары дополнений
          ^ 0x55555555;        // и обменов

    // Операция PP-XOR распространяет информацию о дополнениях
    // и обменах слева направо, парами (шаг cs ^= cs >> 1
    // отсутствует), что дает в результате две операции
    // параллельного префикса над двумя чередующимися
    // множествами из 16 бит каждое

    cs = cs ^ (cs >> 2);
    cs = cs ^ (cs >> 4);
    cs = cs ^ (cs >> 8);
    cs = cs ^ (cs >> 16);
    swap = cs & 0x55555555; // Разделяем биты обмена
    comp = (cs >> 1) & 0x55555555; // и дополнения

    t = (s & swap) ^ comp; // Вычисляем x и y в
    s = s ^ sr ^ t ^ (t << 1); // нечетных и четных
    // позиций соответственно
    s = s & ((1 << 2*n) - 1); // Убираем лишние биты слева

    // Разделяем значения x и y

    t = (s ^ (s >> 1)) & 0x22222222; s = s ^ t ^ (t << 1);
    t = (s ^ (s >> 2)) & 0x0C0C0C0C; s = s ^ t ^ (t << 2);
    t = (s ^ (s >> 4)) & 0x00F000F0; s = s ^ t ^ (t << 4);
    t = (s ^ (s >> 8)) & 0x0000FF00; s = s ^ t ^ (t << 8);

    *xp = s >> 16; // Присваиваем две половины t
    *yp = s & 0xFFFF; // переменным x и y
}
```

Оба сигнала (дополнения и обмена) распространяются одной и той же операцией PP-XOR.

Следующие четыре присвоения транслируют каждую пару битов s в значения (x, y) , где x состоит из битов, находящихся в нечетных (левых) позициях, а y — из битов в четных позициях. Хотя сама логика может показаться и непонятной, не так уж сложно убедиться в том, что каждая пара битов s преобразуется в соответствии с первыми двумя уравнениями на рис. 16.3 (*указание*: рассмотрите отдельно преобразования битов в четных и нечетных позициях, не забывая о том, что в нечетных позициях биты t и sr равны 0).

Оставшаяся часть процедуры очевидна. Всего она требует выполнения 66 базовых RISC-команд (без ветвлений), в то время как код из листинга 16.4 требует в среднем выполнения $19n + 10$ команд (вычислено на основе скомпилированного кода). Таким образом, уже при $n \geq 3$ метод параллельного префикса оказывается более быстрым.

16.3. Преобразование координат в расстояние вдоль кривой Гильберта

Если заданы координаты точки на кривой Гильберта, то расстояние до данной точки вдоль кривой от ее начала можно вычислить с помощью таблицы переходов состояний, подобной табл. 16.2. Для нашей задачи такой таблицей переходов является табл. 16.5.

ТАБЛИЦА 16.5. ТАБЛИЦА ПЕРЕХОДОВ СОСТОЯНИЙ ДЛЯ ВЫЧИСЛЕНИЯ s ПО ЗНАЧЕНИЯМ (x, y)

Если текущее состояние	и следующие справа два бита (x, y) равны	то добавляем к s	и переходим в состояние
A	(0,0)	00	B
A	(0,1)	01	A
A	(1,0)	11	D
A	(1,1)	10	A
B	(0,0)	00	A
B	(0,1)	11	C
B	(1,0)	01	B
B	(1,1)	10	B
C	(0,0)	10	C
C	(0,1)	11	B
C	(1,0)	01	C
C	(1,1)	00	D
D	(0,0)	10	D
D	(0,1)	01	D
D	(1,0)	11	A
D	(1,1)	00	C

Работа с этой таблицей аналогична работе с таблицей переходов состояний в предыдущем разделе. Сначала значения x и y должны быть дополнены слева ведущими нулевыми битами с тем, чтобы длина битовых строк, представляющих эти величины, стала равной n , где n — порядок рассматриваемой кривой Гильберта. Затем биты x и y сканируются слева направо, и значение s строится в том же направлении.

Программа на языке программирования C, реализующая эти действия, представлена в листинге 16.6.

ЛИСТИНГ 16.6. Программа для вычисления s по значениям (x, y)

```
unsigned hil_s_from_xy(unsigned x,
                      unsigned y, int n)
{
    int i;
    unsigned state, s, row;

    state = 0;           // Инициализация
    s = 0;

    for (i = n - 1; i >= 0; i--)
    {
        row = 4*state | 2*((x >> i) & 1) | (y >> i) & 1;
        s = (s << 2) | (0x361E9CB4 >> 2*row) & 3;
        state = (0xFE65831 >> 2*row) & 3;
    }
    return s;
}
```

В [82] приведен алгоритм для вычисления s по значениям (x, y) , похожий на соответствующий алгоритм для вычислений в обратном направлении — (x, y) по значению s . Этот алгоритм основан на сканировании слева направо и показан в табл. 16.6 и листинге 16.7.

ТАБЛИЦА 16.6. Метод вычисления s по значениям (x, y)

Если следующие два бита (x, y) справа равны	то	и добавляем к s
(0,0)	Обмениваем x и y	00
(0,1)	Оставляем x и y без изменений	01
(1,0)	Обмениваем x и y и выполняем операцию дополнения	11
(1,1)	Оставляем x и y без изменений	10

ЛИСТИНГ 16.7. Метод вычисления s по значениям (x, y)

```
unsigned hil_s_from_xy(unsigned x,
                      unsigned y, int n)
{
    int i, xi, yi;
    unsigned s, temp;

    s = 0;           // Инициализация

    for (i = n - 1; i >= 0; i--)
    {
        xi = (x >> i) & 1;           // Получаем  $i$ -й бит  $x$ 
```

```

yi = (y >> i) & 1;          // Получаем i-й бит y
if (yi == 0)
{
    temp = x;                // Обмениваем x и y и,
    x = y^(-xi);            // если xi = 1,
    y = temp^(-xi);         // дополняем их
}
s = 4*s + 2*xi + (xi^yi); // Добавляем два бита к s
}
return s;
}

```

16.4. Увеличение координат кривой Гильберта

Пусть заданы координаты (x, y) точки на кривой Гильберта n -го порядка. Каким образом можно найти координаты следующей точки? Простейший путь, конечно, состоит в преобразовании координат в расстояние вдоль кривой s , увеличении s на 1 и обратном преобразовании в координаты (x, y) с использованием описанных в предыдущих разделах алгоритмов.

Более удачный (но не слишком) путь основывается на том факте, что при перемещении вдоль кривой Гильберта на каждом шаге или x , или y (но не обе величины одновременно) либо увеличиваются, либо уменьшаются на 1. Поэтому можно разработать алгоритм, который сканирует координаты слева направо для определения того, какой тип U-образной кривой представляют собой младшие два бита. Затем на основании этой информации и значения двух младших битов осуществляется увеличение или уменьшение x или y .

В данном алгоритме есть сложность, возникающая на каждом четвертом шаге, когда рассматриваемая точка оказывается в конце U-образной кривой. В этой точке направление определяется *предыдущими* битами x и y и U-образной кривой более высокого порядка, с которой эти биты связаны. Если же точка на U-образной кривой более высокого порядка также оказывается в конце, то необходимо в очередной раз рассмотреть предыдущие биты x и y и U-образную кривую еще более высокого порядка, и т.д.

Этот алгоритм описывается табл. 16.7, где A, B, C и D обозначают U-образные кривые, показанные в табл. 16.1 на с. 386. Для использования данной таблицы сначала надо добавить к x и y ведущие нулевые биты с тем, чтобы получились n -битовые значения (где n — порядок рассматриваемой кривой Гильберта). Начиная с состояния A, сканируем биты x и y слева направо. Первая строка табл. 16.7 означает, что если текущим состоянием является A, а сканированы биты $(0, 0)$, то нужно установить переменную-флаг, указывающую на необходимость увеличения y , и перейти в состояние B. Прочие строки таблицы интерпретируются аналогично; суффикс “минус” означает, что соответствующая координата должна быть уменьшена. Прочерк в третьем столбце таблицы говорит о том, что переменная-флаг, отслеживающая изменение координат, остается на данном шаге неизменной.

По окончании сканирования битов x и y увеличиваем (или уменьшаем) значение одной из координат в соответствии со значением переменной-флага.

ТАБЛИЦА 16.7. ДОБАВЛЕНИЕ ОДНОГО ЗВЕНА КРИВОЙ ГИЛЬБЕРТА

Если текущее состояние	а очередные два бита (x,y) справа	то нужно изменить координату	и перейти в состояние
A	(0,0)	y+	B
A	(0,1)	x+	A
A	(1,0)	—	D
A	(1,1)	y-	A
B	(0,0)	x+	A
B	(0,1)	—	C
B	(1,0)	y+	B
B	(1,1)	x-	B
C	(0,0)	y+	C
C	(0,1)	—	B
C	(1,0)	x-	C
C	(1,1)	y-	D
D	(0,0)	x+	D
D	(0,1)	y-	D
D	(1,0)	—	A
D	(1,1)	x-	C

Программа на языке C, реализующая рассмотренный алгоритм, представлена в листинге 16.8. Переменная `dx` инициализируется таким образом, чтобы при многократных вызовах данной функции в цикле генерировалась одна и та же кривая Гильберта.

ЛИСТИНГ 16.8. Программа для добавления одного звена кривой Гильберта

```
void hil_inc_xy(unsigned *xp, unsigned *yp, int n)
{
    int i;
    unsigned x, y, state, dx, dy, row, dochange;

    x = *xp;
    y = *yp;
    state = 0; // Инициализация
    dx = -((1 << n) - 1); // -(2**n - 1)
    dy = 0;

    for (i = n-1; i >= 0; i--) // Выполняем n раз
    {
        row = 4*state | 2*((x >> i) & 1) | (y >> i) & 1;
        dochange = (0xBDDDB >> row) & 1;
        if (dochange)
        {
            dx = ((0x16451659 >> 2*row) & 3) - 1;
            dy = ((0x51166516 >> 2*row) & 3) - 1;
        }
        state = (0x8FE65831 >> 2*row) & 3;
    }
}
```

```

)
*xp = *xp + dx;
*yр = *yp + dy;
)
    
```

Табл. 16.7 легко реализуется представленной на рис. 16.4 логикой. На этом рисунке использованы приведенные ниже обозначения.

x_i	i -й бит входного значения x
y_i	i -й бит входного значения y
X, Y	Обмененные и дополненные в соответствии со значениями S_{i+1} и C_{i+1} значения x_i и y_i
I	Флаг: если равен 1 — увеличение на 1, если 0 — уменьшение на 1
W	Флаг: если равен 1, увеличивается или уменьшается x , если 0 — увеличивается или уменьшается y
S	Если значение S равно 1, выполняется обмен x_i и y_i
C	Если значение C равно 1, выполняется дополнение x_i и y_i

Пары значений S и C определяют “состояние” в табл. 16.7: пары (S, C) , равные $(0, 0)$, $(0, 1)$, $(1, 0)$ и $(1, 1)$, обозначают соответственно состояния А, В, С и D. Выходные сигналы I_0 и W_0 указывают соответственно, должно ли выполняться уменьшение или увеличение координаты и какой именно. (В дополнение к указанной на рисунке логике требуются схема увеличения/уменьшения с мультиплексорами, указывающими, какое именно значение — x или y — должно быть увеличено или уменьшено, а также схема, которая поместит вычисленное значение назад в регистр, где хранится переменная x или y . В качестве альтернативы можно обойтись двумя схемами увеличения/уменьшения.)

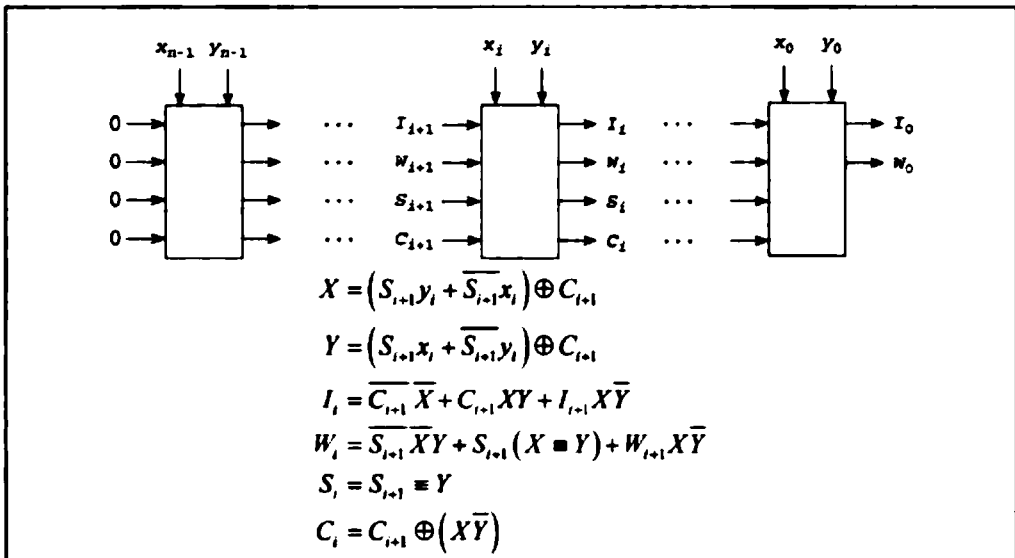


Рис. 16.4. Логическая схема для добавления одного звена кривой Гильберта

16.5. Нерекурсивный алгоритм генерации кривой Гильберта

Алгоритмы, приведенные в табл. 16.2 и 16.7, предоставляют два нерекурсивных алгоритма построения кривой Гильберта любого порядка. Каждый из них без особых сложностей может быть реализован аппаратно. Аппаратное обеспечение для реализации алгоритма на базе табл. 16.2 включает регистр для хранения величины s , увеличивающейся на каждом шаге построения кривой и преобразуемой в координаты (x, y) . Аппаратное обеспечение, реализующее алгоритм на основе табл. 16.7, не требует регистра для хранения s , но используемый при этом алгоритм более сложен.

16.6. Другие кривые, заполняющие пространство

Как упоминалось, Пеано первым открыл в 1890 году кривую, заполняющую пространство. С тех пор открыто множество других кривых, которые часто носят общее название "кривые Пеано". В 1900 году Элиахимом Гастингсом Муром (Eliakim Hastings Moore) была открыта интересная вариация кривой Гильберта. Она "циклична" в том смысле, что ее конечная точка отстоит на один шаг от начальной. На рис. 16.5 показаны кривая Пеано третьего порядка и кривая Мура четвертого порядка. Кривая Мура иррегулярна в том, что кривая первого порядка представляет собой кривую "вверх-вправо-вниз" (Γ), но эта форма не проявляется в кривых более высокого порядка. За этим малым исключением алгоритмы для работы с кривой Мура очень схожи с соответствующими алгоритмами для кривой Гильберта.

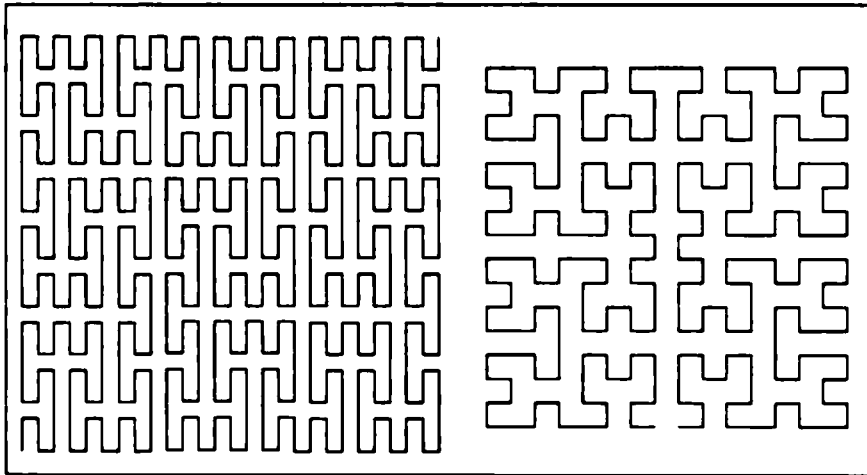
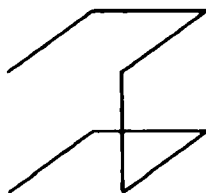


Рис. 16.5. Кривые Пеано (слева) и Мура (справа)

Кривая Гильберта обобщена для произвольных прямоугольников и на три и большее число размерностей. Базовый строительный блок трехмерной кривой Гильберта показан ниже. Он проходит по всем точкам куба $2 \times 2 \times 2$. Эти и множество других кривых, заполняющих пространство, рассматриваются в работе [99].



16.7. Применения

Кривые, заполняющие пространство, находят применение при обработке изображений — сжатии, формировании полутонного изображения и текстурном анализе [82]. Еще одно приложение состоит в повышении производительности трассировки лучей и графической визуализации. Условно сцена сканируется лучом, который перемещается в соответствии с обычным растром — слева направо через весь экран с перемещением сверху вниз. При этом, когда луч попадает на объект в базе данных моделируемой сцены, определяются его цвет и прочие свойства и изменяется вид (цвет и яркость) соответствующего пикселя на экране. (Конечно, это сверхупрощенное описание, но для наших целей его вполне достаточно.) Одна из проблем состоит в том, что зачастую такая база данных достаточно велика и данные по каждому объекту должны загружаться в память, когда сканирующий луч попадает на него. При построчном сканировании обычна ситуация, когда при проходе по строке приходится загружать в память те же объекты, которые загружались при предыдущем сканировании. Количество загрузок в память можно резко снизить, если сканирование будет обладать свойством локализации, например сканирование по квадрантам (когда переход к сканированию очередного квадранта происходит только по завершении сканирования предыдущего) зачастую позволяет существенно снизить количество загрузок в память информации об одном объекте.

Кривая Гильберта, пожалуй, обладает искомым свойством локализации. Она рекурсивно сканирует квадрант за квадрантом и не делает длинных переходов из одного квадранта в другой.

Дуглас Вурхис (Douglas Voorhies) [107] смоделировал поведение загрузки информации об объектах в память для обычного однонаправленного сканирования, кривой Пеано и кривой Гильберта. Его метод состоял в случайном размещении на экране окружностей заданного размера. Когда путь сканирования попадает в окружность, представляющую новый объект, последний загружается в память. Когда путь сканирования покидает объект, информация о нем не выгружается из памяти до тех пор, пока путь сканирования не отойдет от объекта на расстояние, равное удвоенному радиусу. Таким образом, если путь сканирования удаляется недалеко от объекта и вновь возвращается к нему, операции выгрузки и загрузки в память не происходит. Вурхис проводил эксперименты для окружностей разного радиуса на экране размером 1024×1024.

Считаем, что каждый раз, когда путь сканирования попадает в объект и покидает его, происходит одна операция загрузки в память. Тогда очевидно, что обычное построчное сканирование одной окружности диаметром D приводит к D операциям загрузки в па-

мять. Интересным результатом моделирования Вурхиса оказалось то, что для кривой Пеано количество операций загрузки в память в среднем составляет 2.7 и не зависит от диаметра окружности. Для кривой Гильберта среднее количество операций загрузки в память также не зависит от диаметра окружности, но составляет около 1.4. Таким образом, экспериментально доказано, что в смысле локализации при сканировании кривая Гильберта превосходит кривую Пеано и обе они существенно превосходят построчное сканирование в смысле сокращения операций загрузки в память.

Кривая Гильберта применяется при распределении заданий процессорам в случае, когда процессоры соединены друг с другом двух- или трехмерной прямоугольной решеткой [18]. Программное обеспечение распределения процессоров по заданиям использует линейный список процессоров, соответствующий кривой Гильберта. При планировании задания, для выполнения которого требуется несколько процессоров, последние выбираются из линейного списка, наподобие того как система распределения памяти выделяет свободную память. Выделенные таким образом процессоры имеют тенденцию к пространственной локализации, что положительно отражается на параметрах межпроцессорного взаимодействия.

Упражнения

1. Простой способ покрытия решетки $n \times n$ без слишком большого количества переходов и прохода по каждой точке только один раз, заключается в том, чтобы иметь $2n$ -битовую переменную s , которая на каждом шаге увеличивается, и формировать x из первого и всех прочих битов s , а y — из второго и всех прочих битов s . Это эквивалентно вычислению идеального внешнего упорядочения s , после чего x и y представляют собой просто левую и правую половины получившегося результата. Исследуйте свойства локальности этой кривой, набросав ее внешний вид для $n = 3$.
2. Вариантом метода из упр. 1 является первоначальное преобразование s в $\text{Gray}(s)$ (см. с. 339), с последующим получением x и y так, как описано в упр. 1. Изобразите внешний вид такой кривой для $n = 3$. Улучшились ли свойства локальности в результате внесенных изменений?
3. Как бы вы строили трехмерный аналог кривой из упр. 1?

ГЛАВА 17

ЧИСЛА С ПЛАВАЮЩЕЙ ТОЧКОЙ

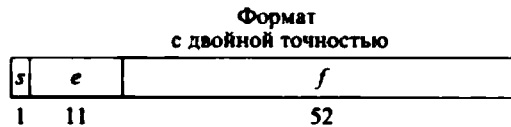
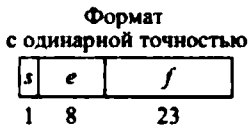
*Бог создал целые числа;
все остальные — дело рук человеческих.*
Леопольд Кронекер (Leopold Kronecker)

Выполнение операций над числами с плавающей точкой с помощью целочисленной арифметики и логических команд — предложение далеко не лучшее. Особенно это относится к правилам и форматам стандарта *IEEE Standard for Binary Floating-Point Arithmetic*, IEEE Std. 754-2008, общеизвестного как “IEEE-арифметика”. В этой арифметике имеются число NaN (not a number — не число) и бесконечности, которые представляют собой частные случаи практически для всех операций. В ней есть плюс ноль и минус ноль, которые должны быть эквивалентны при сравнении, кстати, имеющем в IEEE-арифметике четвертый возможный результат — “неупорядочено”. Старший значащий бит дроби в “нормальных” числах явно не участвует, но участвует в “субнормальных” числах. Дробные части представлены в виде обычных чисел со знаком, в то время как показатели степени — в смещенном виде. На все это, конечно, есть свои причины, но в результате получаются программы, заполненные условными переходами, которые очень сложно эффективно реализовать.

В данной главе предполагается, что читатель немного знаком со стандартом IEEE, так что о нем будет сказано предельно коротко.

17.1. Формат IEEE

Стандарт 2008 года включает три бинарных и два десятичных формата. Мы ограничимся рассмотрением чисел одинарной и двойной точности (32 и 64 бита), показанных ниже.



Знаковый бит s равен 0 для положительных чисел и 1 — для отрицательных. Смещенный показатель степени e и дробная часть f представляют собой величины, старший бит которых находится слева. Числа с плавающей точкой кодируются в соответствии со стандартом IEEE так, как показано ниже.

Одинарная точность			Двойная точность		
e	f	Значение	e	f	Значение
0	0	± 0	0	0	± 0
0	$\neq 0$	$\pm 2^{-126} (0.f)$	0	$\neq 0$	$\pm 2^{-1022} (0.f)$
От 1 до 254	–	$\pm 2^{e-127} (1.f)$	От 1 до 2046	–	$\pm 2^{e-1023} (1.f)$
255	0	$\pm \infty$	2047	0	$\pm \infty$
255	$\neq 0$	NaN	2047	$\neq 0$	NaN

В качестве примера рассмотрим кодирование числа π в формате с одинарной точностью. В соответствии с [66] в двоичном представлении

$$\pi \approx 11.0010\ 0100\ 0011\ 1111\ 0110\ 1010\ 1000\ 1000\ 1000\ 0101\ 1010\ 0011\ 0000\ 10\dots$$

Эта величина находится в диапазоне “нормальных” чисел, представленном третьей строкой приведенной выше таблицы. Старший единичный бит числа π опускается, поскольку ведущий единичный бит не хранится при кодировании нормальных чисел. Чтобы двоичная точка размещалась в правильном месте, показатель степени $e-127$ должен быть равен 1, так что $e = 128$. Таким образом, представление числа π выглядит как

$$0\ 10000000\ 10010010000111111011011$$

или в шестнадцатеричной записи как

$$40490FDB,$$

где дробь округлена до ближайшего представимого числа.

Числа, для которых $1 \leq e \leq 254$, называются “нормальными”. Они “нормализованы” в том смысле, что старший бит таких чисел равен 1 и не хранится явно. Ненулевые числа с $e = 0$ называются “субнормальными”, и их старший бит *хранится* в представлении числа в явном виде. Такая схема иногда называется схемой с “постепенной потерей значимости” (gradual underflow). В табл. 17.1 приведены некоторые крайние значения из различных диапазонов чисел с плавающей точкой. “Максимальное целое” в приведенной таблице означает наибольшее целое число, такое, что все целые числа, не превосходящие его по абсолютному значению, представимы в виде чисел с плавающей точкой точно; следующее за ним число уже должно быть округлено.

Для нормальных чисел единица в последней позиции представляет собой относительное значение в диапазоне от $1/2^{24}$ до $1/2^{23}$ (примерно от 5.96×10^{-8} до 1.19×10^{-7}) для одинарной точности и от $1/2^{53}$ до $1/2^{52}$ (примерно от 1.11×10^{-16} до 2.22×10^{-16}) для двойной точности. Максимальная “относительная ошибка” при округлении к ближайшему числу составляет половину приведенной величины.

ТАБЛИЦА 17.1. НЕКОТОРЫЕ ПРЕДЕЛЬНЫЕ ЗНАЧЕНИЯ ЧИСЕЛ С ПЛАВАЮЩЕЙ ТОЧКОЙ

	Шестнадцатеричное представление	Точное значение	Приближенное значение
Однинарная точность			
Наименьшее субнормальное	0000 0001	2^{-149}	1.401×10^{-45}
Наибольшее субнормальное	007F FFFF	$2^{-126}(1-2^{-23})$	1.175×10^{-38}
Наименьшее нормальное	0080 0000	2^{-126}	1.175×10^{-38}
1.0	3F80 0000	1	1
Максимальное целое	4B80 0000	2^{24}	1.677×10^7
Наибольшее нормальное	7F7F FFFF	$2^{128}(1-2^{-24})$	3.403×10^{38}
∞	7F80 0000	∞	∞
Двойная точность			
Наименьшее субнормальное	0...0001	2^{-1074}	4.941×10^{-324}
Наибольшее субнормальное	000F...F	$2^{-1023}(1-2^{-52})$	2.225×10^{-308}
Наименьшее нормальное	0010...0	2^{-1022}	2.225×10^{-308}
1.0	3FF0...0	1	1
Максимальное целое	4340...0	2^{53}	9.007×10^{15}
Наибольшее нормальное	7BEF...F	$2^{1024}(1-2^{-53})$	1.798×10^{308}
∞	7FF0...0	∞	∞

Диапазон целых чисел, точно представимых в формате с плавающей точкой, простирается от -2^{24} до $+2^{24}$ (т.е. от -16777216 до $+16777216$) для одинарной точности и от -2^{53} до $+2^{53}$ (от -9007199254740992 до $+9007199254740992$) для двойной точности. Конечно, некоторые целые числа вне этого диапазона также могут быть представлены точно, например большие степени 2; однако указанные диапазоны являются максимальными диапазонами, все числа которых представимы точно.

Замена деления на константу умножением может быть выполнена с полной (IEEE) точностью только для чисел, для которых соответствующие множители могут быть точно представлены в формате с плавающей точкой. Это степени 2 от -127 до 127 для одинарной точности и от -1023 до 1023 для двойной. Заметим, что числа 2^{-127} и 2^{-1023} являются субнормальными, так что их использования лучше избегать на машинах, которые недостаточно эффективно реализуют операции над субнормальными числами.

17.2. Преобразование чисел с плавающей точкой в целые и обратно

В табл. 17.2 приведены некоторые формулы для преобразования чисел в формате IEEE с плавающей точкой в целые и обратно. Эти методы краткие и быстрые, но не дают точного результата во всем диапазоне входных значений. В таблице указаны диапазоны значений, в которых эти методы работают точно. Все они дают корректные результаты для ± 0.0 и для субнормальных чисел в указанных диапазонах. Большинство из методов не дает приемлемых результатов для NaN или бесконечности. Эти формулы могут подходить для непосредственного применения в некоторых приложениях или в библиотечных подпрограммах для быстрой обработки распространенных ситуаций.

ТАБЛИЦА 17.2. ПРЕОБРАЗОВАНИЯ ЧИСЕЛ С ПЛАВАЮЩЕЙ ТОЧКОЙ

Тип	Режим	Формула	Диапазон	Примечание
double в int64, n	n	$\left(x + c_{521}\right)^d - c_{521}$	От -2^{51} до $2^{51} + 0.5$	1
double в int64, d	d	$\left(x + c_{521}\right)^d - c_{521}$	От $-2^{51} - 0.5$ до $2^{51} + 0.5$	1
double в int64, u	u	$\left(x + c_{521}\right)^d - c_{521}$	От -2^{51} до $2^{51} + 1$	1
double в int64, z	d или z	if ($x \geq 0.0$) $\left(x + c_{52}\right)^d - c_{52}$ else $c_{52} - \left(c_{52} - x\right)^d$	От -2^{52} до 2^{52}	1
double в uint64, n	n	$\left(x + c_{52}\right)^d - c_{52}$	От -0.25 до 2^{52}	1
double в uint64, d	d	$\left(x + c_{52}\right)^d - c_{52}$	От 0 до 2^{52}	1
double в uint64, u	u	$\left(x + c_{52}\right)^d - c_{52}$	От $-0.5 + \text{ulp}$ до $2^{52} + 1$	1
double в int32 или uint32, n	n	$\text{low32}\left(x + c_{521}\right)^d$	От $-2^{31} - 0.5$ до $2^{31} - 0.5 - \text{ulp}$ или от -0.5 до $2^{32} - 0.5 - \text{ulp}$	1

Продолжение табл. 17.2

Тип	Режим	Формула	Диапазон	Примечание
double в int32 или uint32, d	d	$\text{low32}\left(x^d + c_{321}\right)$	От -2^{31} до $2^{31} - \text{ulp}$ или от 0 до $2^{32} - \text{ulp}$	1
double в int32 или uint32, u	u	$\text{low32}\left(x^d + c_{321}\right)$	От $-2^{31} - 1 + \text{ulp}$ до $2^{31} - 1$ или от $-1 + \text{ulp}$ до $2^{32} - 1$	1
double в int32 или uint32, z	d или z	if ($x \geq 0.0$) $\text{low32}\left(x^d + c_{321}\right)$ else $-\text{low32}\left(c_{321} - x\right)$	От $-2^{31} - 1 + \text{ulp}$ до $2^{31} - \text{ulp}$ или от $-1 + \text{ulp}$ до $2^{32} - \text{ulp}$	1
float в int32, n	n	$\left(x^d + c_{231}\right) - c_{231}$	От -2^{22} до $2^{22} + 0.5$	
float в int32, d	d	$\left(x^d + c_{231}\right) - c_{231}$	От $-2^{22} - 0.5$ до $2^{22} + 0.5$	
float в int32, u	u	$\left(x^d + c_{231}\right) - c_{231}$	От -2^{22} до $2^{22} + 1$	
float в int32, z	d или z	if ($x \geq 0.0$) $\left(x^d + c_{23}\right) - c_{23}$ else $c_{23} - \left(c_{23} - x\right)$	От -2^{23} до 2^{23}	
float в uint32, n	n	$\left(x^d + c_{23}\right) - c_{23}$	От -0.25 до 2^{23}	
float в uint32, d	d	$\left(x^d + c_{23}\right) - c_{23}$	От 0 до 2^{23}	
float в uint32, u	u	$\left(x^d + c_{23}\right) - c_{23}$	От $-0.5 + \text{ulp}$ до $2^{23} + 1$	
Округление double к ближайшему	n	$\left(x^d + c_{321}\right)^d - c_{321}$	От -2^{51} до $2^{51} + 0.5$	1

Окончание табл. 17.2

Тип	Режим	Формула	Диапазон	Примечание
Округление неотрицательного double к ближайшему	n	$(x + c_{32})^d - c_{32}$	От -0.25 до 2^{32}	1, 3
Округление float к ближайшему	n	$(x + c_{23})^d - c_{23}$	От -2^{22} до $2^{22} + 0.5$	2
Округление отрицательного float к ближайшему	n	$(x + c_{23})^z - c_{23}$	От -0.25 до 2^{23}	2, 3
int64 в double	-	$(x + c_{321})^d - c_{321}$	От -2^{31} до 2^{31}	4
uint64 в double	-	$(x + c_{32})^d - c_{32}$	От 0 до $2^{32} - 1$	4
int32 в float	-	$(x + c_{231})^d - c_{231}$	От -2^{22} до 2^{22}	
uint32 в float	-	$(x + c_{23})^d - c_{23}$	От 0 до 2^{23}	

Константы

$$\begin{aligned}
 c_{321} &= 0x4338\ 0000\ 0000\ 0000 &= 2^{32} + 2^{31} \\
 c_{32} &= 0x4330\ 0000\ 0000\ 0000 &= 2^{32} \\
 c_{231} &= 0x4B40\ 0000 &= 2^{23} + 2^{22} \\
 c_{23} &= 0x4B00\ 0000 &= 2^{23}
 \end{aligned}$$

Примечания

1. Операции над числами с плавающей точкой должны выполняться с двойной точностью IEEE (точность 53 бита) и не более. Большинство машин на базе Intel по умолчанию в этом режиме не работают. На таких машинах необходимо установить точность (поле PC в управляющем слове FPU) равной двойной точности.
2. Операции над числами с плавающей точкой должны выполняться с одинарной точностью IEEE (точность 24 бита) и не более. Большинство машин на базе Intel по умолчанию в этом режиме не работают. На таких машинах необходимо установить точность (поле PC в управляющем слове FPU) равной одинарной точности.
3. "Неотрицательный" означает -0.0 или больше или равно 0.0 .
4. Для преобразования 32-битового знакового или беззнакового целого числа в double, требуется знаково расширить 32-битовое целое число до 64-битового и воспользоваться соответствующей формулой из приведенных.

В столбце "Тип" описан тип требуемого преобразования, включая режим округления: n для приведения к ближайшему четному, d — для округления вниз, к меньшему значению, u — для округления вверх, к большему значению, и z — для округления по направлению к нулю. В столбце "Режим" описан режим округления, в котором должна нахо-

даться машина для того, чтобы формула дала корректный результат. (На некоторых машинах, таких как IA-32, режим округления может быть указан в самой команде, а не в регистре режима.)

“double” соответствует двойной точности IEEE длиной 64 бита; “float” соответствует одинарной точности IEEE длиной 32 бита.

Обозначение “ulp” означает одну единицу в последней позиции. Например, $1.0 - \text{ulp}$ означает число в формате IEEE, ближайшее к 1.0, но меньшее 1.0, что-то наподобие 0.99999.... Обозначение “int64” описывает знаковое 64-битовое целое число (дополненное до 2), а “int32” — знаковое 32-битовое целое число. “uint64” и “uint32” имеют аналогичный беззнаковый смысл.

Функция $\text{low32}(x)$ выделяет младшие 32 бита числа x .

Операции $\overset{\circ}{+}$ и $\overset{\circ}{-}$ означают сложение чисел с плавающей точкой с двойной и одинарной точностью соответственно. Аналогично операторы $\overset{\circ}{-}$ и $\overset{\circ}{+}$ означают вычитание с двойной и одинарной точностью.

Может показаться интересным, что на большинстве машин Intel преобразование числа с плавающей точкой двойной точности в целое требует, чтобы режим точности был снижен до 53 бит, в то время как при преобразовании числа с плавающей точкой одинарной точности в целое такое снижение точности не требуется: корректный результат получается при работе машины в режиме расширенной точности (64 бита). Дело в том, что, когда к числу с двойной точностью прибавляется константа, дробная часть может быть сдвинута вправо на расстояние до 52 бит, что может привести к сдвигу единичных битов за пределы 64-битовых границ, а следовательно, к их потере. Таким образом, выполняются два округления — сначала до 64 бит, а затем до 53 бит. С другой стороны, при добавлении константы к числу с одинарной точностью максимальный сдвиг — 23 бита. При таком малом сдвиге ни один бит не сдвигается за пределы 64-битовой границы, так что выполняется только одна операция округления. Преобразования из числа с плавающей точкой одинарной точности в целое дают корректный результат на машинах Intel во всех трех режимах точности.

На машинах Intel в режимах расширенной точности преобразования чисел с плавающей точкой в int64 и uint64 выполняются без изменений в режиме точности путем применения других константных значений и одной дополнительной операции с плавающей точкой. Вычисление имеет вид

$$\left(\left(x \overset{\circ}{+} c_1 \right) \overset{\circ}{-} c_2 \right) - c_3,$$

где $\overset{\circ}{+}$ и $\overset{\circ}{-}$ означают соответственно сложение и вычитание с расширенной точностью. (Результат сложения должен оставаться в 80-битовом регистре для использования операцией вычитания с расширенной точностью.)

Для преобразования double в int64

$$\begin{aligned}c_1 &= 0x43E0\ 0300\ 0000\ 0000 = 2^{63} + 2^{52} + 2^{51} \\c_2 &= 0x43E0\ 0000\ 0000\ 0000 = 2^{63} \\c_3 &= 0x4338\ 0000\ 0000\ 0000 = 2^{52} + 2^{51}\end{aligned}$$

Для преобразования `double` в `uint64`

$$\begin{aligned}c_1 &= 0x43E0\ 0200\ 0000\ 0000 = 2^{63} + 2^{52} \\c_2 &= 0x43E0\ 0000\ 0000\ 0000 = 2^{63} \\c_3 &= 0x4330\ 0000\ 0000\ 0000 = 2^{52}\end{aligned}$$

Используя эти константы, можно получить аналогичные выражения для операций преобразования и округления, показанных в табл. 17.2, к которым относится примечание 1. Диапазоны применимости при этом близки к показанным в таблице.

Однако при округлении `double` к ближайшему, если сначала выполняется вычитание, а затем сложение, т.е.

$$\left(\left(x - c_1 \right) + c_2 \right) + c_3$$

(с использованием первого из приведенных выше наборов констант), то диапазон, для которого будет получаться корректный результат, — от $-2^{51} - 0.5$ до ∞ , но не для NaN.

17.3. Сравнение чисел с плавающей точкой с использованием целых операций

Одной из особенностей IEEE-кодирования чисел с плавающей точкой является то, что все не NaN-значения корректно упорядочены, если рассматривать их как знаковые целые числа.

Для того чтобы запрограммировать сравнение чисел с плавающей точкой с использованием целых операций, необходимо отказаться от результата сравнения “неупорядочено”. В IEEE 754 такой результат получается в том случае, если один или оба аргумента оператора сравнения представляют собой значение NaN. Описываемый далее метод трактует значения NaN как числа, превышающие бесконечное значение.

Сравнение становится существенно проще, если считать, что величина -0.0 строго меньше $+0.0$ (что не согласуется со стандартом IEEE 754). Считая описанные допущения приемлемыми и используя для сравнений с плавающей точкой обозначения $\overset{f}{<}$, $\overset{f}{\leq}$ и $\overset{f}{=}$, а символ \approx — как напоминание о том, что приведенные формулы не совсем корректно работают со значениями ± 0.0 , можно использовать приведенные ниже формулы. Это те же сравнения, что и предикат “глобального упорядочения” IEEE 754-2008.

$$\begin{aligned}a \overset{f}{=} b &\approx (a = b) \\ a \overset{f}{<} b &\approx \left(a \geq 0 \ \& \ a < b \right) \vee \left(a < 0 \ \& \ a \overset{f}{>} b \right)\end{aligned}$$

$$a \overset{f}{\leq} b \approx (a \geq 0 \& a \leq b) | (a < 0 \& a \overset{a}{\geq} b)$$

Если же значение -0.0 в обязательном порядке должно быть равно значению $+0.0$, то приведенные формулы существенно усложняются (хотя и не настолько, как могло бы показаться необходимым на первый взгляд).

$$\begin{aligned} a \overset{f}{=} b &\equiv (a = b) | (-a = a \& -b = b) \\ &\equiv (a = b) | ((a | b) = 0x80000000) \\ &\equiv (a = b) | (((a | b) \& 0x7FFFFFFF) = 0) \end{aligned}$$

$$a \overset{f}{<} b \equiv \left((a \geq 0 \& a < b) | (a < 0 \& a \overset{a}{>} b) \right) \& \left((a | b) \neq 0x80000000 \right)$$

$$a \overset{f}{\leq} b \equiv (a \geq 0 \& a \leq b) | (a < 0 \& a \overset{a}{\geq} b) | \left((a | b) = 0x80000000 \right)$$

В некоторых приложениях более эффективным может оказаться путь, состоящий в некотором преобразовании чисел и последующем сравнении с плавающей точкой с использованием единственной команды. Например, при сортировке n чисел преобразование придется выполнить для каждого числа только один раз, в то время как количество сравнений составляет как минимум $\lceil n \log_2 n \rceil$ (в смысле минимакса).

В табл. 17.3 приведены четыре таких преобразования. Для преобразований в левом столбце -0.0 эквивалентно $+0.0$, в правом столбце $-0.0 < +0.0$. В любом случае смысл сравнения при преобразовании не изменяется. Переменная n — знаковая, t — беззнаковая, а c может быть как знаковой, так и беззнаковой.

ТАБЛИЦА 17.3. Подготовка чисел с плавающей точкой к целым сравнениям

$-0.0 = +0.0$ (IEEE)	$-0.0 < +0.0$ (He-IEEE)
<pre>if (n >= 0) n = n+0x80000000; else n = ~n; Используется беззнаковое сравнение</pre>	<pre>if (n >= 0) n = n+0x80000000; else n = ~n; Используется беззнаковое сравнение</pre>
<pre>c = 0x7FFFFFFF; if (n < 0) n = (n ^ c) + 1; Используется знаковое сравнение</pre>	<pre>c = 0x7FFFFFFF; if (n < 0) n = n ^ c; Используется знаковое сравнение</pre>
<pre>c = 0x80000000; if (n < 0) n = c - n; Используется знаковое сравнение</pre>	<pre>c = 0x7FFFFFFF; if (n < 0) n = c - n; Используется знаковое сравнение</pre>
<pre>t = n >> 31; n = (n ^ (t >> 1)) - t; Используется знаковое сравнение</pre>	<pre>t = (unsigned)(n >> 30) >> 1; n = n ^ t; Используется знаковое сравнение</pre>

В последней строке приведен код без ветвлений, который может быть реализован на RISC-компьютере с базовым набором команд, с использованием четырех команд в левом столбце, и трех — в правом (эти команды должны быть выполнены для каждого из аргументов операции сравнения).

17.4. Программа приближенного вычисления обратного к квадратному корню

В начале 2000-х годов в программистских кругах наблюдалось определенное оживление, связанное с найденной удивительной подпрограммой вычисления приближенного значения, обратного к квадратному корню значения в формате числа с плавающей точкой одинарной точности IEEE. Такая подпрограмма весьма полезна в графических приложениях, например, для нормализации вектора путем умножения его компонент x , y и z на $1/\sqrt{x^2 + y^2 + z^2}$. Код этой функции на языке программирования C показан в листинге 17.1 [106].

Листинг 17.1. Приближенное обратное к квадратному корню

```
float rsqrt(float x0)
{
    union
    {
        int ix;
        float x;
    };

    x = x0; // x можно рассматривать как int
    float xhalf = 0.5f*x;
    ix = 0x5f375a82 - (ix >> 1); // Начальная оценка
    x = x*(1.5f - xhalf*x*x); // Шаг алгоритма Ньютона
    return x;
}
```

Относительная ошибка полученного результата находится в диапазоне от 0 до -0.00176 для всех нормальных чисел с единичной точностью. В случае NaN-аргумента получается корректный результат — NaN. Однако в случае, когда аргумент представляет собой $\pm\infty$, в результате получается отрицательное число — -0 . Если аргумент равен $+0$ или положительному субнормальному числу, результат отличается от того, каким он должен быть, но представляет собой большое число (большее, чем 9×10^{18}), которое может быть вполне приемлемым в некоторых приложениях.

Абсолютное значение относительной ошибки может быть снижено до диапазона ± 0.000892 , если заменить константу 1.5 в шаге алгоритма Ньютона значением 1.5008908.

Еще одно возможное улучшение заключается в замене умножения на 0.5 вычитанием 1 из показателя степени x , т.е. определение `xhalf` заменяется следующим.

```
union {int ihalf; float xhalf;};
ihalf = ix - 0x00800000;
```

Однако в этом случае функция дает неточные результаты (хотя и большие, чем 6×10^{18}) для нормальных чисел x , меньших примерно 2.34×10^{-38} , и NaN для x , являющегося субнормальным числом. Для $x = 0$ результат равен $+\infty$ (что верно).

Шаг алгоритма Ньютона представляет собой стандартное вычисление Ньютона-Рафсона для функции обратного к квадратному корню (см. приложение Б). Простое повторение этого шага уменьшает относительную ошибку до диапазона от 0 до -0.0000047 . Оптимальной константой для этого случая является `0x5F37599E`.

С другой стороны, удаление шага алгоритма Ньютона приводит к существенно более быстрой функции с относительной ошибкой в диапазоне ± 0.035 при применении константы `0x5F37642F`. Она состоит из всего лишь двух целочисленных команд, плюс код для загрузки константы (переменную `half` можно удалить).

Чтобы разобраться, почему этот метод работает, предположим, что $x = 2^n (1 + f)$, где n — несмещенный показатель степени, а f — дробная часть ($0 \leq f < 1$). Тогда

$$\frac{1}{\sqrt{x}} = 2^{-n/2} (1 + f)^{-1/2}.$$

Игнорируя дробную часть, мы видим, что должны изменить смещенный показатель степени с $127 + n$ на $127 - n/2$. Если $e = 127 + n$, то $127 - n/2 = 127 - (e - 127)/2 = 190.5 - e/2$. Следовательно, похоже, что вычисление наподобие сдвига x на одну позицию вправо и вычитание его из 190 в позиции показателя степени должно дать грубое приближение $1/\sqrt{x}$. На языке программирования C это можно выразить следующим образом.¹

```
union {int ix; float x;}; // Делаем ix и x перекрывающимися
...
0x5F000000 - (ix >> 1); // Обращаемся к x как к целому ix
```

Поиск путем анализа значения, лучшего, чем `0x5F000000`, достаточно сложен. Следует проанализировать четыре случая: случай, когда нулевой или единичный бит смещается из поля показателя степени в поле дробной части, и случаи, когда вычитание генерирует или не генерирует заем, распространяющийся на поле показателя степени. Такой анализ проведен в [79]. Здесь мы сделаем только некоторые простые наблюдения.

Используя обозначение $\text{гер}(x)$ для представления числа с плавающей точкой x в формате IEEE с одинарной точностью, мы хотим получить формулу вида

$$\text{гер}(1/\sqrt{x}) \approx k - \left(\text{гер}(x) \gg 1 \right)$$

для некоторой константы k . (Выполняется ли знаковый или беззнаковый сдвиг, значения не имеет, так как мы исключаем отрицательные значения x и -0.0 .) Как вариант можно грубо представить k в виде

¹ Это код на не совсем официальном стандарте языка программирования C, но он вполне работоспособен на большинстве компиляторов.

$$k \approx \text{гер}(1/\sqrt{x}) + \left(\text{гер}(x) \gg 1 \right)$$

и попробовать подставить несколько значений x . Результаты эксперимента представлены в табл. 17.4 (в шестнадцатеричной записи).

ТАБЛИЦА 17.4. ОПРЕДЕЛЕНИЕ КОНСТАНТЫ

Испытываемое x	$\text{гер}(x)$	$\text{гер}(1/\sqrt{x})$	k
1.0	3F80 0000	3F80 0000	54F0 0000
1.5	3FC0 0000	3F51 05EC	5F31 05EC
2.0	4000 0000	3735 04F3	5F35 04F3
2.5	4020 0000	3F21 E89B	5F31 E89B
3.0	4040 0000	3F13 CD3A	5F33 CD3A
3.5	4060 0000	3F08 D677	5F38 D677
4.0	4080 0000	3F00 0000	5F40 0000

Похоже на то, что k приближенно можно считать константой. Заметим, что для $x = 1.0$ и $x = 4.0$ получается одно и то же значение. На самом деле одни и те же значения получаются для любых x и $4x$ (если оба этих числа — нормальные). Это связано с тем, что в формуле для k при учетверении x член $\text{гер}(1/\sqrt{x})$ уменьшается на 1 в поле показателя степени, а член $\text{гер}(x) \gg 1$ увеличивается на 1 в поле показателя степени.

Что более важно, относительные ошибки для x и $4x$ одинаковы при условии, что оба значения являются нормальными числами. Чтобы убедиться, что это так, можно показать, что если аргумент x функции `sqrt` учетверяется, то результат функции уменьшается ровно в два раза, при этом не имеет значения, сколько шагов алгоритма Ньютона выполняется. Конечно, значение $1/\sqrt{x}$ также уменьшается в два раза, а следовательно, относительная ошибка не изменяется.

Это важно, поскольку означает, что если мы найдем оптимальное значение константы (в соответствии с некоторым критерием, таким как минимизация максимального абсолютного значения ошибки) для значений x от 1.0 до 4.0, то это же значение k будет оптимальным для всех нормальных чисел.

После этого перед нами встает простая задача написания программы, которая для данного значения k вычисляет истинное значение $1/\sqrt{x}$ (с помощью точной библиотечной подпрограммы) и его приближенное значение по рассмотренной формуле, скажем, для некоторых 10000 x в диапазоне от 1.0 до 4.0 и вычисляет максимальную ошибку. Оптимальное значение k можно определить, например, вручную. Это поистине удивительно, что имеется такая константа, для которой ошибка не превышает $\pm 3.5\%$ в функ-

ции, состоящей всего лишь из двух целочисленных операций, и без какого-бы то ни было табличного поиска!

17.5. Распределение ведущих цифр

Когда IBM выпустила в 1964 году свою ЭВМ System/360, все были озабочены снижением точности арифметики с плавающей точкой. Предыдущая линия ЭВМ, семейство 704–709–7090, имела 36-битовое слово. Представление числа с плавающей точкой с одинарной точностью состояло из 9-битового поля знака и показателя степени, за которым следовали 27 бит двоичного представления дроби. Старший бит дроби явно включался в это представление числа, так что точность представления чисел с плавающей точкой была равна 27 битам.

ЭВМ System/360 имела 32-битовое слово. Для хранения чисел с плавающей точкой одинарной точности IBM был выбран формат с 8-битовым полем знака и показателя степени и 24-битовым полем дробной части. Это снижало точность представления с 27 до 24 бит, что само по себе достаточно плохо, но на самом деле все обстояло еще хуже. Для того чтобы обеспечить большой диапазон показателя степени, каждая единица в 7-битовом показателе приравнивалась к 16. Это привело, по сути, к представлению дробной части в шестнадцатеричном виде, где первая цифра числа может принимать любое двоичное значение от 0001 до 1111 (от 1 до 15). Числа, старшая цифра которых равна 1, таким образом, имели точность всего лишь 21 бит (поскольку три старших бита такого числа нулевые), и таких чисел должно быть примерно $1/15$, или 6.7% от всех чисел.

Но на этом беды не заканчиваются, и на самом деле все обстоит еще хуже. Как тут же было показано теоретическим анализом и эмпирическими экспериментами, ведущие цифры чисел распределены *не* равномерно и среди чисел с плавающей точкой в шестнадцатеричной системе счисления следует ожидать, что единица будет ведущей цифрой у четверти всех чисел.

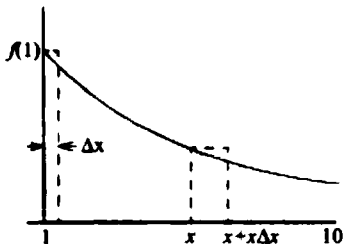
Рассмотрим распределение ведущих цифр у десятичных чисел. Предположим, что у нас есть большое множество чисел, представляющих различные физические величины — длины, объемы, массы и так далее, выраженные в “научной” записи (например, 6.022×10^{23}). Если старшая цифра у таких чисел имеет некоторое распределение, то это распределение не должно зависеть от используемых единиц измерений, например сантиметров или дюймов, фунтов или килограммов и т.д. Таким образом, если умножить все числа данного множества на некоторую константу, распределение старших цифр должно остаться тем же. Например, рассматривая умножение на 2, мы должны сделать вывод о том, что чисел с ведущей единицей (от 1.0 до 1.999..., умноженных на некоторую степень 10) должно быть столько же, сколько и чисел с ведущими цифрами 2 или 3 (от 2.0 до 3.999..., умноженных на некоторую степень 10), поскольку не должно иметь значения, измеряется ли длина в метрах или в полуметрах, масса — в килограммах или полукилограммах и т.п.

Пусть $f(x)$, $1 \leq x < 10$, — функция плотности вероятности для ведущих цифр множества чисел с физической размерностью. Эта функция обладает тем свойством, что

$$\int_a^b f(x) dx$$

пропорционально количеству чисел, старшая цифра которых находится в диапазоне от a до b . Для малых приращений Δx должно выполняться соотношение (см. приведенный ниже рисунок)

$$f(1) \cdot \Delta x = f(x) \cdot x \Delta x,$$



поскольку $f(1) \cdot \Delta x$ примерно пропорционально количеству чисел в диапазоне от 1 до $1 + \Delta x$ (игнорируя множитель степени 10), а $f(x) \cdot x \Delta x$ приблизительно пропорционально количеству чисел в диапазоне от x до $x + x \Delta x$. Так как последнее множество представляет собой первое, умноженное на x , коэффициенты пропорциональности должны быть одинаковы, а следовательно, функция плотности вероятности имеет простейший вид:

$$f(x) = f(1)/x.$$

Поскольку площадь под кривой от $x=1$ до $x=10$ должна быть равна 1 (все числа со старшими цифрами лежат в диапазоне от 1.000... до 9.999...), легко показать, что

$$f(1) = 1/\ln 10.$$

Числа со старшей цифрой в диапазоне от a до b ($1 \leq a \leq b < 10$) составляют часть общего количества всех чисел, выражаемую формулой

$$\int_a^b \frac{dx}{x \ln 10} = \frac{\ln x}{\ln 10} \Big|_a^b = \frac{\ln b/a}{\ln 10} = \log_{10} \frac{b}{a}.$$

Таким образом, в десятичной системе счисления количество чисел со старшей цифрой 1 составляет $\log_{10}(2/1) \approx 30.103\%$ от всех чисел, а чисел со старшей цифрой девять — только $\log_{10}(10/9) \approx 4.58\%$.

В шестнадцатеричной системе счисления доля чисел, старшей цифрой которых является цифра из диапазона $1 \leq a \leq b < 16$, равна $\log_{16}(b/a)$. Следовательно, в шестнадцатеричной системе счисления доля чисел со старшей цифрой, равной 1, равна $\log_{16}(2/1) = 1/\log_2 16 = 0.25$.

17.6. Таблица различных значений

В табл. 17.5 приведены IEEE-представления различных чисел, имеющие практический интерес. Приведенные значения не являются точными и округлены до ближайшего представимого числа.

ТАБЛИЦА 17.5. IEEE-ПРЕДСТАВЛЕНИЯ РАЗЛИЧНЫХ ЗНАЧЕНИЙ

Десятичное значение	Однимарная точность	Двойная точность
	(шестнадцатеричное число)	(шестнадцатеричное число)
-∞	FF80 0000	FFF0 0000 0000 0000
-2.0	C000 0000	C000 0000 0000 0000
-1.0	BF80 0000	BFF0 0000 0000 0000
-0.5	BF00 0000	BFE0 0000 0000 0000
-0.0	8000 0000	8000 0000 0000 0000
+0.0	0000 0000	0000 0000 0000 0000
Наименьшее положительное субнормальное	0000 0001	0000 0000 0000 0001
Наибольшее субнормальное	007F FFFF	000F FFFF FFFF FFFF
Наименьшее положительное нормальное	0080 0000	0010 0000 0000 0000
$\pi/180$ (0.01745...)	3C8E FA35	3F91 DF46 A252 9D39
0.1	3DCC CCCD	3FB9 9999 9999 999A
$\log_{10} 2$ (0.3010...)	3E9A 209B	3FD3 4413 509F 79FF
$1/e$ (0.3678...)	3EBC 5AB2	3FD7 8B56 362C EF38
$1/\ln 10$ (0.4342...)	3EDE 5BD9	3FDB CB7B 1526 E50E
0.5	3F00 0000	3FE0 0000 0000 0000
$\ln 2$ (0.6931...)	3F31 7218	3FE6 2E42 FEFA 39EF
$1/\sqrt{2}$ (0.7071...)	3F35 04F3	3FE6 A09E 667F 3BCD
$1/\ln 3$ (0.9102...)	3F69 0570	3FED 20AE 03BC C153
1.0	3F80 0000	3FF0 0000 0000 0000
$\ln 3$ (1.0986...)	3F8C 9F54	3FF1 93EA 7AAD 030B
$\sqrt{2}$ (1.414...)	3FB5 04F3	3FF6 A09E 667F 3BCD
$1/\ln 2$ (1.442...)	3FB8 AA3B	3FF7 1547 652B 82FE
$\sqrt{3}$ (1.732...)	3FDD B3D7	3FFB B67A E858 4CAA
2.0	4000 0000	4000 0000 0000 0000
$\ln 10$ (2.302...)	4013 5D8E	4002 6BB1 BBB5 5516
e (2.718...)	402D F854	4005 BF0A 8B14 5769
3.0	4040 0000	4008 0000 0000 0000
π (3.141...)	4049 0FDB	4009 21FB 5444 2D18
$\sqrt{10}$ (3.162...)	404A 62C2	4009 4C58 3ADA 5B53

Окончание табл. 17.5

Десятичное значение	Одинарная точность	Двойная точность
	(шестнадцатеричное число)	(шестнадцатеричное число)
$\log_2 10$ (3.321...)	4054 9A78	400A 934F 0979 A371
4.0	4080 0000	4010 0000 0000 0000
5.0	40A0 0000	4014 0000 0000 0000
6.0	40C0 0000	4018 0000 0000 0000
2π (6.283...)	40C9 0FDB	4019 21FB 5444 2D18
7.0	40E0 0000	401C 0000 0000 0000
8.0	4100 0000	4020 0000 0000 0000
9.0	4110 0000	4022 0000 0000 0000
10.0	4120 0000	4024 0000 0000 0000
11.0	4130 0000	4026 0000 0000 0000
12.0	4140 0000	4028 0000 0000 0000
13.0	4150 0000	402A 0000 0000 0000
14.0	4160 0000	402C 0000 0000 0000
15.0	4170 0000	402E 0000 0000 0000
16.0	4180 0000	4030 0000 0000 0000
$180/\pi$ (57.295...)	4265 2EE1	404C A5DC 1A63 C1F8
$2^{23} - 1$	4AFF FFFE	415F FFFF C000 0000
2^{23}	4B00 0000	4160 0000 0000 0000
$2^{24} - 1$	4B7F FFFF	416F FFFF E000 0000
2^{24}	4B80 0000	4170 0000 0000 0000
$2^{31} - 1$	4F00 0000	41DF FFFF FFC0 0000
2^{31}	4F00 0000	41E0 0000 0000 0000
$2^{32} - 1$	4F80 0000	41EF FFFF FFE0 0000
2^{32}	4F80 0000	41F0 0000 0000 0000
2^{32}	5980 0000	4330 0000 0000 0000
2^{43}	5F00 0000	43E0 0000 0000 0000
2^{44}	5F80 0000	43F0 0000 0000 0000
Наибольшее нормальное	7F7F FFFF	7FEF FFFF FFFF FFFF
∞	7F80 0000	7FF0 0000 0000 0000
"Наименьшее" SNaN	7F80 0001	7FF0 0000 0000 0001
"Наибольшее" SNaN	7FBF FFFF	7FF7 FFFF FFFF FFFF
"Наименьшее" QNaN	7FC0 0000	7FF8 0000 0000 0000
"Наибольшее" QNaN	7FFF FFFF	7FFF FFFF FFFF FFFF

Стандарт IEEE 754 не оговаривает, каким образом различаются сигнализирующие нечисла (SNaN) и несигнализирующие нечисла (QNaN). В табл. 17.5 использовано соглашение, применяющееся на платформах PowerPC, AMD 29050, Intel x86 и i860.

SPARC и семействе ARM: старший значащий бит дробной части равен 0 для SNaN и 1 — для QNaN. На некоторых (в основном старых) машинах используется тот же бит, но с обратным значением (0 — QNaN, 1 — SNaN).

Упражнения

1. Какие числа имеют одно и то же представление (не считая завершающих нулей) как в случае одинарной точности, так и в случае двойной?
2. Имеется ли программа, подобная рассмотренной программе для вычисления обратного к квадратному корню, вычисляющая приближенное значение квадратного корня?
3. Имеется ли такая программа для приближенного вычисления значения кубического корня?
4. Имеется ли программа, подобная рассмотренной программе для вычисления обратного к квадратному корню, но работающая с числами с плавающей точкой двойной точности? Считаем, что у нас 64-разрядная машина, на которой доступен целочисленный 64-битовый тип `long long`.

ГЛАВА 18

ФОРМУЛЫ ДЛЯ ПРОСТЫХ ЧИСЕЛ

18.1. Введение

Как и все молодые студенты, я в свое время был очарован простыми числами и пытался найти формулу для их поиска. Я не знал точно, какие операции должны использоваться в этой “формуле”, и даже не очень представлял себе, какую именно функцию я ищу: функцию, которая по заданному n будет возвращать n -е простое число, или по заданному простому числу будет находить следующее простое число, или будет просто давать простые числа, хотя и не все из них. Сейчас, наперекор всем этим неоднозначностям, я хочу рассмотреть, что же все же известно нам по этой проблеме? В результате вы узнаете, что а) формулы для простых чисел существуют, но б) все они не очень удовлетворительны.

Начнем наше рассмотрение с краткой истории вопроса.

В 1640 году Ферма предположил, что формула

$$F_n = 2^{2^n} + 1$$

всегда дает простые числа, и числа этого вида получили название “числа Ферма”. Предположение Ферма справедливо для n от 0 до 4, но в 1732 году Эйлер обнаружил, что

$$F_5 = 2^{2^5} + 1 = 641 \cdot 6\,700\,417.$$

(Мы уже имели дело с этими множителями, когда рассматривали вопросы деления на константу на 32-битовом компьютере.) Затем в 1880 году Ландри показал, что

$$F_6 = 2^{2^6} + 1 = 274\,177 \cdot 67\,280\,421\,310\,721.$$

В настоящее время известно, что F_n — составные числа для многих больших значений n , в частности для n от 7 до 16 включительно. В настоящее время неизвестно ни одно значение $n > 4$, для которого число Ферма было бы простым [58]. Увы, гипотеза оказалась слишком поспешной...¹

Почему же Ферма использовал двойное возведение в степень? Он знал, что если m имеет нечетный множитель, отличный от 1, то $2^m + 1$ — составное число. Если $m = ab$, где b нечетно и не равно 1, то

$$2^{ab} + 1 = (2^a + 1)(2^{a(b-1)} - 2^{a(b-2)} + 2^{a(b-3)} - \dots + 1).$$

¹ Справедливости ради следует заметить, что это единственная известная ошибочная гипотеза Ферма [110].

Зная это, Ферма заинтересовался числами $2^m + 1$, такими, где m не имеет ни одного нечетного множителя, отличного от 1, т.е. $m = 2^n$. Он испробовал несколько значений n и убедился, что они дают простые числа $2^{2^n} + 1$.

Практически все согласны, что полином вполне можно считать "формулой". В 1772 году Эйлером был открыт удивительный полином

$$f(n) = n^2 + n + 41,$$

обладающий тем свойством, что он дает простые числа для всех n от 0 до 39. Этот результат может быть расширен. Так как

$$f(-n) = n^2 - n + 41 = f(n-1),$$

$f(-n)$ дает простые числа для каждого n из диапазона от 1 до 40, т.е. $f(n)$ дает простые числа для всех n от -1 до -40 . Таким образом, полином

$$f(n-40) = (n-40)^2 + (n-40) + 41 = n^2 - 79n + 1601$$

дает простые числа для всех n от 0 до 79. (Однако эстетическая привлекательность формулы при этом теряется, так как она дает повторяющиеся и немонотонные результаты: для $n = 0, 1, \dots, 79$ формула дает числа 1601, 1523, 1447, ..., 43, 41, 41, 43, ..., 1447, 1523, 1601.)

Несмотря на всю привлекательность полинома, открытого Эйлером, известно, что не существует полиномиальной формулы $f(n)$, которая давала бы простое число для каждого n (не считая тривиального случая константного полинома типа $f(n) = 5$). Более того, справедлива следующая теорема [58].

ТЕОРЕМА. Если $f(n) = P(n, 2^n, 3^n, \dots, k^n)$ представляет собой полином с целочисленными коэффициентами от своих аргументов и $f(n) \rightarrow \infty$ при $n \rightarrow \infty$, то $f(n)$ является составным значением для бесконечного множества значений n .

Следовательно, формула типа $n^2 \cdot 2^n + 2n^3 + 2n + 5$ должна давать бесконечное количество составных чисел. С другой стороны, теорема ничего не говорит о формулах, содержащих члены наподобие 2^{2^n} , n^n или $n!$.

Формула для n -го целого простого числа может быть получена с помощью функции получения максимального целого, не превосходящего данное число ("пол"), и магического значения

$$a = 0.203\,005\,000\,700\,011\,000\,013\dots$$

Число a в десятичной системе счисления представляет собой первое простое число, записанное в первой позиции после десятичной точки, второе простое число, записанное в следующих двух позициях, третье, записанное в следующих трех позициях, и т.д. При этом для очередного простого числа всегда найдется достаточно свободного места, так как $p_n < 10^n$. Не доказывая этого, просто заметим: поскольку известно, что между n и $2n$ ($n \geq 2$) всегда есть простое число, значит, между n и $10n$ оно есть наверняка, а отсюда следует, что $p_n < 10^n$. Формула для n -го простого числа выглядит следующим образом:

$$p_n = \left[10^{\frac{n^2+n}{2}} a \right] - 10^n \left[10^{\frac{n^2-n}{2}} a \right],$$

где было использовано соотношение

$$\sum_{i=1}^n n = \frac{n^2 + n}{2}.$$

Например:

$$p_3 = [10^6 a] - 10^3 [10^3 a] = 203\,005 - 203\,000 = 5.$$

Это довольно дешевый трюк, который требует знания окончательного результата для корректного определения a . Формула такого рода может представлять интерес, только если имеется некоторый путь определения a независимо от знания всех простых чисел, но, увы, никто такого способа определения a пока что не знает.

Понятно, что такая формула может служить для многих различных последовательностей, но сейчас нас это не интересует.

18.2. Формулы Вилланса

Первая формула

К.П. Вилланс (C.P. Willans) дал следующую формулу для n -го простого числа [111]:

$$p_n = 1 + \sum_{m=1}^{\infty} \left[\sqrt[n]{n} \left(\sum_{x=1}^m \cos^2 \pi \frac{(x-1)!+1}{x} \right) \right]^{-V_n}.$$

Ее вывод начинается с теоремы Вильсона, которая гласит, что p является простым или равно 1 тогда и только тогда, когда $(p-1)! \equiv -1 \pmod{p}$. Таким образом,

$$\frac{(x-1)!+1}{x}$$

будет целым для простого x (или $x=1$) и дробным для всех составных x . Следовательно,

$$F(x) = \left[\cos^2 \pi \frac{(x-1)!+1}{x} \right] = \begin{cases} 1, & x \text{ простое или } 1, \\ 0, & x \text{ составное.} \end{cases} \quad (1)$$

Таким образом, если $\pi(m)$ обозначает² количество простых чисел, не превышающих m , то

$$\pi(m) = -1 + \sum_{x=1}^m F(x). \quad (2)$$

² Автор приносит извинения за разное по смыслу применение символа π в столь тесной близости одно к другому, но это стандартное обозначение, использование которого не должно привести к каким-либо трудностям при чтении книги.

Заметим, что $\pi(p_n) = n$, а кроме того,

$$\begin{aligned} \pi(m) < n & \text{ для } m < p_n \text{ и} \\ \pi(m) \geq n & \text{ для } m \geq p_n. \end{aligned}$$

Таким образом, количество значений m от 1 до ∞ , для которых $\pi(m) < n$, равно $p_n - 1$, т.е.

$$p_n = 1 + \sum_{m=1}^{\infty} (\pi(m) < n), \quad (3)$$

где под знаком суммы стоит выражение-предикат, значение которого может быть равно 0 или 1.

Поскольку у нас есть формула для $\pi(m)$, уравнение (3) представляет собой формулу для n -го простого числа как функцию n . Однако в ней есть две вещи, которые делают ее неприемлемой: бесконечная сумма и использование выражения-предиката, которое обычно в математике не применяется.

Доказано, что для $n \geq 1$ имеется как минимум одно простое число между n и $2n$. Следовательно, количество простых чисел, не превышающих 2^n , по меньшей мере равно n , т.е. $\pi(2^n) \geq n$. Таким образом, предикат $\pi(m) < n$ для $m \geq 2^n$ равен 0, так что верхний предел суммирования можно заменить на 2^n .

Вилланс использовал более искусную замену для выражения-предиката. Пусть

$$LT(x, y) = \left\lfloor \sqrt[y]{\frac{y}{1+x}} \right\rfloor \text{ для } x = 0, 1, 2, \dots; y = 1, 2, \dots$$

Тогда, если $x < y$, то $1 \leq y/(1+x) \leq y$, так что $1 \leq \sqrt[y]{y/(1+x)} \leq \sqrt[y]{y} < 2$. Кроме того, если $x \geq y$, то $0 < y/(1+x) < 1$, так что $0 \leq \sqrt[y]{y/(1+x)} < 1$. Применяя функцию наибольшего целого, не превосходящего заданное число ("пол"), получим

$$LT(x, y) = \begin{cases} 1, & x < y, \\ 0, & x \geq y, \end{cases}$$

так что $LT(x, y)$ представляет собой предикат $x < y$ (для x и y из указанных диапазонов).

Подставив полученное выражение в (3), получим

$$p_n = 1 + \sum_{m=1}^{\infty} LT(\pi(m), n) = 1 + \sum_{m=1}^{\infty} \left\lfloor \sqrt[n]{\frac{n}{1+\pi(m)}} \right\rfloor.$$

Дальнейшая подстановка уравнения (2) для выражения $\pi(m)$ через $F(x)$ и уравнения (1) для $F(x)$ дает нам формулу, приведенную в начале этого раздела.

Вторая формула

Вилланс дал еще одну формулу для n -го простого числа:

$$p_n = \sum_{m=1}^{\infty} m F(m) \left\lfloor 2^{-\pi(m)-n} \right\rfloor.$$

Здесь F и π — функции, использовавшиеся в первой формуле. Таким образом, $mF(m)$ равно m , если m — простое число или единица, и равно 0, если m — составное число. Третий множитель представляет собой предикат $\pi(m) = n$. Все члены суммы равны 0, кроме одного, равного n -му простому числу. Например:

$$p_n = 1 \cdot 1 \cdot 0 + 2 \cdot 1 \cdot 0 + 3 \cdot 1 \cdot 0 + 4 \cdot 0 \cdot 0 + 5 \cdot 1 \cdot 0 + 6 \cdot 0 \cdot 0 + 7 \cdot 1 \cdot 1 + 8 \cdot 0 \cdot 1 + 9 \cdot 0 \cdot 1 + 10 \cdot 0 \cdot 1 + 11 \cdot 1 \cdot 0 + \dots + 16 \cdot 0 \cdot 0 = 7.$$

Третья формула

На этом Вилланс не остановился и предложил еще одну формулу для n -го простого числа, которая не использует ни одной “неаналитической”³ функции типа “пол” или абсолютное значение величины. Он заметил, что для $x = 2, 3, \dots$, функция

$$\frac{((x-1)!)^2}{x} = \begin{cases} \text{целое} + \frac{1}{x}, & \text{если } x \text{ — простое число,} \\ \text{целое,} & \text{если } x \text{ — составное число или } 1. \end{cases}$$

Первая часть следует из того, что

$$\frac{((x-1)!)^2}{x} = \frac{((x-1)!+1) \cdot ((x-1)!-1)}{x} + \frac{1}{x}$$

и $(x-1)!+1$ делится на x в соответствии с теоремой Вильсона. Таким образом, предикат “ x — простое число” для $x \geq 2$ задается следующим образом:

$$H(x) = \frac{\sin^2 \pi \frac{((x-1)!)^2}{x}}{\sin^2 \frac{\pi}{x}}.$$

Из этого следует

$$\pi(m) = \sum_{x=2}^m H(x), \quad m = 2, 3, \dots$$

Эту формулу нельзя преобразовать в формулу для p_n методами, использованными в первых двух формулах Вилланса, так как в них применяется функция “пол”. Вместо этого Вилланс предложил следующую формулу⁴ для предиката $x < y$ при $x, y \geq 1$:

$$LT(x, y) = \sin \left(\frac{\pi}{2} \cdot 2^y \right),$$

³ Это термин автора книги, а не Вилланса.

⁴ Мы немного упростили эту формулу.

где

$$e = \prod_{i=0}^{x-1} (x-i).$$

Таким образом, если $x < y$, то $e = x(x-1)\dots(0)(-1)\dots(x-(y-1)) = 0$, так что $LT(x, y) = \sin(\pi/2) = 1$. Если $x \geq y$, произведение не включает 0, так что $e \geq 1$ и $LT(x, y) = \sin((\pi/2) \cdot (\text{четное число})) = 0$.

Наконец, как и в первой формуле Вилланса,

$$p_n = 2 + \sum_{m=2}^{\infty} LT(\pi(m), n).$$

Объединив все вместе, получим следующую “жуть”:

$$p_n = 2 + \sum_{m=2}^{\infty} \sin \left(\frac{\pi}{2} \cdot 2 \left(\prod_{i=0}^{m-1} \left(\sum_{j=0}^{m-i} \frac{((x-1)!)^2}{m^{\frac{2}{j}}} \right) \right) \right).$$

Четвертая формула

После этого Вилланс приводит еще одну формулу, которая выражает простое число p_{n+1} через p_n :

$$p_{n+1} = 1 + p_n + \sum_{i=1}^{2p_n} \prod_{j=1}^i f(p_n + j),$$

где $f(x)$ — предикат “ x — составное число”, $x \geq 2$, т.е.

$$f(x) = \left[\cos^2 \pi \frac{((x-1)!)^2}{x} \right].$$

В качестве альтернативы можно использовать соотношение $f(x) = 1 - H(x)$, что позволит избежать применения функции “пол”.

Рассмотрим пример работы этой формулы для $p_n = 7$. Тогда

$$\begin{aligned} p_{n+1} &= 1 + 7 + f(8) + f(8)f(9) + f(8)f(9)f(10) + \\ &\quad + f(8)f(9)f(10)f(11) + \dots + f(8)f(9)\dots f(14) \\ &= 1 + 7 + 1 + 1 \cdot 1 + 1 \cdot 1 \cdot 1 + 1 \cdot 1 \cdot 1 \cdot 0 + \dots + 1 \cdot 1 \cdot 1 \cdot 0 \cdot 1 \cdot 0 \cdot 1 = 11. \end{aligned}$$

18.3. Формула Вормелла

К.П. Вормелл (С.Р. Wormell) [113] усовершенствовал формулу Вилланса, убрав из нее как тригонометрические функции, так и функцию “пол”. Вычисления по формуле

Вормелла в принципе могут быть проведены с помощью простой компьютерной программы с использованием только целочисленной арифметики. Ее вывод не использует теорему Вильсона. Вормелл начал с того, что для $x \geq 2$

$$B(x) = \prod_{a=2}^x \prod_{b=2}^x (x-ab)^2 = \begin{cases} \text{положительное целое, если } x \text{ — простое число,} \\ 0, \text{ если } x \text{ — составное число.} \end{cases}$$

Таким образом, количество простых чисел, не превосходящих m , задается соотношением

$$\pi(m) = \sum_{r=2}^m \frac{1 + (-1)^{2^{\pi(r)}}}{2},$$

потому что под знаком суммы находится предикат “ x — простое число”.

Заметим, что для $n \geq 1, a \geq 0$

$$\prod_{r=1}^n (1-r+a)^2 = \begin{cases} 0, \text{ если } a < n, \\ \text{положительное целое, если } a \geq n. \end{cases}$$

Повторя описанный выше прием, находим, что предикат $a < n$ можно выразить следующим образом:

$$(a < n) = \frac{1 - (-1)^{2^{\sum_{r=1}^n (1-r+a)^2}}}{2}.$$

Поскольку

$$p_n = 2 + \sum_{m=2}^n (\pi(m) < n),$$

после вынесения постоянных множителей за знак суммы получаем

$$p_n = \frac{3}{2} + 2^{n-1} - \frac{1}{2} \sum_{m=2}^n (-1)^2 \left[\prod_{r=1}^n \left(\frac{(n-r+1)^2}{2} - \frac{1}{2} \sum_{s=1}^n (1-r+s)^2 \right) \right].$$

Как и было обещано, формула Вормелла не использует тригонометрических функций. Однако, как указал Вормелл, они появятся вновь, если заменить степени -1 , используя соотношение $(-1)^n = \cos \pi n$.

18.4. Формулы для других сложных функций

Еще раз внимательно посмотрим на то, что сделали Вилланс и Вормелл. Ниже постулированы правила, определяющие, что именно мы подразумеваем под классом функций, которые могут быть представлены “формулами” и которые мы назовем “формульными функциями”. Здесь \bar{x} означает сокращение x_1, x_2, \dots, x_n для любого $n \geq 1$. Область значений представляет собой целые числа $\dots -2, -1, 0, 1, 2, \dots$

1. Константы ... -1, 0, 1, ... являются формульными функциями.
2. Функции проекции $f(\bar{x}) = x_i$, для $1 \leq i \leq n$ являются формульными.
3. $x + y$, $x - y$ и xy являются формульными функциями, если таковыми являются x и y .
4. Класс формульных функций замкнут по отношению к суперпозиции функций (подстановке), т.е. $f(g_1(\bar{x}), g_2(\bar{x}), \dots, g_m(\bar{x}))$ является формульной функцией, если f и g_i являются таковыми для $i = 1, \dots, m$.
5. Ограниченные суммы и произведения

$$\sum_{i=a(\bar{x})}^{b(\bar{x})} f(i, \bar{x}) \quad \text{и} \quad \prod_{i=a(\bar{x})}^{b(\bar{x})} f(i, \bar{x})$$

представляют собой формульные функции, если a , b и f являются таковыми и $a(\bar{x}) \leq b(\bar{x})$.

Требование ограниченности сумм и произведений сохраняет вычислительный характер формул, т.е. формулы могут быть вычислены подстановкой значений аргументов и проведением конечного числа операций. Причина наличия штриха у символов Σ и Π поясняется далее в этой главе.

При формировании новых формульных функций с использованием суперпозиции в случае необходимости применяются круглые скобки в соответствии со стандартными соглашениями об их применении.

Заметим, что деление в приведенный выше список не включено. Однако даже при этом приведенный список нельзя считать минимальным. Конечно, интересно найти минимально необходимый список правил, но не будем останавливаться на этом вопросе в нашей книге.

Определение "формульной функции" близко к определению "элементарной функции", приведенному в [19]. Однако область значений, использованная в [19], представляет собой неотрицательные целые числа (как и в обычной теории рекурсивных функций). Кроме того, в [19] требуется, чтобы границы итерационных сумм и произведений были 0 и $x-1$ (где x — переменная) и допускали пустые диапазоны (в этом случае сумма считается равной 0, а произведение — 1).

Далее покажем, что класс формульных функций существенно шире и включает большинство функций, встречающихся в математике. Тем не менее он не включает в себя все функции, которые легко определяются и имеют элементарный характер.

По сравнению с теорией рекурсивных функций наши выводы несколько усложнены, так как здесь переменные могут иметь отрицательные значения. Однако то, что некоторые значения могут быть отрицательными, легко исправить, используя возведение соответствующего выражения в квадрат. Еще одним усложнением является настойчивое требование, чтобы диапазоны сумм и произведений были непустыми.

В нашем рассмотрении под "предикатом" понимается функция, которая возвращает значение, равное 0 или 1, в то время как в теории рекурсивных функций предикат представляет собой функцию, которая возвращает значение "истинно/ложно", и каждый пре-

дикат имеет связанную с ним “характеристическую функцию”, которая возвращает значения 0/1. Мы ассоциируем значение 1 с истиной, а 0 с ложью, как обычно делается в языках программирования и в вычислительных машинах в целом (в плане работы команд *и* и *или*). Однако в теории логических и рекурсивных функций встречается и противоположная ассоциация логических и числовых значений.

Приведем ряд примеров формульных функций.

1. $a^2 = aa, a^3 = aaa$ и т.д.
2. Предикат $a = b$:

$$(a = b) = \prod_{j=0}^{(a-b)^2} (1-j).$$

3. $(a \neq b) = 1 - (a = b)$.
4. Предикат $a \geq b$:

$$(a \geq b) = \sum_{i=0}^{(a-b)^2} ((a-b) = i) = \sum_{i=0}^{(a-b)^2} \prod_{j=0}^{(a-b)-i} (1-j).$$

Теперь можно объяснить, почему не используется соглашение, в соответствии с которым итеративная сумма/произведение принимает при пустом диапазоне суммирования/произведения значения 0/1. Если поступить подобным образом, то можно получить такие “жульничества”, как

$$(a = b) = \sum_{i=0}^{-(a-b)^2} 1 \quad \text{и} \quad (a \geq b) = \prod_{i=0}^{b-1} 0.$$

Предикаты сравнений являются ключом ко всему последующему материалу, поэтому не желательно, чтобы они были основаны на искусственных построениях такого рода.

5. $(a > b) = (a \geq b + 1)$.
6. $(a \leq b) = (b \geq a)$.
7. $(a < b) = (b > a)$.
8. $|a| = (2(a \geq 0) - 1)a$.
9. $\max(a, b) = (a \geq b)(a - b) + b$.
10. $\min(a, b) = (a \geq b)(b - a) + a$.

Теперь можно откорректировать итерационные суммы и произведения таким образом, чтобы они давали корректный результат и при пустом диапазоне суммирования или произведения.

$$11. \sum_{i=a(\bar{x})}^{b(\bar{x})} f(i, \bar{x}) = (b(\bar{x}) \geq a(\bar{x})) \sum_{i=a(\bar{x})}^{b(\bar{x})} f(i, \bar{x}).$$

$$12. \prod_{i=a(\bar{x})}^{b(\bar{x})} f(i, \bar{x}) = 1 + \left(b(\bar{x}) \geq a(\bar{x}) \right) \left(-1 + \prod_{i=a(\bar{x})}^{\max(a(\bar{x}), b(\bar{x}))} f(i, \bar{x}) \right).$$

Начиная с этого момента, символы суммы и произведения используются без штриха. Таким образом, все дальнейшие определения функций корректны для любых значений аргументов.

$$13. n! = \prod_{i=1}^n i. \text{ Это определение дает нам } n! = 1 \text{ для } n \leq 0.$$

В следующих формулах P и Q обозначают предикаты.

$$14. \neg P(\bar{x}) = 1 - P(\bar{x}).$$

$$15. P(\bar{x}) \& Q(\bar{x}) = P(\bar{x})Q(\bar{x}).$$

$$16. P(\bar{x}) | Q(\bar{x}) = 1 - (1 - P(\bar{x}))(1 - Q(\bar{x})).$$

$$17. P(\bar{x}) \oplus Q(\bar{x}) = (P(\bar{x}) - Q(\bar{x}))^2.$$

$$18. \text{if } P(\bar{x}) \text{ then } f(\bar{y}) \text{ else } g(\bar{z}) = P(\bar{x})f(\bar{y}) + (1 - P(\bar{x}))g(\bar{z}).$$

$$19. a^n = \text{if } n \geq 0 \text{ then } \prod_{i=1}^n a \text{ else } 0.$$

Эта формула дает результат 0 для $n < 0$ и 1 для 0^0 , что в ряде случаев может оказаться некорректным решением.

$$20. (m \leq \forall x \leq n) P(x, \bar{y}) = \prod_{x=m}^n P(x, \bar{y}).$$

$$21. (m \leq \exists x \leq n) P(x, \bar{y}) = 1 - \prod_{x=m}^n (1 - P(x, \bar{y})).$$

Пустое \forall истинно, пустое \exists ложно.

$$22. (m \leq \min x \leq n) P(x, \bar{y}) = m + \sum_{i=m}^n \prod_{j=i}^n (1 - P(j, \bar{y})).$$

Значение этого выражения есть наименьшее x в диапазоне от m до n , такое, что предикат для него оказывается истинным, либо m в случае пустого диапазона, либо $n+1$, если предикат ложен для всего (непустого) диапазона. Эта операция называется ограниченной минимизацией и является весьма мощным инструментом при разработке новых формульных функций. Вычисление такой минимизации, представляющей собой тип функциональной инверсии, предложено Гудштейном (Goodstein) [41].

$$23. \lfloor \sqrt{n} \rfloor = (0 \leq \min k \leq |n|) \left((k+1)^2 > n \right).$$

Эта функция представляет собой "целочисленный квадратный корень", который для обобщенности равен 0 при $n < 0$.

24. $d | n = (-|n| \leq \exists q \leq |n|)(n = qd)$.

Это предикат " d является делителем n ", в соответствии с которым $0 | 0$, но $\neg(0 | n)$ при $n \neq 0$.

25. $n + d = \text{if } n \geq 0 \text{ then } (-n \leq \min q \leq n)(0 \leq \exists r \leq |d| - 1)(n = qd + r)$
 else $(n \leq \min q \leq -n)(-|d| + 1 \leq \exists r \leq 0)(n = qd + r)$.

Это формула для отсекающего целочисленного деления. Для $d = 0$ формула произвольно дает результат $|n| + 1$.

26. $\text{rem}(n, d) = n - (n + d)d$.

Это обычная функция получения остатка при делении. Если $\text{rem}(n, d)$ имеет ненулевое значение, его знак равен знаку числителя n . Если $d = 0$, остаток равен n .

27. $\text{isprime}(n) = n \geq 2 \ \& \ \neg(2 \leq \exists d \leq |n| - 1)(d | n)$.

28. $\pi(n) = \sum_{i=1}^n \text{isprime}(i)$ (количество простых чисел, не превосходящих n).

29. $p_n = (1 \leq \min k \leq 2^n)(\pi(k) = n)$.

30. $\text{exponent}(p, n) = (0 \leq \min x \leq |n|) \neg(p^{x+1} | n)$.

Это формула показателя степени простого множителя p данного числа $n \geq 1$.

31. Для $n \geq 0$:

$$2^n = \prod_{i=1}^n 2, \quad 2^{2^n} = \prod_{i=1}^{2^n} 2, \quad 2^{2^{2^n}} = \prod_{i=1}^{2^{2^n}} 2 \text{ и т.д.}$$

32. n -я цифра десятичного представления $\sqrt{2}$:

$$\text{rem}\left(\left\lfloor \sqrt{2 \cdot 10^{2n}} \right\rfloor, 10\right).$$

Таким образом, класс формульных функций весьма большой. Тем не менее он ограничен как минимум следующей теоремой.

ТЕОРЕМА. Если f — формульная функция, то существует константа k , такая, что

$$f(\bar{x}) \leq 2^{2^{\dots^{2^k}}},$$

где количество двоек равно k .

Эту теорему можно доказать, показав, что каждое из правил 1–5 на с. 426 сохраняет справедливость теоремы. Например, если $f(\bar{x}) = c$ (правило 1), то для некоторого h

$$f(\bar{x}) \leq 2^{2^{h+1}} h,$$

где имеется h двоек. Таким образом,

$$f(\bar{x}) \leq 2^{2^{2^{h+1} + h}} \} h + 2,$$

потому что $\max(|x_1|, \dots, |x_n|) \geq 0$.

Для $f(\bar{x}) = x_i$ (правило 2) $f(\bar{x}) \leq \max(|x_1|, \dots, |x_n|)$, так что теорема выполняется с $k = 0$.

Рассмотрим правило 3. Пусть

$$f(\bar{x}) \leq 2^{2^{2^{h+1} + k_1}} \} k_1 \text{ и } g(\bar{x}) \leq 2^{2^{2^{h+1} + k_2}} \} k_2.$$

Тогда очевидно, что

$$\begin{aligned} f(\bar{x}) \pm g(\bar{x}) &\leq 2 \cdot 2^{2^{2^{h+1} + \max(k_1, k_2)}} \} \max(k_1, k_2) \\ &\leq 2^{2^{2^{h+1} + \max(k_1, k_2) + 1}} \} \max(k_1, k_2) + 1. \end{aligned}$$

Аналогично можно показать, что теорема выполняется для $f(x, y) = xy$.

Доказательства того, что правила 4 и 5 сохраняют справедливость теоремы, утомительны и сложны, так что здесь они не приводятся.

Из теоремы следует, что функция

$$f(x) = 2^{2^{x^2}} \} x \tag{4}$$

не является формульной, поскольку всегда существует достаточно большое x , для которого значение из (4) превышает значение выражения с любым фиксированным количеством двоек k .

Для тех, кто интересуется теорией рекурсивных функций, укажем, что (4) представляет собой примитивную рекурсию. Кроме того, можно легко показать непосредственно из определения примитивной рекурсии, что формульные функции представляют собой примитивную рекурсию. Таким образом, класс формульных функций представляет собой истинное подмножество примитивных рекурсивных функций. Заинтересовавшегося этим вопросом читателю можно посоветовать обратиться к [19].

В данной главе представлена не только формула для n -го простого числа, состоящая из элементарных функций, но и многие другие функции, встречающиеся в математике. Кроме того, наши "формульные функции" не используют тригонометрические функции, функцию "пол", абсолютное значение, степени значения -1 и даже деление. Они очень ловко используют тот факт, что произведение большого количества чисел равно 0, если хотя бы одно из них равно 0 (см. формулу для предиката $a = b$).

Впрочем, после знакомства с приведенными функциями для простых чисел они перестают казаться столь интересными, и вопрос об “интересных” формулах для простых чисел остается открытым. Например, в [96] приведена удивительная теорема У.Г. Миллса (W.H. Mills) о том, что существует θ , такое, что выражение

$$\lfloor \theta^{3^n} \rfloor$$

дает простые числа для всех $n \geq 1$. На самом деле имеется бесконечное число таких значений (например, 1.3063778838+ и 1.4537508625483+). Кроме того, нет ничего выдающегося и в числе 3: теорема остается верна, если заменить тройку любым большим ее целым числом (разумеется, для других значений θ). Более того, тройку можно заменить даже двойкой, если справедливо недоказанное до сих пор утверждение о том, что между n^2 и $(n+1)^2$ всегда имеется по крайней мере одно простое число. Более того... Впрочем, заинтересовавшийся читатель сам найдет в [96] и [28] множество формул такого типа.

Упражнения

1. Докажите, что для любого неконстантного полинома $f(x)$ с целыми коэффициентами $\lfloor f(x) \rfloor$ является составным для бесконечного количества значений x . *Указание:* если $f(x_0) = k$, рассмотрите $f(x_0 + rk)$, где r — целое число, большее 1.

2. Докажите теорему Вильсона (Wilson): целое число $p > 1$ простое тогда и только тогда, когда

$$(p-1)! \equiv -1 \pmod{p}.$$

Указание: чтобы показать, что если p — простое, то $(p-1)! \equiv -1 \pmod{p}$, сгруппируйте члены факториала в пары (a, b) , такие, что $ab \equiv 1 \pmod{p}$. Воспользуйтесь теоремой М1 из раздела 10.16 на с. 267.

3. Покажите, что если n — целое составное число, большее 4, то

$$(n-1)! \equiv 0 \pmod{n}.$$

4. Вычислите оценку значения θ , удовлетворяющего теореме Миллса, а в процессе вычислений дайте неформальное доказательство этой теоремы. Считайте, что для $n > 1$ между n^3 и $(n+1)^3$ имеется простое число. (Утверждение опирается на гипотезу Римана, хотя для больших n оно доказано независимо от нее.)

5. Рассмотрим множество чисел вида $a + b\sqrt{-5}$, где a и b — целые числа. Докажите, что 2 и 3 в этом множестве являются простыми, т.е. что их нельзя разложить на множители из данного множества, если только один из множителей не равен ± 1 (“единице”). Найдите в этом множестве число, имеющее два разных разложения на произведения простых чисел. (“Фундаментальная теорема арифметики” гласит, что разложение на простые сомножители является единственным, за исключением единиц и степеней множителей. Единственность для данного множества чисел при использовании сложений и умножений комплексных чисел не выполняется. Это пример “кольца”).

ОТВЕТЫ К УПРАЖНЕНИЯМ

Глава 1. Введение

1. Вот как выглядит точный перевод.

```
e1;  
while (e2) {  
    оператор  
    e3;  
}
```

Если условие e_2 в цикле `for` отсутствует, вместо него в приведенном коде используется константа `1` (что приводит к бесконечному циклу, если только *оператор* не содержит инструкцию, завершающую работу цикла).

Выражение цикла `for` с помощью цикла `do` достаточно запутанное, поскольку тело цикла `do` всегда выполняется по крайней мере один раз, в то время как тело цикла `for` может не выполняться ни разу — в зависимости от того, что собой представляют e_1 и e_2 . Тем не менее цикл `for` можно записать следующим образом.

```
e1;  
if (e2) {  
    do { оператор; e3; } while (e2);  
}
```

Здесь также при отсутствии условия e_2 в цикле `for` в приведенном коде вместо него используется константа `1`.

2. Если вы запишете код как

```
for (i = 0; i <= 0xFFFFFFFF; i++) {...},
```

то получите бесконечный цикл. Правильным решением является код

```
i = 0xFFFFFFFF;  
do {i = i + 1; ...} while (i < 0xFFFFFFFF);
```

3. В тексте упоминается *умножение*, которое в случае умножения $32 \times 32 \Rightarrow 64$ требует двух выходных регистров.

Там же упоминается *деление*. Обычная реализация этой команды генерирует как частное, так и остаток, так что во многих программах, в которых требуются оба значения, так можно сократить время выполнения.

В действительности наиболее естественная операция машинного деления получает двойное слово делимого и одно слово делителя, и после выполнения возвращает частное и остаток. При этом используются три регистра-источника и два целевых регистра.

Для эффективной работы с полями регистра многие машины предоставляют команды *извлечения* (*extract*) и *вставки* (*insert*). Обобщенная форма команды *извлечения* требует трех регистров-источников и одного целевого. Регистры-источ-

ники — это регистры, которые содержат извлекаемое поле, номер начального бита и номер конечного бита. Результат выравнивается вправо, расширяется (либо знаково, либо нулем) и помещается в целевой регистр. Некоторые машины предоставляют эту команду только в форме, в которой длина поля является непосредственным значением, что представляет собой разумный компромисс в силу распространенности этого случая.

Обобщенная команда *вставки* считывает четыре исходных регистра и записывает один целевой. В обычной реализации исходными являются регистр, содержащий вставляемые биты (они берутся из младших битов регистра-источника), стартовая позиция вставки в целевой регистр и длина вставляемой области. В дополнение к чтению этих трех регистров команда должна считать целевой регистр, скомбинировать его со вставляемыми битами и записать результат в целевой регистр. Как и в случае команды *извлечения*, длина поля может представлять собой непосредственное значение, и в этом случае команда выполняет три чтения и одну запись. Некоторые машины предоставляют семейство команд *выбора* (*select*).

SELcc RT, RA, RB, RC

При выполнении команды тестируется регистр RC, и если он удовлетворяет условию, задаваемому кодом операции (показанному выше как cc, и может быть EQ, GT, GE и др.), то выбирается регистр RA; в противном случае выбирается регистр RB. Выбранный регистр копируется в целевой регистр RT.

Хотя команда *выбор бита*, или *мультиплексор*, и не так распространена, она также может пригодиться на практике.

MUX RT, RA, RB, RC

Здесь RC содержит маску. Когда бит маски равен 1, выбирается соответствующий бит регистра RA, а когда он равен 0, выбирается соответствующий бит регистра RB, т.е. выполняется следующая операция.

RT <-- RA & RC | RB & ~RC

Двойной сдвиг вправо/влево: иногда полезной оказывается такая команда.

SHLD RT, RA, RB, RC

Она соединяет регистры RA и RB, рассматривая их как регистр двойной длины, и сдвигает его влево (или вправо) на количество позиций, задаваемых регистром RC. Регистр RT получает часть результата, в которой имеются биты и из RA, и из RB. Такие команды полезны в арифметике больших чисел и в некоторых иных ситуациях.

В обработке сигналов и других приложениях полезно иметь команду, вычисляющую $A * B + C$ (как для целочисленных данных, так и для чисел с плавающей точкой).

Конечно, имеются также команды *множественной загрузки* и *множественного сохранения*, которые должны читать и записывать множество регистров. Хотя эти команды имеются на многих RISC-машинах, обычно как RISC-команды они не рассматриваются.

Глава 2. Основы

1. (Вывод Дэвида де Клота.) Ясно, что тело цикла `while` выполняется столько раз, сколько в x имеется завершающих нулевых битов. k единичных битов разбивают n -битовое слово на $k+1$ сегмент, каждый из которых содержит 0 или более нулевых битов. Количество нулевых битов в каждом слове равно $n-k$. Если количество слов равно N ($N = \binom{n}{k}$, но здесь это неважно), то общее количество нулевых битов во всех словах составляет $N(n-k)$. Из соображений симметрии количество нулевых битов в каждом сегменте, просуммированное по всем словам, составляет $N(n-k)/(k+1)$, так что среднее количество нулевых битов в сегменте равно $(n-k)/(k+1)$, и это применимо и к последнему сегменту завершающих нулевых битов.

Если взять в качестве примера $n = 32$ и $k = 3$, то цикл `while` в среднем выполняется 7.25 раз. На многих машинах цикл `while` можно реализовать всего лишь тремя командами (*и, сдвиг вправо* и *условное ветвление*), которым требуется для выполнения всего лишь четыре такта процессора. При этих условиях цикл `while` выполняется в среднем за $4 \cdot 7.25 = 29$ тактов. Это меньше, чем время деления на большинстве 32-разрядных машин, что делает алгоритм де Клота более быстрым, чем алгоритм Госпера. Для больших значений k алгоритм де Клота остается более предпочтительным.

2. Команда `и с 1` делает величину сдвига независимой от всех битов x , за исключением крайнего справа. Таким образом, зная только крайний справа бит величины сдвига, можно выяснить, чему равен результат, — x или $x \ll 1$. Поскольку x и $x \ll 1$ вычислимы справа налево, то таковым же является и выбор одного из этих значений, основанный на значении крайнего справа бита. Кстати, функция $x \ll (x \& 2)$ не является вычислимой справа налево, в отличие от функции $(x \& -2) \ll (x \& 2)$.

Другим примером является функция x^n , где принимается, что $x^0 = 1$. Эта функция не является вычислимой справа налево, так как если x четно, то крайний справа бит результата зависит от выполнения условия $x = 0$ и, таким образом, является функцией от битов, находящихся слева от крайней справа позиции. Но если нам известно *априори*, что переменная n равна либо 0, либо 1, то значение x^n вычислимо справа налево. Аналогично $x^{n \& 1}$ вычислимо справа налево, например, исходя из того, что

$$x^{n \& 1} = (x \& -(n \& 1)) + 1 - (n \& 1) = \begin{cases} 1, & n \text{ четно} \\ x, & n \text{ нечетно} \end{cases}$$

Заметим, что функция x^n похожа на функцию сдвига влево в том, что x^n вычислима справа налево для любого конкретного значения n , либо если n представляет собой переменную, ограниченную значениями 0 и 1, но не в случае, когда n — переменная, принимающая произвольное значение.

3. Практически очевидная формула сложения приведена на с. 36, п. (ж).

$$x + y = (x \oplus y) + 2(x \& y)$$

Деление каждой части на 2 дает нам формулу Дитца. Сложение в формуле Дитца не может привести к переполнению, поскольку среднее двух представимых чисел является представимым.

Заметим, что если начать с п. (и) на с. 36, то можно получить формулу, приведенную в разделе для вычисления потолка среднего двух беззнаковых чисел

$$\left\lceil \frac{x+y}{2} \right\rceil = (x|y) - \left((x \oplus y) \gg 1 \right)$$

4. Метод заключается в вычислении среднего с применением пола a и b , а также c и d с использованием формулы Дитца. Затем вычисляется среднее с применением пола x и y и выполняется приведенная ниже коррекция.

$$x = (a \& b) + \left((a \oplus b) \gg 1 \right)$$

$$y = (c \& d) + \left((c \oplus d) \gg 1 \right)$$

$$r = (x \& y) + \left((x \oplus y) \gg 1 \right)$$

$$r = r + ((a \oplus b) \& (c \oplus d) \& (x \oplus y) \& 1)$$

Шаг коррекции в действительности содержит четыре операции, а не семь, поскольку члены с *исключаящим* или были вычислены ранее. Этот метод работает по следующей причине: вычисленное значение x может быть меньше истинного (действительного значения) среднего на $1/2$, и эта ошибка получается, когда a нечетно, а b четно или наоборот. Эта ошибка составляет $1/4$ после усреднения x и y . Если имеется единственная ошибка округления, то первое значение, вычисляемое для получения r , корректно, так как в этом случае истинное значение равно целому числу плюс $1/4$, а так как нам нужен пол среднего, мы в любом случае отбрасываем эту $1/4$. Аналогично округление отбрасыванием при вычислении y может сделать вычисленное значение среднего меньшим истинного на $1/4$. Таким образом, первоначально вычисленное значение r может оказаться меньшим, чем истинное значение среднего x и y , на $1/2$. Эти ошибки накапливаются. Если сумма ошибок меньше 1, их можно игнорировать, так как мы все равно отбрасываем дробную часть истинного среднего. Но если произойдут все три ошибки, то их сумма достигнет $1/4 + 1/4 + 1/2 = 1$, и требуется коррекция путем прибавления 1 к r . Это и делает последняя строка: если одно из значений a и b нечетно и одно из значений c и d нечетно, и одно из значений x и y нечетно, то мы прибавляем к r единицу.

5. Упрощаемое выражение для $x \dot{\leq} y$ таково.

$$(\neg x | y) \& ((x \oplus y) | \neg(y - x))$$

В логических операциях в этом выражении имеют значение только 31-е биты x и y . Поскольку $y_{31} = 0$, выражение тут же упрощается.

$$\neg x \& (x | \neg(y - x))$$

Применение закона дистрибутивности дает следующее.

$$\neg x \& \neg(y - x)$$

После этого применение правила де Моргана упрощает выражение до трех элементарных команд.

$$\neg(x | (y - x))$$

(Удаление оператора дополнения дает решение из двух команд для предиката $x \stackrel{>}{\neq} y$.)

Если y представляет собой константу, можно воспользоваться тождеством $\neg u = -1 - u$, чтобы переписать выражение, полученное из закона дистрибутивности.

$$\neg x \& (x - (y + 1))$$

Иначе говоря, можно получить три команды, поскольку прибавление 1 к y можно выполнить перед вычислением выражения. Эта форма предпочтительна в случае, когда y представляет собой малую константу, потому что при этом можно использовать команду *добавления непосредственного значения*. (Задача предложена Джорджем Тиммсом (George Timms).)

6. Чтобы получить перенос при втором сложении, перенос из первого сложения должен быть равен 1, а все младшие 32 бита первой суммы должны быть единичными, т.е. первая сумма должна быть равна как минимум $2^{33} - 1$. Но каждый операнд не превышает $2^{32} - 1$, так что их сумма не может превышать $2^{33} - 2$.
7. Для простоты рассмотрим 4-битовую машину. Пусть x и y представляют собой целые значения 4-битовых величин, рассматриваемые как беззнаковые целые. Пусть $f(x, y)$ обозначает целочисленный результат применения обычного бинарного сложения с переносом значений x и y , с 4-битовым сумматором и 4-битовым результатом. Тогда

$$f(x, y) = \text{mod} \left(x + y + \left\lfloor \frac{x + y}{16} \right\rfloor, 16 \right).$$

В приведенной на следующей странице таблице показана интерпретация 4-битовых бинарных слов в формате дополнения до единицы. Эта интерпретация слова, обычная бинарная интерпретация которого — x , дается следующей формулой.

$$\text{ones}(x) = \begin{cases} x, & 0 \leq x \leq 7 \\ x - 15, & 8 \leq x \leq 15 \end{cases}$$

Нужно показать, что $f(x, y)$, в случае рассмотрения результата как целой величины в формате дополнения до единицы, представляет собой сумму x и y , рассматриваемых в том же формате, т.е. мы должны показать, что

$$\text{ones}(x) + \text{ones}(y) = \text{ones}(f(x, y)).$$

Нас интересуют только случаи без переполнения (т.е. когда сумма может быть выражена как целое в формате дополнения до единицы).

Случай 0, $0 \leq x, y \leq 7$. Тогда $\text{ones}(x) + \text{ones}(y) = x + y$, и

$$f(x, y) = \text{mod}(x + y + 0, 16) = x + y.$$

Чтобы не было переполнения, результат в формате дополнения до единицы должен находиться в диапазоне от 0 до 7, и из таблицы видно, что должно выполняться $x + y \leq 7$. Следовательно, $\text{ones}(x + y) = x + y$.

Случай 1, $0 \leq x \leq 7, 8 \leq y \leq 15$. Переполнение не может произойти, поскольку $\text{ones}(x) \geq 0$ и $\text{ones}(y) \leq 0$. В этом случае $\text{ones}(x) + \text{ones}(y) = x + y - 15$. Если $x + y < 16$,

$$f(x, y) = \text{mod}(x + y + 0, 16) = x + y.$$

В этом случае $x + y$ должно быть не менее 8, так что $\text{ones}(x + y) = x + y - 15$. С другой стороны, если $x + y \geq 16$,

$$f(x, y) = \text{mod}(x + y + 1, 16) = x + y + 1 - 16 = x + y - 15.$$

Поскольку $x + y$ не превышает 22 и не менее 16, $1 \leq x + y - 15 \leq 7$, так что $\text{ones}(x + y - 15) = x + y - 15$.

Случай 2, $8 \leq x \leq 15, 0 \leq y \leq 7$. Этот случай аналогичен рассмотренному выше случаю 1.

Случай 3, $8 \leq x \leq 15, 8 \leq y \leq 15$. Тогда $\text{ones}(x) + \text{ones}(y) = x - 15 + y - 15 = x + y - 30$, и

$$f(x, y) = \text{mod}(x + y + 1, 16) = x + y + 1 - 16 = x + y - 15.$$

В силу указанных границ x и y $16 \leq x + y \leq 30$. Как указано в таблице, чтобы избежать переполнения, мы должны иметь $x + y \geq 23$. В терминах дополнения до единицы мы можем сложить без переполнения -6 с -1 , или -6 с -0 , но не -6 с -2 . Таким образом, $23 \leq x + y \leq 30$, а следовательно, $8 \leq x + y - 15 \leq 15$, так что $\text{ones}(x + y - 15) = x + y - 30$.

Что касается вопроса о распространении переноса, в случае дополнения до единицы наимудший случай имеет вид наподобие следующего.

111...1111	
+ 000...0100	

000...0011	
+ 1 (циклический перенос)	

000...0100	

x	$\text{ones}(x)$
0000	0
0001	1
0010	2
0011	3
0100	4
0101	5
0110	6
0111	7
1000	-7
1001	-6
1010	-5
1011	-4
1100	-3
1101	-2
1110	-1
1111	-0

В этом случае перенос распространяется на n позиций, где n — размер слова. В случае дополнения до 2 в наихудшем случае распространение составляет $n-1$ позиций в предположении, что перенос из позиции старшего бита отбрасывается.

Интересно выполнить следующие сравнения (с использованием для иллюстрации 4-битовых величин). В простой бинарной (беззнаковой) арифметике или арифметике дополнения до двух сумма двух чисел всегда (даже при переполнении) корректна по модулю 16. В арифметике дополнения до единицы сумма всегда корректна по модулю 15. Если x_n обозначает n -й бит x , то в обозначениях дополнения до двух $x = -8x_3 + 4x_2 + 2x_1 + x_0$. При дополнении до единицы $x = -7x_3 + 4x_2 + 2x_1 + x_0$.

8. $((x \oplus y) \& m) \oplus y$.
9. $x \oplus y = (x | y) \& \neg(x \& y)$.
10. [5] Переменная t равна 1, если биты отличаются (шесть команд).

$$t \leftarrow \left(\left(x \gg i \right) \oplus \left(x \gg j \right) \right) \& 1$$

$$x \leftarrow x \oplus (t \ll j)$$

Добавление строки $x \leftarrow x \oplus (t \ll i)$ позволяет обменять биты i и j .

11. Как описано в разделе, любая булева функция $f(x_1, x_2, \dots, x_n)$ может быть разложена в виде $g(x_1, x_2, \dots, x_{n-1}) \oplus x_n h(x_1, x_2, \dots, x_{n-1})$. Пусть $c(n)$ — количество команд, необходимых для декомпозиции булевой функции от n переменных в бинарные команды ($n \geq 2$). Тогда

$$c_{n+1} = 2c_n + 2,$$

где $c_2 = 1$. Это рекуррентное соотношение имеет следующее решение.

$$c_n = 3 \cdot 2^{n-2} - 2$$

(Наименьшая верхняя граница гораздо меньше.)

12. (а)

$$\begin{aligned} f(x, y, z) &= \bar{z}f_0(x, y) + zf_1(x, y) \\ &= \bar{z}f_0(x, y) \oplus zf_1(x, y) \\ &= \bar{z}f_0(x, y) \oplus (1 \oplus \bar{z})f_1(x, y) \\ &= \bar{z}f_0(x, y) \oplus f_1(x, y) \oplus \bar{z}f_1(x, y) \\ &= f_1(x, y) \oplus \bar{z}(f_0(x, y) \oplus f_1(x, y)), \end{aligned}$$

что представляет собой вид, требующийся в условии.

- (б) Из части (а)

$$\begin{aligned}
 f(x, y, z) &= f_1(x, y) \oplus \bar{z}(f_0(x, y) \oplus f_1(x, y)) \\
 &= \overline{f_1(x, y)} \oplus \overline{\bar{z}(f_0(x, y) \oplus f_1(x, y))} \\
 &= \overline{f_1(x, y)} \oplus \left(z + \overline{(f_0(x, y) \oplus f_1(x, y))} \right),
 \end{aligned}$$

что представляет собой вид, требующийся в условии.

13. Используя обозначения из табл. 2.3 на с. 26, отсутствующие функции можно получить из $0000 = \text{and}(x, x)$, $0011 = \text{and}(x, x)$, $0100 = \text{and}(y, x)$, $0101 = \text{and}(y, y)$, $1010 = \text{nand}(y, y)$, $1011 = \text{cor}(y, x)$, $1100 = \text{nand}(x, x)$ и $1111 = \text{cor}(x, x)$.
14. Нет. Десять истинно бинарных функций в числовом виде представляют собой

0001 0010 0100 0110 0111
1000 1001 1011 1101 1110

При наличии реализации функции 0010 путем обмена операндов можно получить функцию 0100; аналогично функция 1011 дает 1101. Это все, чего можно добиться обменом операндов, поскольку другие функции коммутативны. Очевидно, что использование в качестве обоих операндов одной переменной сводит функцию к константе или унарной функции. Следовательно, требуется восемь типов команд.

15. В приведенной ниже таблице показаны шесть типов команд, которые решают поставленную задачу. Здесь x означает содержимое операнда-регистра, а k — содержимое поля непосредственного значения.

НАБОР ИЗ ШЕСТИ R-I БУЛЕВЫХ КОМАНД

Значения функции	Формула	Команда
0001	xk	<i>and</i>
0111	$x + k$	<i>or</i>
0110	$x \oplus k$	<i>xor</i>
1110	\overline{xk}	<i>nand</i>
1000	$\overline{x + k}$	<i>nor</i>
0101	k	<i>const</i>

Отсутствующие функции могут быть получены из $0000 = \text{and}(x, 0)$, $0010 = \text{and}(x, \bar{k})$, $0011 = \text{or}(x, 0)$, $0100 = \text{nor}(x, \bar{k})$, $1001 = \text{xor}(x, \bar{k})$, $1010 = \text{const}(x, \bar{k})$, $1011 = \text{or}(x, \bar{k})$, $1100 = \text{nor}(x, 0)$, $1101 = \text{nand}(x, \bar{k})$ и $1111 = \text{nand}(x, 0)$.

16. Автор не знает “аналитического” способа решения поставленной задачи. Однако нетрудно написать программу, которая генерирует все булевы функции от трех переменных, которые могут быть реализованы тремя бинарными командами. Такая программа на языке программирования С приведена ниже. Она написана максимально просто, чтобы дать убедительный ответ на поставленный вопрос.


```

        fun[5] = boole(c3, fun[k1], fun[k2]);
        found[fun[5]] = 1;
    })
})
})
printf(" 0 1 2 3 4 5 6 7 8 9 A B C D E F\n");
for (i = 0; i < 16; i++)
{
    printf("%X", i);
    for (j = 0; j < 16; j++)
    {
        printf("%2d", found[16*i + j]);
    }
    printf("\n");
}
return 0;
}

```

Программа представляет функцию как 8-битовую строку, которая является таблицей истинности функции, со значениями x , y и z , записанными обычным способом. Каждый раз при генерации функции соответствующий байт в векторе `found` устанавливается равным единице. Этот вектор имеет длину 256 байт и изначально все байты в нем — нулевые.

Таблица истинности, с которой работает программа, показана ниже.

ТАБЛИЦА ИСТИННОСТИ ДЛЯ ТРЕХ ПЕРЕМЕННЫХ

$f_0 = x$	$f_1 = y$	$f_2 = z$	f_3	f_4	f_5
0	0	0			
0	0	1			
0	1	0			
0	1	1			
1	0	0			
1	0	1			
1	1	0			
1	1	1			

Шесть столбцов таблицы истинности хранятся в векторе `fun`. Первые три позиции `fun` содержат столбцы x , y и z таблицы истинности. Эти столбцы содержат шестнадцатеричные значения 0F, 33 и 55, представляющие тривиальные функции $f(x, y, z) = x$, $f(x, y, z) = y$ и $f(x, y, z) = z$. Следующие три позиции будут хранить столбцы для функций, сгенерированных для текущей пробы одной, двумя и тремя бинарными командами соответственно.

Концептуально программа состоит из трех вложенных циклов, по одному для каждой испытываемой в настоящее время команды. Внешний цикл выполняет итерации по всем 16 бинарным булевым операциям, работая со всеми парами из x , y и z ($16 \cdot 3 \cdot 3 = 144$ итерации). На каждой итерации результат применения операции параллельно ко всем восьми битам x , y и/или z помещается в элемент `fun[3]`.

В циклах следующего уровня аналогично выполняются итерации по всем 16 бинарным булевым операциям над всеми парами, составляемыми из x , y , z и результата внешнего цикла ($16 \cdot 4 \cdot 4 = 256$ итераций). Для каждой итерации результат помещается в элемент `fun[4]`.

Наиболее глубоко вложенные циклы аналогично итерируют все 16 булевых операций над всеми возможными парами из x , y , z и результатов двух внешних циклов ($16 \cdot 5 \cdot 5 = 400$ итераций). Для каждой из этих вычисленных функций соответствующий байт `found` устанавливается равным 1.

В конце программа записывает вектор `found` в виде 16 строк по 16 элементов. Некоторые элементы вектора `found` нулевые, что свидетельствует о том, что трех бинарных булевых команд недостаточно для того, чтобы реализовать все 256 булевых функций от трех переменных. Первая такая функция — под номером $0x16$, или бинарная 00010110 , которая представляет собой $\bar{x}y + x\bar{y} + xy\bar{z}$.

Чтобы сократить количество итераций, можно воспользоваться множеством симметрий. Например, для данной операции op и операндов x и y необязательно вычислять и $op(x, y)$, и $op(y, x)$, поскольку если вычислено $op(x, y)$, то $op(y, x)$ будет представлять собой результат $op'(x, y)$, где op' — другая из 16 бинарных операций. Аналогично нет смысла в вычислении $op(x, x)$, поскольку это значение также будет равно $op'(x, y)$ для некоторой другой функции op' . Таким образом, внешние циклы, выбирающие комбинации операндов для испытаний, можно переписать так.

```
for (i1 = 0; i1 < 3; i1++) {
for (i2 = i1 + 1; i2 < 3; i2++) {
```

То же самое можно сделать и для прочих циклов.

Еще одно усовершенствование вытекает из наблюдения, что нет необходимости включать в таблицу все 16 бинарных булевых операций. Операции под номерами 0, 3, 5, 10, 12 и 15 можно опустить, сокращая цикл по операциям с 16 до 10 итераций. Доказательство правомочности этого решения достаточно длинное и здесь не приводится.

Программу легко изменить для проведения эксперимента с набором команд меньшего размера или большего их количества, или для большего числа переменных. Но учтите: время работы программы катастрофически растет с увеличением разрешенного количества команд, поскольку оно определяет уровень вложенности в основной программе. На практике вам придется ограничиться не более чем пятью командами.

Глава 3. Округление к степени 2

1. (а) $(x+4) \& -8$.
- (б) $(x+3) \& -8$.
- (в) $\left(x+3+\left(\left(x \gg 3\right) \& 1\right)\right) \& -8$.

Часть (в) можно выполнить четырьмя командами, если доступна команда *извлечения*, позволяющая выполнить $\left(x \gg 3\right) \& 1$ одной командой.

Примечание: несмещенное округление сохраняет среднее значение большого множества случайных чисел.

2. Стандартный способ выполнения части (а) — $\left((x+5) \div 10\right) * 10$. Если легко доступна функция получения остатка, задание из части (а) можно выполнить как $x+5 - \text{rem}(x+5, 10)$, что экономит умножение ценой дополнительного сложения.

Часть (б) аналогична, с заменой 5 на 4 в ответе к части (а) задания.

Часть (в): используем тот факт, что целое число является нечетным кратным 10 тогда и только тогда, когда оно является нечетным кратным 2.

```
r = x % 10;
y = x - r;
if (r > 5 | (r == 5 & (y & 2) != 0))
    y = y + 10;
```

Альтернативное решение (должно выполняться $x \leq 2^{32} - 6$) таково.

```
r = (x + 5) % 10;
y = x + 5 - r;
if (r == 0 & (y & 2) != 0)
    y = y - 10;
```

3. Возможная реализация на языке программирования C показана ниже.

```
int loadUnaligned(int *a) {
    int *alo, *ahi;
    int xlo, xhi, shift;

    alo = (int *) ((int)a & -4);
    ahi = (int *) (((int)a + 3) & -4);
    xlo = *alo;
    xhi = *ahi;
    shift = ((int)a & 3) << 3;
    return ((unsigned)xlo >> shift) | (xhi << (32-shift));
}
```

Глава 4. Арифметические границы

1. Для $a = c = 0$ неравенства (5) принимают следующий вид.

$$0 \leq x - y \leq 2^{32} - 1, \text{ если } -d < 0 \text{ и } b \geq 0$$

$$-d \leq x - y \leq b \quad \text{в противном случае}$$

Поскольку величины беззнаковые, $-d < 0$ эквивалентно $d \neq 0$, а $b \geq 0$ истинно. Следовательно, неравенства упрощаются.

$$\begin{aligned} 0 \leq x - y \leq 2^{32} - 1, & \text{ если } d \neq 0 \\ -d \leq x - y \leq b, & \text{ если } d = 0 \end{aligned}$$

Легко заметить, что когда $d = 0$, то $y = 0$, так что тривиально получается $0 \leq x - y \leq b$. С другой стороны, если $d \neq 0$, то разность может достичь значения 0 путем выбора $x = y = 0$ и максимального беззнакового числа путем выбора $x = 0$ и $y = 1$.

2. Если $a = 0$, то проверка `if (temp >= a)` всегда истинна. Следовательно, когда найдена первая позиция слева, в которой биты b и d равны 1, программа устанавливает бит b равным 0, а следующие биты равными 1 и возвращает значение, полученное с помощью операции *или* с d . Это можно выполнить более просто с помощью следующей замены тела процедуры. Оператор `if` требуется только на машинах, на которых сдвиги выполняются по модулю 32, таких как семейство Intel x86.

```
temp = nlz(c & d);
if (temp == 0) return 0xFFFFFFFF;
m = 1 << (32 - temp);
return b | d | (m - 1);
```

Предположим, например, что

$$\begin{aligned} 0 \leq x \leq 0b01001000 \text{ и} \\ 0b00000011 \leq y \leq 0b00101010. \end{aligned}$$

Тогда, чтобы найти максимальное значение $x | y$, процедура выполняет сканирование слева направо в поисках бита, который и в b , и в d равен 1. Максимальное значение равно $c | d$ для битов слева от этой позиции и единиц для битов справа от нее. В нашем примере это $0b01001000 | 0b00101010 | 0b00001111 = 0b01101111$.

Глава 5. Подсчет битов

1. Версия Норберта Джаффа (Norbert Juffa).

```
int ntz (unsigned int n)
{
    static unsigned char tab[32] =
    {
        0, 1, 2, 24, 3, 19, 6, 25,
        22, 4, 20, 10, 16, 7, 12, 26,
        31, 23, 18, 5, 21, 9, 15, 11,
        30, 17, 8, 14, 29, 13, 28, 27
    };
    unsigned int k;
    n = n & (-n); /* Отделение младшего бита */
    #if defined(SLOW_MUL)
    k = (n << 11) - n;
    k = (k << 2) + k;
    k = (k << 8) + n;
```

```

    k = (k << 5) - k;
#else
    k = n * 0x4d7651f;
#endif
    return n ? tab[k>>27] : 32;
)

```

2. $x \dot{\gg} (x \& -x)$. Применяется в функции `snoob` (с. 36).
3. Обозначим применение операции параллельного префикса к x как $PP\text{-}XOR(x)$.

Тогда, если $y = PP\text{-}XOR(x)$, то $x = y \oplus (y \dot{\gg} 1)$. Чтобы убедиться в этом, пусть x представляет собой 4-битовую величину $abcd$ (где каждая буква обозначает отдельный бит). Тогда

$$y = PP\text{-}XOR(x) = (a)(a \oplus b)(a \oplus b \oplus c)(a \oplus b \oplus c \oplus d),$$

$$y \dot{\gg} 1 = (0)(a)(a \oplus b)(a \oplus b \oplus c), \text{ так что}$$

$$y \oplus (y \dot{\gg} 1) = (a)(b)(c)(d).$$

Для операции параллельного суффикса, если $y = PS\text{-}XOR(x)$, то, как вы можете догадаться, $x = y \oplus (y \ll 1)$.

Глава 6. Поиск в слове

1. Длина и позиция самой длинной строки из единиц (Норберт Джатфа (Norbert Juffa)).

```

int fmaxstr1(unsigned x, int *apos)
{
    int k;
    unsigned oldx;

    oldx = 0;
    for (k = 0; x != 0; k++)
    {
        oldx = x;
        x &= 2*x;
    }
    *apos = nlz(oldx);
    return k;
}

```

2. Как сказано в разделе, это можно сделать путем распространения нулевых битов x влево на $n-1$ позиций с последующим поиском кратчайшей строки из единичных битов во вновь полученном x . Выполнить распространение можно с помощью кода из листинга 6.5 на с. 149, имеющего логарифмическое время работы. (Однако вторая часть алгоритма линейна относительно длины кратчайшей строки единиц в полученном x .) Код приведен ниже. В нем предполагается, что $1 \leq n \leq 32$. В случае, когда искомая подстрока не найдена, функция возвращает позицию `apos = 32`. В этом случае длина должна рассматриваться как неопределенная, хотя функция и возвращает значение длины, равное $n-1$.

```

int bestfit(unsigned x, int n, int *apos)
{
    int m, s;

    m = n;
    while (m > 1)
    {
        s = m >> 1;
        x = x & (x << s);
        m = m - s;
    }
    return fminstr1(x, apos) + n - 1;
}

```

3. В следующем коде используется выражение на с. 31 для обнуления крайней справа непрерывной последовательности единичных битов.

```

int fminstr1(unsigned x, int *apos)
{
    int k, kmin, y0, y;
    unsigned int x0, xmin;

    kmin = 32;
    y0 = pop(x);
    x0 = x;
    do {
        x = ((x & -x) + x) & x; // Обнуление крайней справа
        y = pop(x); // строки единиц
        k = y0 - y; // k = длина обнуленной
        if (k <= kmin) // строки
        {
            kmin = k; // Сохраняем кратчайшую
            xmin = x; // найденную длину и строку
        }
        y0 = y;
    } while (x != 0);
    *apos = nlz(x0 ^ xmin);
    return kmin;
}

```

Функция выполняет $5 + 11n$ команд, где n — количество строк из единиц в x , при $n \geq 1$ (т.е. для $x \neq 0$). Предполагается также, что в конечном счете проверка *if* выполняется в половине случаев, а функции $\text{pop}(x)$ и $\text{nlz}(x)$ учитываются как одна команда каждая. При изменении смысла проверки "*if* ($k \leq kmin$)" и значения инициализации $kmin$ можно находить наидлиннейшую строку единиц (либо крайнюю слева, либо крайнюю справа в случае наличия строк одинаковой длины). Легко также модифицировать этот код для вычисления функции "наилучшего размещения".

4. Первый бит x будет равен 1 (а следовательно, отмечать начало строки из единиц) с вероятностью 0.5. Любой другой бит отмечает начало строки из единичных битов с вероятностью 0.25 (он должен быть единичным, а бит слева от него — нулевым). Следовательно, среднее количество строк единиц составляет $0.5 + 31 \cdot 0.25 = 8.25$.

5. Можно ожидать, что подавляющее большинство слов, если они достаточно длинны, содержат строку из единичных битов единичной длины. Если слово начинается с 10 или заканчивается 01, или содержит строку 010, то кратчайшая содержащаяся в нем строка из единичных битов имеет единичную длину. Следовательно, средняя длина, вероятно, лишь немного больше 1.

Исчерпывающая проверка всех 2^{32} слов показывает, что средняя длина кратчайшей строки из единичных битов составляет около 1.011795.

6. (Решение Джона Ганнелса (John Gunnels)). Эта задача на удивление сложная, но используемый метод решения достаточно известный. Решение основано на рекурсии, подсчитывающей количество слов в каждом из четырех множеств, показанных в таблице ниже. В этой таблице “синглтон” означает строку из единичных битов единичной длины, “ппп” означает строку длиной ≥ 0 , не содержащую синглтона, а “sss” означает строку длиной ≥ 1 , содержащую синглтон. Троекотия означают 0 или большее количество повторений предыдущего бита. Каждое бинарное слово принадлежит одному и только одному из этих четырех множеств.

Множество	Слова вида	Описание
<i>A</i>	ппп0... или пустое	Не содержит синглтон, но может получить его на следующем шаге
<i>B</i>	ппп01 или 1	Содержит синглтон, но может потерять его на следующем шаге
<i>C</i>	ппп011... или 11...	Не содержит синглтон и не может получить его на следующем шаге
<i>D</i>	sss0 или sss01...	Содержит синглтон и будет содержать его на следующем шаге

На каждом шаге к слову справа добавляется бит. После этого слово перемещается из одного множества в другое, как показано ниже. Оно перемещается в первое множество, если добавленный бит нулевой, и в правое, если единичный.

$$A \Rightarrow A \text{ или } B$$

$$B \Rightarrow D \text{ или } C$$

$$C \Rightarrow A \text{ или } C$$

$$D \Rightarrow D \text{ или } D$$

Например, слово 1101 находится в множестве *B*. Если к нему добавить 0, оно превращается в 11010, находящееся в множестве *D*. Если же добавить 1, то оно становится равным 11011, принадлежащим множеству *C*.

Пусть a_n , b_n , c_n и d_n обозначают размеры множеств *A*, *B*, *C* и *D* соответственно после n шагов (когда слова имеют длину n). Тогда

$$a_{n+1} = a_n + c_n$$

$$b_{n+1} = a_n$$

$$c_{n+1} = b_n + c_n$$

$$d_{n+1} = b_n + 2d_n$$

Эти формулы отражают тот факт, что множество A на шаге $n+1$ содержит все члены множества A на шаге n с добавленным 0, а также все члены множества C на шаге n с добавленным 0. Множество B на шаге $n+1$ содержит только все члены множества A на шаге n , к которым добавлены единичные биты, и т.д.

Начальные условия следующие: $a_0 = 1$ и $b_0 = c_0 = d_0 = 0$.

Достаточно легко провести вычисления по указанным формулам с помощью программы или даже вручную. Для $n = 32$ результат имеет вид

$$a_{32} = 26\,931\,732$$

$$b_{32} = 15\,346\,786$$

$$c_{32} = 20\,330\,163$$

$$d_{32} = 4\,232\,358\,615$$

$$b_{32} + d_{32} = 4\,247\,705\,401$$

Последняя строка дает интересующую нас величину — количество слов, в которых кратчайшая из имеющихся строк из единичных битов имеет длину 1. Это около 98.9% от общего количества 32-битовых чисел, равного 2^{32} .

Но нельзя ли получить решение в аналитическом виде? Оказывается, это достаточно сложно. Набросаем здесь “костяк” такого решения.

Пусть $e_n = b_n + d_n$ — величина, которую требуется найти. Тогда из разностных уравнений и используя тот факт, что $a_n + b_n + c_n + d_n = 2^n$, получаем

$$\begin{aligned} e_n &= b_n + d_n \\ &= 2^n - a_n - c_n \\ &= 2^n - a_{n+1}. \end{aligned}$$

Таким образом, если мы найдем аналитическое решение для a_n , тем самым мы найдем его и для e_n .

Разностное уравнение с одной переменной для a_n можно найти следующим образом.

$$\begin{aligned} a_n &= a_n + c_{n-1} \\ &= a_{n-1} + b_{n-2} + c_{n-2} \\ &= a_{n-1} + a_{n-3} + c_{n-2} \\ &= a_{n-1} + a_{n-3} + a_{n-1} - a_{n-2} \\ &= 2a_{n-1} - a_{n-2} + a_{n-3} \end{aligned}$$

Это разностное уравнение можно решить хорошо известными методами. Процесс этот достаточно длинный и неприятный, так что здесь мы его не приводим. В частности, он включает решение кубического уравнения с двумя комплексными корнями. По окончании работы для e_n мы получим следующее приближенное выражение.

$$e_n \approx 2^n - 0.41150 \cdot 1.7549^{n+1} \\ - (0.29425 - 0.13811i)(0.12256 + 0.74486i)^{n+1} \\ - (0.29425 + 0.13811i)(0.12256 - 0.74486i)^{n+1}$$

Если n — целое число, мнимые части сокращаются (это не так трудно доказать). *Указание:* если x и y — комплексно сопряженные, то это же утверждение справедливо и для x^n и y^n).

Можно получить формулу, включающую только действительные числа. Действительная часть второго члена приведенной выше формулы, определенно, меньше

$$|0.29425 - 0.13811i| \cdot |0.12256 + 0.74486i|^{n+1},$$

что для $n = 0$ составляет

$$0.32505 \cdot 0.75488 \approx 0.24573,$$

и уменьшается с ростом n . То же самое справедливо и для последнего члена в формуле для e_n . Следовательно, действительная часть двух последних членов не превышает 0.5. Поскольку априори известно, что e_n — целое число, это означает, что e_n определяется первым членом, округленным к ближайшему целому, т.е.

$$e_n \approx \left[2^n - 0.41150 \cdot 1.7549^{n+1} + 0.5 \right].$$

7. Вкратце: эта задача может быть решена с использованием 10 множеств слов, описанных ниже. В этой таблице “*nnn*” означает строку длиной ≥ 0 , кратчайшая подстрока из единиц которой имеет либо нулевую длину, либо длину ≥ 3 ; “*ddd*” означает строку длиной ≥ 2 , кратчайшая подстрока из единиц которой имеет длину 2, а “*sss*” означает строку длиной ≥ 1 , чья кратчайшая подстрока из единиц которой имеет длину 1. (Эти множества отслеживают слова, которые содержат синглтон в позиции, отличающейся от крайней справа, поскольку такие слова никогда не будут иметь кратчайших подстрок из единиц длиной 2.) Тросточия означают 0 или большее количество повторов предыдущего бита.

Множество	Слова вида	Множество	Слова вида
<i>A</i>	<i>nnn</i> 0... или пустое	<i>F</i>	<i>ddd</i> 01
<i>B</i>	<i>nnn</i> 01 или 1	<i>G</i>	<i>ddd</i> 011
<i>C</i>	<i>nnn</i> 011 или 11	<i>H</i>	<i>ddd</i> 0111
<i>D</i>	<i>nnn</i> 0111... или 111...	<i>I</i>	<i>sss</i> 0
<i>E</i>	<i>ddd</i> 0	<i>J</i>	<i>sss</i> 01...

На каждом шаге к слову справа добавляется бит. После этого слово перемещается из одного множества в другое, как показано ниже. Оно перемещается в первое множество, если добавленный бит нулевой, и в правое, если единичный.

$$\begin{array}{ll}
 A \Rightarrow A \text{ или } B & F \Rightarrow I \text{ или } G \\
 B \Rightarrow I \text{ или } C & G \Rightarrow E \text{ или } H \\
 C \Rightarrow E \text{ или } D & H \Rightarrow E \text{ или } H \\
 D \Rightarrow A \text{ или } D & I \Rightarrow I \text{ или } J \\
 E \Rightarrow E \text{ или } F & J \Rightarrow I \text{ или } J
 \end{array}$$

Пусть a_n, b_n, \dots, j_n обозначают размеры множеств A, B, \dots, J соответственно после n шагов (когда слова имеют длину n). Тогда

$$\begin{array}{ll}
 a_{n+1} = a_n + d_n & f_{n+1} = e_n \\
 b_{n+1} = a_n & g_{n+1} = f_n \\
 c_{n+1} = b_n & h_{n+1} = g_n + h_n \\
 d_{n+1} = c_n + d_n & i_{n+1} = b_n + f_n + i_n + j_n \\
 e_{n+1} = c_n + e_n + g_n + h_n & j_{n+1} = i_n + j_n
 \end{array}$$

Начальные условия следующие: $a_0 = 1$, а все прочие переменные равны 0.

Интересующее нас значение, количество слов, наименьшая подстрока из единиц в которых имеет длину 2, задается суммой $c_n + e_n + g_n + h_n$. Для $n = 32$ разностные уравнения дают значение 44410452, что составляет около 1.034% от общего количества 32-битовых слов. В качестве дополнительного результата мы получаем, что количество слов, наименьшая подстрока из единиц в которых имеет единичную длину, задается суммой $b_n + f_n + i_n + j_n$ и для $n = 32$ равно 4247705401, что согласуется с результатом предыдущего упражнения.

Глава 7. Перестановка битов и байтов

1. Обычное целое число можно увеличить на 1 путем выполнения операции дополнения над некоторым количеством последовательных младших битов.¹ Например, чтобы добавить 1 к $0x321F$, достаточно применить операцию *исключающего или* к этому числу с маской $0x003F$. Аналогично для увеличения целого числа с обратным порядком битов достаточно выполнить операцию дополнения над некоторым количеством последовательных старших битов, воспользовавшись маской, которая состоит из строки единичных битов, за которой следуют нулевые биты. Формула Мёбиуса вычисляет эту маску и применяет ее к числу с обратным порядком битов. (Метод, описанный в разделе и использующий операцию `nlz`, также выполняет эти действия.)

¹ Этим свойством обладает также система счисления по основанию -2 , но не по основанию $-1+i$.

Для обычного целого числа маска состоит из нулей, за которыми следуют единицы, начиная с крайнего справа нулевого бита и заканчивая младшим битом слова. Целое число, которое состоит из единичного бита в позиции крайнего справа нулевого бита в i получается с помощью выражения $-i \& (i + 1)$ (см. раздел 2.1). Для увеличения на 1 обычного целого числа x следует вычислить маску путем распространения вправо единичного бита, а затем выполнить *исключающее или* получившегося результата с числом x . Для увеличения на 1 целого числа с обратным порядком битов нужно выполнить реверс битов этой маски (получить ее отражение). Для однобитовой (представляющей собой степень двойки) величины $-i \& (i + 1)$ ее отражение можно получить путем деления на него $m/2$ (это ключевой шаг алгоритма). Например, в случае 4-битовых целых чисел $m/2 = 8$, $8/1 = 8$, $8/2 = 4$, $8/4 = 2$ и $8/8 = 1$. Для вычисления маски необходимо только распространение влево единичного бита частного, что выполняется путем вычитания частного из m . Наконец выполняем *исключающее или* маски и отраженного целого числа, что дает следующее отраженное число.

В качестве примера предположим, что целые числа имеют длину, равную 8 бит, так что $m = 256$. Пусть $i = 19$ (бинарное 00010011), так что $rev i$ в бинарном виде равно 11001000. Тогда $-i \& (i + 1)$ в бинарном виде записывается как 00000100 (десятичная 4). Деление $m/2$ на него дает частное, равное 32 (бинарное 00100000). Вычитание этого значения из m дает бинарное число 11100000. Наконец *исключающее или*, примененное к маске и $rev i$, дает бинарное число 00101000, которое представляет собой десятичное 20 с обращенными битами.

2. Заметим, что

$$\begin{aligned} m_0 &= 2 \cdot 0x11111111, \\ m_1 &= 0x \cdot C0x01010101, \\ m_2 &= 0xF0 \cdot 0x00010001 \text{ и} \\ m_3 &= 0xFF00 \cdot 0x00000001. \end{aligned}$$

Заметим также, что

$$\begin{aligned} 0x11111111 &= \lfloor 2^{32}/15 \rfloor, \\ 0x01010101 &= \lfloor 2^{32}/255 \rfloor, \\ 0x00010001 &= \lfloor 2^{32}/(2^{16} - 1) \rfloor \text{ и} \\ 0x00000001 &= \lfloor 2^{32}/(2^{32} - 1) \rfloor. \end{aligned}$$

Таким образом, мы получаем формулы

$$m_0 = (2-1)2 \lfloor 2^{32}/(2^4-1) \rfloor,$$

$$m_1 = (2^2-1)2^2 \lfloor 2^{32}/(2^8-1) \rfloor,$$

$$m_2 = (2^4-1)2^4 \lfloor 2^{32}/(2^{16}-1) \rfloor \text{ и}$$

$$m_3 = (2^8-1)2^8 \lfloor 2^{32}/(2^{32}-1) \rfloor.$$

В общем случае

$$m_i = (2^{2^i} - 1)2^{2^i} \lfloor 2^W / (2^{2^{i+1}} - 1) \rfloor,$$

где W — длина перемешиваемого слова, которая должна быть степенью 2.

3. Необходимо всего лишь заменить две строки

```
s = s + b;
x = x >> 1;
```

строками

```
s = s + 1;
x = x >> b;
```

4. Любой корректный алгоритм LRU должен записывать полный порядок обращений к n строкам кеша в множестве. Поскольку имеется $n!$ способов перестановки n предметов, любая реализация LRU должна использовать как минимум $\lceil \log_2 n! \rceil$ бит памяти. В таблице, приведенной ниже, это число сравнивается с количеством битов, требуемых методом матрицы обращений.

Степень ассоциативности	Теоретический минимум	Метод матрицы обращений
2	1	1
4	5	6
8	16	28
16	45	120
32	118	496

Глава 8. Умножение

1. Как показано в разделе 8.3, если x и y являются операндами операции умножения, рассматриваемыми как знаковые целые числа, то их произведение, рассматриваемое как беззнаковая операция, представляет собой

$$(x + 2^{32} x_{31})(y + 2^{32} y_{31}) = xy + 2^{32}(x_{31}y + y_{31}x) + 2^{64} x_{31}y_{31}.$$

где x_{31} и y_{31} — знаковые биты x и y соответственно, которые интерпретируются как целые числа, равные 0 или 1. Поскольку произведение отличается от xu на кратное 2^{32} , младшие 32 бита произведения оказываются теми же самыми.

2. Метод 1. Скорее всего, машина имеет команду умножения, которая дает младшие 32 бита произведения двух 32-битовых целых чисел, т.е.

$$\text{low} = u * v;$$

Метод 2. Непосредственно перед инструкцией `return` вставить

$$\text{low} = (w1 \ll 16) + (w0 \& 0xFFFF);$$

Метод 3. Сохранить значения произведений $u1 * v0$ и $u0 * v1$ во временных переменных $t1$ и $t2$. Затем

$$\text{low} = ((t1 + t2) \ll 16) + w0;$$

Методы 2 и 3 используют по три команды базового набора RISC и работают как с `mulhs`, так и с его беззнаковым двойником (и могут оказаться быстрее метода 1).

3. Разобьем 32-битовые операнды u и v на 16-битовые беззнаковые компоненты a , b , c и d , так что

$$u = 2^{16} a + b \text{ и}$$

$$v = 2^{16} c + d,$$

где $0 \leq a, b, c, d \leq 2^{16} - 1$. Пусть

$$p = ac,$$

$$q = bd \text{ и}$$

$$r = (-a + b)(c - d).$$

Тогда $uv = 2^{32} p + 2^{16} (r + p + q) + q$, в чем легко убедиться.

Теперь $0 \leq p, q \leq 2^{32} - 2^{17} + 1$, так что p и q могут быть представлены 32-битовыми беззнаковыми числами. Однако легко вычислить, что

$$-2^{32} + 2^{17} - 1 \leq r \leq 2^{32} - 2^{17} + 1,$$

так что r является знаковой 33-битовой величиной. Ее удобнее представить знаковым 64-битовым числом, все старшие 32 бита которого либо нулевые, либо единичные. Машинная команда умножения вычислит 32 младших бита r , а старшие 32 бита можно получить из значений $-a + b$ и $c - d$. Это 17-битовые знаковые целые числа. Если они имеют противоположные знаки и не равны нулю, то r отрицательно, и, следовательно, все старшие 32 бита имеют значения 1. Если же у них одинаковые знаки или одно из них равно нулю, то r неотрицательно, а следовательно, все старшие 32 бита имеют нулевые значения. Проверку того, что одно из значений $-a + b$ и $c - d$ равно нулю, можно выполнить, проверяя только младшие 32 бита r . Если они равны нулю, то один из множителей — нулевой, так как $|r| < 2^{32}$.

Эти рассуждения приводят к следующей функции для вычисления старших 32 бит произведения u и v .

```
unsigned mulhu(unsigned u, unsigned v)
{
    unsigned a, b, c, d, p, q, rlow, rhigh;

    a = u >> 16;      b = u & 0xFFFF;
    c = v >> 16;      d = v & 0xFFFF;

    p = a*c;
    q = b*d;
    rlow = (-a + b)*(c - d);
    rhigh = (int)((-a + b)^(c - d)) >> 31;
    if (rlow == 0)
    {
        rhigh = 0;    // Коррекция
    }
    q = q + (q >> 16); // Переполнения здесь быть не может
    rlow = rlow + p;
    if (rlow < p)
    {
        rhigh = rhigh + 1;
    }
    rlow = rlow + q;
    if (rlow < q)
    {
        rhigh = rhigh + 1;
    }

    return p + (rlow >> 16) + (rhigh << 16);
}
```

После вычисления p , q , $rlow$ и $rhigh$ функция выполняет следующее сложение.

```

|.....p.....|
|....rhigh....||.....rlow.....|
                |.....p.....|
                |.....q.....|
                |.....q.....|
```

Инструкция "if (rlow<p) rhigh=rhigh+1" добавляет к $rhigh$ единицу, если имеется перенос при добавлении p к $rlow$ в предыдущей инструкции.

Младшие 32 бита произведения можно получить из следующего выражения, вставленного непосредственно после шага "коррекции" выше.

```
q + ((p + q + rlow) << 16)
```

Далее приводится версия кода без ветвлений.

```
unsigned mulhu(unsigned u, unsigned v)
{
    unsigned a, b, c, d, p, q, x, y, rlow, rhigh, t;
    a = u >> 16;      b = u & 0xFFFF;
    c = v >> 16;      d = v & 0xFFFF;
```



```

p = a*c;
q = b*d;
x = -a + b;
y = c - d;
rlow = x*y;
rhigh = (x ^ y) & (rlow | -rlow);
rhigh = (int)rhigh >> 31;

q = q + (q >> 16); // Переполнения здесь быть не может
t = (rlow & 0xFFFF) + (p & 0xFFFF) + (q & 0xFFFF);
p += (t >> 16) + (rlow >> 16) + (p >> 16) + (q >> 16);
p += (rhigh << 16);
return p;
}

```

Эти функции имеют большие накладные расходы по сравнению с функцией с четырьмя умножениями из листинга 8.2 на с. 200, и будут превосходить ее только на машинах с очень медленными по сравнению с современными машинами командами умножения. В случае арифметики “больших чисел” (арифметики с “многословными” целыми числами) время, требующееся для умножения, существенно больше времени, требуемого для сложения двух целых чисел того же размера. В таких приложениях метод, известный как умножение Карацубы (Karatsuba multiplication) [62], рекурсивно применяет схему трех умножений и оказывается более быстрым, чем непосредственное применение схемы четырех умножений к достаточно большим числам. Умножение Карацубы, как оно обычно описывается, использует

$$\begin{aligned}
 p &= ac, \\
 q &= bd, \\
 r &= (a+b)(c+d) \text{ и} \\
 uv &= 2^{32}p + 2^{16}(r - p - q) + q.
 \end{aligned}$$

Для нашего приложения этот метод работает не очень хорошо, поскольку r может достигать 2^{34} , а, похоже, простого способа вычисления двух старших битов 34-битового значения r не существует.

Знаковая версия приведенных выше функций имеет проблемы с переполнением. Вероятно, следует использовать беззнаковую функцию и откорректировать ее, как описано в разделе 8.3 на с. 200.

Глава 9. Целочисленное деление

1. Пусть $x = x_0 + \delta$, где x_0 — целое число, и $0 \leq \delta < 1$. Тогда $\lceil -x \rceil = \lceil -x_0 - \delta \rceil = -x_0$ в соответствии с определением функции “потолок” как ближайшего целого числа, большего или равного аргументу. Следовательно, $-\lceil -x \rceil = x_0$, что равно $\lfloor x \rfloor$.
2. Пусть n/d обозначает частное знакового целочисленного деления с отсечением. Тогда мы должны вычислить

$$\begin{aligned} n/d, & \quad \text{если } n \geq 0, d > 0, \\ n/d-1, & \quad \text{если } n < 0, d > 0, \\ n/d, & \quad \text{если } n \geq 0, d < 0, \\ n/d+1, & \quad \text{если } n < 0, d < 0. \end{aligned}$$

(Если $d = 0$, результат не существует.) Все это можно вычислить как $n/d + c$, где

$$c = \left(\left(n \dot{\gg} 31 \right) \oplus \left(d \dot{\gg} 31 \right) \right) - \left(d \dot{\gg} 31 \right),$$

так что для вычисления c требуется четыре команды (член $d \dot{\gg} 31$ — общий). Другой способ вычисления c с помощью четырех команд, но с беззнаковым сдвигом, заключается в вычислении

$$c = \left(d \ddot{\gg} 31 \right) - \left((n \oplus d) \ddot{\gg} 31 \right).$$

Если ваша машина оснащена командами сдвига по модулю 32, c можно вычислить с помощью трех команд:

$$c = \left(n \dot{\gg} 31 \right) \ddot{\gg} \left(d \dot{\gg} 31 \right).$$

Что касается остатка, то обозначим через $\text{rem}(n, d)$ остаток при делении знакового целого числа n на знаковое целое число d с использованием деления с отсечением. Тогда мы должны вычислить

$$\begin{aligned} \text{rem}(n, d), & \quad \text{если } n \geq 0, d > 0, \\ \text{rem}(n, d) + d, & \quad \text{если } n < 0, d > 0, \\ \text{rem}(n, d), & \quad \text{если } n \geq 0, d < 0, \\ \text{rem}(n, d) - d, & \quad \text{если } n < 0, d < 0. \end{aligned}$$

Величина, которую следует прибавить к $\text{rem}(n, d)$, либо равна нулю, либо имеет абсолютное значение d . Вычисления можно выполнить следующим образом:

$$\begin{aligned} |d| &= \left(d \oplus \left(d \dot{\gg} 31 \right) \right) - \left(d \dot{\gg} 31 \right), \\ c &= |d| \& \left(n \dot{\gg} 31 \right) \end{aligned}$$

(с помощью пяти команд для вычисления c). Если ваша машина оснащена командами сдвига по модулю 32 и вы используете команды умножения, c можно вычислить с помощью четырех команд (детали опущены).

3. Чтобы получить частное при делении с округлением к меньшему значению, необходимо только лишь вычесть 1 из частного деления с отсечением, если делимое и делитель имеют разные знаки.

$$n/d - \left((n \oplus d) \gg 31 \right)$$

Что касается остатка, то необходимо только добавить делитель к остатку от деления с отсечением, если делимое и делитель имеют разные знаки.

$$\text{rem}(n, d) + \left(\left((n \oplus d) \gg 31 \right) \& d \right)$$

4. Обычный метод, скорее всего, состоит в вычислении $\lfloor (n+d-1)/d \rfloor$. Проблема заключается в том, что $n+d-1$ может вызвать переполнение (рассмотрите вычисление $\lceil 12/5 \rceil$ на четырехразрядной машине).

Другой стандартный метод заключается в том, чтобы вычислить $q = \lfloor n/d \rfloor$ с помощью машинной команды *деления*, вычислить остаток как $r = n - qd$ и, если r имеет ненулевое значение, добавить 1 к q . (Или, как альтернативный вариант, добавить 1, если $n \neq qd$.) Это дает корректный результат для всех n и $d \neq 0$, но иногда оказывается дорогостоящим решением из-за применения умножения, вычитания и прибавления 1. С другой стороны, если ваша машина оснащена командой *деления*, которая в качестве побочного результата дает остаток, а в особенности если она может эффективно вычислять $q = q + (r \neq 0)$, то это достаточно хороший способ.

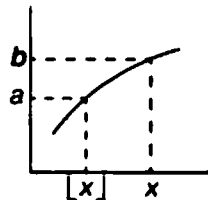
Еще один способ заключается в вычислении $q = \lfloor (n-1)/d \rfloor + 1$. К сожалению, он даст неверный результат при $n = 0$. Эту ошибку можно исправить, если машина позволяет легко вычислить предикат $x \neq 0$, такой, как команда *сравнения*, которая позволяет устанавливать значение регистра общего назначения равным 1 или 0 (см. также раздел 2.12 на с. 43). Тогда можно вычислить

$$c \leftarrow (x \neq 0),$$

$$q \leftarrow \lfloor (n-c)/d \rfloor + c.$$

Наконец можно вычислить $q = \lfloor (n-1)/d \rfloor + 1$, а затем, если $n = 0$, сделать результат равным 0, например, командой условного перемещения или выбора.

5. Пусть $f(\lfloor x \rfloor) = a$ и $f(x) = b$, как проиллюстрировано ниже.



Если b — целое число, то согласно (в) то же самое относится и к x , так что $\lfloor x \rfloor = x$, и доказывать ничего не требуется. Следовательно, при дальнейших рассуждениях предполагаем, что b — не целое, а a может быть как целым, так и не целым.

Не может быть целого числа k , такого, что $a < k \leq b$, поскольку если бы оно было, то имелось бы целое число между $\lfloor x \rfloor$ и x (в соответствии со свойствами (а), (б) и (в)), что невозможно. Следовательно, $\lfloor a \rfloor = b$, т.е. $\lfloor f(\lfloor x \rfloor) \rfloor = \lfloor f(x) \rfloor$.

Примеры применения доказанного утверждения (для целых a и b) приведены ниже.

$$\begin{aligned} \left\lfloor \frac{\lfloor x \rfloor + a}{b} \right\rfloor &= \left\lfloor \frac{x+a}{b} \right\rfloor \\ \left\lfloor \sqrt{\lfloor x \rfloor} \right\rfloor &= \left\lfloor \sqrt{x} \right\rfloor \\ \left\lfloor \log_2(\lfloor x \rfloor) \right\rfloor &= \left\lfloor \log_2(x) \right\rfloor \end{aligned}$$

Аналогично можно показать, что если $f(x)$ обладает свойствами (а), (б) и (в), то

$$\lceil f(\lceil x \rceil) \rceil = \lceil f(x) \rceil.$$

Глава 10. Целое деление на константы

1. (а) Если делитель четный, то младший бит делимого не влияет на частное (при делении с округлением к меньшему значению); если этот бит равен единице, это делает остаток нечетным. После сброса этого бита в нулевое значение остаток от деления оказывается четным числом. Следовательно, для четного делителя d остаток не превышает $d-2$. Это небольшое изменение максимального возможного остатка, как мы сейчас увидим, приводит к тому, что максимальный множитель m является W -битовым, а не $(W+1)$ -битовым числом (а следовательно, команда `shlxi` не нужна). На самом деле мы изучим, к каким упрощениям приводит завершение делителя z нулевыми битами, т.е. при его умножении на 2^z , $z \geq 0$. В этом случае z младших битов делимого можно сбросить в нуль без влияния на частное, а после этого сброса максимальный остаток равен $d-2^z$.

Следуя выводу из раздела 10.9 на с. 256, но внося изменения, так, чтобы максимальный остаток был равен $d-2^z$, мы имеем $n_c = 2^W - \text{rem}(2^W, d) - 2^z$, и неравенство (24а) превращается в

$$2^W - d \leq n_c \leq 2^W - 2^z.$$

Неравенство (25) принимает вид

$$\frac{2^p}{d} \leq m < \frac{2^p}{d} \frac{n_c + 2^z}{n_c}.$$

Уравнение (26) остается неизменным, а неравенство (27) превращается в

$$2^p > \frac{n_c}{2^z} (d-1 - \text{rem}(2^p - 1, d)). \quad (27')$$

Неравенство (28) принимает вид

$$1 \leq 2^p \leq \frac{2n_c}{2^z} (d-1) + 1.$$

В случае, когда p не обязано быть равным W , объединение этих неравенств дает

$$\begin{aligned} \frac{1}{d} \leq m &< \frac{2n_c(d-1) + 2^z n_c + 2^z}{2^z d n_c}, \\ 1 \leq m &< \frac{2d - 2 + 2^z/n_c}{2^z d} (n_c + 2^z), \\ 1 \leq m &< \frac{2}{2^z} (n_c + 2^z) \leq \frac{2}{2^z} 2^W. \end{aligned}$$

Таким образом, если $z \geq 1$, то $m < 2^W$, так что m укладывается в W -битовое слово. Тот же результат получается и в случае, когда p должно быть равно W .

Для вычисления множителя для данного делителя вычисляем n_c , как показано выше, затем находим наименьшее $p \geq W$, которое удовлетворяет неравенству (27'), и вычисляем m из (26). В качестве примера для $d = 14$ и $W = 32$ мы имеем $n_c = 2^{32} - \text{rem}(2^{32}, 14) - 2 = 0\text{x}\text{FFFFFFFA}$. Повторное применение (27') дает $p = 35$, откуда (26) дает $m = (2^{35} + 14 - 1 - 3)/14 = 0\text{x}92492493$. Таким образом, код для деления на 14 имеет следующий вид.

```
ins  n, R0, 0, 1      Сброс младшего бита n
li   m, 0x92492493   Загрузка магического числа
mulhu q, M, n        q = floor(M*n/2**32)
shri q, q, 3         q = q/8
```

(б) Как и ранее, если делитель кратен 2^z , то младшие z бит делимого не влияют на значение частного. Следовательно, можно очистить младшие z бит делимого и поделить делитель на 2^z без изменения значения частного. (Деление делителя выполняется во время компиляции.)

Используя преобразованные значения n и d , каждое из которых меньше 2^{W-z} , (24а) можно записать как

$$2^{W-z} - d \leq n_c \leq 2^{W-z} - 1.$$

Уравнение (26) и неравенство (27) не меняются, но они также работают с преобразованными значениями n_c и d . Мы опустим здесь доказательство того, что множитель будет меньше 2^W , и вновь рассмотрим пример для $d = 14$ и $W = 32$. В уравнениях мы используем $d = 7$. Таким образом, имеем $n_c = 2^{31} - \text{rem}(2^{31}, 7) - 1 = 0\text{x}7\text{FFFFFFF}$.

Повторное применение (27) дает $p = 34$, откуда из (26) получаем $m = (2^M + 5)/7 = 0x92492493$, и код для деления на 14 имеет следующий вид.

shri	n, n, 1	Уменьшаем делимое вдвое
li	M, 0x92492493	Загружаем магическое число
mulhu	q, M, n	$q = \text{floor}(M \cdot n / 2^{**} 32)$.
shri	q, q, 2	$q = q/4$.

Эти методы *не* должны применяться *всегда*, когда делитель представляет собой четное число. Например, при делении на 10, 12, 18 или 22 лучше использовать метод, описанный в разделе, поскольку можно обойтись без команды очистки младших битов делимого или сдвига делимого вправо. Вместо этого следует использовать алгоритм из листинга 10.3 на с. 262, и если он дает значение индикатора "add", равное единице, а делитель четен, то тогда на большинстве машин можно прибегнуть к одному из описанных выше методов. Среди делителей, не превышающих 100, эти методы пригодны для 14, 28, 38, 42, 54, 56, 62, 70, 74, 76, 78, 84 и 90.

Какой же из методов лучше, (а) или (б)? Эксперименты показывают, что метод (б) предпочтительнее в плане количества выполняемых команд, поскольку, похоже, ему всегда требуется либо то же количество команд, что и методу (а), либо на одну меньше. Однако в некоторых случаях методам (а) и (б) требуется одинаковое количество команд, но при этом в методе (а) получается меньший множитель. Некоторые репрезентативные случаи показаны ниже. Метод "Книга" представляет собой код, получаемый с помощью листинга 10.3. Мы предполагаем, что команда компьютера *сложение с непосредственным значением* выполняет распространение старшего бита непосредственного значения.

$d = 6$

Книга		(а)	(б)		
li	M, 0xaaaaaaab	andi	n, n, -2	shri	n, n, 1
mulhu	q, M, n	li	M, 0x2aaaaaab	li	M, 0x55555556
shri	q, q, 2	mulhu	q, M, n	mulhu	q, M, n

$d = 28$

Книга		(а)	(б)		
li	M, 0x24924925	andi	n, n, -4	shri	n, n, 2
mulhu	q, M, n	li	M, 0x24924925	li	M, 0x24924925
add	q, q, n	mulhu	q, M, n	mulhu	q, M, n
shrxl	q, q, 5	shri	q, q, 2		

Эти методы не слишком пригодны для знакового деления. В этом случае разность между наилучшим и наихудшим кодом составляет только две команды (как показано в примерах кодов для деления на 3 и на 7 в разделе 10.3). Корректирующий код для метода (а) может потребовать прибавления единицы к делимому, если оно отрицательное и нечетное, и вычитания единицы, если оно неотрицательное и нечетное, что в результате потребует более двух команд. В случае метода (б) коррек-

тирующий код состоит в делении делимого на 2, что требует трех команд из базового набора RISC (см. раздел 10.1 на с. 231), так что этот метод также не является победителем в соревновании.

2. Ниже приведен код на языке программирования Python.

```
def magicg(nmax, d):
    nc = (nmax//d)*d - 1
    nbits = int(log(nmax, 2)) + 1
    for p in range(0, 2*nbits - 1):
        if 2**p > nc*(d - (2**p)%d):
            m = (2**p + d - (2**p)%d)//d
            return (m, p)
    print "Невозможно найти p, ошибка".
    sys.exit(1)
```

3. Поскольку $81 = 3^4$, нам нужно начальное значение — мультипликативное обратное к d по модулю 3. Это просто остаток от деления d на 3, поскольку $1 \cdot 1 \equiv 1 \pmod{3}$ и $2 \cdot 2 \equiv 1 \pmod{3}$ (и если остаток нулевой, то мультипликативное обратное не существует). Для $d = 146$ вычисления выполняются следующим образом.

$$\begin{aligned}x_0 &= 146 \bmod 3 = 2 \\x_1 &= 2(2 - 146 \cdot 2) = -580 \equiv 68 \pmod{81} \\x_2 &= 68(2 - 146 \cdot 68) = 647968 \equiv 5 \pmod{81} \\x_3 &= 5(2 - 146 \cdot 5) = -3640 \equiv 5 \pmod{81}\end{aligned}$$

При этом достигнута фиксированная точка, так что мультипликативным обратным к 146 по модулю 81 является 5. Проверим: $146 \cdot 5 = 730 \equiv 1 \pmod{81}$. В действительности *априори* известно, что достаточно двух итераций.

Глава 11. Некоторые элементарные функции

1. Да. Результат корректен, несмотря на двойное отсечение. Предположим, что $\lfloor \sqrt{x} \rfloor = a$. Тогда по определению данной операции a представляет собой целое число, такое, что $a^2 \leq x$ и $(a+1)^2 > x$.

Пусть $\lfloor \sqrt{a} \rfloor = b$. Тогда $b^2 \leq a$ и $(b+1)^2 > a$. Таким образом, $b^4 \leq a^2$ и, поскольку $a^2 \leq x$, $b^4 \leq x$.

Так как $(b+1)^2 > a$, $(b+1)^2 \geq a+1$, так что $(b+1)^4 \geq (a+1)^2$. Поскольку $(a+1)^2 > x$, $(b+1)^4 > x$. Следовательно, b представляет собой целочисленный корень четвертой степени из x .

Более просто это можно получить из упр. 5 к главе 9.

2. Простейший код имеет следующий вид.

```
int icbrt64(unsigned long long x)
{
    int s;
    unsigned long long y, b, bs;
```

```

y = 0;
for (s = 63; s >= 0; s = s - 3)
{
    y = 2*y;
    b = 3*y*(y + 1) + 1;
    bs = b << s;
    if (x >= bs && b == (bs >> s))
    {
        x = x - bs;
        y = y + 1;
    }
}
return y;
}

```

Переполнение b (bs в приведенном коде) может произойти только на второй итерации цикла. Следовательно, другой способ обработки переполнения состоит в раскрытии первых двух итераций цикла и выполнении цикла только начиная с $s = 57$ и далее с удаленным фрагментом “&& $b == (bs \gg s)$ ”.

Методом подбора выяснено, что результат первых двух итераций цикла следующий.

Если $x \geq 2^{63}$, устанавливается $x = x - 2^{63}$ и $y = 2$.

Если $2^{60} \leq x < 2^{63}$, устанавливается $x = x - 2^{60}$ и $y = 1$.

Если $x < 2^{60}$, x не изменяется, а y устанавливается равным нулю.

Таким образом, начало подпрограммы можно закодировать так, как показано ниже.

```

y = 0;
if (x >= 0x1000000000000000LL) {
    if (x >= 0x8000000000000000LL) {
        x = x - 0x8000000000000000LL;
        y = 2;
    } else {
        x = x - 0x1000000000000000LL;
        y = 1;
    }
}

for (s = 57; s >= 0; s = s - 3) {
    ...
}

```

И как упоминалось ранее, фрагмент “&& $b == (bs \gg s)$ ” может быть удален.

3. Шесть [67]. Метод бинарного разложения, основанный на $x^{23} = x^{16} \cdot x^4 \cdot x^2 \cdot x$, требует семи умножений. Разложения x^{23} как $(x^{11})^2 \cdot x$ и $((x^5)^2 \cdot x)^2 \cdot x$ также приводят к семи умножениям. Но вычисление степеней x в порядке $x^2, x^3, x^4, x^{10}, x^{13}, x^{23}$, где каждый член является произведением двух предыдущих или x , выполняется с помощью шести умножений.

4. (а) x округляется в сторону уменьшения до целой степени 2. (б) x округляется до целой степени 2 (в обоих случаях x остается неизменным, если является целочисленной степенью 2).

Глава 12. Системы счисления с необычными основаниями

1. Если B представляет собой бинарное число, а N — его эквивалент по основанию -2 , то

$$B \leftarrow 0x55555555 - (N \oplus 0x55555555) \text{ и}$$

$$N \leftarrow (0x55555555 - B) \oplus 0x55555555.$$

2. Простой способ выполнения поставленной задачи — преобразовать число x из системы по основанию -2 в бинарное, добавить единицу и выполнить обратное преобразование к основанию -2 . Применяя формулу Шреппеля и упрощая, получим

$$((x \oplus 0xAAAAAAAA) + 1) \oplus 0xAAAAAAAA \text{ или}$$

$$((x \oplus 0x55555555) - 1) \oplus 0x55555555.$$

3. Как и в упр. 1, можно выполнить преобразование числа x из системы по основанию -2 в бинарное, выполнить операцию $\&$ с $0xFFFFF0$, и выполнить обратное преобразование к основанию -2 . Все это выполняется с помощью пяти операций. Однако при применении любой из приведенных ниже формул можно обойтись четырьмя операциями.²

$$(((x \oplus 0xAAAAAAAA) - 10) \oplus 0xAAAAAAAA) \& -16$$

$$(((x \oplus 0x55555555) + 10) \oplus 0x55555555) \& -16$$

Формулы, приведенные далее, выполняют округление в сторону *увеличения* к ближайшей большей степени 16.

$$(((x \oplus 0xAAAAAAAA) + 5) \oplus 0xAAAAAAAA) \& -16$$

$$(((x \oplus 0x55555555) - 5) \oplus 0x55555555) \& -16$$

Имеются аналогичные формулы для округления в сторону уменьшения и увеличения для других степеней 2.

4. Это очень простая программа на языке программирования Python, так как данный язык программирования поддерживает комплексные числа.

```
import sys
import cmath

num = sys.argv[1:]
```

² Эти формулы найдены с помощью программы исчерпывающего поиска Aha! (A Hacker's Assistant).

```

if len(num) == 0:
    print "Преобразование в систему счисления с основанием"
    print "-1 + 1j числа, заданного в десятичной или"
    print "шестнадцатеричной системе счисления в виде"
    print "a + bj, где a и b - действительные числа".
    sys.exit()
num = eval(num[0])
r = 0
weight = 1
while num > 0:
    if num & 1:
        r = r + weight;
    weight = (-1 + 1j)*weight
    num = num >> 1;
print 'r =', r
    
```

5. Для изменения знака в системе счисления с основанием $-1+i$ нужно либо вычесть его из 0, либо умножить на -1 (11101) с применением правил арифметики с основанием $-1+i$.

Чтобы выделить действительную часть числа x , к нему следует добавить его мнимую часть с обратным знаком. Будем обрабатывать биты x группами по четыре начиная с правого (младшего) конца. Пронумеруем биты в каждой группе как 0, 1, 2 и 3 справа налево. Тогда получим следующее.

Если бит 1 равен единице, добавить $-i$ (0111) в позиции текущей группы.

Если бит 2 равен единице, добавить $2i$ (1110100) в позиции текущей группы.

Если бит 3 равен единице, добавить $-2i$ (0100) в позиции текущей группы.

Бит 1 имеет вес $-1+i$, так что добавление $-i$ сокращает мнимую часть. То же самое замечание применимо и к битам 2 и 3. Делать что-либо с битом 0 нет необходимости, поскольку он не имеет мнимой части. Каждая группа из четырех битов имеет вес, равный -4 , умноженному на вес предыдущей группы, поскольку 10000 в системе счисления с основанием $-1+i$ представляет собой -4 в десятичной системе счисления. Таким образом, вес бита n числа x равен действительному числу $(-4)^n$, умноженному на вес бита $n-4$.

Пример ниже показывает выделение действительной части числа 101101101 в системе счисления с основанием $-1+i$.

```

1 0110 1101   x
   111 0100   2i, добавленное для бита 2
      0100   -2i, добавленное для бита 3
     0111   -i(-4), добавленное для бита 5
111 0100   2i(-4), добавленное для бита 6
-----
1100 1101 1101   Сумма
    
```

Читатель может самостоятельно убедиться в том, что $x = 23 + 4i$, а сумма представляет собой 23. При выполнении суммирования генерируется множество непоказанных здесь переносов. Возможны определенные сокращения: если единичными являются одновременно и бит 2, и бит 3, то ничего добавлять не требуется,

так как это приводит к добавлению одновременно $2i$ и $-2i$. Если группа заканчивается битами 11, то эти биты можно просто отбросить, поскольку они представляют чисто мнимое число (i). Аналогично бит 2 также можно просто отбросить как имеющий мнимый вес $-2i$.

В пределе метод, применяющий все возможные сокращения, будет преобразовывать каждую группу из четырех битов в их действительную часть независимо. В некоторых случаях будет генерироваться перенос, и такие переносы будут добавляться к преобразуемому числу. Чтобы проиллюстрировать сказанное, представим каждую группу из четырех битов в шестнадцатеричном виде. Преобразование имеет следующий вид.

$0 \Rightarrow 0$	$4 \Rightarrow 0$	$8 \Rightarrow C$	$C \Rightarrow C$
$1 \Rightarrow 1$	$5 \Rightarrow 1$	$9 \Rightarrow D$	$D \Rightarrow 6$
$2 \Rightarrow 1D$	$6 \Rightarrow 1D$	$A \Rightarrow 1$	$E \Rightarrow 1$
$3 \Rightarrow 0$	$7 \Rightarrow 0$	$B \Rightarrow C$	$F \Rightarrow C$

Цифры 2 и 6 имеют действительную часть -1 , которая в системе счисления с основанием $-1+i$ записывается как 1D, т.е. при единичных указанных битах исходная цифра заменяется цифрой D с переносом 1. Переносы можно прибавлять с использованием арифметики в системе счисления с основанием $-1+i$, но в случае работы вручную имеется более подходящий способ. После преобразования имеются только четыре возможные цифры: 0, 1, C и D, как видно из приведенной выше таблицы. Правила прибавления 1 к этим цифрам показаны в левом столбце ниже.

$0+1=1$	$0+1D=1D$
$1+1=C$	$1+1D=0$
$C+1=D$	$C+1D=1$
$D+1=1D0$	$D+1D=C$

Прибавление 1 к D генерирует перенос 1D (так как $3+1=4$). Мы переносим обе цифры в один и тот же столбец. Как обрабатывать перенос 1D, показано в правом столбце. При выполнении сложения можно получить перенос и 1, и D в одном и том же столбце (первый перенос при преобразовании, а второй — при сложении). В этом случае переносы сокращаются друг с другом, так как 1D представляет собой -1 в системе счисления с основанием $-1+i$. Получить два переноса единиц или два переноса 1D в одном и том же столбце невозможно.

Приведенный ниже пример иллюстрирует использование этого метода для выделения действительной части из числа EA26 (записанного в шестнадцатеричном виде) в системе счисления с основанием $-1+i$.

EA26	x
11	Переносы при преобразовании
11DD	x с преобразованными шестнадцатеричными цифрами

110D	Сумма

Читатель может самостоятельно убедиться, что $x = -45 + 21i$, а сумма равна -45 .

Кстати, число в системе счисления с основанием $-1+i$ является действительным тогда и только тогда, когда, будучи записанным в шестнадцатеричном виде, оно состоит только из цифр 0, 1, C и D.

Для выделения мнимой части из x можно, конечно, выделить и вычесть из x действительную часть. Чтобы сделать то же непосредственно “сокращенным” методом, можно воспользоваться приведенной ниже таблицей преобразования каждой шестнадцатеричной цифры в чисто мнимую часть.

$0 \Rightarrow 0$	$4 \Rightarrow 4$	$8 \Rightarrow 74$	$C \Rightarrow 0$
$1 \Rightarrow 0$	$5 \Rightarrow 4$	$9 \Rightarrow 74$	$D \Rightarrow 0$
$2 \Rightarrow 3$	$6 \Rightarrow 7$	$A \Rightarrow 77$	$E \Rightarrow 3$
$3 \Rightarrow 3$	$7 \Rightarrow 7$	$B \Rightarrow 77$	$F \Rightarrow 3$

Таким образом, может произойти перенос 7, так что нам требуются правила прибавления 7 к четырем возможным преобразованным цифрам 0, 3, 4 и 7. Они показаны в левом столбце ниже.

$0+7=7$	$0+3=3$
$3+7=0$	$3+3=74$
$4+7=33$	$4+3=7$
$7+7=4$	$7+3=0$

Теперь может произойти перенос 3; как с ним поступить — показано в правом столбце.

Приведенный ниже пример иллюстрирует использование этого метода для выделения действительной части из числа 568A (записанного в шестнадцатеричном виде) в системе счисления с основанием $-1+i$.

568A	x
77	Переносы при преобразовании
4747	x с преобразованными шестнадцатеричными цифрами

4737	Сумма

Читатель может самостоятельно убедиться, что $x = -87 + 107i$, а сумма равна $107i$.

Число в системе счисления с основанием $-1+i$ является мнимым тогда и только тогда, когда, будучи записанным в шестнадцатеричном виде, оно состоит только из цифр 0, 3, 4 и 7.

Для преобразования числа в комплексно сопряженное нужно дважды вычесть из него мнимую часть. Можно прибегнуть, как и ранее, к таблице преобразования, но в этом случае преобразование может генерировать большее количество переносов, а преобразованное число состоять из любых шестнадцатеричных цифр. Таблица преобразования показана ниже.

$0 \Rightarrow 0$	$4 \Rightarrow 74$	$8 \Rightarrow 38$	$C \Rightarrow C$
$1 \Rightarrow 1$	$5 \Rightarrow 75$	$9 \Rightarrow 39$	$D \Rightarrow D$
$2 \Rightarrow 6$	$6 \Rightarrow 2$	$A \Rightarrow 3E$	$E \Rightarrow 3A$
$3 \Rightarrow 7$	$7 \Rightarrow 3$	$B \Rightarrow 3F$	$F \Rightarrow 3B$

Переносы можно прибавлять с использованием арифметики в системе счисления с основанием $-1+i$ либо можно разработать таблицу для выполнения сложения шестнадцатеричных цифр. Эта таблица больше предыдущих, поскольку переносы могут прибавляться ко всем возможным шестнадцатеричным цифрам.

Глава 13. Код Грея

1. набросок доказательства 1: это очевидно из способа построения отраженного бинарного кода Грея.

Набросок доказательства 2: из формулы $G(x) = x \oplus \left(x \gg 1 \right)$ видно, что $G(x)$ имеет 1 в позиции i , когда из позиции i к биту слева от нее имеется переход от 0 к 1 или от 1 к 0, а в противном случае имеет 0. Если x четно, имеется четное количество переходов, а если x нечетно, то и количество переходов нечетно.

Набросок доказательства 3: по индукции по длине x с применением приведенной выше формулы. Утверждение истинно для однобитовых слов 0 и 1. Пусть x — бинарное слово длиной n , и предположим, что доказываемое утверждение справедливо для x . Если к x спереди добавляется нулевой бит, то нулевой бит добавляется спереди и к $G(x)$, а остальные биты представляют собой $G(x)$. Если же к x спереди добавляется единичный бит, то единичный бит добавляется спереди и к $G(x)$, а следующий по старшинству бит инвертируется. Все прочие биты остаются неизменными. Следовательно, количество единичных битов в $G(x)$ либо увеличивается на 2, либо остается неизменным.

Набросок доказательства 3: по индукции по длине x с применением приведенной выше формулы. Утверждение истинно для однобитовых слов 0 и 1. Пусть x — бинарное слово длиной n , и предположим, что доказываемое утверждение справедливо для x . Если к x спереди добавляется нулевой бит, то нулевой бит добавляется спереди и к $G(x)$, а остальные биты представляют собой $G(x)$. Если же к x спереди добавляется единичный бит, то единичный бит добавляется спереди и к $G(x)$, а следующий по старшинству бит инвертируется. Все прочие биты остаются неизменными. Следовательно, количество единичных битов в $G(x)$ либо увеличивается на 2, либо остается неизменным.

Таким образом, можно построить генератор случайных чисел, которые генерирует целые числа с четным (или нечетным) количеством единичных битов с помощью генератора равномерно распределенных целых чисел, устанавливая младший бит равным 0 (или 1) и преобразуя результат в код Грея [5].

2. (а) Поскольку каждый столбец представляет собой циклический сдвиг первого столбца, мы тут же получаем требуемый результат.
(б) Такой код не существует. В этом нетрудно убедиться путем перечисления всех возможных кодов Грея для $n = 3$. Без потери общности можно начать с

000
001
011

поскольку любой код Грея можно сделать начинающимся таким образом путем дополнения и перестановки столбцов. Следствие: не существует кода Грея для одной дорожки для $n = 3$, в котором имеется восемь кодовых слов.

3. Представленный ниже код был получен путем отражения первых пяти кодовых слов отраженного бинарного кода Грея.

0000
0001
0011
0010
0110
1110
1010
1011
1001
1000

Другой код можно получить, если взять бинарно закодированные десятичные числа (BCD), увеличенные на 3, и преобразовать их в код Грея. Оказывается, получающийся в результате код цикличен. Такой код для кодирования десятичных чисел обладает тем свойством, что сложение кодовых слов генерирует переносы только тогда, когда их генерирует сложение десятичных цифр.

Код Грея с увеличением на 3

Десятичная цифра	Код с увеличением на 3	Эквивалентный код Грея
0	0011	0010
1	0100	0110
2	0101	0111
3	0110	0101
4	0111	0100
5	1000	1100
6	1001	1101
7	1010	1111
8	1011	1110
9	1100	1010

4. Достаточно просто вывести код Грея со “смешанным основанием”, используя принцип отражения. Для числа с разложением на простые множители вида $2^a 3^b 5^c \dots$ столбцы кода Грея должны иметь основания $e_1 + 1, e_2 + 1, e_3 + 1, \dots$. Например, для числа $72 = 2^3 \cdot 3^2$ приведенный ниже список демонстрирует код Грея “основание 4 — основание 3”; во втором столбце показаны делители 72, которые представляют кодовые слова.

00 1
01 3
02 9
12 18

11	6
10	2
20	4
21	12
22	36
32	72
31	24
30	8

Ясно, что каждый делитель получается из предыдущего с помощью единственного умножения или деления на простое число.

И даже проще: бинарный код Грея может использоваться для итерации по всем подмножествам множества таким образом, что на каждом шаге добавляется или удаляется только один член множества.

Глава 14. Циклический избыточный код

- Из раздела следует, что полином сообщения M и полином генератора G удовлетворяют условию $Mx' = QG + R$, где R представляет собой полином контрольной суммы. Пусть M' — полином сообщения, который отличается от M в члене x^e (т.е. бинарное сообщение отличается битом в позиции e). Тогда $M' = M + x^e$, и

$$M'x' = (M + x^e)x' = Mx' + x^{e+1} = QG + R + x^{e+1}.$$

Член x^{e+1} не делится на G , поскольку G состоит из двух или большего количества членов. (Делитель x^{e+1} должен иметь вид x^n .) Следовательно, остаток при делении $M'x'$ на G отличается от R , так что ошибка обнаруживается.

- Основной цикл можно закодировать так, как показано ниже, где `word` имеет тип `unsigned int` [22].

```

crc = 0xFFFFFFFF;
while (((word = *(unsigned int *)message) & 0xFF) != 0) {
    crc = crc ^ word;
    crc = (crc >> 8) ^ table[crc & 0xFF];
    crc = (crc >> 8) ^ table[crc & 0xFF];
    crc = (crc >> 8) ^ table[crc & 0xFF];
    crc = (crc >> 8) ^ table[crc & 0xFF];
    message = message + 4;
}

```

По сравнению с кодом в листинге 14.3 на с. 355 мы экономим три команды загрузки байта и три команды *исключающего или* для каждого слова `message`. Кроме того, выполняется меньше управляющих команд цикла.

Глава 15. Коды с коррекцией ошибок

- Ваша таблица должна выглядеть так, как табл. 15.1 на с. 359, с удаленными крайним справа столбцом и строками с нечетными номерами.
- В первом случае, если ошибка произошла в проверочном бите, получатель не может знать этого и выполнит неверную “коррекцию” информационных битов.

Во втором случае, если ошибка произошла в проверочном бите, синдром будет одним из k различных значений $100\dots 0$, $010\dots 0$, ..., $000\dots 1$. Следовательно, k должно быть достаточно велико, чтобы закодировать как эти k значений, так и m значений для кодирования одной ошибки в одном из m информационных битов, а также значение для случая отсутствия ошибок. Так что в результате выполняется правило Хэмминга.

3. Рассматривая k и m как действительные числа, мы получим быструю сходимость следующих итераций:

$$k_0 = 0,$$

$$k_{i+1} = \lg(k_i + m + 1), \quad i = 0, 1, \dots,$$

где $\lg(x)$ обозначает логарифм x по основанию 2. Корректный результат получается как $\text{ceil}(k_2)$, т.е. для всех $m \geq 0$ достаточно всего двух итераций.

Подходя с другой стороны, нетрудно доказать, что для $m \geq 0$

$$\text{bitsize}(m) \leq k \leq \text{bitsize}(m) + 1.$$

Здесь $\text{bitsize}(m)$ представляет размер m в битах ($\text{bitsize}(3) = 2$, $\text{bitsize}(4) = 3$ и так далее) (эти значения отличаются от функции с тем же именем, описанной в разделе 5.3, которая предназначена для знаковых целых чисел). *Указание:* $\text{bitsize}(m) = \lceil \lg(m+1) \rceil = \lfloor \lg(m) \rfloor + 1$, где мы принимаем $\lg(0) = -1$. Таким образом, можно взять $k = \text{bitsize}(m)$, испытать его и, если это значение слишком мало, просто добавить к нему единицу. При использовании функции для вычисления количества ведущих нулевых битов можно выполнить вычисления следующим образом:

$$k \leftarrow W - \text{nlz}(m),$$

$$k \leftarrow k + \left(((1 \ll k) - 1 - k) \ll m \right),$$

где W — размер машинного слова и $0 \leq m \leq 2^W - 1$.

4. Ответ: если $d(x, z) > d(x, y) + d(y, z)$, то это должно выполняться по крайней мере для одной битовой позиции i , которая вносит 1 в $d(x, z)$ и 0 в $d(x, y) + d(y, z)$. Отсюда вытекает, что $x_i \neq z_i$, но $x_i = y_i$ и $y_i = z_i$, т.е. очевидное противоречие.
5. В заданном коде длиной n и с минимальным расстоянием d просто удвоим каждую единицу и каждый ноль в каждом кодовом слове. В результате мы получим код длиной $2n$ с минимальным расстоянием $2d$ и того же размера, что и ранее.
6. Будем рассматривать заданный код длиной n с минимальным расстоянием d и размера $A(n, d)$, как если бы он был изображен в табл. 15.1 на с. 359. Удалим произвольные $d-1$ столбцов. Получившиеся в результате кодовые слова длиной $n - (d-1)$ имеют

минимальное расстояние, не меньше 1, т.е. все они различны. Следовательно, их количество не может превышать $2^{n-(d-1)}$, так что $A(n, d) \leq 2^{n-(d-1)}$.

7. Правило Хэмминга применимо к случаю, когда $d = 3$ и код содержит 2^m кодовых слов, где m — количество информационных битов. Правая часть неравенства (6), с $A(n, d) = 2^m$ и $d = 3$, имеет вид

$$2^m \leq \frac{2^n}{\binom{n}{0} + \binom{n}{1}} = \frac{2^n}{1+n}.$$

Заменяя n на $m+k$, получаем неравенство

$$2^m \leq \frac{2^{m+k}}{1+m+k},$$

которое после сокращения 2^m в обеих частях становится неравенством (1).

8. Код должен состоять из произвольной битовой строки и ее дополнения до единицы, так что размер кода равен 2. В том, что такие коды являются идеальными при нечетных n , можно убедиться, показав, что они достигают верхней границы в неравенстве (6). Вот набросок такого доказательства. n -битовое бинарное целое число можно рассматривать как представляющее единственный выбор из n объектов, где единичный бит означает выбор соответствующего объекта, а ноль — отказ от него. Таким образом, имеется 2^n способов выбора от 0 до n объектов из множества n объектов, т.е. $\sum_{i=0}^n \binom{n}{i} = 2^n$. Если n нечетно, i в диапазоне от 0 до $(n-1)/2$ охватывает половину членов этой суммы, а так как в силу симметрии $\binom{n}{i} = \binom{n}{n-i}$, суммирование дает половину всей суммы. Следовательно, $\sum_{i=0}^{(n-1)/2} \binom{n}{i} = 2^{n-1}$, так что верхняя граница (6) равна 2. Таким образом, рассматриваемый код достигает верхней границы (6).
9. Для простоты воспользуемся понятием эквивалентности кодов. Ясно, что код не изменяется существенным образом при перестановке его столбцов (как изображено в табл. 15.1 на с. 359) или при выполнении операции дополнения над любым его столбцом. Если один код может быть получен из другого путем таких преобразований, такие коды называются эквивалентными. Поскольку код представляет собой неупорядоченное множество кодовых слов, порядок вывода его кодовых слов значения не имеет. Путем выполнения операции дополнения столбцов любой код можно преобразовать в эквивалентный код, имеющий кодовое слово, состоящее только из нулевых битов.

Также для простоты проиллюстрируем доказательство с помощью случая $n = 9$ и $d = 6$.

Без потери общности пусть кодовым словом 0 (первым, которое мы обозначим как sw_0) является 000000000. Тогда все прочие кодовые слова должны иметь как минимум шесть единиц, чтобы отличаться от sw_0 как минимум в шести местах.

Предположим (что будет показано позже), что код имеет как минимум три кодовых слова. Тогда ни одно кодовое слово не может содержать семь или большее количество единиц. Если бы такое слово имелось, то другое кодовое слово (в котором обязательно имеется не менее шести единиц) имело бы как минимум четыре единицы в тех же позициях, что и слово с семью или большим количеством единиц. Это означает, что данные кодовые слова были бы равны не менее чем в четырех позициях, так что отличались бы они не более чем в пяти позициях ($9 - 4$), тем самым нарушая требование $d = 6$. Следовательно, все кодовые слова, отличные от первого, должны иметь ровно шесть единиц.

Без потери общности переставим столбцы так, чтобы первыми двумя кодовыми словами были

$$sw_0 : 000\ 000\ 000$$

$$sw_1 : 111111\ 000$$

Следующее кодовое слово, sw_2 , не может иметь четыре или более своих единиц в шести левых столбцах, поскольку тогда оно будет совпадать с sw_1 не менее чем в четырех позициях, так что отличаться от sw_1 оно будет не более чем в пяти позициях. Следовательно, в левых шести столбцах у этого слова находятся ровно три единицы. Переставим шесть левых столбцов (всех трех кодовых слов) так, чтобы sw_2 имело следующий вид:

$$sw_2 : 111\ 000\ 111$$

Аналогично очередное кодовое слово sw_3 не может иметь четыре своих единицы одновременно среди трех левых и трех правых позиций, потому что тогда оно будет совпадать с sw_2 в четырех позициях. Следовательно, оно имеет три или менее единиц в левой и правой тройке, так что в средних трех позициях оно должно иметь три единицы. Путем аналогичного сравнения с sw_1 мы заключаем, что оно должно иметь три единицы в трех правых позициях, т.е.

$$sw_3 : 000\ 111\ 111$$

Сравнивая очередное кодовое слово, если таковое возможно, с sw_1 , мы заключаем, что оно должно иметь три единицы в трех правых позициях, а сравнение с sw_2 приводит к требованию трех единиц в трех средних позициях. Таким образом, это слово должно иметь вид 000111111, т.е. совпадать с sw_3 . Следовательно, пятое кодовое слово невозможно. Полученные методом подбора четыре приведенных выше кодовых слова удовлетворяют условию $d = 6$, так что $A(9,6) = 4$.

10. Очевидно, что $A(n,d)$ имеет значение, не меньшее 2, поскольку два кодовых слова могут быть состоящими из одних только нулевых или только единичных битов. Рас-

суждая, как и в предыдущем упражнении, положим одно кодовое слово, sw_0 , состоящим из одних нулей. Тогда каждое прочее кодовое слово должно иметь более $2n/3$ единиц. Если код включает три или большее количество кодовых слов, то любые два кодовых слова, отличные от sw_0 , должны иметь единицы в одних и тех же позициях более чем в $2n/3 - n/3 = n/3$ случаях, как показано на рисунке ниже.

$$\begin{array}{c} 1111\dots11110\dots0 \\ > 2n/3 < n/3 \end{array}$$

(Здесь представлено кодовое слово sw_1 со всеми единицами, сдвинутыми влево. Теперь представьте размещение более чем $2n/3$ единиц кодового слова sw_2 так, чтобы минимизировать перекрытие единиц.) Поскольку sw_1 и sw_2 перекрываются более чем в $n/3$ позициях, они могут отличаться менее чем в $n - n/3 = 2n/3$ позициях, так что получающееся в результате расстояние меньше $n/3$.

11. Это SEC-DED код, поскольку минимальное расстояние между кодовыми словами равно 4. Чтобы убедиться в этом, предположим, что первые два кодовых слова отличаются одним информационным битом. Тогда в дополнение к информационному биту в этих двух кодовых словах будут отличаться четность строк, четность столбцов и угловой проверочный бит, так что расстояние между кодовыми словами оказывается равным 4. Если информационные слова отличаются двумя битами, находящимися в разных строках, то в этих двух кодовых словах будет одинаковая четность строк, но четность столбцов будет различна в двух битах. Следовательно, расстояние между этими словами также равно 4. Тот же результат получается, если отличающиеся биты находятся в одном столбце. Если же эти биты находятся в разных строках и в разных столбцах, расстояние между кодовыми словами становится равным 6. Наконец, если информационные слова различаются тремя битами, легко убедиться, что, независимо от распределения последних по строкам и столбцам, будет отличаться как минимум один бит четности. Следовательно, в этом случае расстояние не меньше 4.

Если угловой бит не используется, минимальное расстояние равно 3. Следовательно, в этом случае получается просто SEC-код, но не SEC-DED.

Независимо от того, представляет ли угловой бит сумму строк или сумму столбцов, он является суммой всех 64 информационных битов по модулю 2, так что его значение одинаково в обоих случаях.

Этот код использует 17 проверочных битов, в то время как коду Хэмминга требуются только 8 (см. табл. 15.3 на с. 362), так что в этом отношении данный код не очень эффективен.

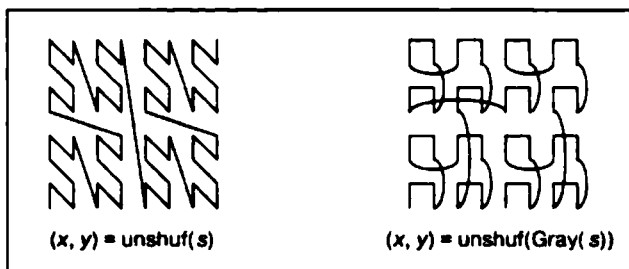
Однако он *эффективен* при обнаружении пакетных ошибок. Представим, что массив 9×9 передается через последовательный битовый канал построчно, начиная с первой строки. Тогда любая последовательность из десяти или меньшего количества ошибок находится в одной или двух строках не более чем с одним перекры-

вающимся битом. Следовательно, если ошибки при передаче сосредоточены в подмножестве из десяти последовательных битов, в большинстве случаев ошибка будет обнаружена при проверке четностей столбцов либо четностей строк, если ошибочны только первый и десятый биты.

Ошибка, которая *не* обнаруживается кодом, — четыре поврежденных бита, расположенных прямоугольником.

Глава 16. Кривая Гильберта

1 и 2.



Среднее расстояние перемещения для обхода, показанного в левой части, составляет примерно 1.46. То же значение для обхода справа составляет около 1.33. Следовательно, применение кода Грея, похоже, повышает локальность, по крайней мере с точки зрения данного параметра. (В случае кривой Гильберта все переходы выполняются на единичное расстояние.)

По предложению Эдсгера Дейкстры (Edsger Dijkstra) алгоритм тасования использовался в раннем компиляторе языка программирования Algol при отображении матриц на магнитные диски. Цель заключалась в снижении количества страничных операций при обращении матрицы. Похоже, что этот алгоритм был независимо открыт многими исследователями.

3. Воспользуйтесь каждым третьим битом s .

Глава 17. Числа с плавающей точкой

1. ± 0 , ± 2.0 и некоторые NaN.
2. Да! Программу легко получить, если заметить, что если $x = 2^n(1+f)$, то

$$\sqrt{x} = 2^{n/2}(1+f)^{1/2}.$$

Игнорируя дробное значение, мы видим, что должны заменить смещенный показатель степени $127+n$ на $127+n/2$, что представляет собой $(127+n)/2+127/2$. Таким образом, похоже, что грубое приближение \sqrt{x} получается путем сдвига $\text{гер}(x)$ на одну позицию вправо с добавлением 63 в позиции показателя степени, что равно $0x1F800000$. Это приближение,

$$\text{гер}(\sqrt{x}) \approx k + \left(\text{гер}(x) \gg 1 \right),$$

также обладает тем свойством, что если мы найдем оптимальное значение k для значений x от 1.0 до 4.0, то то же самое значение k будет оптимальным для всех нормальных чисел. После уточнения значения k с помощью программы, которая находит максимальную и минимальную ошибку для заданного значения k мы получаем приведенную ниже программу, которая включает один шаг итерации алгоритма Ньютона-Рафсона.

```
float asqrt(float x0)
{
    union {int ix; float x;};

    x = x0;                // x можно рассматривать как int
    ix = 0x1fbb67a8 +
        (ix >> 1);        // Начальная оценка
    x = 0.5f*(x + x0/x);  // Шаг алгоритма Ньютона
    return x;
}
```

Для нормальных чисел относительная ошибка лежит в диапазоне от 0 до 0.000601. Эта функция дает точные результаты для $x = \text{inf}$ и $x = \text{NaN}$ (inf и NaN соответственно). Для $x = 0$ результат приблизительно равен $4.0 \cdot 10^{-20}$. При $x = -0$ результат оказывается бесполезным: $-1.35 \cdot 10^{19}$. Для положительных денормализованных значений x результат либо находится в указанных ранее границах, либо представляет собой положительное число, меньшее 10^{-19} .

Шаг алгоритма Ньютона использует деление, так что на большинстве машин программа не такая быстрая, как вычисление обратного к квадратному корню.

Если добавить второй шаг алгоритма Ньютона, относительная ошибка для нормальных чисел оказывается в диапазоне от 0 до 0.00000023. Оптимальной является константа 0x1FBB3F80. Если не использовать алгоритм Ньютона вовсе, относительная ошибка немного меньше ± 0.035 при использовании константы 0x1FBB4F2E. Это примерно та же относительная ошибка, что и в программе вычисления обратного к квадратному корню без шага Ньютона, и так же, как и она, использует только две целочисленные операции.

3. Да, вычислить кубический корень положительного нормального числа можно практически таким же способом. Ключевым является первое приближение

```
i = 0x2a51067f + i/3;    // Первое приближение
```

Таким образом можно вычислить кубический корень с относительной ошибкой около ± 0.0316 .

Деление на 3 можно приближенно выполнить как

$$\frac{i}{3} \approx \frac{i}{4} + \frac{i}{16} + \frac{i}{64} + \dots + \frac{i}{65536}$$

(где деление на степень 2 реализуется с помощью правого сдвига), для чего требуется семь команд. Можно также немного повысить точность так, как показано в приведенной ниже программе. (Этот фокус с делением рассматривался в разделе 10.18.)

```
float acbrt(float x0)
{
    union
    {
        int ix;
        float x;
    };

    x = x0; // x можно рассматривать как int
    ix = ix/4 + ix/16; // Приближенное деление на 3
    ix = ix + ix/16;
    ix = ix + ix/256;
    ix = 0x2a5137a0 + ix; // Первое приближение
    x = 0.33333333f* // Шаг алгоритма Ньютона
        (2.0f*x + x0/(x*x));
    return x;
}
```

Хотя мы и избежали деления на 3 (ценой семи элементарных целочисленных команд), в шаге алгоритма Ньютона имеются как деление, так и еще четыре другие команды. Относительная ошибка находится в диапазоне от 0 до примерно +0.00103. Таким образом, этот метод не столь успешен, как вычисление обратного к квадратному корню и вычисление квадратного корня, но может быть полезен в некоторых ситуациях.

Если добавить еще один шаг алгоритма Ньютона при использовании той же константы, относительная ошибка окажется в диапазоне от 0 до примерно +0.00000116.

4. Да. Приведенная ниже программа вычисляет обратное к квадратному корню для чисел с плавающей точкой двойной точности, с точностью около $\pm 3.5\%$. Эту точность можно увеличить с помощью одного или двух шагов алгоритма Ньютона-Рафсона. Использование константы `0x5FE80...0` дает относительную ошибку в диапазоне от 0 до примерно +0.887, а константа `0x5FE618FDF80...0` — относительную ошибку в диапазоне от 0 до примерно -0.0613.

```
double rsqrtdd(double x0)
{
    union
    {
        long long ix;
        double x;
    };

    x = x0;
    ix = 0x5fe6ec85e8000000LL - (ix >> 1);
    return x;
}
```

Глава 18. Формулы для простых чисел

1. Пусть $f(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_0$. Такой полином монотонно стремится к бесконечности при x , стремящемся к бесконечности. (При достаточно больших значениях x первый член превосходит по величине сумму остальных.)

Пусть x_0 — целое число, такое, что $|f(x)| \geq 2$ для всех $x > x_0$. Пусть $f(x_0) = k$ и пусть r — произвольное положительное целое число. Тогда $|k| \geq 2$, и

$$\begin{aligned} |f(x_0 + rk)| &= |a_n (x_0 + rk)^n + a_{n-1} (x_0 + rk)^{n-1} + \dots + a_0| \\ &= |f(x_0) + \text{кратное } rk| \\ &= |k + \text{кратное } rk|. \end{aligned}$$

Таким образом, с ростом r $|f(x_0 + rk)|$ пробегает по множеству составных чисел с увеличивающимися значениями, а следовательно, различными. Следовательно, $f(x)$ принимает бесконечное множество составных значений.

Другая формулировка теоремы: не существует неконстантного полинома от одной переменной, который принимает только простые значения, даже при достаточно больших значениях своего аргумента.

Пример. Пусть $f(x) = x^2 + x + 41$. Тогда $f(1) = 43$ и

$$\begin{aligned} f(1+43r) &= (1+43r)^2 + (1+43r) + 41 \\ &= (1+86+43^2 r^2) + (1+43r) + 41 \\ &= 1+1+43+86r+43^2 r^2 + 43r \\ &= 43 + (2+43r+1) \cdot 43r, \end{aligned}$$

что со всей очевидностью с увеличением r дает еще быстрее возрастающие числа, кратные 43.

2. Предположим, что p — составное. Запишем тождество как

$$(p-1)! = pk - 1$$

для некоторого целого значения k . Пусть a — собственный делитель p . Тогда a делит левую часть, но не правую, так что равенство не может быть выполнено.

Легко убедиться в справедливости теоремы для $p = 1, 2$ и 3 . Предположим, что p — простое число, большее 3. Тогда в факториале

$$(p-1)! = (p-1)(p-2) \cdots (3)(2)$$

первый член, $p-1$, сравним с -1 по модулю p . Каждый из остальных членов взаимно простой с p и, таким образом, имеет мультипликативное обратное по модулю p (см. раздел 10.16), а кроме того, это обратное является единственным и не равно самому члену.

Чтобы убедиться в том, что мультипликативное обратное по простому модулю не равно исходному значению (за исключением 1 и $p-1$), предположим, что $a^2 \equiv 1 \pmod{p}$. Тогда $a^2 - 1 \equiv 0 \pmod{p}$, так что $(a+1)(a-1) \equiv 0 \pmod{p}$. Поскольку p — простое, либо $a-1$, либо $a+1$ сравнимо с 0 по модулю p . В первом случае $a \equiv 1 \pmod{p}$, а в последнем $a \equiv -1 \equiv p-1 \pmod{p}$.

Следовательно, целые числа $p-2, p-3, \dots, 2$ можно разобрать попарно так, что произведение каждой пары сравнимо с 1 по модулю p , т.е.

$$(p-1)! = (p-1)(ab)(cd)\dots,$$

где a и b мультипликативно обратны одно другому, то же самое относится к c и d , и т.д. Таким образом,

$$(p-1)! \equiv (-1)(1)(1)\dots \equiv -1 \pmod{p}.$$

Например, при $p=11$

$$\begin{aligned} 10! \pmod{11} &\equiv 10 \cdot 9 \cdot 8 \cdot 7 \cdot 6 \cdot 5 \cdot 4 \cdot 3 \cdot 2 \pmod{11} \\ &\equiv 10 \cdot (9 \cdot 5)(8 \cdot 7)(6 \cdot 2)(4 \cdot 3) \pmod{11} \\ &\equiv (-1)(1)(1)(1)(1) \pmod{11} \\ &\equiv -1 \pmod{11}. \end{aligned}$$

Теорема носит имя Джона Вильсона (John Wilson), ученика английского математика Эдуарда Уоринга (Edward Waring). Уоринг сформулировал ее без доказательства в 1770 году. Первым доказательство этой теоремы в 1773 году опубликовал Лагранж (Lagrange). Теорема была известна в средневековой Европе еще около 1000 года н.э.

3. Если $n = ab$, где a и b различны и ни одно из них не равно ни 1, ни n , то очевидно, что a и b меньше n , а значит, являются сомножителями $(n-1)!$. Следовательно, n делит $(n-1)!$.

Если $n = a^2$, то для $a > 2$ справедливо $a^2 = n > 2a$, так что a и $2a$ являются сомножителями $(n-1)!$, а значит, a^2 делит $(n-1)!$.

4. Вероятно, это один из тех случаев, когда вычисления привносят в математику больше, чем формальные доказательства.

Согласно теореме Миллса существует действительное число θ , такое, что $\lfloor \theta^n \rfloor$ простое для всех целых $n \geq 1$. Исследуем возможность того, что при $n=1$ это простое число 2. Тогда

$$\lfloor \theta^1 \rfloor = 2,$$

так что

$$\begin{aligned}
 2 \leq \theta^3 < 3, \text{ или} & \quad (1) \\
 2^{V^3} \leq \theta < 3^{V^3}, \text{ или} \\
 1.2599... \leq \theta < 1.4422....
 \end{aligned}$$

Возведение неравенства (1) в куб дает

$$8 \leq \theta^3 < 27. \quad (2)$$

В этом диапазоне имеется простое число. (Из нашего исходного предположения между 2^3 и $(2+1)^3$ имеется простое число.) Выберем в качестве второго простого числа 11. Тогда мы получим $\lfloor \theta^3 \rfloor = 11$ при ограничении (2) до

$$11 \leq \lfloor \theta^3 \rfloor < 12. \quad (3)$$

Продолжим работу, возводя (3) в куб и получая

$$1331 \leq \theta^{3^2} < 1728. \quad (4)$$

Мы знаем, что между 1331 и 1728 имеется простое число. Выберем наименьшее — 1361. Тогда (4) ограничивается до

$$1361 \leq \theta^{3^2} < 1362.$$

Продолжая, мы должны показать, что существует действительное число θ , такое, что $\lfloor \theta^{3^n} \rfloor$ является простым числом при $n = 1, 2$ и 3 , и, вычисляя корень 27-й степени из $\lfloor 1361$ и 1362 , находим, что θ находится между 1.30637 и 1.30642.

Очевидно, что процесс можно продолжить. Можно показать, что существует предельное значение θ , но в действительности это необязательно. Если в пределе θ представляет собой произвольное число из некоторого конечного диапазона, это также подтверждает теорему Миллса.

Проведенные вычисления показывают, что теорема Миллса несколько “натянута” — будучи формулой для простых чисел, она требует предварительного знания этих простых чисел для вычисления θ . Она похожа на формулу для простых чисел на с. 421, в которую входит константа

$$a = 0.203005000700011000013....$$

Ясно, что эта теорема имеет малое отношение к простым числам. Аналогичная теорема справедлива для любой возрастающей последовательности при условии достаточной ее плотности.

Приведенные выше шаги вычисляют наименьшее значение θ , удовлетворяющее теореме Миллса. Оно иногда называется константой Миллса, и для нее вычислено более 6850 десятичных цифр [15].

5. Предположим, что существуют такие целые числа a, b, c и d , что

$$(a + b\sqrt{-5})(c + d\sqrt{-5}) = 2. \quad (5)$$

Приравнивая действительные и мнимые части, находим

$$ac - 5bd = 2 \text{ и} \quad (6)$$

$$ad + bc = 0. \quad (7)$$

Ясно, что $c \neq 0$, поскольку если $c = 0$, то из (6) $-5bd = 2$, что не имеет решения в целых числах.

Также $b \neq 0$, поскольку при $b = 0$ из (7) вытекает, что либо a , либо d равно нулю. Но $a = 0$ не позволяет удовлетворить уравнение (5). Следовательно, $d = 0$, но тогда (5) превращается в $ac = 2$, так что один из множителей в (5) должен быть равен единице, что делает разложение неприемлемым.

Из (7) $abd + b^2c = 0$. Из (6) $a^2c - 5abd = 2a$. Объединяя, получим $a^2c + 5b^2c = 2a$, или

$$a^2 + 5b^2 = 2a/c \quad (8)$$

(вспомним, что $c \neq 0$). Левая часть (8) имеет значение, не меньшее, чем $a^2 + 5$, что превышает $2a/c$, какими бы ни были значения a и c .

Чтобы увидеть, что 3 — простое, можно аналогично вывести уравнение

$$a^2 + 5b^2 = 3a/c,$$

в котором $b \neq 0$ и $c \neq 0$. Оно также не может быть удовлетворено ни при каких целых числах.

Число 6 имеет два различных разложения на простые сомножители:

$$6 = 2 \cdot 3 = (1 + \sqrt{-5})(1 - \sqrt{-5}).$$

Мы не показали, что $1 \pm \sqrt{-5}$ — простые. Это можно показать с помощью рассуждений, подобных приведенным выше (хотя и более длинных), но на самом деле в этом нет необходимости, чтобы продемонстрировать, что разложения на простые сомножители в этом кольце не единственные. Это связано с тем, что даже если указанные числа можно разложить на простые сомножители, общее разложение все равно не будет иметь вид $2 \cdot 3$.

ПРИЛОЖЕНИЕ А

АРИФМЕТИЧЕСКИЕ ТАБЛИЦЫ ДЛЯ 4-БИТОВОЙ МАШИНЫ

Во всех таблицах данного приложения подчеркивание обозначает знаковое переполнение. Например, в табл. А.1 $7+1=8$, что не представимо в виде знакового целого числа на 4-разрядной машине, так что происходит знаковое переполнение.

ТАБЛИЦА А.1. СЛОЖЕНИЕ

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
1	1	2	3	4	5	6	7	<u>8</u>	9	A	B	C	D	E	F	10
2	2	3	4	5	6	7	<u>8</u>	<u>9</u>	A	B	C	D	E	F	10	11
3	3	4	5	6	7	<u>8</u>	<u>9</u>	<u>A</u>	B	C	D	E	F	10	11	12
4	4	5	6	7	<u>8</u>	<u>9</u>	<u>A</u>	<u>B</u>	C	D	E	F	10	11	12	13
5	5	6	7	<u>8</u>	<u>9</u>	<u>A</u>	<u>B</u>	<u>C</u>	D	E	F	10	11	12	13	14
6	6	7	<u>8</u>	<u>9</u>	<u>A</u>	<u>B</u>	<u>C</u>	<u>D</u>	E	F	10	11	12	13	14	15
7	7	<u>8</u>	<u>9</u>	<u>A</u>	<u>B</u>	<u>C</u>	<u>D</u>	<u>E</u>	F	10	11	12	13	14	15	16
-8 8	8	9	A	B	C	D	E	F	<u>10</u>	<u>11</u>	<u>12</u>	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>	<u>17</u>
-7 9	9	A	B	C	D	E	F	10	<u>11</u>	<u>12</u>	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>	<u>17</u>	18
-6 A	A	B	C	D	E	F	10	11	<u>12</u>	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>	<u>17</u>	18	19
-5 B	B	C	D	E	F	10	11	12	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>	<u>17</u>	18	19	1A
-4 C	C	D	E	F	10	11	12	13	<u>14</u>	<u>15</u>	<u>16</u>	<u>17</u>	18	19	1A	1B
-3 D	D	E	F	10	11	12	13	14	<u>15</u>	<u>16</u>	<u>17</u>	18	19	1A	1B	1C
-2 E	E	F	10	11	12	13	14	15	<u>16</u>	<u>17</u>	18	19	1A	1B	1C	1D
-1 F	F	10	11	12	13	14	15	16	<u>17</u>	18	19	1A	1B	1C	1D	1E

В таблице для вычитания (табл. А.2) предполагается, что бит переноса для $a-b$ устанавливается как и для $a+\bar{b}+1$, так что перенос эквивалентен отсутствию заема.

ТАБЛИЦА А.2. ВЫЧИТАНИЕ (СТРОКА – СТОЛБЕЦ)

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	10	F	E	D	C	B	A	9	8	7	6	5	4	3	2	1
1	11	10	F	E	D	C	B	A	<u>9</u>	<u>8</u>	7	6	5	4	3	2
2	12	11	10	F	E	D	C	B	<u>8</u>	<u>9</u>	7	6	5	4	3	2
3	13	12	11	10	F	E	D	C	B	<u>8</u>	<u>9</u>	7	6	5	4	3
4	14	13	12	11	10	F	E	D	C	B	<u>8</u>	<u>9</u>	7	6	5	4
5	15	14	13	12	11	10	F	E	D	C	B	<u>8</u>	<u>9</u>	7	6	5
6	16	15	14	13	12	11	10	F	E	D	C	B	<u>8</u>	<u>9</u>	7	6
7	17	16	15	14	13	12	11	10	F	E	D	C	B	<u>8</u>	<u>9</u>	7
-8 8	18	<u>17</u>	<u>16</u>	<u>15</u>	<u>14</u>	<u>13</u>	<u>12</u>	<u>11</u>	10	F	E	D	C	B	A	9
-7 9	19	18	<u>17</u>	<u>16</u>	<u>15</u>	<u>14</u>	<u>13</u>	<u>12</u>	11	10	F	E	D	C	B	A
-6 A	1A	19	18	<u>17</u>	<u>16</u>	<u>15</u>	<u>14</u>	<u>13</u>	12	11	10	F	E	D	C	B
-5 B	1B	1A	19	18	<u>17</u>	<u>16</u>	<u>15</u>	<u>14</u>	13	12	11	10	F	E	D	C
-4 C	1C	1B	1A	19	18	<u>17</u>	<u>16</u>	<u>15</u>	14	13	12	11	10	F	E	D
-3 D	1D	1C	1B	1A	19	18	<u>17</u>	<u>16</u>	15	14	13	12	11	10	F	E
-2 E	1E	1D	1C	1B	1A	19	18	<u>17</u>	16	15	14	13	12	11	10	F
-1 F	1F	1E	1D	1C	1B	1A	19	18	17	16	15	14	13	12	11	10

В случае умножения (табл. А.3 и А.4) переполнение означает, что результат не может быть выражен 4-битовой величиной. В случае знакового умножения (см. табл. А.3) это эквивалентно тому, что первые 5 бит 8-битового результата не равны пяти единицам или пяти нулям.

ТАБЛИЦА А.3. ЗНАКОВОЕ УМНОЖЕНИЕ

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	1	2	3	4	5	6	7	F8	F9	FA	FB	FC	FD	FE	FF
2	0	2	4	6	8	A	C	E	F0	F2	F4	F6	F8	FA	FC	FE
3	0	3	6	9	C	F	12	15	E8	EB	EE	F1	F4	F7	FA	FD
4	0	4	8	C	10	14	18	1C	E0	E4	E8	EC	F0	F4	F8	FC
5	0	5	A	F	14	19	1E	23	D8	DD	E2	E7	EC	F1	F6	FB
6	0	6	C	12	18	1E	24	2A	D0	D6	DC	E2	E8	EE	F4	FA
7	0	7	E	15	1C	23	2A	31	C8	CF	D6	DD	E4	EB	F2	F9
-8	8	0	F8	F0	E8	E0	D8	D0	C8	40	38	30	28	20	18	10
-7	9	0	F9	F2	EB	E4	DD	D6	CF	38	31	2A	23	1C	15	E
-6	A	0	FA	F4	EE	E8	E2	DC	D6	30	2A	24	1E	18	12	C
-5	B	0	FB	F6	F1	EC	E7	E2	DD	28	23	1E	19	14	F	A
-4	C	0	FC	F8	F4	F0	EC	E8	E4	20	1C	18	14	10	C	B
-3	D	0	FD	FA	F7	F4	F1	EE	EB	18	15	12	F	C	9	6
-2	E	0	FE	FC	FA	F8	F6	F4	E2	10	E	C	A	8	6	4
-1	F	0	FF	FE	FD	FC	FB	FA	F9	8	7	6	5	4	3	2

ТАБЛИЦА А.4. БЕЗЗНАКОВОЕ УМНОЖЕНИЕ

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
2	0	2	4	6	8	A	C	E	10	12	14	16	18	1A	1C	1E
3	0	3	6	9	C	F	12	15	18	1B	1E	21	24	27	2A	2D
4	0	4	8	C	10	14	18	1C	20	24	28	2C	30	34	38	3C
5	0	5	A	F	14	19	1E	23	28	2D	32	37	3C	41	46	4B
6	0	6	C	12	18	1E	24	2A	30	36	3C	42	48	4E	54	5A
7	0	7	E	15	1C	23	2A	31	38	3F	46	4D	54	5B	62	69
8	0	8	10	18	20	28	30	38	40	48	50	58	60	68	70	78
9	0	9	12	1B	24	2D	36	3F	48	51	5A	63	6C	75	7E	87
A	0	A	14	1E	28	32	3C	46	50	5A	64	6E	78	82	8C	96
B	0	B	16	21	2C	37	42	4D	58	63	6E	79	84	8F	9A	A5
C	0	C	18	24	30	3C	48	54	60	6C	78	84	90	9C	A8	B4
D	0	D	1A	27	34	41	4E	5B	68	75	82	8F	9C	A9	B6	C3
E	0	E	1C	2A	38	46	54	62	70	7E	8C	9A	A8	B6	C4	D2
F	0	F	1E	2D	3C	4B	5A	69	78	87	96	A5	B4	C3	D2	E1

В табл. А.7 и А.8 содержатся остатки при делении с отсечением. В табл. А.7 для случая деления максимального по модулю отрицательного числа на -1 показан результат, равный 0 с переполнением, хотя на большинстве машин такая операция либо будет запрещена, либо ее результат будет неопределенным.

ТАБЛИЦА А.7. ОСТАТОК ПРИ ЗНАКОВОМ КОРОТКОМ ДЕЛЕНИИ (СТРОКА + СТОЛБЕЦ)

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	-	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	-	0	1	1	1	1	1	1	1	1	1	1	1	1	1	0
2	-	0	0	2	2	2	2	2	2	2	2	2	2	2	2	0
3	-	0	1	0	3	3	3	3	3	3	3	3	3	3	0	1
4	-	0	0	1	0	4	4	4	4	4	4	4	4	0	1	0
5	-	0	1	2	1	0	5	5	5	5	5	0	1	2	1	0
6	-	0	0	0	2	1	0	6	6	6	0	1	2	0	0	0
7	-	0	1	1	3	2	1	0	7	0	1	2	3	1	1	0
-8 8	-	0	0	E	0	D	E	F	0	F	E	D	0	E	0	0
-7 9	-	0	F	F	D	E	F	0	9	0	F	E	D	F	F	0
-6 A	-	0	0	0	E	F	0	A	A	A	0	F	E	0	0	0
-5 B	-	0	F	E	F	0	B	B	B	B	B	0	F	E	F	0
-4 C	-	0	0	F	0	C	C	C	C	C	C	C	0	F	0	0
-3 D	-	0	F	0	D	D	D	D	D	D	D	D	D	0	F	0
-2 E	-	0	0	E	E	E	E	E	E	E	E	E	E	E	0	0
-1 F	-	0	F	F	F	F	F	F	F	F	F	F	F	F	F	0

ТАБЛИЦА А.8. ОСТАТОК ПРИ БЕЗЗНАКОВОМ КОРОТКОМ ДЕЛЕНИИ (СТРОКА + СТОЛБЕЦ)

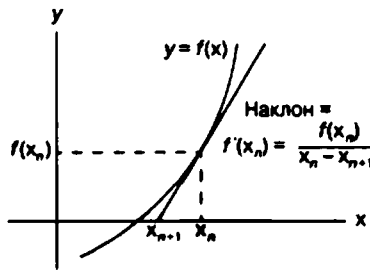
	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	-	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	-	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	-	0	0	2	2	2	2	2	2	2	2	2	2	2	2	2
3	-	0	1	0	3	3	3	3	3	3	3	3	3	3	3	3
4	-	0	0	1	0	4	4	4	4	4	4	4	4	4	4	4
5	-	0	1	2	1	0	5	5	5	5	5	5	5	5	5	5
6	-	0	0	0	2	1	0	6	6	6	6	6	6	6	6	6
7	-	0	1	1	3	2	1	0	7	7	7	7	7	7	7	7
8	-	0	0	2	0	3	2	1	0	8	8	8	8	8	8	8
9	-	0	1	0	1	4	3	2	1	0	9	9	9	9	9	9
A	-	0	0	1	2	0	4	3	2	1	0	A	A	A	A	A
B	-	0	1	2	3	1	5	4	3	2	1	0	B	B	B	B
C	-	0	0	0	0	2	0	5	4	3	2	1	0	C	C	C
D	-	0	1	1	1	3	1	6	5	4	3	2	1	0	D	D
E	-	0	0	2	2	4	2	0	6	5	4	3	2	1	0	E
F	-	0	1	0	3	0	3	1	7	6	5	4	3	2	1	0

ПРИЛОЖЕНИЕ Б МЕТОД НЬЮТОНА

Для краткого рассмотрения метода Ньютона предположим, что имеется дифференцируемая функция f от действительной переменной x и нужно решить уравнение $f(x) = 0$ относительно x . Метод Ньютона по данному приближению x_n корня f дает при определенных условиях улучшенное приближение x_{n+1} по формуле

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}.$$

Здесь $f'(x_n)$ означает производную f в точке $x = x_n$. Вывод этой формулы можно пояснить приведенным ниже рисунком (решая показанное уравнение относительно x_{n+1}).



Этот метод хорошо работает для простых доброкачественных функций, например полиномов, если первое приближение достаточно близко к решению. Если приближенное значение достаточно близко к точному, метод обладает квадратичной сходимостью. Это означает, что если r — точное значение корня, а x_n — достаточно близкая оценка, то

$$|x_{n+1} - r| \leq (x_n - r)^2.$$

Таким образом, количество точных цифр решения удваивается при каждой итерации (например, если $|x_n - r| \leq 0.001$, то $|x_{n+1} - r| \leq 0.000001$).

Если первое приближение далеко от значения корня, то это может привести к медленной сходимости итераций, расходимости в бесконечность, сходимости к другому корню (не являющемуся ближайшим к первому приближению) или привести к бесконечному циклу с некоторым набором значений.

Приведенное рассмотрение метода достаточно туманно из-за наличия фраз наподобие “определенные условия”, “доброкачественная функция” и “достаточно близко”. Для того чтобы познакомиться с методом Ньютона более строго, обратитесь практически к любому учебнику по вычислительным методам.

Несмотря на предупреждения об области применения данного метода, зачастую он очень полезен в области целых чисел. Для того чтобы определить применимость этого метода для той или иной функции, следует провести соответствующий анализ наподобие описанного в разделе 11.1, "Целочисленный квадратный корень", на с. 305.

В табл. Б.1 приведены некоторые итеративные формулы, выведенные из метода Ньютона. В первом столбце приведены искомые значения, во втором — функции, для которых эти значения являются корнями, и в третьем — правые части формул Ньютона, соответствующие этим функциям.

ТАБЛИЦА Б.1. Метод Ньютона для вычисления некоторых значений

Вычисляемое значение	Функция	Итеративная формула
\sqrt{a}	$x^2 - a$	$\frac{1}{2}\left(x_n + \frac{a}{x_n}\right)$
$\sqrt[3]{a}$	$x^3 - a$	$\frac{1}{3}\left(2x_n + \frac{a}{x_n^2}\right)$
$\frac{1}{\sqrt{a}}$	$x^{-2} - a$	$\frac{x_n}{2}(3 - ax_n^2)$
$\frac{1}{a}$	$x^{-1} - a$	$x_n(2 - ax_n)$
$\log_2 a$	$2^x - a$	$x_n + \frac{1}{\ln 2}\left(\frac{a}{2^{x_n}} - 1\right)$

К сожалению, не всегда просто найти подходящую функцию. Разумеется, существует множество функций, корнем которых является интересующее нас значение, но только немногие из них приводят к практичным итеративным формулам. Обычно функция представляет собой некоторый тип обратных вычислений по отношению к вычислению искомого значения. Например, для поиска \sqrt{a} используется $f(x) = x^2 - a$, для поиска $\log_2 a$ — $f(x) = 2^x - a$ и т.д.¹

Итеративная формула для $\log_2 a$ сходится (к $\log_2 a$), даже если заменить множитель $1/\ln 2$ некоторым другим значением (например, 1 или 2). Однако скорость сходимости при этом снизится. Для ряда приложений вместо множителя $1/\ln 2$ можно использовать величину $3/2$ или $23/16$ ($1/\ln 2 \approx 1.4427$).

¹ Метод Ньютона для частного случая квадратного корня был известен еще в Древнем Вавилоне около 4000 лет назад.

ПРИЛОЖЕНИЕ В

ГРАФИКИ ДИСКРЕТНЫХ ФУНКЦИЙ

В этом приложении приведены графики ряда дискретных функций, полученных с помощью Mathematica. Для каждой функции приведены два графика: один для размера слова 3 бита, и второй — для размера слова 5 бит. Этот материал представлен Гаем Стилом (Guy Steele).

В.1. Графики логических операций над целыми числами

В этом разделе приведены трехмерные графики $\text{and}(x, y)$, $\text{or}(x, y)$ и $\text{xor}(x, y)$ как функций от целых чисел x и y (рис. В.1–В.3).

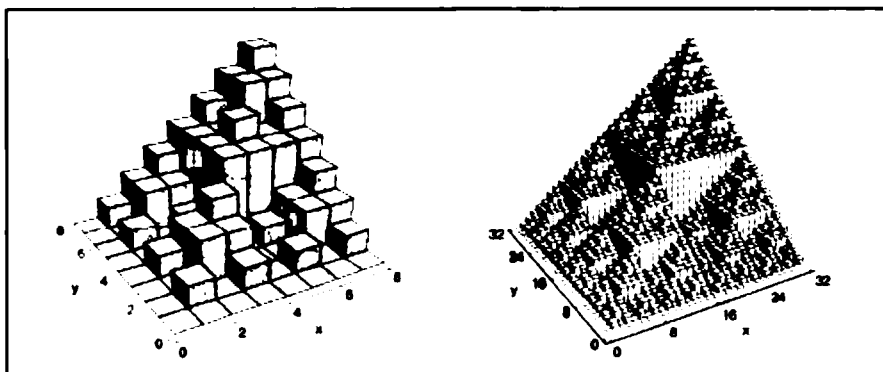


Рис. В.1. График логической функции and

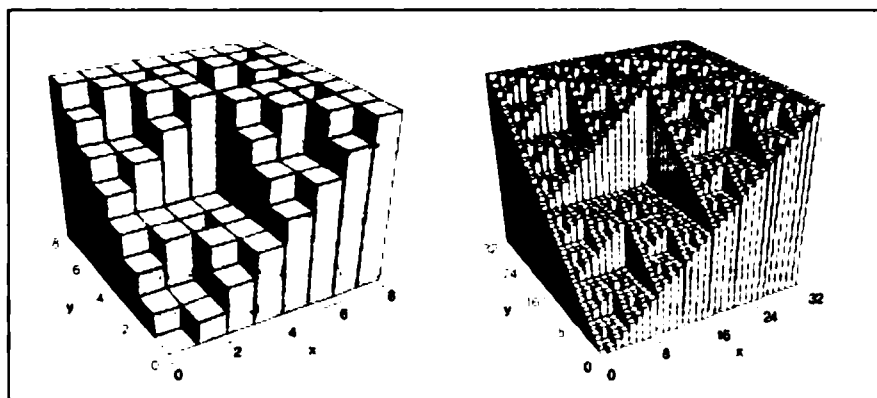


Рис. В.2. График логической функции or

На рис. В.3 половина точек скрыта за диагональной плоскостью $x = \bar{y}$.

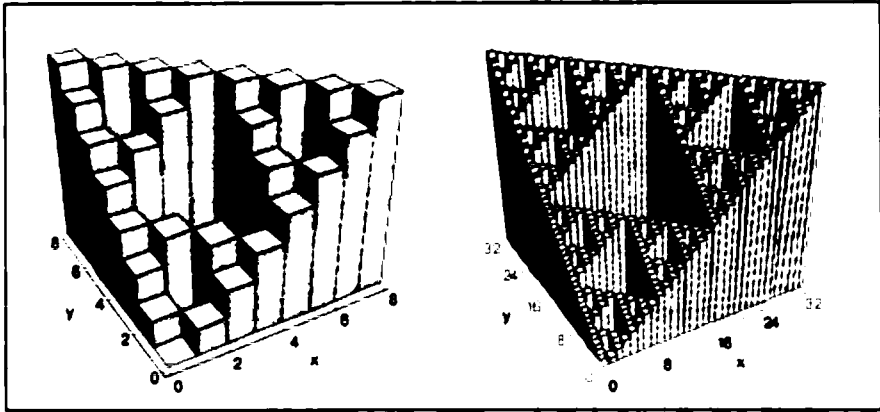


Рис. В.3. График логической функции xor

В графике функции $\text{and}(x, y)$ просматривается определенный самоподобный, или фрактальный, узор из треугольников. Если рассматривать рисунок вдоль оси y и перейти в пределе к большим целым числам, то картина будет такой, как показано на рис. В.4.

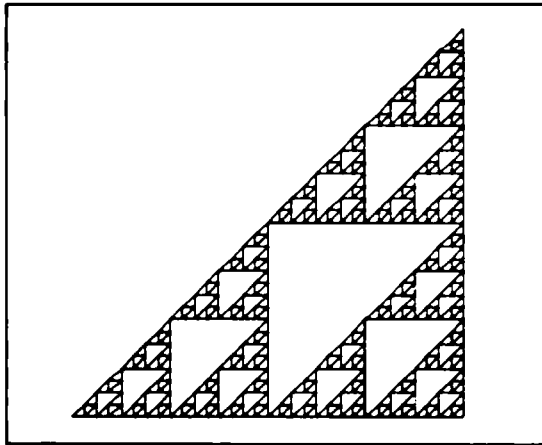


Рис. В.4. Самоподобная структура функции $\text{and}(x, y)$

Она очень похожа на треугольник Серпинского [99] с тем отличием, что на рис. В.4 изображен прямоугольный треугольник, в то время как Серпинский использовал равносторонний треугольник. Если присмотреться к рис. В.3, то можно увидеть, что в наклонной плоскости получен точный неискаженный треугольник Серпинского.

В.2. Графики для сложения, вычитания и умножения

В этом разделе на рис. В.5–В.9 показаны трехмерные графики для сложения, вычитания и трех видов умножения беззнаковых чисел с помощью “компьютерной арифметики”. Обратите внимание, что на графике для операции сложения начало координат находится в дальнем левом углу.

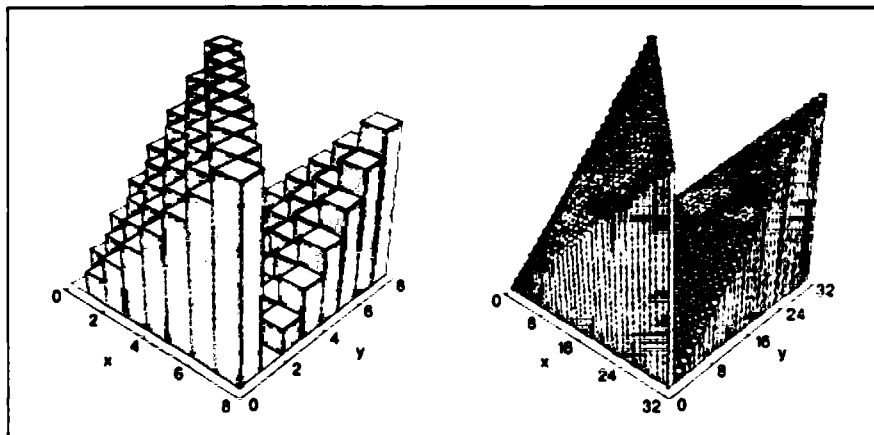


Рис. В.5. График $x + y$ (компьютерная арифметика)

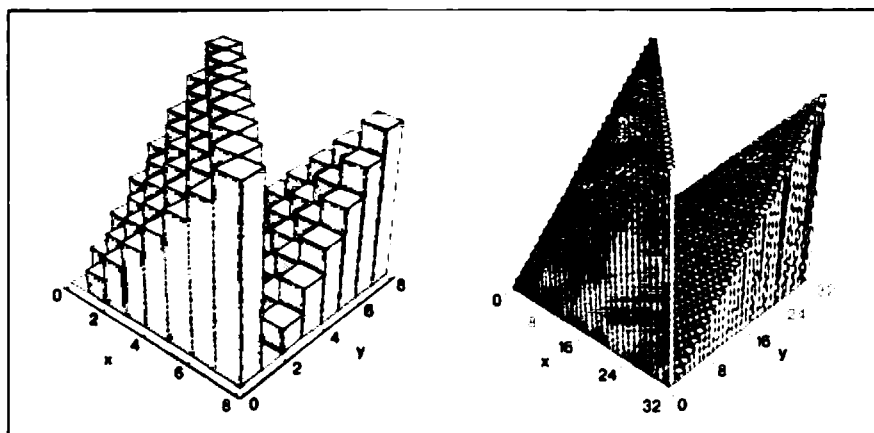


Рис. В.6. График $x + y$ (компьютерная арифметика)

На рис. В.7 изображение сжато по вертикали; наивысший пик в левой части рисунка имеет высоту $7 \cdot 7 = 49$.

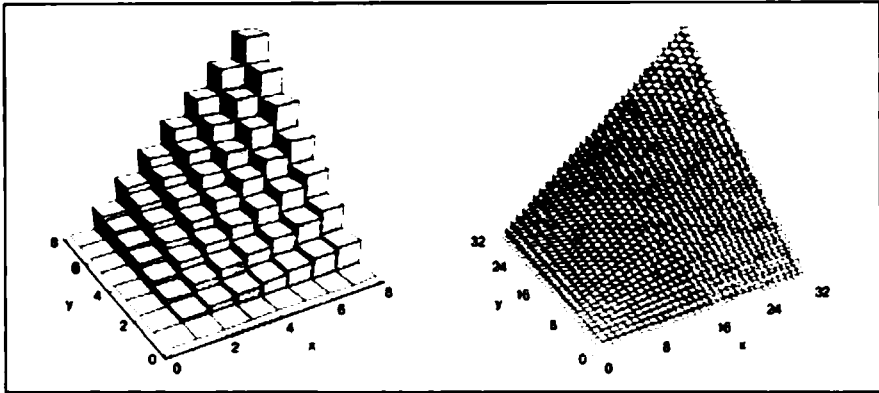


Рис. В.7. График беззнакового произведения x и y

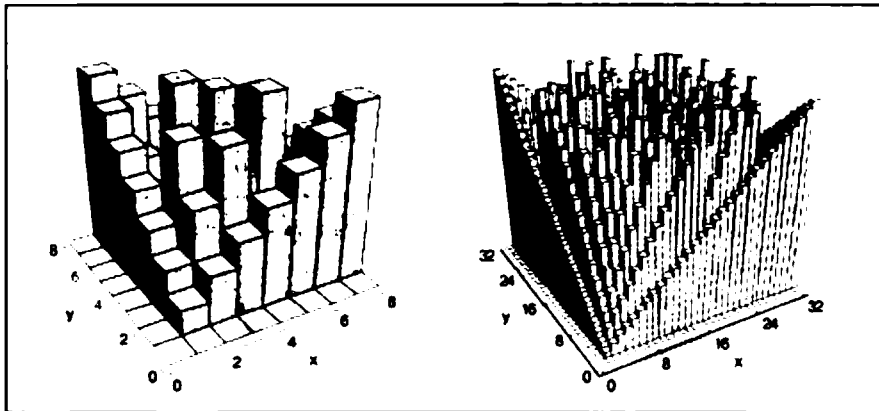


Рис. В.8. График младшей половины беззнакового произведения x и y

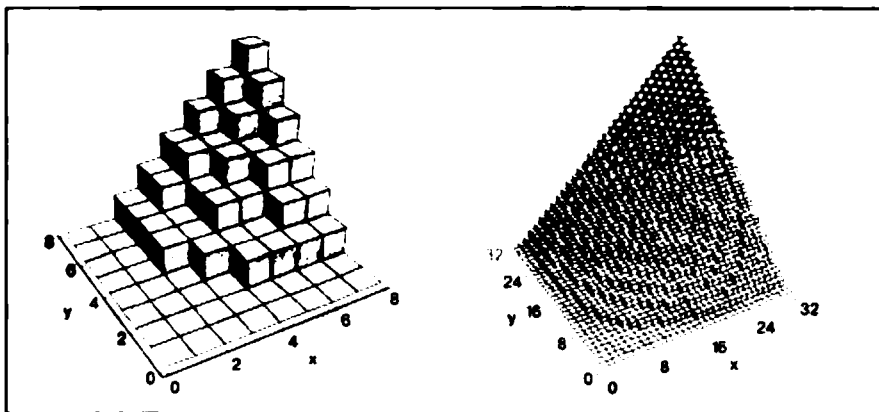


Рис. В.9. График старшей половины беззнакового произведения x и y

В.3. Графики функций, включающих деление

В этом разделе показаны трехмерные графики для частного, остатка, наибольшего общего делителя и наименьшего общего кратного для неотрицательных целых чисел x и y (на рис. В.10, В.11, В.12 и В.13 соответственно). Обратите внимание, что на рис. В.10 начало координат находится в дальнем правом углу.

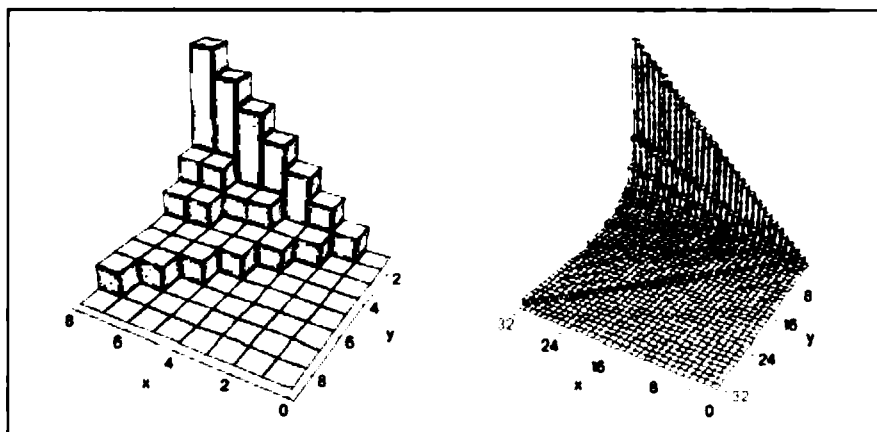


Рис. В.10. График функции целочисленного частного $x + y$

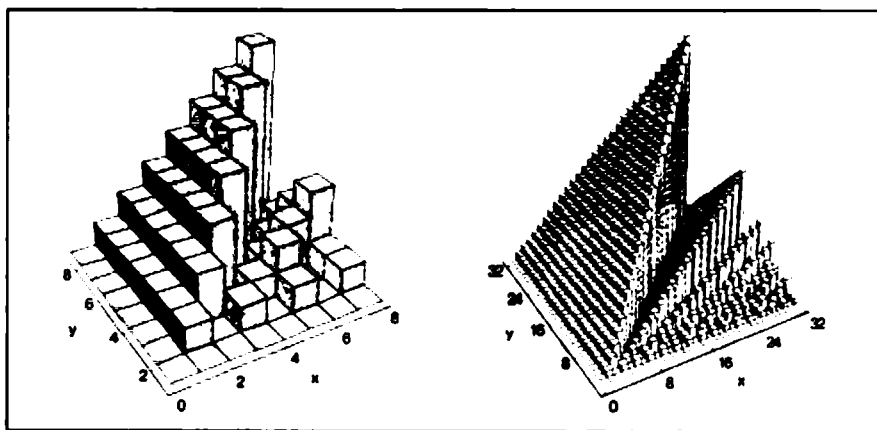


Рис. В.11. График функции остатка $\text{rem}(x, y)$

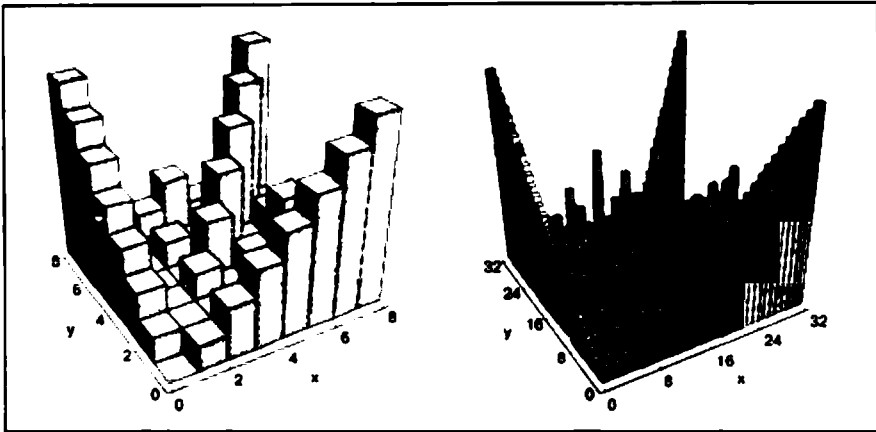


Рис. В.12. График функции наибольшего общего делителя $\text{gcd}(x, y)$

На рис. В.13 изображение сжато по вертикали; наивысший пик в левой части рисунка имеет высоту $\text{lcm}(6, 7) = 42$.

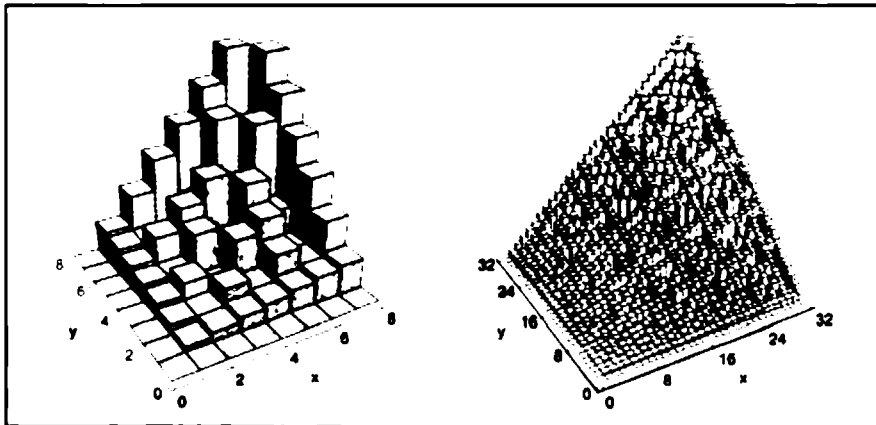


Рис. В.13. График функции наименьшего общего кратного $\text{lcm}(x, y)$

В.4. Графики функций сжатия, обобщенного упорядочения и циклического сдвига влево

В этом разделе показаны трехмерные графики для функций $\text{compress}(x, m)$, $\text{SAG}(x, m)$ и $x \ll r$ как функций от целых чисел x , m и r (на рис. В.14, В.15 и В.16 соответственно).

Для функций compress и SAG m представляет собой маску. При сжатии биты x , выбранные маской m , извлекаются и сжимаются вправо с дополнением нулевыми битами слева. В случае функции SAG биты x , выбранные маской m , сжимаются влево, а невыбранные — вправо.

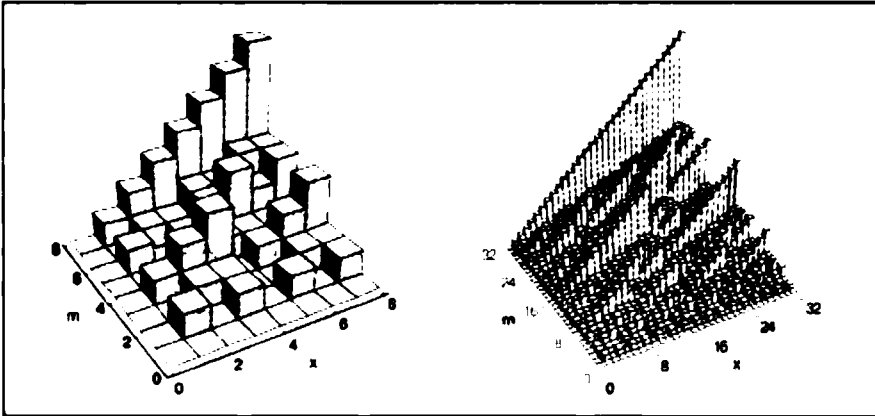


Рис. В.14. График функции $\text{compress}(x, m)$

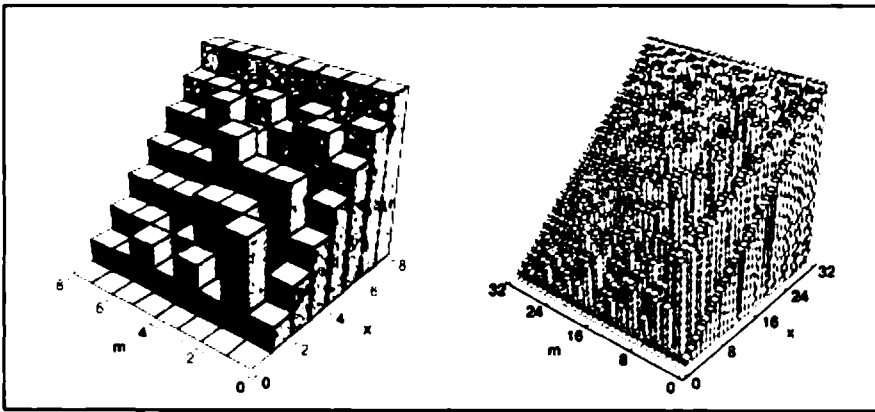


Рис. В.15. График функции $\text{SAG}(x, m)$

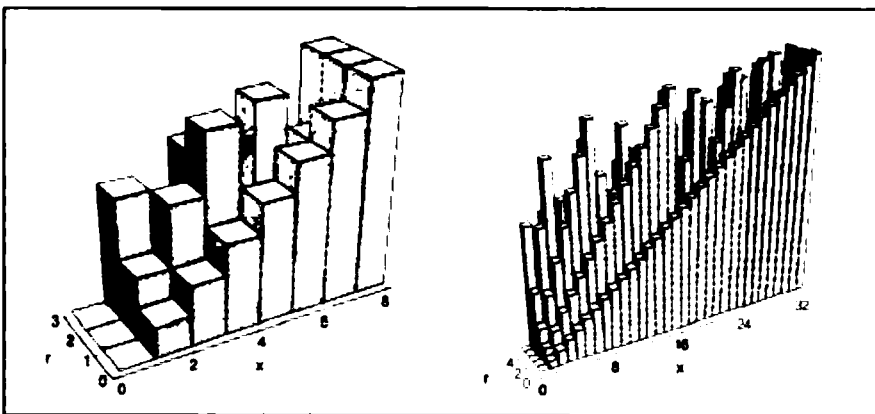


Рис. В.16. График функции $x \ll r$

В.5. Графики некоторых унарных функций

На рис. В.17–В.21 показаны двухмерные графики некоторых унарных функций от битовых строк, которые рассматриваются как функции от целых чисел. Двухмерные графики, так же, как и трехмерные, получены с помощью Mathematica. Для большинства функций приведены два графика: для слова размером 4 бита и для слова размером 7 бит.

Под “функциями кода Грея” подразумеваются функции, отображающие целое число, которое представляет величину смещения или циклического сдвига, в код Грея для этого смещения или циклического сдвига. Обратная функция, соответственно, отображает код Грея на величину смещения или циклического сдвига. (См. рис. 13.1 на с. 338.)

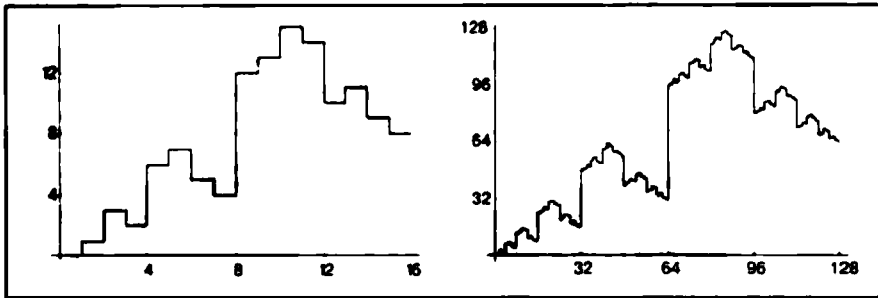


Рис. В.17. Графики функции кода Грея

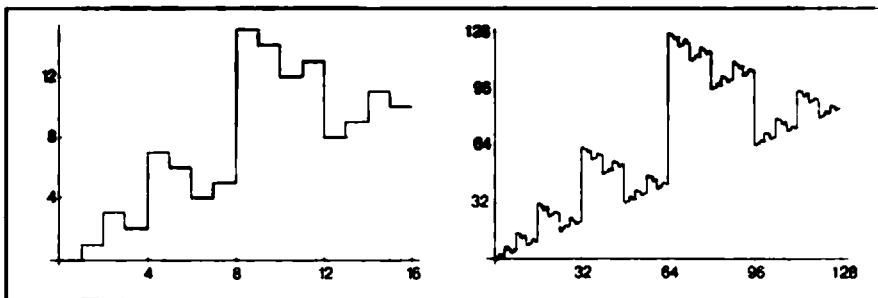


Рис. В.18. Графики обратной функции кода Грея

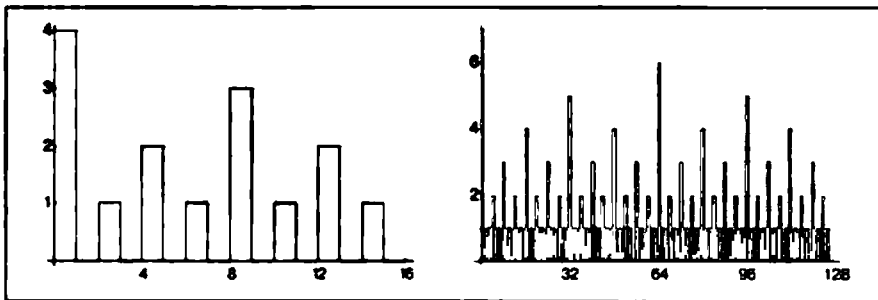


Рис. В.19. Графики линейной функции (количество завершающих нулей)

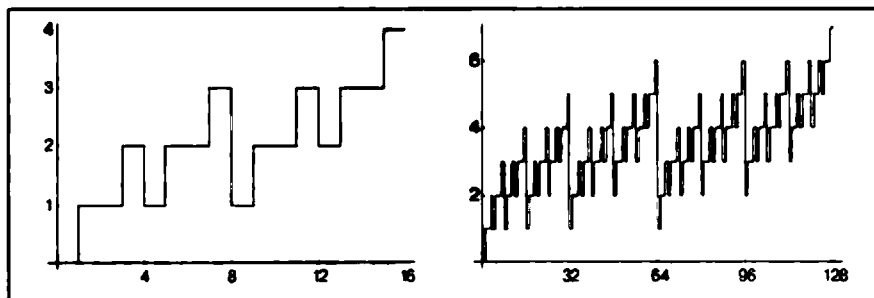


Рис. В.20. Графики количества единичных битов

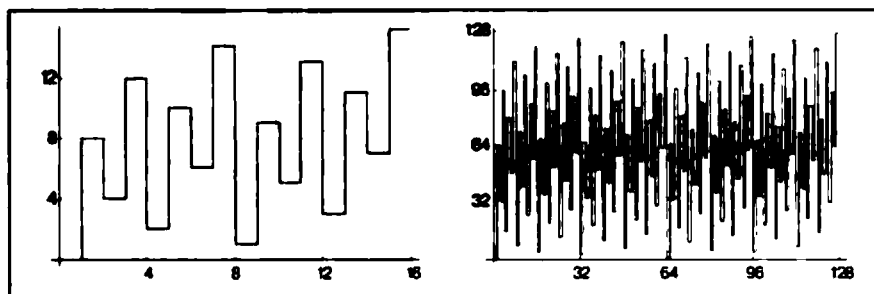


Рис. В.21. Графики функции обращения порядка битов в слове

На рис. В.22 показано, что произойдет с колодой из 16 карт при выполнении одного, двух и трех внешних идеальных перемешиваний (при которых первая и последняя карты не перемещаются). Координата x указывает исходное положение карты, а координата y — окончательное положение карты после одного, двух и трех перемешиваний. На рис. В.23 показаны те же результаты для *внутренних* перемешиваний. На рис. В.24 и В.25 показаны графики для обратных операций.

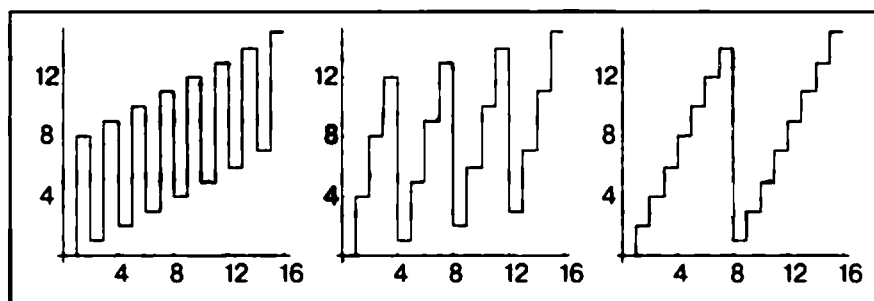


Рис. В.22. Графики функции внешнего идеального перемешивания

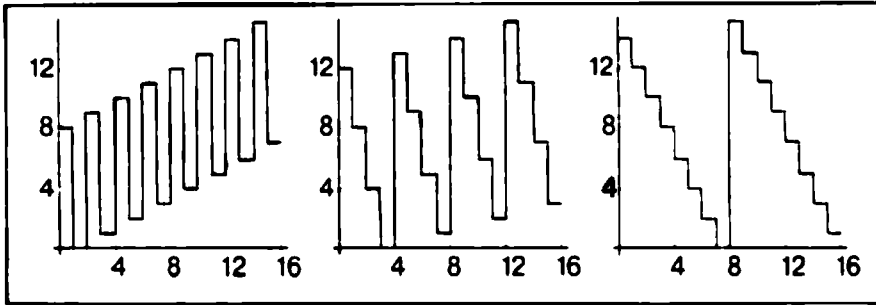


Рис. В.23. Графики функции внутреннего идеального перемешивания

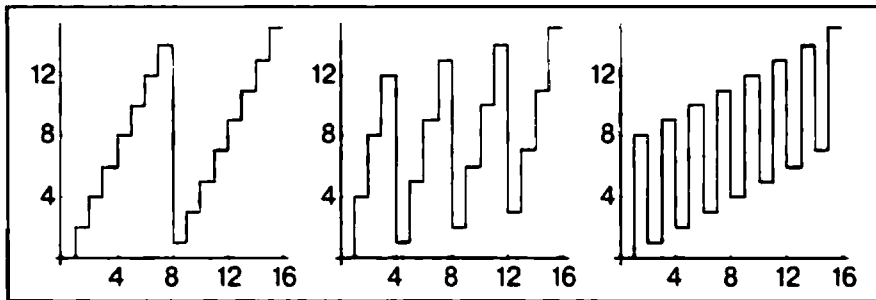


Рис. В.24. Графики функции внешнего идеального упорядочения

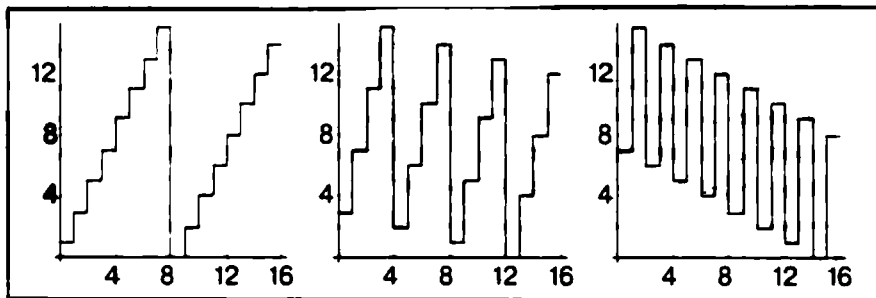


Рис. В.25. Графики функции внутреннего идеального упорядочения

На рис. В.26 и В.27 показано отображение, получающееся в результате перемешивания битов целого числа длиной 4 и 8 бит. Говоря неформально,

$$\text{shuffleBits}(x) = \text{asInteger}(\text{shuffle}(\text{bits}(x))).$$

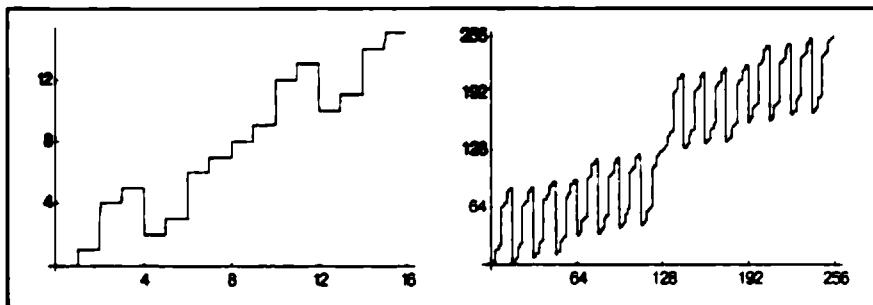


Рис. В.26. Графики функции битов внешнего идеального перемешивания

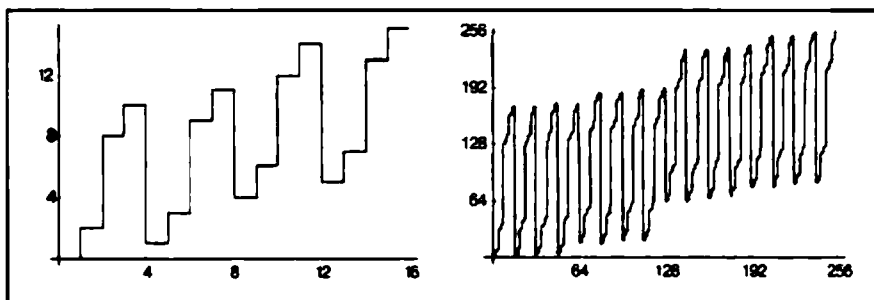


Рис. В.27. Графики функции битов внутреннего идеального перемешивания

СПИСОК ЛИТЕРАТУРЫ

1. *Advanced Encryption Standard (AES)*, National Institute of Standards and Technology, FIPS PUB 197 (November 2001). Доступно по адресу <http://csrc.nist.gov/publications/fips/fips197/fips-197.pdf>.
2. Agrell, Erik. <http://webfiles.portal.chalmers.se/s2/research/kit/bounds/>, последнее обновление — июль 2004.
3. Allen, Joseph H. Частное сообщение.
4. Alverson, Robert. "Integer Division Using Reciprocals". In *Proceedings IEEE 10th Symposium on Computer Arithmetic*, June 26–28, 1991, Grenoble, France, 186–190.
5. Arndt, Jörg. *Matters Computational: Ideas, Algorithms, Source Code*. Springer-Verlag, 2010. Также доступно по адресу <http://www.jjj.de/fxt/#fxtbook>.
6. Найдено в подпрограмме интерпретатора REXX, написанной Марком Аусландером (Marc A. Auslander).
7. Auslander, Marc A. Частное сообщение.
8. Д. Кнут приписывает тернарный метод неопубликованной работе Брюса Баумгарта (Bruce Baumgart) середины 1970-х годов, в которой сравнивается около 20 различных методов обращения порядка битов на PDP10.
9. Bernstein, Robert. "Multiplication by Integer Constants". *Software — Practice and Experience* 16, 7 (July 1986), 641–652.
10. Burks, Arthur W., Goldstine, Herman H., and von Neumann, John. "Preliminary Discussion of the Logical Design of an Electronic Computing Instrument, Second Edition" (1947). In *Papers of John von Neumann on Computing and Computing Theory*, Volume 12 in the Charles Babbage Institute Reprint Series for the History of Computing, MIT Press, 1987.
11. Black, Richard. Веб-сайт www.cl.cam.ac.uk/Research/SRG/bluebook/21/crc/crc.html. University of Cambridge Computer Laboratory Systems Research Group, February 1994.
12. Bonzini, Paolo. Частное сообщение.
13. Brouwer, Andries E. <http://www.win.tue.nl/~aeb/codes/binary-1.html>, последнее обновление — январь 2012.
14. Cavagnino, D. and Werbrouck, A. E. "Efficient Algorithms for Integer Division by Constants Using Multiplication". *The Computer Journal* 51, 4 (2008), 470–480.
15. Caldwell, Chris K. and Cheng, Yuanyou. "Determining Mills' Constant and a Note on Honaker's Problem". *Journal of Integer Sequences* 8, 4 (2005), article 05.4.1, 9 pp. Также доступно по адресу <http://www.cs.uwaterloo.ca/journals/JIS/VOL8/Caldwell/caldwell178.pdf>.
16. Stephenson, Christopher J. Частное сообщение.
17. На эти правила было указано Норманом Кохеном (Norman H. Cohen).
18. Leung, Vitus J., et. al. "Processor Allocation on Cplnt: Achieving General Processor Locality Using One-Dimensional Allocation Strategies". In *Proceedings 4th IEEE International Conference on Cluster Computing*, September 2002, 296–304.
19. Cutland, Nigel J. *Computability: An Introduction to Recursive Function Theory*. Cambridge University Press, 1980.
20. Hoxey, Karim, Hay, and Warren (Editors). *The PowerPC Compiler Writer's Guide*. Warthman Associates, 1996.
21. Dalton, Michael. Частное сообщение.
22. Dannemiller, Christopher M. Частное сообщение. Он приписывает данный код Linux Source base, www.gelato.unsw.edu.au/lxr/source/lib/crc32.c, строки 105–111.

23. *Data Encryption Standard (DES)*, National Institute of Standards and Technology, FIPS PUB 46-2 (December 1993). Доступно по адресу <http://www.itl.nist.gov/fipspubs/fip46-2.htm>.
24. Dewdney, A. K. *The Turing Omnibus*. Computer Science Press, 1989.
25. Dietz, Henry G. <http://aggregate.org/MAGIC/>.
26. Ditlow, Gary S. Частное сообщение.
27. Dubé, Danny. Newsgroup comp.compression.research, October 3, 1997.
28. Dudley, Underwood. "History of a Formula for Primes". *American Mathematics Monthly* 76 (1969), 23–28.
29. Ercegovac, Miloš D. and Lang, Tomás. *Division and Square Root: Digit-Recurrence Algorithms and Implementations*. Kluwer Academic Publishers, 1994.
30. Etzion, Tuvi. "Constructions for Perfect 2-Burst-Correcting Codes", *IEEE Transactions on Information Theory* 47, 6 (September 2001), 2553–2555.
31. Floyd, Robert W. "Permuting Information in Idealized Two-Level Storage". In *Complexity of Computer Computations* (Conference proceedings), Plenum Press, 1972, 105–109. Это самое раннее известное автору описание этого метода транспонирования матриц.
32. Gardner, Martin. "Mathematical Games" column in *Scientific American* 227, 2 (August 1972), 106–109.
33. Gaudet, Dean. Частное сообщение.
34. Gregoire, Dennis G., Groves, Randall D., and Schmookler, Martin S. *Single Cycle Merge/Logic Unit*, US Patent No. 4,903,228, February 20, 1990.
35. Granlund, Torbjörn and Kenner, Richard. "Eliminating Branches Using a Superoptimizer and the GNU C Compiler". In *Proceedings of the 5th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, July 1992, 341–352.
36. Graham, Ronald L., Knuth, Donald E., and Patashnik, Oren. *Concrete Mathematics: A Foundation for Computer Science, Second Edition*. Addison-Wesley, 1994.
37. Steele, Guy L., Jr. Частное сообщение.
38. Steele, Guy L., Jr. "Arithmetic Shifting Considered Harmful". AI Memo 378, MIT Artificial Intelligence Laboratory (September 1976); also in *SIGPLAN Notices* 12, 11 (November 1977), 61–69.
39. Granlund, Torbjörn and Montgomery, Peter L. "Division by Invariant Integers Using Multiplication". In *Proceedings of the ACM SIGPLAN '94 Conference on Programming Language Design and Implementation (PLDI)*, August 1994, 61–72.
40. Второе выражение принадлежит Ричарду Гольдбергу (Richard Goldberg).
41. Goodstein, Prof. R. L. "Formulae for Primes". *The Mathematical Gazette* 51 (1967), 35–36.
42. Goryavsky, Julius. Частное сообщение.
43. Найдено GNU Superoptimizer.
44. Beeler, M., Gosper, R. W., and Schroepel, R. *НАКМЕМ*, MIT Artificial Intelligence Laboratory AIM 239, February 1972.
45. Hamming, Richard W., "Error Detecting and Error Correcting Codes", *The Bell System Technical Journal* 26, 2 (April 1950), 147–160.
46. Harley, Robert. Newsgroup comp.arch, July 12, 1996.
47. Hay, R. W. Частное сообщение.
48. Первое выражение найдено в подпрограмме компилятора, написанного Р. Хэем (R. W. Hay).
49. Hilbert, David. "Ueber die stetige Abbildung einer Linie auf ein Flächenstück". *Mathematischen Annalen* 38 (1891), 459–460.
50. Hill, Raymond. *A First Course in Coding Theory*. Clarendon Press, 1986.

51. Hiltgen, Alain P. and Paterson, Kenneth G. "Single-Track Circuit Codes". *IEEE Transactions on Information Theory* 47, 6 (2001), 2587–2595.
52. Hopkins, Martin E. Частное сообщение.
53. Hillis, W. Daniel and Steele, Guy L., Jr. "Data Parallel Algorithms". *Comm. ACM* 29, 12 (December 1986), 1170–1183.
54. Hsieh, Paul. Newsgroup comp.lang.c, April 29, 2005.
55. Hueffner, Falk. Частное сообщение.
56. Hennessy, John L. and Patterson, David A. *Computer Architecture: A Quantitative Approach*. Morgan Kaufmann, 1990.
57. Harbison, Samuel P. and Steele, Guy L., Jr. *C: A Reference Manual*, Fourth Edition. Prentice-Hall, 1995.
58. Hardy, G. H. and Wright, E. M. *An Introduction to the Theory of Numbers*, Fourth Edition. Oxford University Press, 1960.
59. Из курса программирования IBM, 1961.
60. Irvine, M. M. "Early Digital Computers at Bell Telephone Laboratories". *IEEE Annals of the History of Computing* 23, 3 (July–September 2001), 22–42.
61. von Neumann, John. "First Draft of a Report on the EDVAC". In *Papers of John von Neumann on Computing and Computing Theory*, Volume 12 in the Charles Babbage Institute Reprint Series for the History of Computing, MIT Press, 1987.
62. Karatsuba, A. and Ofman, Yu. "Multiplication of multidigit numbers on automata." *Soviet Physics-Doklady* 7, 7 (January 1963), 595–596. Авторами работы показано, что умножение m -битовых целых чисел оказывается более громоздким, чем метод на основе схемы трех умножений Гаусса для комплексных чисел.
63. Karvonen, Vesa. Найдено на веб-странице Found at "The Assembly Gems", www.af.lth.se/~john_e/tr_gems.html.
64. Keane, Joe. Newsgroup sci.math.num-analysis, July 9, 1995.
65. Найдено в компиляторе GNU C, перенесенном на платформу RS/6000 Ричардом Кеннером (Richard Kenner). Он ссылается на статью, написанную им в соавторстве с Торбьёрном Гранлундом (Torbjörn Granlund) и представленную в 1992 году на конференции PLDI.
66. Knuth, Donald E. *The Art of Computer Programming, Volume 1, Third Edition: Fundamental Algorithms*. Addison-Wesley, 1997.
67. Knuth, Donald E. *The Art of Computer Programming, Volume 2, Third Edition: Seminumerical Algorithms*. Addison-Wesley, 1998.
68. Идея использовать в качестве основания системы счисления отрицательные числа была независимо открыта различными людьми. Наиболее раннее упоминание этой идеи (по Кнуту) — в 1885 году Витторио Грюнвальдом (Vittorio Grünwald). Более детально этот вопрос рассмотрен в томе 2 книги *Искусство программирования* [67].
69. Knuth, Donald E. *The Art of Computer Programming, Volume 4A: Combinatorial Algorithms, Part 1*, Section 7.1.1. Addison-Wesley, 2011.
70. *Ibid*, Section 7.1.3. Кнут приписывает отношение равенства У.Ч. Линчу (W. C. Lynch).
71. *Ibid*, Section 7.2.1.1, Exercise 80.
72. Knuth, Donald E. *The Art of Computer Programming, Volume 1, Fascicle 1: MMIX—A RISC Computer for the New Millennium*. Addison-Wesley, 2005.
73. Knuth, Donald E. Частное сообщение.
74. Kruskal, Clyde P., Rudolph, Larry, and Snir, Marc. "The Power of Parallel Prefix". *IEEE Transactions on Computers* C-34, 10 (October 1985), 965–968.
75. Этот рисунок предложен Говри Кумаром (Gowri Kumar). Частное сообщение.
76. Lamport, Leslie. "Multiple Byte Processing with Full-Word Instructions". *Communications of the ACM* 18, 8 (August 1975), 471–475.

77. Langdon, Glen G. Jr., "Subtraction by Minuend Complementation", *IEEE Transactions on Computers C-18*, 1 (January 1969), 74–76.
78. Lin, Shu and Costello, Daniel J., Jr. *Error Control Coding: Fundamentals and Applications*. Prentice-Hall, 1983.
79. Lomont, Chris. *Fast Inverse Square Root*.
www.lomont.org/Math/Papers/2003/InvSqrt.pdf.
80. Leiserson, Charles E., Prokop, Harald, and Randall, Keith H. *Using de Bruijn Sequences to Index a 1 in a Computer Word*. MIT Laboratory for Computer Science, July 7, 1998. Доступно также по адресу <http://supertech.csail.mit.edu/papers/debruijn.pdf>.
81. Lee, Ruby B., Shi, Zhijie, and Yang, Xiao. "Efficient Permutation Instructions for Fast Software Cryptography". *IEEE Micro* 21, 6 (November–December 2001), 56–69.
82. Lam, Warren M. and Shapiro, Jerome M. "A Class of Fast Algorithms for the Peano-Hilbert Space-Filling Curve". In *Proceedings ICIP 94*, 1 (1994), 638–641.
83. Denneau, Monty M. Частное сообщение.
84. Kane, Gerry and Heinrich, Joe. *MIPS RISC Architecture*. Prentice-Hall, 1992.
85. Morton, Mike. "Quibbles & Bits." *Computer Language* 7, 12 (December 1990), 45–55.
86. Möbius, Stefan K. Частное сообщение.
87. MacWilliams, Florence J. and Sloane, Neil J. A. *The Theory of Error-Correcting Codes, Part II*. North-Holland, 1977.
88. Mycroft, Alan. Newsgroup comp.arch, April 8, 1987.
89. Neumann, Jasper L. Частное сообщение.
90. Niven, Ivan, Zuckerman, Herbert S., and Montgomery, Hugh L. *An Introduction to the Theory of Numbers, Fifth Edition*. John Wiley & Sons, Inc., 1991.
91. Peterson, W. W. and Brown, D. T. "Cyclic Codes for Error Detection". In *Proceedings of the IRE*, 1 (January 1961), 228–235.
92. Oden, Peter H. Частное сообщение.
93. Я научился этому трюку у компилятора PL.8.
94. Purdom, Paul Walton Jr., and Brown, Cynthia A. *The Analysis of Algorithms*. Holt, Rinehart and Winston, 1985.
95. Reiser, John. Newsgroup comp.arch.arithmetic, December 11, 1998.
96. Ribenboim, Paulo. *The Little Book of Big Primes*. Springer-Verlag, 1991.
97. Reingold, Edward M., Nievergelt, Jürg, and Deo, Narsingh. *Combinatorial Algorithms: Theory and Practice*. Prentice-Hall, 1977.
98. Roman, Steven. *Coding and Information Theory*. Springer-Verlag, 1992.
99. Sagan, Hans. *Space-Filling Curves*. Springer-Verlag, 1994. Чудесная книга, рекомендуемая всем интересующимся данной темой.
100. Seal, David. Newsgroup comp.arch.arithmetic, May 13, 1997.
101. Seal, David. Newsgroup comp.sys.acorn.tech, February 16, 1994.
102. Shepherd, Arvin D. Частное сообщение.
103. Stallman, Richard M. *Using and Porting GNU CC*. Free Software Foundation, 1998.
104. Strachey, Christopher. "Bitwise Operations." *Communications of the ACM* 4, 3 (March 1961), 146. В этой работе содержится другая статья, в которой имеются два метода обращения порядка битов ("Two Methods for Word Inversion on the IBM 709," by Robert A. Price and Paul Des Jardins).
105. Tanenbaum, Andrew S. *Computer Networks*, Second Edition. Prentice Hall, 1988.
106. Похоже, автор этой программы навсегда потерян для истории. Дополнительную информацию можно почерпнуть по адресу <http://www.beyond3d.com/content/articles/8/>.

107. Voorhies, Douglas. "Space-Filling Curves and a Measure of Coherence". *Graphics Gems II*, AP Professional (1991).
108. Warren, H. S., Jr. "Functions Realizable with Word-Parallel Logical and Two's-Complement Addition Instructions". *Communications of the ACM* 20, 6 (June 1977), 439–441.
109. Наиболее ранняя известная мне ссылка — Wegner, P. A. "A Technique for Counting Ones in a Binary Computer". *Communications of the ACM* 3, 5 (May 1960), 322.
110. Wells, David. *The Penguin Dictionary of Curious and Interesting Numbers*. Penguin Books, 1997.
111. Willans, C. P. "On Formulae for the n th Prime Number." *The Mathematical Gazette* 48 (1964), 413–415.
112. Woodrum, Luther. Частное сообщение. Вторая формула не использует литералов и хорошо работает на IBM System/370.
113. Wormell, C. P. "Formulae for Primes." *The Mathematical Gazette* 51 (1967), 36–38.
114. Zadeck, F. Kenneth. Частное сообщение.

ПРЕДМЕТНЫЙ УКАЗАТЕЛЬ

С

CRC. См. Циклический избыточный код

Е

ECC. См. Коды с коррекцией ошибок

Н

НАКМЕМ, 13

Р

PDP-10, 13

R

RISC

базовый набор команд, 23

дополнительный набор команд, 25

А

Абсолютное значение, 39

многобайтовой величины, 63

Алгоритм Евклида, 268

Б

Базовый набор команд RISC, 23

Бинарный поиск, 309

Бит четности, 345

В

Ведущие цифры чисел, 413

Вектор кодовый, 345

Возведение в степень, 314

Выбор среди множества значений, 70

Вычисление отношений, 48

Г

Гильберта кривая. См. Кривая Гильберта

Границы логических выражений, 96

Границы суммы и разности, 92

Грея код. См. Код Грея

Д

Деление

больших чисел, 210

больших чисел знаковое, 215

длинное беззнаковое, 219

короткое беззнаковое, 215

на 3, 234, 253

на 5, 235

на 7, 235, 254

на делитель, не превышающий -2 , 243

на известную степень 2, 231

на константу, 231–304, 236, 256

на константу точное, 266

проверка кратности константе, 274

целочисленное, 207–29

Дополнительный набор команд RISC, 25

З

Закон де Моргана, 33; 99

Знаковый сдвиг вправо, 40

К

Код Грея, 337

Инкремент, 339

Отрицательно-двоичный, 341

Циклический, 337

Код условия, 58

Кодовое слово, 357

Коды с коррекцией ошибок, 357

граница Гильберта–Варшавова, 376

граница Хэмминга, 375

идеальный код, 359

код Хэмминга, 358–364

кодовая скорость, 370

размер кода, 370

расширенный код Хэмминга, 360

синдром, 360

скорость кода, 370

Кривая Гильберта, 381
 нерекурсивное построение, 398
 преобразование координат
 в расстояние, 393
 преобразование расстояния
 в координаты, 385
 рекурсивный алгоритм построения, 383

Кривая Пеано, 381; 398

Л

Логические операции, 38

М

Манипуляции младшими битами, 31

Метод Ньютона, 271, 305, 487

Монус, 64

Мультипликативное обратное число, 268

Н

Неравенства, 38

О

Обмен полей одного регистра, 69

Обмен полей регистров, 68

Обмен содержимого регистров, 68

Обнаружение переполнения, 49

Обобщенная перестановка битов, 187

Обобщенное извлечение битов, 176

Обобщенное упорядочение битов, 188

Обратный порядок байтов, 155

Округление к ближайшей степени 2, 82

Округление к кратному степени 2, 81

П

Параллельный префикс, 119, 339, 391

Параллельный суффикс, 120, 177

Перемешивание битов, 165

Перенос знака, 42

Переполнение, 49

Подсчет ведущих нулевых битов, 122

Подсчет единичных битов, 103

в массиве, 111

Подсчет завершающих нулевых битов, 130

Поиск первого байта из заданного
 диапазона, 146

Поиск первого нулевого байта, 141

Поиск строки единичных битов заданной
 длины, 147

Поиск числа с тем же количеством битов, 35

Поле Галуа, 347

Порядок байтов, 155

Правило Хэмминга, 358

Предикаты сравнения, 43

Проверка границ, 89

логических выражений, 96

суммы и разности, 92

Прямой порядок байтов, 155

Р

Разреженный массив, 117

Распространение знака, 40

Расстояние Хэмминга, 117, 370

Реверс байтов, 155

Реверс битов, 155

С

Сдвиг двойного слова, 61

Сжатие битов, 176

Синдром, 360

Система счисления

По основанию $-1+i$, 332

По основанию -2 , 325

Эффективность, 335

Среднее двух целых чисел, 39

Т

Тасование, 165

Теорема Вильсона, 421

Теорема Миллса, 431

Транспонирование битовой матрицы, 167

Трехзначная функция сравнения, 42

У

Умножение больших чисел, 197

Умножение на константу, 201

Условный обмен, 70

Ф

Формат IEEE, 401
 Формула Вормелла для n -го простого
 числа, 424
 Формулы Вилланса для n -го простого
 числа, 421
 Функция
 bitsize(), 130
 ceil(), 82, 209
 clp2(), 82
 cmp(), 42
 dist(), 117
 doz(), 63
 floor(), 82, 209
 fp2(), 82
 ilog10(), 317
 ilog2(), 317
 ISIGN, 42
 nlz(), 122
 ntz(), 130
 pop(), 104
 pow2(), 315
 sign(), 41
 zbytel(), 141
 zbyter(), 141

перенос знака, 42
 сравнения трехзначная, 42
 формульная, 425

Х

Хакер, 16

Ц

Целочисленный квадратный корень, 305
 Целочисленный кубический корень, 313
 Целочисленный логарифм, 316
 Циклический избыточный код, 345
 Аппаратная реализация, 351
 Полиномы, 349
 Программная реализация, 353–56
 Стандарты, 350
 Циклический сдвиг, 59

Ч

Четность, 118
 Числа с плавающей точкой
 Диапазон точно представимых целых
 чисел, 403
 Сравнение, 408
 Формат IEEE, 401
 Числа Ферма, 419