



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

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

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

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

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

George Boole  
1815 - 1864

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

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

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

$$A(n, d) = A(n-1, d-1), d \text{ even} \quad -\bar{x} = x + 1$$

# Hacker's Delight

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

## SECOND EDITION

$$1111^7 = 11100001$$

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

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

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

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

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

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

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

HENRY S. WARREN, JR.

FREE SAMPLE CHAPTER



SHARE WITH OTHERS

# Hacker's Delight

*This page intentionally left blank*

# 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

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

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

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

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

For sales outside the United States, please contact:

International Sales  
international@pearsoned.com

Visit us on the Web: [informit.com/aw](http://informit.com/aw)

*Library of Congress Cataloging-in-Publication Data*

Warren, Henry S.

Hacker's delight / Henry S. Warren, Jr. -- 2nd ed.  
p. cm.

Includes bibliographical references and index.

ISBN 0-321-84268-5 (hardcover : alk. paper)

1. Computer programming. I. Title.

QA76.6.W375 2013

005.1—dc23

2012026011

Copyright © 2013 Pearson Education, Inc.

All rights reserved. Printed in the United States of America. 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. To obtain permission to use material from this work, please submit a written request to Pearson Education, Inc., Permissions Department, One Lake Street, Upper Saddle River, New Jersey 07458, or you may fax your request to (201) 236-3290.

ISBN-13: 978-0-321-84268-8

ISBN-10: 0-321-84268-5

Text printed in the United States on recycled paper at Courier in Westford, Massachusetts.

First printing, September 2012

*To Joseph W. Gauld,  
my high school algebra teacher,  
for sparking in me a delight  
in the simple things in mathematics*

*This page intentionally left blank*

# CONTENTS

<i>Foreword</i> .....	<i>xiii</i>
<i>Preface</i> .....	<i>xv</i>
CHAPTER 1. INTRODUCTION .....	1
1-1 Notation .....	1
1-2 Instruction Set and Execution Time Model .....	5
CHAPTER 2. BASICS .....	11
2-1 Manipulating Rightmost Bits .....	11
2-2 Addition Combined with Logical Operations .....	16
2-3 Inequalities among Logical and Arithmetic Expressions .....	17
2-4 <i>Absolute Value</i> Function .....	18
2-5 Average of Two Integers .....	19
2-6 Sign Extension .....	19
2-7 Shift Right Signed from Unsigned .....	20
2-8 <i>Sign</i> Function .....	20
2-9 <i>Three-Valued Compare</i> Function .....	21
2-10 <i>Transfer of Sign</i> Function .....	22
2-11 Decoding a “Zero Means $2^{*n}$ ” Field .....	22
2-12 Comparison Predicates .....	23
2-13 Overflow Detection .....	28
2-14 Condition Code Result of <i>Add</i> , <i>Subtract</i> , and <i>Multiply</i> .....	36
2-15 Rotate Shifts .....	37
2-16 Double-Length Add/Subtract .....	38
2-17 Double-Length Shifts .....	39
2-18 Multibyte <i>Add</i> , <i>Subtract</i> , <i>Absolute Value</i> .....	40
2-19 Doz, Max, Min .....	41
2-20 Exchanging Registers .....	45
2-21 Alternating among Two or More Values .....	48
2-22 A Boolean Decomposition Formula .....	51
2-23 Implementing Instructions for all 16 Binary Boolean Operations .....	53
CHAPTER 3. POWER-OF-2 BOUNDARIES .....	59
3-1 Rounding Up/Down to a Multiple of a Known Power of 2 .....	59
3-2 Rounding Up/Down to the Next Power of 2 .....	60
3-3 Detecting a Power-of-2 Boundary Crossing .....	63



CHAPTER 4. ARITHMETIC BOUNDS .....	67
4-1 Checking Bounds of Integers .....	67
4-2 Propagating Bounds through <i>Add</i> 's and <i>Subtract</i> 's .....	70
4-3 Propagating Bounds through Logical Operations .....	73
CHAPTER 5. COUNTING BITS .....	81
5-1 Counting 1-Bits .....	81
5-2 Parity .....	96
5-3 Counting Leading 0's .....	99
5-4 Counting Trailing 0's .....	107
CHAPTER 6. SEARCHING WORDS .....	117
6-1 Find First 0-Byte .....	117
6-2 Find First String of 1-Bits of a Given Length .....	123
6-3 Find Longest String of 1-Bits .....	125
6-4 Find Shortest String of 1-Bits .....	126
CHAPTER 7. REARRANGING BITS AND BYTES .....	129
7-1 Reversing Bits and Bytes .....	129
7-2 Shuffling Bits .....	139
7-3 Transposing a Bit Matrix .....	141
7-4 <i>Compress</i> , or <i>Generalized Extract</i> .....	150
7-5 <i>Expand</i> , or <i>Generalized Insert</i> .....	156
7-6 Hardware Algorithms for Compress and Expand .....	157
7-7 General Permutations, Sheep and Goats Operation .....	161
7-8 Rearrangements and Index Transformations .....	165
7-9 An LRU Algorithm .....	166
CHAPTER 8. MULTIPLICATION .....	171
8-1 Multiword Multiplication .....	171
8-2 High-Order Half of 64-Bit Product .....	173
8-3 High-Order Product Signed from/to Unsigned .....	174
8-4 Multiplication by Constants .....	175
CHAPTER 9. INTEGER DIVISION .....	181
9-1 Preliminaries .....	181
9-2 Multiword Division .....	184
9-3 Unsigned Short Division from Signed Division .....	189

9-4	Unsigned Long Division	192
9-5	Doubleword Division from Long Division	197
CHAPTER 10. INTEGER DIVISION BY CONSTANTS		205
10-1	Signed Division by a Known Power of 2	205
10-2	Signed Remainder from Division by a Known Power of 2	206
10-3	Signed Division and Remainder by Non-Powers of 2	207
10-4	Signed Division by Divisors $\geq 2$	210
10-5	Signed Division by Divisors $\leq -2$	218
10-6	Incorporation into a Compiler	220
10-7	Miscellaneous Topics	223
10-8	Unsigned Division	227
10-9	Unsigned Division by Divisors $\geq 1$	230
10-10	Incorporation into a Compiler (Unsigned)	232
10-11	Miscellaneous Topics (Unsigned)	234
10-12	Applicability to Modulus and Floor Division	237
10-13	Similar Methods	237
10-14	Sample Magic Numbers	238
10-15	Simple Code in Python	240
10-16	Exact Division by Constants	240
10-17	Test for Zero Remainder after Division by a Constant	248
10-18	Methods Not Using Multiply High	251
10-19	Remainder by Summing Digits	262
10-20	Remainder by Multiplication and Shifting Right	268
10-21	Converting to Exact Division	274
10-22	A Timing Test	276
10-23	A Circuit for Dividing by 3	276
CHAPTER 11. SOME ELEMENTARY FUNCTIONS		279
11-1	Integer Square Root	279
11-2	Integer Cube Root	287
11-3	Integer Exponentiation	288
11-4	Integer Logarithm	291
CHAPTER 12. UNUSUAL BASES FOR NUMBER SYSTEMS		299
12-1	Base $-2$	299
12-2	Base $-1 + i$	306
12-3	Other Bases	308
12-4	What Is the Most Efficient Base?	309

CHAPTER 13. GRAY CODE .....	311
13-1 Gray Code .....	311
13-2 Incrementing a Gray-Coded Integer .....	313
13-3 Negabinary Gray Code .....	315
13-4 Brief History and Applications .....	315
CHAPTER 14. CYCLIC REDUNDANCY CHECK .....	319
14-1 Introduction .....	319
14-2 Theory .....	320
14-3 Practice .....	323
CHAPTER 15. ERROR-CORRECTING CODES .....	331
15-1 Introduction .....	331
15-2 The Hamming Code .....	332
15-3 Software for SEC-DED on 32 Information Bits .....	337
15-4 Error Correction Considered More Generally .....	342
CHAPTER 16. HILBERT'S CURVE .....	355
16-1 A Recursive Algorithm for Generating the Hilbert Curve ....	356
16-2 Coordinates from Distance along the Hilbert Curve .....	358
16-3 Distance from Coordinates on the Hilbert Curve .....	366
16-4 Incrementing the Coordinates on the Hilbert Curve .....	368
16-5 Non-Recursive Generating Algorithms .....	371
16-6 Other Space-Filling Curves .....	371
16-7 Applications .....	372
CHAPTER 17. FLOATING-POINT .....	375
17-1 IEEE Format .....	375
17-2 Floating-Point To/From Integer Conversions .....	377
17-3 Comparing Floating-Point Numbers Using Integer Operations	381
17-4 An Approximate Reciprocal Square Root Routine .....	383
17-5 The Distribution of Leading Digits .....	385
17-6 Table of Miscellaneous Values .....	387
CHAPTER 18. FORMULAS FOR PRIMES .....	391
18-1 Introduction .....	391
18-2 Willans's Formulas .....	393
18-3 Wormell's Formula .....	397
18-4 Formulas for Other Difficult Functions .....	398

ANSWERS TO EXERCISES .....	405
APPENDIX A. ARITHMETIC TABLES FOR A 4-BIT MACHINE .....	453
APPENDIX B. NEWTON'S METHOD .....	457
APPENDIX C. A GALLERY OF GRAPHS OF DISCRETE FUNCTIONS .....	459
C-1 Plots of Logical Operations on Integers .....	459
C-2 Plots of Addition, Subtraction, and Multiplication .....	461
C-3 Plots of Functions Involving Division .....	463
C-4 Plots of the Compress, SAG, and Rotate Left Functions .....	464
C-5 2D Plots of Some Unary Functions .....	466
<i>Bibliography</i> .....	471
<i>Index</i> .....	481

*This page intentionally left blank*

## FOREWORD

### Foreword from the First Edition

When I first got a summer job at MIT's Project MAC almost 30 years ago, I was delighted to be able to work with the DEC PDP-10 computer, which was more fun to program in assembly language than any other computer, bar none, because of its rich yet tractable set of instructions for performing bit tests, bit masking, field manipulation, and operations on integers. Though the PDP-10 has not been manufactured for quite some years, there remains a thriving cult of enthusiasts who keep old PDP-10 hardware running and who run old PDP-10 software—entire operating systems and their applications—by using personal computers to simulate the PDP-10 instruction set. They even write new software; there is now at least one Web site with pages that are served up by a simulated PDP-10. (Come on, stop laughing—it's no sillier than keeping antique cars running.)

I also enjoyed, in that summer of 1972, reading a brand-new MIT research memo called HAKMEM, a bizarre and eclectic potpourri of technical trivia.<sup>1</sup> The subject matter ranged from electrical circuits to number theory, but what intrigued me most was its small catalog of ingenious little programming tricks. Each such gem would typically describe some plausible yet unusual operation on integers or bit strings (such as counting the 1-bits in a word) that could easily be programmed using either a longish fixed sequence of machine instructions or a loop, and then show how the same thing might be done much more cleverly, using just four or three or two carefully chosen instructions whose interactions are not at all obvious until explained or fathomed. For me, devouring these little programming nuggets was like eating peanuts, or rather bonbons—I just couldn't stop—and there was a certain richness to them, a certain intellectual depth, elegance, even poetry.

"Surely," I thought, "there must be more of these," and indeed over the years I collected, and in some cases discovered, a few more. "There ought to be a book of them."

I was genuinely thrilled when I saw Hank Warren's manuscript. He has systematically collected these little programming tricks, organized them thematically, and explained them clearly. While some of them may be described in terms of machine instructions, this is not a book only for assembly language programmers. The subject matter is basic structural relationships among integers and bit strings

---

1. Why "HAKMEM"? Short for "hacks memo"; one 36-bit PDP-10 word could hold six 6-bit characters, so a lot of the names PDP-10 hackers worked with were limited to six characters. We were used to glancing at a six-character abbreviated name and instantly decoding the contractions. So naming the memo "HAKMEM" made sense at the time—at least to the hackers.

in a computer and efficient techniques for performing useful operations on them. These techniques are just as useful in the C or Java programming languages as they are in assembly language.

Many books on algorithms and data structures teach complicated techniques for sorting and searching, for maintaining hash tables and binary trees, for dealing with records and pointers. They overlook what can be done with very tiny pieces of data—bits and arrays of bits. It is amazing what can be done with just binary addition and subtraction and maybe some bitwise operations; the fact that the carry chain allows a single bit to affect all the bits to its left makes addition a peculiarly powerful data manipulation operation in ways that are not widely appreciated.

Yes, there ought to be a book about these techniques. Now it is in your hands, and it's terrific. If you write optimizing compilers or high-performance code, you must read this book. You otherwise might not use this bag of tricks every single day—but if you find yourself stuck in some situation where you apparently need to loop over the bits in a word, or to perform some operation on integers and it just seems harder to code than it ought, or you really need the inner loop of some integer or bit-fiddly computation to run twice as fast, then this is the place to look. Or maybe you'll just find yourself reading it straight through out of sheer pleasure.

Guy L. Steele, Jr.  
Burlington, Massachusetts  
April 2002

## PREFACE

*Caveat Emptor: The cost of software maintenance increases with the square of the programmer's creativity.*

First Law of Programmer Creativity,  
Robert D. Bliss, 1992

This is a collection of small programming tricks that I have come across over many years. Most of them will work only on computers that represent integers in two's-complement form. Although a 32-bit machine is assumed when the register length is relevant, most of the tricks are easily adapted to machines with other register sizes.

This book does not deal with large tricks such as sophisticated sorting and compiler optimization techniques. Rather, it deals with small tricks that usually involve individual computer words or instructions, such as counting the number of 1-bits in a word. Such tricks often use a mixture of arithmetic and logical instructions.

It is assumed throughout that integer overflow interrupts have been masked off, so they cannot occur. C, Fortran, and even Java programs run in this environment, but Pascal and Ada users beware!

The presentation is informal. Proofs are given only when the algorithm is not obvious, and sometimes not even then. The methods use computer arithmetic, “floor” functions, mixtures of arithmetic and logical operations, and so on. Proofs in this domain are often difficult and awkward to express.

To reduce typographical errors and oversights, many of the algorithms have been executed. This is why they are given in a real programming language, even though, like every computer language, it has some ugly features. C is used for the high-level language because it is widely known, it allows the straightforward mixture of integer and bit-string operations, and C compilers that produce high-quality object code are available.

Occasionally, machine language is used, employing a three-address format, mainly for ease of readability. The assembly language used is that of a fictitious machine that is representative of today's RISC computers.

Branch-free code is favored, because on many computers, branches slow down instruction fetching and inhibit executing instructions in parallel. Another problem with branches is that they can inhibit compiler optimizations such as instruction scheduling, commoning, and register allocation. That is, the compiler may be more effective at these optimizations with a program that consists of a few large basic blocks rather than many small ones.



The code sequences also tend to favor small immediate values, comparisons to zero (rather than to some other number), and instruction-level parallelism. Although much of the code would become more concise by using table lookups (from memory), this is not often mentioned. This is because loads are becoming more expensive relative to arithmetic instructions, and the table lookup methods are often not very interesting (although they *are* often practical). But there are exceptional cases.

Finally, I should mention that the term “hacker” in the title is meant in the original sense of an aficionado of computers—someone who enjoys making computers do new things, or do old things in a new and clever way. The hacker is usually quite good at his craft, but may very well not be a professional computer programmer or designer. The hacker’s work may be useful or may be just a game. As an example of the latter, more than one determined hacker has written a program which, when executed, writes out an exact copy of itself.<sup>1</sup> This is the sense in which we use the term “hacker.” If you’re looking for tips on how to break into someone else’s computer, you won’t find them here.

## Acknowledgments

First, I want to thank Bruce Shriver and Dennis Allison for encouraging me to publish this book. I am indebted to many colleagues at IBM, several of whom are cited in the Bibliography. One deserves special mention: Martin E. Hopkins, whom I think of as “Mr. Compiler” at IBM, has been relentless in his drive to make every cycle count, and I’m sure some of his spirit has rubbed off on me. Addison-Wesley’s reviewers have improved the book immensely. Most of their names are unknown to me, but the review by one whose name I did learn was truly outstanding: Guy L. Steele, Jr., completed a 50-page review that included new subject areas to address, such as bit shuffling and unshuffling, the sheep and goats operation, and many others. He suggested algorithms that beat the ones I used. He was extremely thorough. For example, I had erroneously written that the hexadecimal number AAAAAAAAA factors as  $2 \cdot 3 \cdot 17 \cdot 257 \cdot 65537$ ; Guy pointed out that the 3 should be a 5. He suggested improvements to style and did not shirk from mentioning minutiae. Wherever you see “parallel prefix” in this book, the material is due to Guy.

See [www.HackersDelight.org](http://www.HackersDelight.org) for additional material related to this book.

H. S. Warren, Jr.  
Yorktown, New York  
June 2012

---

1. One such program, written in C, is:

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

## CHAPTER 2

### BASICS

#### 2-1 Manipulating Rightmost Bits

Some of the formulas in this section find application in later chapters.

Use the following formula to turn off the rightmost 1-bit in a word, producing 0 if none (e.g., 01011000  $\Rightarrow$  01010000):

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

This can be used to determine if an unsigned integer is a power of 2 or is 0: apply the formula followed by a 0-test on the result.

Use the following formula to turn on the rightmost 0-bit in a word, producing all 1's if none (e.g., 10100111  $\Rightarrow$  10101111):

$$x \mid (x + 1)$$

Use the following formula to turn off the trailing 1's in a word, producing  $x$  if none (e.g., 10100111  $\Rightarrow$  10100000):

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

This can be used to determine if an unsigned integer is of the form  $2^n - 1$ , 0, or all 1's: apply the formula followed by a 0-test on the result.

Use the following formula to turn on the trailing 0's in a word, producing  $x$  if none (e.g., 10101000  $\Rightarrow$  10101111):

$$x \mid (x - 1)$$

Use the following formula to create a word with a single 1-bit at the position of the rightmost 0-bit in  $x$ , producing 0 if none (e.g., 10100111  $\Rightarrow$  00001000):

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

Use the following formula to create a word with a single 0-bit at the position of the rightmost 1-bit in  $x$ , producing all 1's if none (e.g., 10101000  $\Rightarrow$  11110111):

$$\neg x \mid (x - 1)$$

Use one of the following formulas to create a word with 1's at the positions of the trailing 0's in  $x$ , and 0's elsewhere, producing 0 if none (e.g., 01011000  $\Rightarrow$  00000111):

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

The first formula has some instruction-level parallelism.

Use the following formula to create a word with 0's at the positions of the trailing 1's in  $x$ , and 0's elsewhere, producing all 1's if none (e.g., 10100111  $\Rightarrow$  11111000):

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

Use the following formula to isolate the rightmost 1-bit, producing 0 if none (e.g., 01011000  $\Rightarrow$  00001000):

$$x \& (\neg x)$$

Use the following formula to create a word with 1's at the positions of the rightmost 1-bit and the trailing 0's in  $x$ , producing all 1's if no 1-bit, and the integer 1 if no trailing 0's (e.g., 01011000  $\Rightarrow$  00001111):

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

Use the following formula to create a word with 1's at the positions of the rightmost 0-bit and the trailing 1's in  $x$ , producing all 1's if no 0-bit, and the integer 1 if no trailing 1's (e.g., 01010111  $\Rightarrow$  00001111):

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

Use either of the following formulas to turn off the rightmost contiguous string of 1's (e.g., 01011100  $\Rightarrow$  01000000) [Wood]:

$$\begin{aligned} &(((x \mid (x - 1)) + 1) \& x), \text{ or} \\ &((x \& \neg x) + x) \& x \end{aligned}$$

These can be used to determine if a nonnegative integer is of the form  $2^j - 2^k$  for some  $j \geq k \geq 0$ : apply the formula followed by a 0-test on the result.

### De Morgan's Laws Extended

The logical identities known as De Morgan's laws can be thought of as distributing, or "multiplying in," the *not* sign. This idea can be extended to apply to the expressions of this section, and a few more, as shown here. (The first two are De Morgan's laws.)

$$\begin{aligned}
\neg(x \& y) &= \neg x \mid \neg y \\
\neg(x \mid y) &= \neg x \& \neg y \\
\neg(x + 1) &= \neg x - 1 \\
\neg(x - 1) &= \neg x + 1 \\
\neg x &= x - 1 \\
\neg(x \oplus y) &= \neg x \oplus y = x \equiv y \\
\neg(x \equiv y) &= \neg x \equiv y = x \oplus y \\
\neg(x + y) &= \neg x - y \\
\neg(x - y) &= \neg x + y
\end{aligned}$$

As an example of the application of these formulas,  $\neg(x \mid \neg(x + 1)) = \neg x \& \neg\neg(x + 1) = \neg x \& ((x + 1) - 1) = \neg x \& x = 0$ .

### Right-to-Left Computability Test

There is a simple test to determine whether or not a given function can be implemented with a sequence of *add*'s, *subtract*'s, *and*'s, *or*'s, and *not*'s [War]. We can, of course, expand the list with other instructions that can be composed from the basic list, such as *shift left* by a fixed amount (which is equivalent to a sequence of *add*'s), or *multiply*. However, we exclude instructions that cannot be composed from the list. The test is contained in the following theorem.

*THEOREM. A function mapping words to words can be implemented with word-parallel add, subtract, and, or, and not instructions if and only if each bit of the result depends only on bits at and to the right of each input operand.*

That is, imagine trying to compute the rightmost bit of the result by looking only at the rightmost bit of each input operand. Then, try to compute the next bit to the left by looking only at the rightmost two bits of each input operand, and continue in this way. If you are successful in this, then the function can be computed with a sequence of *add*'s, *and*'s, and so on. If the function cannot be computed in this right-to-left manner, then it cannot be implemented with a sequence of such instructions.

The interesting part of this is the latter statement, and it is simply the contrapositive of the observation that the functions *add*, *subtract*, *and*, *or*, and *not* can all be computed in the right-to-left manner, so any combination of them must have this property.

To see the "if" part of the theorem, we need a construction that is a little awkward to explain. We illustrate it with a specific example. Suppose that a function of two variables  $x$  and  $y$  has the right-to-left computability property, and suppose that bit 2 of the result  $r$  is given by

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

We number bits from right to left, 0 to 31. Because bit 2 of the result is a function of bits at and to the right of bit 2 of the input operands, bit 2 of the result is “right-to-left computable.”

Arrange the computer words  $x$ ,  $x$  shifted left two, and  $y$  shifted left one, as shown below. Also, add a mask that isolates bit 2.

$$\begin{array}{cccccccc}
 x_{31} & x_{30} & \cdots & x_3 & x_2 & x_1 & x_0 & \\
 x_{29} & x_{28} & \cdots & x_1 & x_0 & 0 & 0 & \\
 y_{30} & y_{29} & \cdots & y_2 & y_1 & y_0 & 0 & \\
 0 & 0 & \cdots & 0 & 1 & 0 & 0 & \\
 0 & 0 & \cdots & 0 & r_2 & 0 & 0 & 
 \end{array}$$

Now, form the word-parallel *and* of lines 2 and 3, *or* the result with row 1 (following Equation (1)), and *and* the result with the mask (row 4 above). The result is a word of all 0's except for the desired result bit in position 2. Perform similar computations for the other bits of the result, *or* the 32 resulting words together, and the result is the desired function.

This construction does not yield an efficient program; rather, it merely shows that it can be done with instructions in the basic list.

Using the theorem, we immediately see that there is no sequence of such instructions that turns off the leftmost 1-bit in a word, because to see if a certain 1-bit should be turned off, we must look to the left to see if it is the leftmost one. Similarly, there can be no such sequence for performing a right shift, or a rotate shift, or a left shift by a variable amount, or for counting the number of trailing 0's in a word (to count trailing 0's, the rightmost bit of the result will be 1 if there are an odd number of trailing 0's, and we must look to the left of the rightmost position to determine that).

### A Novel Application

An application of the sort of bit twiddling discussed above is the problem of finding the next higher number after a given number that has the same number of 1-bits. You might very well wonder why anyone would want to compute that. It has application where bit strings are used to represent subsets. The possible members of a set are listed in a linear array, and a subset is represented by a word or sequence of words in which bit  $i$  is on if member  $i$  is in the subset. Set unions are computed by the logical *or* of the bit strings, intersections by *and*'s, and so on.

You might want to iterate through all the subsets of a given size. This is easily done if you have a function that maps a given subset to the next higher number (interpreting the subset string as an integer) with the same number of 1-bits.

A concise algorithm for this operation was devised by R. W. Gosper [HAK, item 175].<sup>1</sup> Given a word  $x$  that represents a subset, the idea is to find the

---

1. A variation of this algorithm appears in [H&S] sec. 7.6.7.

rightmost contiguous group of 1's in  $x$  and the following 0's, and "increment" that quantity to the next value that has the same number of 1's. For example, the string `xxx0 1111 0000`, where `xxx` represents arbitrary bits, becomes `xxx1 0000 0111`. The algorithm first identifies the "smallest" 1-bit in  $x$ , with  $s = x \& -x$ , giving `0000 0001 0000`. This is added to  $x$ , giving  $r = xxx1\ 0000\ 0000$ . The 1-bit here is one bit of the result. For the other bits, we need to produce a right-adjusted string of  $n - 1$  1's, where  $n$  is the size of the rightmost group of 1's in  $x$ . This can be done by first forming the *exclusive or* of  $r$  and  $x$ , which gives `0001 1111 0000` in our example.

This has two too many 1's and needs to be right-adjusted. This can be accomplished by dividing it by  $s$ , which right-adjusts it ( $s$  is a power of 2), and shifting it right two more positions to discard the two unwanted bits. The final result is the *or* of this and  $r$ .

In computer algebra notation, the result is  $y$  in

$$\begin{aligned} s &\leftarrow x \& -x \\ r &\leftarrow s + x \\ y &\leftarrow r \mid (((x \oplus r) \ggg 2) \div s) \end{aligned} \tag{2}$$

A complete C procedure is given in Figure 2-1. It executes in seven basic RISC instructions, one of which is division. (Do not use this procedure with  $x = 0$ ; that causes division by 0.)

If division is slow but you have a fast way to compute the *number of trailing zeros* function  $\text{ntz}(x)$ , the *number of leading zeros* function  $\text{nlz}(x)$ , or *population count* ( $\text{pop}(x)$  is the number of 1-bits in  $x$ ), then the last line of Equation (2) can be replaced with one of the following formulas. (The first two methods can fail on a machine that has modulo 32 shifts.)

$$\begin{aligned} y &\leftarrow r \mid ((x \oplus r) \ggg (2 + \text{ntz}(x))) \\ y &\leftarrow r \mid ((x \oplus r) \ggg (33 - \text{nlz}(s))) \\ y &\leftarrow r \mid ((1 \ll (\text{pop}(x \oplus r) - 2)) - 1) \end{aligned}$$

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

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

FIGURE 2-1. Next higher number with same number of 1-bits.

## 2-2 Addition Combined with Logical Operations

We assume the reader is familiar with the elementary identities of ordinary algebra and Boolean algebra. Below is a selection of similar identities involving addition and subtraction combined with logical operations.

- a.  $-x = \neg x + 1$
- b.  $\quad \quad = \neg(x - 1)$
- c.  $\quad \quad -x = -x - 1$
- d.  $- \neg x = x + 1$
- e.  $\neg -x = x - 1$
- f.  $x + y = x - \neg y - 1$
- g.  $\quad \quad = (x \oplus y) + 2(x \& y)$
- h.  $\quad \quad = (x \mid y) + (x \& y)$
- i.  $\quad \quad = 2(x \mid y) - (x \oplus y)$
- j.  $x - y = x + \neg y + 1$
- k.  $\quad \quad = (x \oplus y) - 2(\neg x \& y)$
- l.  $\quad \quad = (x \& \neg y) - (\neg x \& y)$
- m.  $\quad \quad = 2(x \& \neg y) - (x \oplus y)$
- n.  $x \oplus y = (x \mid y) - (x \& y)$
- o.  $x \& \neg y = (x \mid y) - y$
- p.  $\quad \quad = x - (x \& y)$
- q.  $\neg(x - y) = y - x - 1$
- r.  $\quad \quad = \neg x + y$
- s.  $x \equiv y = (x \& y) - (x \mid y) - 1$
- t.  $\quad \quad = (x \& y) + \neg(x \mid y)$
- u.  $x \mid y = (x \& \neg y) + y$
- v.  $x \& y = (\neg x \mid y) - \neg x$

Equation (d) can be applied to itself repeatedly, giving  $- \neg - \neg x = x + 2$ , and so on. Similarly, from (e) we have  $\neg - \neg - x = x - 2$ . So we can add or subtract any constant using only the two forms of complementation.

Equation (f) is the dual of (j), where (j) is the well-known relation that shows how to build a subtracter from an adder.

Equations (g) and (h) are from HAKMEM memo [HAK, item 23]. Equation (g) forms a sum by first computing the sum with carries ignored ( $x \oplus y$ ), and then adding in the carries. Equation (h) is simply modifying the addition operands so that the combination  $0 + 1$  never occurs at any bit position; it is replaced with  $1 + 0$ .

It can be shown that in the ordinary addition of binary numbers with each bit independently equally likely to be 0 or 1, a carry occurs at each position with probability about 0.5. However, for an adder built by preconditioning the inputs using (g), the probability is about 0.25. This observation is probably not of value in building an adder, because for that purpose the important characteristic is the maximum number of logic circuits the carry must pass through, and using (g) reduces the number of stages the carry propagates through by only one.

Equations (k) and (l) are duals of (g) and (h), for subtraction. That is, (k) has the interpretation of first forming the difference ignoring the borrows ( $\mathbf{x} \oplus \mathbf{y}$ ), and then subtracting the borrows. Similarly, Equation (l) is simply modifying the subtraction operands so that the combination  $1 - 1$  never occurs at any bit position; it is replaced with  $0 - 0$ .

Equation (n) shows how to implement *exclusive or* in only three instructions on a basic RISC. Using only *and-or-not* logic requires four instructions ( $(x \mid y) \& \neg(x \& y)$ ). Similarly, (u) and (v) show how to implement *and* and *or* in three other elementary instructions, whereas using DeMorgan's laws requires four.

### 2–3 Inequalities among Logical and Arithmetic Expressions

Inequalities among binary logical expressions whose values are interpreted as unsigned integers are nearly trivial to derive. Here are two examples:

$$(x \oplus y) \leq^u (x \mid y), \text{ and}$$

$$(x \& y) \leq^u (x \equiv y).$$

These can be derived from a list of all binary logical operations, shown in Table 2–1.

Let  $f(x, y)$  and  $g(x, y)$  represent two columns in Table 2-1. If for each row in which  $f(x, y)$  is 1,  $g(x, y)$  also is 1, then for all  $(x, y)$ ,  $f(x, y) \leq g(x, y)$ . Clearly, this extends to word-parallel logical operations. One can easily read off such relations (most of which are trivial) as  $(x \& y) \leq x \leq (x \mid \neg y)$ , and so on. Furthermore, if two columns have a row in which one entry is 0 and the other is 1,

TABLE 2-1. THE 16 BINARY LOGICAL OPERATIONS

[illegible]



and another row in which the entries are 1 and 0, respectively, then no inequality relation exists between the corresponding logical expressions. So the question of whether or not  $f(x, y) \stackrel{u}{\leq} g(x, y)$  is completely and easily solved for all binary logical functions  $f$  and  $g$ .

Use caution when manipulating these relations. For example, for ordinary arithmetic, if  $x + y \leq a$  and  $z \leq x$ , then  $z + y \leq a$ , but this inference is not valid if “+” is replaced with *or*.

Inequalities involving mixed logical and arithmetic expressions are more interesting. Below is a small selection.

- a.  $(x \mid y) \stackrel{u}{\geq} \max(x, y)$
- b.  $(x \& y) \stackrel{u}{\leq} \min(x, y)$
- c.  $(x \mid y) \stackrel{u}{\leq} x + y$  if the addition does not overflow
- d.  $(x \mid y) \stackrel{u}{\geq} x + y$  if the addition overflows
- e.  $|x - y| \stackrel{u}{\leq} (x \oplus y)$

The proofs of these are quite simple, except possibly for the relation  $|x - y| \stackrel{u}{\leq} (x \oplus y)$ . By  $|x - y|$  we mean the absolute value of  $x - y$ , which can be computed within the domain of unsigned numbers as  $\max(x, y) - \min(x, y)$ . This relation can be proven by induction on the length of  $x$  and  $y$  (the proof is a little easier if you extend them on the left rather than on the right).

## 2-4 Absolute Value Function

If your machine does not have an instruction for computing the absolute value, this computation can usually be done in three or four branch-free instructions. First, compute  $y \leftarrow x \ggg 31$ , and then one of the following:

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

By “ $2x$ ” we mean, of course,  $x + x$  or  $x \ll 1$ .

If you have fast multiplication by a variable whose value is  $\pm 1$ , the following will do:

$$((x \ggg 30) \mid 1) * x$$

## 2-5 Average of Two Integers

The following formula can be used to compute the average of two unsigned integers,  $\lfloor (x + y)/2 \rfloor$ , without causing overflow [Dietz]:

$$(x \& y) + ((x \oplus y) \ggg 1) \quad (3)$$

The formula below computes  $\lceil (x + y)/2 \rceil$  for unsigned integers:

$$(x \mid y) - ((x \oplus y) \ggg 1)$$

To compute the same quantities (“floor and ceiling averages”) for signed integers, use the same formulas, but with the unsigned shift replaced with a signed shift.

For signed integers, one might also want the average with the division by 2 rounded toward 0. Computing this “truncated average” (without causing overflow) is a little more difficult. It can be done by computing the floor average and then correcting it. The correction is to add 1 if, arithmetically,  $x + y$  is negative and odd. But  $x + y$  is negative if and only if the result of (3), with the unsigned shift replaced with a signed shift, is negative. This leads to the following method (seven instructions on the basic RISC, after commoning the subexpression  $x \oplus y$ ):

$$\begin{aligned} t &\leftarrow (x \& y) + ((x \oplus y) \ggg 1); \\ t &+ ((t \ggg 31) \& (x \oplus y)) \end{aligned}$$

Some common special cases can be done more efficiently. If  $x$  and  $y$  are signed integers and known to be nonnegative, then the average can be computed as simply  $(x + y) \ggg 1$ . The sum can overflow, but the overflow bit is retained in the register that holds the sum, so that the unsigned shift moves the overflow bit to the proper position and supplies a zero sign bit.

If  $x$  and  $y$  are unsigned integers and  $x \leq y$ , or if  $x$  and  $y$  are signed integers and  $x \leq y$  (signed comparison), then the average is given by  $x + ((y - x) \ggg 1)$ . These are floor averages, for example, the average of  $-1$  and  $0$  is  $-1$ .

## 2-6 Sign Extension

By “sign extension,” we mean to consider a certain bit position in a word to be the sign bit, and we wish to propagate that to the left, ignoring any other bits present. The standard way to do this is with *shift left logical* followed by *shift right signed*. However, if these instructions are slow or nonexistent on your machine, it can be

done with one of the following, where we illustrate by propagating bit position 7 to the left:

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

The “+” above can also be “−” or “ $\oplus$ .” The second formula is particularly useful if you know that the unwanted high-order bits are all 0’s, because then the *and* can be omitted.

## 2-7 Shift Right Signed from Unsigned

If your machine does not have the *shift right signed* instruction, it can be computed using the formulas shown below. The first formula is from [GM], and the second is based on the same idea. These formulas hold for  $0 \leq n \leq 31$  and, if the machine has mod-64 shifts, the last holds for  $0 \leq n \leq 63$ . The last formula holds for any  $n$  if by “holds” we mean “treats the shift amount to the same modulus as does the logical shift.”

When  $n$  is a variable, each formula requires five or six instructions on a basic RISC.

$$\begin{aligned} & ((x + 0x80000000) \ggg n) - (0x80000000 \ggg n) \\ & t \leftarrow 0x80000000 \ggg n; \quad ((x \ggg n) \oplus t) - t \\ & t \leftarrow (x \& 0x80000000) \ggg n; (x \ggg n) - (t + t) \\ & (x \ggg n) \mid (-(x \ggg 31) \ll 31 - n) \\ & t \leftarrow -(x \ggg 31); \quad ((x \oplus t) \ggg n) \oplus t \end{aligned}$$

In the first two formulas, an alternative for the expression  $0x80000000 \ggg n$  is  $1 \ll 31 - n$ .

If  $n$  is a constant, the first two formulas require only three instructions on many machines. If  $n = 31$ , the function can be done in two instructions with  $-(x \ggg 31)$ .

## 2-8 Sign Function

The *sign*, or *signum*, function is defined by

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

It can be calculated with four instructions on most machines [Hop]:

$$(x \overset{s}{\gg} 31) \mid (-x \overset{u}{\gg} 31)$$

If you don't have *shift right signed*, then use the substitute noted at the end of Section 2-7, giving the following nicely symmetric formula (five instructions):

$$-(x \overset{u}{\gg} 31) \mid (-x \overset{u}{\gg} 31)$$

Comparison predicate instructions permit a three-instruction solution, with either

$$\begin{aligned} (x > 0) - (x < 0), \text{ or} \\ (x \geq 0) - (x \leq 0). \end{aligned} \tag{4}$$

Finally, we note that the formula  $(-x \overset{u}{\gg} 31) - (x \overset{u}{\gg} 31)$  almost works; it fails only for  $x = -2^{31}$ .

## 2-9 Three-Valued Compare Function

The *three-valued compare* function, a slight generalization of the *sign* function, is defined by

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

There are both signed and unsigned versions, and unless otherwise specified, this section applies to both.

Comparison predicate instructions permit a three-instruction solution, an obvious generalization of Equations in (4):

$$\begin{aligned} (x > y) - (x < y), \text{ or} \\ (x \geq y) - (x \leq y). \end{aligned}$$

A solution for unsigned integers on PowerPC is shown below [CWG]. On this machine, “carry” is “not borrow.”

```
subf  R5,Ry,Rx    # R5 <-- Rx - Ry.
subfc R6,Rx,Ry    # R6 <-- Ry - Rx, set carry.
subfe R7,Ry,Rx    # R7 <-- Rx - Ry + carry, set carry.
subfe R8,R7,R5    # R8 <-- R5 - R7 + carry, (set carry).
```

If limited to the instructions of the basic RISC, there does not seem to be any particularly good way to compute this function. The comparison predicates  $x < y$ ,  $x \leq y$ , and so on, require about five instructions (see Section 2-12), leading to a solution in about 12 instructions (using a small amount of commonality in computing  $x < y$  and  $x > y$ ). On the basic RISC it's probably preferable to use compares and branches (six instructions executed worst case if compares can be commoned).

## 2-10 Transfer of Sign Function

The *transfer of sign* function, called ISIGN in Fortran, is defined by

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

This function can be calculated (modulo  $2^{32}$ ) with four instructions on most machines:

$$\begin{array}{ll} t \leftarrow y \ggg 31; & t \leftarrow (x \oplus y) \ggg 31; \\ \text{ISIGN}(x, y) = (\text{abs}(x) \oplus t) - t & \text{ISIGN}(x, y) = (x \oplus t) - t \\ = (\text{abs}(x) + t) \oplus t & = (x + t) \oplus t \end{array}$$

## 2-11 Decoding a “Zero Means $2^n$ ” Field

Sometimes a 0 or negative value does not make much sense for a quantity, so it is encoded in an  $n$ -bit field with a 0 value being understood to mean  $2^n$ , and a non-zero value having its normal binary interpretation. An example is the length field of PowerPC's *load string word immediate* (lswi) instruction, which occupies five bits. It is not useful to have an instruction that loads zero bytes when the length is an immediate quantity, but it is definitely useful to be able to load 32 bytes. The length field could be encoded with values from 0 to 31 denoting lengths from 1 to 32, but the “zero means 32” convention results in simpler logic when the processor must also support a corresponding instruction with a variable (in-register) length that employs straight binary encoding (e.g., PowerPC's lswx instruction).

It is trivial to encode an integer in the range 1 to  $2^n$  into the “zero means  $2^n$ ” encoding—simply mask the integer with  $2^n - 1$ . To do the decoding without a test-and-branch is not quite as simple, but here are some possibilities, illustrated for a 3-bit field. They all require three instructions, not counting possible loads of constants.

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

## 2-12 Comparison Predicates

A “comparison predicate” is a function that compares two quantities, producing a single bit result of 1 if the comparison is **true**, and 0 if the comparison is **false**. Below we show branch-free expressions to evaluate the result into the sign position. To produce the 1/0 value used by some languages (e.g., C), follow the code with a *shift right* of 31. To produce the  $-1/0$  result used by some other languages (e.g., Basic), follow the code with a *shift right signed* of 31.

These formulas are, of course, not of interest on machines such as MIPS and our model RISC, which have comparison instructions that compute many of these predicates directly, placing a 0/1-valued result in a general purpose register.

$$\begin{array}{ll}
x = y: & \begin{array}{l} \text{abs}(x - y) - 1 \\ \text{abs}(x - y + 0x80000000) \\ \text{nlz}(x - y) \ll 26 \\ -(\text{nlz}(x - y) \stackrel{u}{\gg} 5) \\ \neg(x - y \mid y - x) \end{array} \\
x \neq y: & \begin{array}{l} \text{nabs}(x - y) \\ \text{nlz}(x - y) - 32 \\ x - y \mid y - x \end{array} \\
x < y: & \begin{array}{l} (x - y) \oplus [(x \oplus y) \& ((x - y) \oplus x)] \\ (x \& \neg y) \mid ((x \equiv y) \& (x - y)) \\ \text{nabs}(\text{doz}(y, x)) \end{array} \quad [\text{GSO}] \\
x \leq y: & \begin{array}{l} (x \mid \neg y) \& ((x \oplus y) \mid \neg(y - x)) \\ ((x \equiv y) \stackrel{s}{\gg} 1) + (x \& \neg y) \end{array} \quad [\text{GSO}] \\
x \stackrel{u}{\leq} y: & \begin{array}{l} (\neg x \& y) \mid ((x \equiv y) \& (x - y)) \\ (\neg x \& y) \mid ((\neg x \mid y) \& (x - y)) \end{array} \\
x \stackrel{u}{\leq} y: & (\neg x \mid y) \& ((x \oplus y) \mid \neg(y - x))
\end{array}$$

A machine instruction that computes the negative of the absolute value is handy here. We show this function as “nabs.” Unlike absolute value, it is well defined in that it never overflows. Machines that do not have nabs, but have the more usual abs, can use  $-\text{abs}(x)$  for  $\text{nabs}(x)$ . If  $x$  is the maximum negative

number, this overflows twice, but the result is correct. (We assume that the absolute value and the negation of the maximum negative number is itself.) Because some machines have neither abs nor nabs, we give an alternative that does not use them.

The “nlz” function is the number of leading 0’s in its argument. The “doz” function (*difference or zero*) is described on page 41. For  $x > y$ ,  $x \geq y$ , and so on, interchange  $x$  and  $y$  in the formulas for  $x < y$ ,  $x \leq y$ , and so on. The *add* of **0x8000 0000** can be replaced with any instruction that inverts the high-order bit (in  $x$ ,  $y$ , or  $x - y$ ).

Another class of formulas can be derived from the observation that the predicate  $x < y$  is given by the sign of  $x/2 - y/2$ , and the subtraction in that expression cannot overflow. The result can be fixed up by subtracting 1 in the cases in which the shifts discard essential information, as follows:

$$\begin{aligned} x < y: & \quad (x \gg^s 1) - (y \gg^s 1) - (\neg x \& y \& 1) \\ x \leq y: & \quad (x \gg^u 1) - (y \gg^u 1) - (\neg x \& y \& 1) \end{aligned}$$

These execute in seven instructions on most machines (six if it has *and not*), which is no better than what we have above (five to seven instructions, depending upon the fullness of the set of logic instructions).

The formulas above involving nlz are due to [Shep], and his formula for the  $x = y$  predicate is particularly useful, because a minor variation of it gets the predicate evaluated to a 1/0-valued result with only three instructions:

$$\text{nlz}(x - y) \gg^u 5.$$

Signed comparisons to 0 are frequent enough to deserve special mention. There are some formulas for these, mostly derived directly from the above. Again, the result is in the sign position.

$$\begin{aligned} x = 0: & \quad \text{abs}(x) - 1 \\ & \quad \text{abs}(x + \mathbf{0x8000\,0000}) \\ & \quad \text{nlz}(x) \ll 26 \\ & \quad -(\text{nlz}(x) \gg^u 5) \\ & \quad \neg(x \mid -x) \\ & \quad \neg x \& (x - 1) \\ x \neq 0: & \quad \text{nabs}(x) \\ & \quad \text{nlz}(x) - 32 \\ & \quad x \mid -x \\ & \quad (x \gg^u 1) - x \quad \quad \quad [\text{CWG}] \end{aligned}$$

$$\begin{aligned}
x < 0: & \quad x \\
x \leq 0: & \quad x \mid (x - 1) \\
& \quad x \mid \neg x \\
x > 0: & \quad x \oplus \text{nabs}(x) \\
& \quad (x \gg^s 1) - x \\
& \quad \neg x \& \neg x \\
x \geq 0: & \quad \neg x
\end{aligned}$$

Signed comparisons can be obtained from their unsigned counterparts by biasing the signed operands upward by  $2^{31}$  and interpreting the results as unsigned integers. The reverse transformation also works.<sup>2</sup> Thus, we have

$$\begin{aligned}
x < y &= x + 2^{31} \stackrel{u}{<} y + 2^{31}, \\
x \stackrel{u}{<} y &= x - 2^{31} < y - 2^{31}.
\end{aligned}$$

Similar relations hold for  $\leq$ ,  $\stackrel{u}{\leq}$ , and so on. In these relations, one can use addition, subtraction, or *exclusive or* with  $2^{31}$ . They are all equivalent, as they simply invert the sign bit. An instruction like the basic RISC's *add immediate shifted* is useful to avoid loading the constant  $2^{31}$ .

Another way to get signed comparisons from unsigned is based on the fact that if  $x$  and  $y$  have the same sign, then  $x < y = x \stackrel{u}{<} y$ , whereas if they have opposite signs, then  $x < y = x \stackrel{u}{>} y$  [Lamp]. Again, the reverse transformation also works, so we have

$$\begin{aligned}
x < y &= (x \stackrel{u}{<} y) \oplus x_{31} \oplus y_{31} \text{ and} \\
x \stackrel{u}{<} y &= (x < y) \oplus x_{31} \oplus y_{31},
\end{aligned}$$

where  $x_{31}$  and  $y_{31}$  are the sign bits of  $x$  and  $y$ , respectively. Similar relations hold for  $\leq$ ,  $\stackrel{u}{\leq}$ , and so on.

Using either of these devices enables computing all the usual comparison predicates other than  $=$  and  $\neq$  in terms of any one of them, with at most three additional instructions on most machines. For example, let us take  $x \stackrel{u}{\leq} y$  as primitive, because it is one of the simplest to implement (it is the carry bit from  $y - x$ ). Then the other predicates can be obtained as follows:

$$\begin{aligned}
x < y &= \neg(y + 2^{31} \stackrel{u}{\leq} x + 2^{31}) \\
x \leq y &= x + 2^{31} \stackrel{u}{\leq} y + 2^{31}
\end{aligned}$$

---

2. This is useful to get unsigned comparisons in Java, which lacks unsigned integers.



$$x > y = \neg(x + 2^{31} \stackrel{u}{\leq} y + 2^{31})$$

$$x \geq y = y + 2^{31} \stackrel{u}{\leq} x + 2^{31}$$

$$x \stackrel{u}{<} y = \neg(y \stackrel{u}{\leq} x)$$

$$x \stackrel{u}{>} y = \neg(x \stackrel{u}{\leq} y)$$

$$x \stackrel{u}{\geq} y = y \stackrel{u}{\leq} x$$

### Comparison Predicates from the Carry Bit

If the machine can easily deliver the carry bit into a general purpose register, this may permit concise code for some of the comparison predicates. Below are several of these relations. The notation  $\text{carry}(\text{expression})$  means the carry bit generated by the outermost operation in *expression*. We assume the carry bit for the subtraction  $x - y$  is what comes out of the adder for  $x + \bar{y} + 1$ , which is the complement of “borrow.”

$x = y$ :	$\text{carry}(0 - (x - y))$ , or $\text{carry}((x + \bar{y}) + 1)$ , or $\text{carry}((x - y - 1) + 1)$
$x \neq y$ :	$\text{carry}((x - y) - 1)$ , i.e., $\text{carry}((x - y) + (-1))$
$x < y$ :	$\neg\text{carry}((x + 2^{31}) - (y + 2^{31}))$ , or $\neg\text{carry}(x - y) \oplus x_{31} \oplus y_{31}$
$x \leq y$ :	$\text{carry}((y + 2^{31}) - (x + 2^{31}))$ , or $\text{carry}(y - x) \oplus x_{31} \oplus y_{31}$
$x \stackrel{u}{<} y$ :	$\neg\text{carry}(x - y)$
$x \stackrel{u}{\leq} y$ :	$\text{carry}(y - x)$
$x = 0$ :	$\text{carry}(0 - x)$ , or $\text{carry}(\bar{x} + 1)$
$x \neq 0$ :	$\text{carry}(x - 1)$ , i.e., $\text{carry}(x + (-1))$
$x < 0$ :	$\text{carry}(x + x)$
$x \leq 0$ :	$\text{carry}(2^{31} - (x + 2^{31}))$

For  $x > y$ , use the complement of the expression for  $x \leq y$ , and similarly for other relations involving “greater than.”

The GNU Superoptimizer has been applied to the problem of computing predicate expressions on the IBM RS/6000 computer and its close relative PowerPC [GK]. The RS/6000 has instructions for  $\text{abs}(x)$ ,  $\text{nabs}(x)$ ,  $\text{doz}(x, y)$ , and a number of forms of *add* and *subtract* that use the carry bit. It was found that the RS/6000 can

compute all the integer predicate expressions with three or fewer elementary (one-cycle) instructions, a result that surprised even the architects of the machine. “All” includes the six two-operand signed comparisons and the four two-operand unsigned comparisons, all of these with the second operand being 0, and all in forms that produce a 1/0 result or a -1/0 result. PowerPC, which lacks  $\text{abs}(x)$ ,  $\text{nabs}(x)$ , and  $\text{doz}(x, y)$ , can compute all the predicate expressions in four or fewer elementary instructions.

### How the Computer Sets the Comparison Predicates

Most computers have a way of evaluating the integer comparison predicates to a 1-bit result. The result bit may be placed in a “condition register” or, for some machines (such as our RISC model), in a general purpose register. In either case, the facility is often implemented by subtracting the comparison operands and then performing a small amount of logic on the result bits to determine the 1-bit comparison result.

Below is the logic for these operations. It is assumed that the machine computes  $x - y$  as  $x + \bar{y} + 1$ , and the following quantities are available in the result:

$C_o$ , the carry out of the high-order position

$C_i$ , the carry into the high-order position

$N$ , the sign bit of the result

$Z$ , which equals 1 if the result, exclusive of  $C_o$ , is all-0, and is otherwise 0

Then we have the following in Boolean algebra notation (juxtaposition denotes *and*,  $+$  denotes *or*):

$V$ :	$C_i \oplus C_o$	(signed overflow)
$x = y$ :	$Z$	
$x \neq y$ :	$\bar{Z}$	
$x < y$ :	$N \oplus V$	
$x \leq y$ :	$(N \oplus V) + Z$	
$x > y$ :	$(N \equiv V)\bar{Z}$	
$x \geq y$ :	$N \equiv V$	
$x \stackrel{u}{<} y$ :	$\overline{C_o}$	
$x \stackrel{u}{\leq} y$ :	$\overline{C_o} + Z$	
$x \stackrel{u}{>} y$ :	$C_o \bar{Z}$	
$x \stackrel{u}{\geq} y$ :	$C_o$	

## 2-13 Overflow Detection

“Overflow” means that the result of an arithmetic operation is too large or too small to be correctly represented in the target register. This section discusses methods that a programmer might use to detect when overflow has occurred, without using the machine’s “status bits” that are often supplied expressly for this purpose. This is important, because some machines do not have such status bits (e.g., MIPS), and even if the machine is so equipped, it is often difficult or impossible to access the bits from a high-level language.

### Signed Add/Subtract

When overflow occurs on integer addition and subtraction, contemporary machines invariably discard the high-order bit of the result and store the low-order bits that the adder naturally produces. Signed integer overflow of addition occurs if and only if the operands have the same sign and the sum has a sign opposite to that of the operands. Surprisingly, this same rule applies even if there is a carry into the adder—that is, if the calculation is  $x + y + 1$ . This is important for the application of adding multiword signed integers, in which the last addition is a signed addition of two fullwords and a carry-in that may be 0 or +1.

To prove the rule for addition, let  $x$  and  $y$  denote the values of the one-word signed integers being added, let  $c$  (carry-in) be 0 or 1, and assume for simplicity a 4-bit machine. Then if the signs of  $x$  and  $y$  are different,

$$\begin{aligned} -8 \leq x \leq -1, \text{ and} \\ 0 \leq y \leq 7, \end{aligned}$$

or similar bounds apply if  $x$  is nonnegative and  $y$  is negative. In either case, by adding these inequalities and optionally adding in 1 for  $c$ ,

$$-8 \leq x + y + c \leq 7.$$

This is representable as a 4-bit signed integer, and thus overflow does not occur when the operands have opposite signs.

Now suppose  $x$  and  $y$  have the same sign. There are two cases:

(a)	(b)
$-8 \leq x \leq -1$	$0 \leq x \leq 7$
$-8 \leq y \leq -1$	$0 \leq y \leq 7$

Thus,

(a)	(b)
$-16 \leq x + y + c \leq -1$	$0 \leq x + y + c \leq 15.$

Overflow occurs if the sum is not representable as a 4-bit signed integer—that is, if

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

In case (a), this is equivalent to the high-order bit of the 4-bit sum being 0, which is opposite to the sign of  $x$  and  $y$ . In case (b), this is equivalent to the high-order bit of the 4-bit sum being 1, which again is opposite to the sign of  $x$  and  $y$ .

For subtraction of multiword integers, the computation of interest is  $x - y - c$ , where again  $c$  is 0 or 1, with a value of 1 representing a borrow-in. From an analysis similar to the above, it can be seen that overflow in the final value of  $x - y - c$  occurs if and only if  $x$  and  $y$  have opposite signs and the sign of  $x - y - c$  is opposite to that of  $x$  (or, equivalently, the same as that of  $y$ ).

This leads to the following expressions for the overflow predicate, with the result being in the sign position. Following these with a *shift right* or *shift right signed* of 31 produces a 1/0- or a -1/0-valued result.

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

By choosing the second alternative in the first column, and the first alternative in the second column (avoiding the *equivalence* operation), our basic RISC can evaluate these tests with three instructions in addition to those required to compute  $x + y + c$  or  $x - y - c$ . A fourth instruction (*branch if negative*) can be added to branch to code where the overflow condition is handled.

If executing with overflow interrupts enabled, the programmer may wish to test to see if a certain addition or subtraction will cause overflow, in a way that does not cause it. One branch-free way to do this is as follows:

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

The assignment to  $z$  in the left column sets  $z = \mathbf{0x80000000}$  if  $x$  and  $y$  have the same sign, and sets  $z = \mathbf{0}$  if they differ. Then, the addition in the second expression is done with  $x \oplus z$  and  $y$  having different signs, so it can't overflow. If  $x$  and  $y$  are nonnegative, the sign bit in the second expression will be 1 if and only if  $(x - 2^{31}) + y + c \geq \mathbf{0}$ —that is, iff  $x + y + c \geq 2^{31}$ , which is the condition for overflow in evaluating  $x + y + c$ . If  $x$  and  $y$  are negative, the sign bit in the second expression will be 1 iff  $(x + 2^{31}) + y + c < \mathbf{0}$ —that is, iff  $x + y + c < -2^{31}$ , which

again is the condition for overflow. The *and* with  $z$  ensures the correct result (0 in the sign position) if  $x$  and  $y$  have opposite signs. Similar remarks apply to the case of subtraction (right column). The code executes in nine instructions on the basic RISC.

It might seem that if the carry from addition is readily available, this might help in computing the signed overflow predicate. This does not seem to be the case; however, one method along these lines is as follows.

If  $x$  is a signed integer, then  $x + 2^{31}$  is correctly represented as an unsigned number and is obtained by inverting the high-order bit of  $x$ . Signed overflow in the positive direction occurs if  $x + y \geq 2^{31}$ —that is, if  $(x + 2^{31}) + (y + 2^{31}) \geq 3 \cdot 2^{31}$ . This latter condition is characterized by carry occurring in the unsigned add (which means that the sum is greater than or equal to  $2^{32}$ ) and the high-order bit of the sum being 1. Similarly, overflow in the negative direction occurs if the carry is 0 and the high-order bit of the sum is also 0.

This gives the following algorithm for detecting overflow for signed addition:

Compute  $(x \oplus 2^{31}) + (y \oplus 2^{31})$ , giving sum  $s$  and carry  $c$ .  
 Overflow occurred iff  $c$  equals the high-order bit of  $s$ .

The sum is the correct sum for the signed addition, because inverting the high-order bits of both operands does not change their sum.

For subtraction, the algorithm is the same except that in the first step a subtraction replaces the addition. We assume that the carry is that which is generated by computing  $x - y$  as  $x + \bar{y} + 1$ . The subtraction is the correct difference for the signed subtraction.

These formulas are perhaps interesting, but on most machines they would not be quite as efficient as the formulas that do not even use the carry bit (e.g., overflow =  $(x \equiv y) \& (s \oplus x)$  for addition, and  $(x \oplus y) \& (d \oplus x)$  for subtraction, where  $s$  and  $d$  are the sum and difference, respectively, of  $x$  and  $y$ ).

### How the Computer Sets Overflow for Signed Add/Subtract

Machines often set “overflow” for signed addition by means of the logic “the carry into the sign position is not equal to the carry out of the sign position.” Curiously, this logic gives the correct overflow indication for both addition and subtraction, assuming the subtraction  $x - y$  is done by  $x + \bar{y} + 1$ . Furthermore, it is correct whether or not there is a carry- or borrow-in. This does not seem to lead to any particularly good methods for computing the signed overflow predicate in software, however, even though it is easy to compute the carry into the sign position. For addition and subtraction, the carry/borrow into the sign position is given by the sign bit after evaluating the following expressions (where  $c$  is 0 or 1):

carry	borrow
$(x + y + c) \oplus x \oplus y$	$(x - y - c) \oplus x \oplus y$

In fact, these expressions give, at each position  $i$ , the carry/borrow into position  $i$ .

### Unsigned Add/Subtract

The following branch-free code can be used to compute the overflow predicate for unsigned add/subtract, with the result being in the sign position. The expressions involving a right shift are probably useful only when it is known that  $c = 0$ . The expressions in brackets compute the carry or borrow generated from the least significant position.

$$\begin{aligned} & x + y + c, \text{ unsigned} \\ & (x \& y) \mid ((x \mid y) \& \neg(x + y + c)) \\ & (x \gg 1) + (y \gg 1) + [((x \& y) \mid ((x \mid y) \& c)) \& 1] \end{aligned}$$

$$\begin{aligned} & x - y - c, \text{ unsigned} \\ & (\neg x \& y) \mid ((x \equiv y) \& (x - y - c)) \\ & (\neg x \& y) \mid ((\neg x \mid y) \& (x - y - c)) \\ & (x \gg 1) - (y \gg 1) - [((\neg x \& y) \mid ((\neg x \mid y) \& c)) \& 1] \end{aligned}$$

For unsigned *add*'s and *subtract*'s, there are much simpler formulas in terms of comparisons [MIPS]. For unsigned addition, overflow (carry) occurs if the sum is less (by unsigned comparison) than either of the operands. This and similar formulas are given below. Unfortunately, there is no way in these formulas to allow for a variable  $c$  that represents the carry- or borrow-in. Instead, the program must test  $c$ , and use a different type of comparison depending upon whether  $c$  is 0 or 1.

$x + y, \text{ unsigned}$	$x + y + 1, \text{ unsigned}$	$x - y, \text{ unsigned}$	$x - y - 1, \text{ unsigned}$
$\neg x \lessdot y$	$\neg x \lessseq y$	$x \lessdot y$	$x \lessseq y$
$x + y \lessdot x$	$x + y + 1 \lessseq x$	$x - y \gtrdot x$	$x - y - 1 \gtrseq x$

The first formula for each case above is evaluated before the add/subtract that may overflow, and it provides a way to do the test without causing overflow. The second formula for each case is evaluated after the add/subtract that may overflow.

There does not seem to be a similar simple device (using comparisons) for computing the signed overflow predicate.

### Multiplication

For multiplication, overflow means that the result cannot be expressed in 32 bits (it can always be expressed in 64 bits, whether signed or unsigned). Checking for overflow is simple if you have access to the high-order 32 bits of the product. Let us denote the two halves of the 64-bit product by  $hi(x \times y)$  and  $lo(x \times y)$ . Then the overflow predicates can be computed as follows [MIPS]:

$$\begin{array}{ll} \mathbf{x} \times \mathbf{y}, \text{ unsigned} & \mathbf{x} \times \mathbf{y}, \text{ signed} \\ \text{hi}(\mathbf{x} \times \mathbf{y}) \neq \mathbf{0} & \text{hi}(\mathbf{x} \times \mathbf{y}) \neq (\text{lo}(\mathbf{x} \times \mathbf{y}) \stackrel{s}{\gg} 31) \end{array}$$

One way to check for overflow of multiplication is to do the multiplication and then check the result by dividing. Care must be taken not to divide by 0, and there is a further complication for signed multiplication. Overflow occurs if the following expressions are **true**:

$$\begin{array}{ll} \text{Unsigned} & \text{Signed} \\ \mathbf{z} \leftarrow \mathbf{x} * \mathbf{y} & \mathbf{z} \leftarrow \mathbf{x} * \mathbf{y} \\ \mathbf{y} \neq \mathbf{0} \ \&\ \mathbf{z} \stackrel{u}{\div} \mathbf{y} \neq \mathbf{x} & (\mathbf{y} < \mathbf{0} \ \&\ \mathbf{x} = -2^{31}) \mid (\mathbf{y} \neq \mathbf{0} \ \&\ \mathbf{z} \div \mathbf{y} \neq \mathbf{x}) \end{array}$$

The complication arises when  $\mathbf{x} = -2^{31}$  and  $\mathbf{y} = -1$ . In this case the multiplication overflows, but the machine may very well give a result of  $-2^{31}$ . This causes the division to overflow, and thus any result is possible (for some machines). Therefore, this case has to be checked separately, which is done by the term  $\mathbf{y} < \mathbf{0} \ \&\ \mathbf{x} = -2^{31}$ . The above expressions use the “conditional *and*” operator to prevent dividing by 0 (in C, use the `&&` operator).

It is also possible to use division to check for overflow of multiplication without doing the multiplication (that is, without causing overflow). For unsigned integers, the product overflows iff  $xy > 2^{32} - 1$ , or  $x > ((2^{32} - 1)/y)$ , or, since  $x$  is an integer,  $x > \lfloor (2^{32} - 1)/y \rfloor$ . Expressed in computer arithmetic, this is

$$\mathbf{y} \neq \mathbf{0} \ \&\ \mathbf{x} > (\mathbf{0x}\mathbf{FFFFFFFF} \stackrel{u}{\div} \mathbf{y}).$$

For signed integers, the determination of overflow of  $\mathbf{x} * \mathbf{y}$  is not so simple. If  $\mathbf{x}$  and  $\mathbf{y}$  have the same sign, then overflow occurs iff  $xy > 2^{31} - 1$ . If they have opposite signs, then overflow occurs iff  $xy < -2^{31}$ . These conditions can be tested as indicated in Table 2-2, which employs signed division. This test is awkward to implement, because of the four cases. It is difficult to unify the expressions very much because of problems with overflow and with not being able to represent the number  $+2^{31}$ .

The test can be simplified if unsigned division is available. We can use the absolute values of  $\mathbf{x}$  and  $\mathbf{y}$ , which are correctly represented under unsigned integer interpretation. The complete test can then be computed as shown below. The variable  $\mathbf{c} = 2^{31} - 1$  if  $\mathbf{x}$  and  $\mathbf{y}$  have the same sign, and  $\mathbf{c} = 2^{31}$  otherwise.

TABLE 2-2. OVERFLOW TEST FOR SIGNED MULTIPLICATION

	$\mathbf{y} > \mathbf{0}$	$\mathbf{y} \leq \mathbf{0}$
$\mathbf{x} > \mathbf{0}$	$\mathbf{x} > \mathbf{0x7FFFFFFF} \div \mathbf{y}$	$\mathbf{y} < \mathbf{0x80000000} \div \mathbf{x}$
$\mathbf{x} \leq \mathbf{0}$	$\mathbf{x} < \mathbf{0x80000000} \div \mathbf{y}$	$\mathbf{x} \neq \mathbf{0} \ \&\ \mathbf{y} < \mathbf{0x7FFFFFFF} \div \mathbf{x}$

$$\begin{aligned}
c &\leftarrow ((x \equiv y) \ggg 31) + 2^{31} \\
x &\leftarrow \text{abs}(x) \\
y &\leftarrow \text{abs}(y) \\
y \neq 0 \ \&\ x >^u (c \div y)
\end{aligned}$$

The *number of leading zeros* instruction can be used to give an estimate of whether or not  $x * y$  will overflow, and the estimate can be refined to give an accurate determination. First, consider the multiplication of unsigned numbers. It is easy to show that if  $x$  and  $y$ , as 32-bit quantities, have  $m$  and  $n$  leading 0's, respectively, then the 64-bit product has either  $m + n$  or  $m + n + 1$  leading 0's (or 64, if either  $x = 0$  or  $y = 0$ ). Overflow occurs if the 64-bit product has fewer than 32 leading 0's. Hence,

$\text{nlz}(x) + \text{nlz}(y) \geq 32$ : Multiplication definitely does not overflow.

$\text{nlz}(x) + \text{nlz}(y) \leq 30$ : Multiplication definitely does overflow.

For  $\text{nlz}(x) + \text{nlz}(y) = 31$ , overflow may or may not occur. In this case, the overflow assessment can be made by evaluating  $t = x \lfloor y/2 \rfloor$ . This will not overflow. Since  $xy$  is  $2t$  or, if  $y$  is odd,  $2t + x$ , the product  $xy$  overflows if  $t \geq 2^{31}$ . These considerations lead to a plan for computing  $xy$ , but branching to “overflow” if the product overflows. This plan is shown in Figure 2-2.

For the multiplication of signed integers, we can make a partial determination of whether or not overflow occurs from the number of leading 0's of nonnegative arguments, and the number of leading 1's of negative arguments. Let

$$\begin{aligned}
m &= \text{nlz}(x) + \text{nlz}(\bar{x}), \text{ and} \\
n &= \text{nlz}(y) + \text{nlz}(\bar{y}).
\end{aligned}$$

```

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;
}
// z is the correct product of x and y.

```

FIGURE 2-2. Determination of overflow of unsigned multiplication.



Then, we have

$m + n \geq 34$ : Multiplication definitely does not overflow.

$m + n \leq 31$ : Multiplication definitely does overflow.

There are two ambiguous cases: 32 and 33. The case  $m + n = 33$  overflows only when both arguments are negative and the true product is exactly  $2^{31}$  (machine result is  $-2^{31}$ ), so it can be recognized by a test that the product has the correct sign (that is, overflow occurred if  $m \oplus n \oplus (m * n) < 0$ ). When  $m + n = 32$ , the distinction is not so easily made.

We will not dwell on this further, except to note that an overflow estimate for signed multiplication can also be made based on  $\text{nlz}(\text{abs}(x)) + \text{nlz}(\text{abs}(y))$ , but again there are two ambiguous cases (a sum of 31 or 32).

### Division

For the signed division  $x \div y$ , overflow occurs if the following expression is **true**:

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

Most machines signal overflow (or trap) for the indeterminate form  $0 \div 0$ .

Straightforward code for evaluating this expression, including a final branch to the overflow handling code, consists of seven instructions, three of which are branches. There do not seem to be any particularly good tricks to improve on this, but here are a few possibilities:

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

That is, evaluate the large expression in brackets, and branch if the result is less than 0. This executes in about nine instructions, counting the load of the constant and the final branch, on a machine that has the indicated instructions and that gets the “compare to 0” for free.

Some other possibilities are to first compute  $z$  from

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

(three instructions on many machines), and then do the test and branch on  $y = 0 \mid z = 0$  in one of the following ways:

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

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

$$((\text{nlz}(y) \mid \text{nlz}(z)) \ggg 5) \neq 0$$

These execute in nine, seven, and eight instructions, respectively, on a machine that has the indicated instructions. The last line represents a good method for PowerPC.

For the unsigned division  $x \stackrel{u}{\div} y$ , overflow occurs if and only if  $y = 0$ .

Some machines have a “long division” instruction (see page 192), and you may want to predict, using elementary instructions, when it would overflow. We will discuss this in terms of an instruction that divides a doubleword by a fullword, producing a fullword quotient and possibly also a fullword remainder.

Such an instruction overflows if either the divisor is 0 or if the quotient cannot be represented in 32 bits. Typically, in these overflow cases both the quotient and remainder are incorrect. The remainder cannot overflow in the sense of being too large to represent in 32 bits (it is less than the divisor in magnitude), so the test that the remainder will be correct is the same as the test that the quotient will be correct.

We assume the machine either has 64-bit general registers or 32-bit registers and there is no problem doing elementary operations (shifts, adds, and so forth) on 64-bit quantities. For example, the compiler might implement a doubleword integer data type.

In the unsigned case the test is trivial: for  $x \div y$  with  $x$  a doubleword and  $y$  a fullword, the division will not overflow if (and only if) either of the following equivalent expressions is true.

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

$$y \neq 0 \ \& \ (x \stackrel{u}{\gg} 32) < y$$

On a 32-bit machine, the shifts need not be done; simply compare  $y$  to the register that contains the high-order half of  $x$ . To ensure correct results on a 64-bit machine, it is also necessary to check that the divisor  $y$  is a 32-bit quantity (e.g., check that  $(y \stackrel{u}{\gg} 32) = 0$ ).

The signed case is more interesting. It is first necessary to check that  $y \neq 0$  and, on a 64-bit machine, that  $y$  is correctly represented in 32 bits (check that  $((y \ll 32) \stackrel{s}{\gg} 32) = y$ ). Assuming these tests have been done, the table that follows shows how the tests might be done to determine precisely whether or not the quotient is representable in 32 bits by considering separately the four cases of the dividend and divisor each being positive or negative. The expressions in the table are in ordinary arithmetic, not computer arithmetic.

In each column, each relation follows from the one above it in an if-and-only-if way. To remove the floor and ceiling functions, some relations from Theorem D1 on page 183 are used.

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

As an example of interpreting this table, consider the leftmost column. It applies to the case in which  $x \geq 0$  and  $y > 0$ . In this case the quotient is  $\lfloor x/y \rfloor$ , and this must be strictly less than  $2^{31}$  to be representable as a 32-bit quantity. From this it follows that the real number  $x/y$  must be less than  $2^{31}$ , or  $x$  must be less than  $2^{31}y$ . This test can be implemented by shifting  $y$  left 31 positions and comparing the result to  $x$ .

When the signs of  $x$  and  $y$  differ, the quotient of conventional division is  $\lceil x/y \rceil$ . Because the quotient is negative, it can be as small as  $-2^{31}$ .

In the bottom row of each column the comparisons are all of the same type (less than). Because of the possibility that  $x$  is the maximum negative number, in the third and fourth columns an unsigned comparison must be used. In the first two columns the quantities being compared begin with a leading 0-bit, so an unsigned comparison can be used there, too.

These tests can, of course, be implemented by using conditional branches to separate out the four cases, doing the indicated arithmetic, and then doing a final compare and branch to the code for the overflow or non-overflow case. However, branching can be reduced by taking advantage of the fact that when  $y$  is negative,  $-y$  is used, and similarly for  $x$ . Hence the tests can be made more uniform by using the absolute values of  $x$  and  $y$ . Also, using a standard device for optionally doing the additions in the second and third columns results in the following scheme:

$$\begin{aligned} x' &= |x| \\ y' &= |y| \\ \delta &= ((x \oplus y) \gg 63) \& y' \\ \text{if } (x' \lessgtr (y' \ll 31) + \delta) &\text{ then \{will not overflow\}} \end{aligned}$$

Using the three-instruction method of computing the absolute value (see page 18), on a 64-bit version of the basic RISC this amounts to 12 instructions, plus a conditional branch.

## 2-14 Condition Code Result of *Add*, *Subtract*, and *Multiply*

Many machines provide a “condition code” that characterizes the result of integer arithmetic operations. Often there is only one *add* instruction, and the characterization reflects the result for both unsigned and signed interpretation of the operands and result (but not for mixed types). The characterization usually consists of the following:

- Whether or not carry occurred (unsigned overflow)
- Whether or not signed overflow occurred
- Whether the 32-bit result, interpreted as a signed two’s-complement integer and ignoring carry and overflow, is negative, 0, or positive

Some older machines give an indication of whether the infinite precision result (that is, 33-bit result for *add*'s and *subtract*'s) is positive, negative, or 0. However, this indication is not easily used by compilers of high-level languages, and so has fallen out of favor.

For addition, only nine of the 12 combinations of these events are possible. The ones that cannot occur are “no carry, overflow, result > 0,” “no carry, overflow, result = 0,” and “carry, overflow, result < 0.” Thus, four bits are, just barely, needed for the condition code. Two of the combinations are unique in the sense that only one value of inputs produces them: Adding 0 to itself is the only way to get “no carry, no overflow, result = 0,” and adding the maximum negative number to itself is the only way to get “carry, overflow, result = 0.” These remarks remain true if there is a “carry in”—that is, if we are computing  $x + y + 1$ .

For subtraction, let us assume that to compute  $x - y$  the machine actually computes  $x + \bar{y} + 1$ , with the carry produced as for an *add* (in this scheme the meaning of “carry” is reversed for subtraction, in that carry = 1 signifies that the result fits in a single word, and carry = 0 signifies that the result does not fit in a single word). Then for subtraction, only seven combinations of events are possible. The ones that cannot occur are the three that cannot occur for addition, plus “no carry, no overflow, result = 0,” and “carry, overflow, result = 0.”

If a machine's multiplier can produce a doubleword result, then two *multiply* instructions are desirable: one for signed and one for unsigned operands. (On a 4-bit machine, in hexadecimal,  $F \times F = 01$  signed, and  $F \times F = E1$  unsigned.) For these instructions, neither carry nor overflow can occur, in the sense that the result will always fit in a doubleword.

For a multiplication instruction that produces a one-word result (the low-order word of the doubleword result), let us take “carry” to mean that the result does not fit in a word with the operands and result interpreted as unsigned integers, and let us take “overflow” to mean that the result does not fit in a word with the operands and result interpreted as signed two's-complement integers. Then again, there are nine possible combinations of results, with the missing ones being “no carry, overflow, result > 0,” “no carry, overflow, result = 0,” and “carry, no overflow, result = 0.” Thus, considering addition, subtraction, and multiplication together, ten combinations can occur.

## 2-15 Rotate Shifts

These are rather trivial. Perhaps surprisingly, this code works for  $n$  ranging from 0 to 32 inclusive, even if the shifts are mod-32.

Rotate left  $n$ :  $y \leftarrow (x \ll n) \mid (x \ggg (32 - n))$

Rotate right  $n$ :  $y \leftarrow (x \ggg n) \mid (x \ll (32 - n))$

If your machine has double-length shifts, they can be used to do rotate shifts. These instructions might be written

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

They treat the concatenation of RA and RB as a single double-length quantity, and shift it left or right by the amount given by the immediate field I. (If the shift amount is in a register, the instructions are awkward to implement on most RISCs because they require reading three registers.) The result of the left shift is the high-order word of the shifted double-length quantity, and the result of the right shift is the low-order word.

Using `shldi`, a rotate left of Rx can be accomplished by

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

and similarly a rotate right shift can be accomplished with `shrdi`.

A rotate left shift of one position can be accomplished by adding the contents of a register to itself with “end-around carry” (adding the carry that results from the addition to the sum in the low-order position). Most machines do not have that instruction, but on many machines it can be accomplished with two instructions: (1) add the contents of the register to itself, generating a carry (into a status register), and (2) add the carry to the sum.

## 2-16 Double-Length Add/Subtract

Using one of the expressions shown on page 31 for overflow of unsigned addition and subtraction, we can easily implement double-length addition and subtraction without accessing the machine’s carry bit. To illustrate with double-length addition, let the operands be  $(x_1, x_0)$  and  $(y_1, y_0)$ , and the result be  $(z_1, z_0)$ . Subscript 1 denotes the most significant half, and subscript 0 the least significant. We assume that all 32 bits of the registers are used. The less significant words are unsigned quantities.

$$\begin{aligned} z_0 &\leftarrow x_0 + y_0 \\ c &\leftarrow [(x_0 \& y_0) \mid ((x_0 \mid y_0) \& \neg z_0)] \ggg 31 \\ z_1 &\leftarrow x_1 + y_1 + c \end{aligned}$$

This executes in nine instructions. The second line can be  $c \leftarrow (z_0 \overset{u}{\prec} x_0)$ , permitting a four-instruction solution on machines that have this comparison operator in a form that gives the result as a **1** or **0** in a register, such as the “SLTU” (*Set on Less Than Unsigned*) instruction on MIPS [MIPS].

Similar code for double-length subtraction  $(x - y)$  is

$$\begin{aligned} z_0 &\leftarrow x_0 - y_0 \\ b &\leftarrow [(\neg x_0 \& y_0) \mid ((x_0 \equiv y_0) \& z_0)] \ggg 31 \\ z_1 &\leftarrow x_1 - y_1 - b \end{aligned}$$

This executes in eight instructions on a machine that has a full set of logical instructions. The second line can be  $b \leftarrow (x_0 \overset{u}{<} y_0)$ , permitting a four-instruction solution on machines that have the “SLTU” instruction.

Double-length addition and subtraction can be done in five instructions on most machines by representing the multiple-length data using only 31 bits of the least significant words, with the high-order bit being 0 except momentarily when it contains a carry or borrow bit.

## 2-17 Double-Length Shifts

Let  $(x_1, x_0)$  be a pair of 32-bit words to be shifted left or right as if they were a single 64-bit quantity, with  $x_1$  being the most significant half. Let  $(y_1, y_0)$  be the result, interpreted similarly. Assume the shift amount  $n$  is a variable ranging from 0 to 63. Assume further that the machine’s shift instructions are modulo 64 or greater. That is, a shift amount in the range 32 to 63 or  $-32$  to  $-1$  results in an all-0 word, unless the shift is a signed right shift, in which case the result is 32 sign bits from the word shifted. (This code will not work on the Intel x86 machines, which have mod-32 shifts.)

Under these assumptions, the *shift left double* operation can be accomplished as follows (eight instructions):

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

The main connective in the first assignment must be *or*, not *plus*, to give the correct result when  $n = 32$ . If it is known that  $0 \leq n \leq 32$ , the last term of the first assignment can be omitted, giving a six-instruction solution.

Similarly, a *shift right double unsigned* operation can be done with

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

*Shift right double signed* is more difficult, because of an unwanted sign propagation in one of the terms. Straightforward code follows:

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

If your machine has the *conditional move* instructions, it is a simple matter to express this in branch-free code, in which form it takes eight instructions. If the conditional move instructions are not available, the operation can be done in ten

instructions by using the familiar device of constructing a mask with the *shift right signed 31* instruction to mask the unwanted sign propagating term:

$$y_0 \leftarrow x_0 \overset{u}{\gg} n \mid x_1 \ll (32 - n) \mid [(x_1 \overset{s}{\gg} (n - 32)) \& ((32 - n) \overset{s}{\gg} 31)]$$

$$y_1 \leftarrow x_1 \overset{s}{\gg} n$$

## 2-18 Multibyte *Add, Subtract, Absolute Value*

Some applications deal with arrays of short integers (usually bytes or halfwords), and often execution is faster if they are operated on a word at a time. For definiteness, the examples here deal with the case of four 1-byte integers packed into a word, but the techniques are easily adapted to other packings, such as a word containing a 12-bit integer and two 10-bit integers, and so on. These techniques are of greater value on 64-bit machines, because more work is done in parallel.

Addition must be done in a way that blocks the carries from one byte into another. This can be accomplished by the following two-step method:

1. Mask out the high-order bit of each byte of each operand and *add* (there will then be no carries across byte boundaries).
2. Fix up the high-order bit of each byte with a 1-bit *add* of the two operands and the carry into that bit.

The carry into the high-order bit of each byte is given by the high-order bit of each byte of the sum computed in step 1. The subsequent similar method works for subtraction:

Addition

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

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

Subtraction

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

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

These execute in eight instructions, counting the load of **0x7F7F7F7F**, on a machine that has a full set of logical instructions. (Change the *and* and *or* of **0x80808080** to *and not* and *or not*, respectively, of **0x7F7F7F7F**.)

There is a different technique for the case in which the word is divided into only two fields. In this case, addition can be done by means of a 32-bit addition followed by subtracting out the unwanted carry. On page 30 we noted that the expression  $(x + y) \oplus x \oplus y$  gives the carries into each position. Using this and similar observations about subtraction gives the following code for adding/subtracting two halfwords modulo  $2^{16}$  (seven instructions):

Addition

$$s \leftarrow x + y$$

$$c \leftarrow (s \oplus x \oplus y) \& 0x00010000$$

$$s \leftarrow s - c$$

Subtraction

$$d \leftarrow x - y$$

$$b \leftarrow (d \oplus x \oplus y) \& 0x00010000$$

$$d \leftarrow d + b$$

Multibyte *absolute value* is easily done by complementing and adding 1 to each byte that contains a negative integer (that is, has its high-order bit on). The following code sets each byte of  $y$  equal to the absolute value of each byte of  $x$  (eight instructions):

$a \leftarrow x \& 0x80808080$	// Isolate signs.
$b \leftarrow a \ggg 7$	// Integer 1 where $x$ is negative.
$m \leftarrow (a - b) \mid a$	// 0xFF where $x$ is negative.
$y \leftarrow (x \oplus m) + b$	// Complement and add 1 where negative.

The third line could as well be  $m \leftarrow a + a - b$ . The addition of  $b$  in the fourth line cannot carry across byte boundaries, because the quantity  $x \oplus m$  has a high-order 0 in each byte.

## 2-19 Doz, Max, Min

The “doz” function is “difference or zero,” defined as follows:

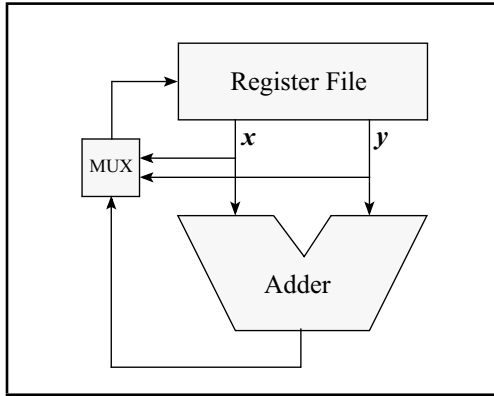
Signed	Unsigned
$\text{doz}(x, y) = \begin{cases} x - y, & x \geq y, \\ 0, & x < y. \end{cases}$	$\text{dozu}(x, y) = \begin{cases} x - y, & x \geq^u y, \\ 0, & x <^u y. \end{cases}$

It has been called “first grade subtraction” because the result is 0 if you try to take away too much.<sup>3</sup> If implemented as a computer instruction, perhaps its most important use is to implement the  $\max(x, y)$  and  $\min(x, y)$  functions (in both signed and unsigned forms) in just two simple instructions, as will be seen. Implementing  $\max(x, y)$  and  $\min(x, y)$  in hardware is difficult because the machine would need paths from the output ports of the register file back to an input port, bypassing the adder. These paths are not normally present. If supplied, they would be in a region that’s often crowded with wiring for register bypasses. The situation is illustrated in Figure 2-3. The adder is used (by the instruction) to do the subtraction  $x - y$ . The high-order bits of the result of the subtraction (sign bit and carries, as described on page 27) define whether  $x \geq y$  or  $x < y$ . The comparison result is fed to a multiplexor

---

3. Mathematicians name the operation *monus* and denote it with  $\dot{-}$ . The terms *positive difference* and *saturated subtraction* are also used.



FIGURE 2-3. Implementing  $\max(x, y)$  and  $\min(x, y)$ .

(MUX) that selects either  $x$  or  $y$  as the result to write into the target register. These paths, from register file outputs  $x$  and  $y$  to the multiplexor, are not normally present and would have little use. The *difference or zero* instructions can be implemented without these paths because it is the output of the adder (or 0) that is fed back to the register file.

Using *difference or zero*,  $\max(x, y)$  and  $\min(x, y)$  can be implemented in two instructions as follows:

Signed	Unsigned
$\max(x, y) = y + \text{doz}(x, y)$	$\maxu(x, y) = y + \text{dozu}(x, y)$
$\min(x, y) = x - \text{doz}(x, y)$	$\minu(x, y) = x - \text{dozu}(x, y)$

In the signed case, the result of the *difference or zero* instruction can be negative. This happens if overflow occurs in the subtraction. Overflow should be ignored; the addition of  $y$  or subtraction from  $x$  will overflow again, and the result will be correct. When  $\text{doz}(x, y)$  is negative, it is actually the correct difference if it is interpreted as an unsigned integer.

Suppose your computer does not have the *difference or zero* instructions, but you want to code  $\text{doz}(x, y)$ ,  $\max(x, y)$ , and so forth, in an efficient branch-free way. In the next few paragraphs we show how these functions might be coded if your machine has the *conditional move* instructions, comparison predicates, efficient access to the carry bit, or none of these.

If your machine has the *conditional move* instructions, it can get  $\text{doz}(x, y)$  in three instructions, and destructive<sup>4</sup>  $\max(x, y)$  and  $\min(x, y)$  in two instructions. For example, on the full RISC,  $z \leftarrow \text{doz}(x, y)$  can be calculated as follows ( $r0$  is a permanent zero register):

---

4. A destructive operation is one that overwrites one or more of its arguments.

<code>sub</code>	<code>z,x,y</code>	Set $z = x - y$ .
<code>cmplt</code>	<code>t,x,y</code>	Set $t = 1$ if $x < y$ , else 0.
<code>movne</code>	<code>z,t,r0</code>	Set $z = 0$ if $x < y$ .

Also on the full RISC,  $x \leftarrow \max(x, y)$  can be calculated as follows:

<code>cmplt</code>	<code>t,x,y</code>	Set $t = 1$ if $x < y$ , else 0.
<code>movne</code>	<code>x,t,y</code>	Set $x = y$ if $x < y$ .

The min function, and the unsigned counterparts, are obtained by changing the comparison conditions.

These functions can be computed in four or five instructions using comparison predicates (three or four if the comparison predicates give a result of -1 for “true”):

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

On some machines, the carry bit may be a useful aid to computing the unsigned versions of these functions. Let  $\text{carry}(x - y)$  denote the bit that comes out of the adder for the operation  $x + \bar{y} + 1$ , moved to a GPR. Thus,  $\text{carry}(x - y) = 1$  iff  $x \geq y$ . Then we have

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

On most machines that have a *subtract* that generates a carry or borrow, and another form of *subtract* that uses that carry or borrow as an input, the expression  $\text{carry}(x - y) - 1$  can be computed in one more instruction after the subtraction of  $y$  from  $x$ . For example, on the Intel x86 machines,  $\text{minu}(x, y)$  can be computed in four instructions as follows:

<code>sub</code>	<code>eax,ecx</code>	; Inputs $x$ and $y$ are in <code>eax</code> and <code>ecx</code> resp.
<code>sbb</code>	<code>edx,edx</code>	; <code>edx</code> = 0 if $x \geq y$ , else -1.
<code>and</code>	<code>eax,edx</code>	; 0 if $x \geq y$ , else $x - y$ .
<code>add</code>	<code>eax,ecx</code>	; Add $y$ , giving $y$ if $x \geq y$ , else $x$ .

In this way, all three of the functions can be computed in four instructions (three instructions for  $\text{dozu}(x, y)$  if the machine has *and with complement*).

A method that applies to nearly any RISC is to use one of the above expressions that employ a comparison predicate, and to substitute for the predicate one of the expressions given on page 23. For example:

$$d \leftarrow x - y$$

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

$$\text{dozu}(x, y) = d \ \& \ \neg[(\neg x \ \& \ y) \mid ((x \equiv y) \ \& \ d)] \ \overset{s}{\gg} 31]$$

These require from seven to ten instructions, depending on the computer's instruction set, plus one more to get max or min.

These operations can be done in four branch-free basic RISC instructions if it is known that  $-2^{31} \leq x - y \leq 2^{31} - 1$  (that is an expression in ordinary arithmetic, not computer arithmetic). The same code works for both signed and unsigned integers, with the same restriction on  $x$  and  $y$ . A sufficient condition for these formulas to be valid is that, for signed integers,  $-2^{30} \leq x, y \leq 2^{30} - 1$ , and for unsigned integers,  $0 \leq x, y \leq 2^{31} - 1$ .

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

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

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

Some uses of the *difference or zero* instruction are given here. In these, the result of  $\text{doz}(x, y)$  must be interpreted as an unsigned integer.

1. It directly implements the Fortran IDIM function.
2. To compute the absolute value of a difference [Knu7]:

$$\begin{aligned} |x - y| &= \text{doz}(x, y) + \text{doz}(y, x), & \text{signed arguments,} \\ &= \text{dozu}(x, y) + \text{dozu}(y, x), & \text{unsigned arguments.} \end{aligned}$$

Corollary:  $|x| = \text{doz}(x, 0) + \text{doz}(0, x)$  (other three-instruction solutions are given on page 18).

3. To clamp the upper limit of the true sum of unsigned integers  $x$  and  $y$  to the maximum positive number  $(2^{32} - 1)$  [Knu7]:

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

4. Some comparison predicates (four instructions each):

$$x > y = (\text{doz}(x, y) \mid \neg \text{doz}(x, y)) \ \overset{u}{\gg} 31,$$

$$x \overset{u}{>} y = (\text{dozu}(x, y) \mid \neg \text{dozu}(x, y)) \ \overset{u}{\gg} 31.$$

5. The carry bit from the addition  $x + y$  (five instructions):

$$\text{carry}(x + y) = x \overset{u}{\succ} \neg y = (\text{dozu}(x, \neg y) \mid -\text{dozu}(x, \neg y)) \overset{u}{\gg} 31.$$

The expression  $\text{doz}(x, \neg y)$ , with the result interpreted as an unsigned integer, is in most cases the true sum  $x + y$  with the lower limit clamped at 0. However, it fails if  $y$  is the maximum negative number.

The IBM RS/6000 computer, and its predecessor the 801, have the signed version of *difference or zero*. Knuth's MMIX computer [Knu7] has the unsigned version (including some varieties that operate on parts of words in parallel). This raises the question of how to get the signed version from the unsigned version, and vice versa. This can be done as follows (where the additions and subtractions simply complement the sign bit):

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

Some other identities that may be useful are:

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

The relation  $\text{doz}(\neg x, \neg y) = \text{doz}(y, x)$  fails if either  $x$  or  $y$ , but not both, is the maximum negative number.

## 2-20 Exchanging Registers

A very old trick is exchanging the contents of two registers without using a third [IBM]:

$$\begin{aligned}x &\leftarrow x \oplus y \\ y &\leftarrow y \oplus x \\ x &\leftarrow x \oplus y\end{aligned}$$

This works well on a two-address machine. The trick also works if  $\oplus$  is replaced by the  $\equiv$  logical operation (complement of *exclusive or*) and can be made to work in various ways with *add*'s and *subtract*'s:

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

Unfortunately, each of these has an instruction that is unsuitable for a two-address machine, unless the machine has “reverse subtract.”

This little trick can actually be useful in the application of double buffering, in which two pointers are swapped. The first instruction can be factored out of the loop in which the swap is done (although this negates the advantage of saving a register):

Outside the loop:  $t \leftarrow x \oplus y$   
 Inside the loop:  $x \leftarrow x \oplus t$   
 $y \leftarrow y \oplus t$

### Exchanging Corresponding Fields of Registers

The problem here is to exchange the contents of two registers  $x$  and  $y$  wherever a mask bit  $m_i = 1$ , and to leave  $x$  and  $y$  unaltered wherever  $m_i = 0$ . By “corresponding” fields, we mean that no shifting is required. The 1-bits of  $m$  need not be contiguous. The straightforward method is as follows:

$$\begin{aligned} x' &\leftarrow (x \& \bar{m}) \mid (y \& m) \\ y &\leftarrow (y \& \bar{m}) \mid (x \& m) \\ x &\leftarrow x' \end{aligned}$$

By using “temporaries” for the four *and* expressions, this can be seen to require seven instructions, assuming that either  $m$  or  $\bar{m}$  can be loaded with a single instruction and the machine has *and* and *not* as a single instruction. If the machine is capable of executing the four (independent) *and* expressions in parallel, the execution time is only three cycles.

A method that is probably better (five instructions, but four cycles on a machine with unlimited instruction-level parallelism) is shown in column (a) below. It is suggested by the “three *exclusive or*” code for exchanging registers.

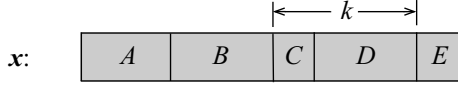
(a)	(b)	(c)
$x \leftarrow x \oplus y$	$x \leftarrow x \equiv y$	$t \leftarrow (x \oplus y) \& m$
$y \leftarrow y \oplus (x \& m)$	$y \leftarrow y \equiv (x \mid \bar{m})$	$x \leftarrow x \oplus t$
$x \leftarrow x \oplus y$	$x \leftarrow x \equiv y$	$y \leftarrow y \oplus t$

The steps in column (b) do the same exchange as that of column (a), but column (b) is useful if  $m$  does not fit in an immediate field, but  $\bar{m}$  does, and the machine has the *equivalence* instruction.

Still another method is shown in column (c) above [GLS1]. It also takes five instructions (again assuming one instruction must be used to load  $m$  into a register), but executes in only three cycles on a machine with sufficient instruction-level parallelism.

### Exchanging Two Fields of the Same Register

Assume a register  $x$  has two fields (of the same length) that are to be swapped, without altering other bits in the register. That is, the object is to swap fields  $B$  and  $D$  without altering fields  $A$ ,  $C$ , and  $E$ , in the computer word illustrated below. The fields are separated by a shift distance  $k$ .



Straightforward code would shift  $D$  and  $B$  to their new positions, and combine the words with *and* and *or* operations, as follows:

$$t_1 = (x \& m) \ll k$$

$$t_2 = (x \gg k) \& m$$

$$x' = (x \& m') \mid t_1 \mid t_2$$

Here,  $m$  is a mask with 1's in field  $D$  (and 0's elsewhere), and  $m'$  is a mask with 1's in fields  $A$ ,  $C$ , and  $E$ . This code requires 11 instructions and six cycles on a machine with unlimited instruction-level parallelism, allowing for four instructions to generate the two masks.

A method that requires only eight instructions and executes in five cycles, under the same assumptions, is shown below [GLS1]. It is similar to the code in column (c) on page 46 for interchanging corresponding fields of two registers. Again,  $m$  is a mask that isolates field  $D$ .

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

$$t_2 = t_1 \ll k$$

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

The idea is that  $t_1$  contains  $B \oplus D$  in position  $D$  (and 0's elsewhere), and  $t_2$  contains  $B \oplus D$  in position  $B$ . This code, and the straightforward code given earlier, work correctly if  $B$  and  $D$  are “split fields”—that is, if the 1-bits of mask  $m$  are not contiguous.

### Conditional Exchange

The exchange methods of the preceding two sections, which are based on *exclusive or*, degenerate into no-operations if the mask  $m$  is 0. Hence, they can perform an exchange of entire registers, or of corresponding fields of two registers, or of two fields of the same register, if  $m$  is set to all 1's if some condition  $c$  is **true**, and to all 0's if  $c$  is **false**. This gives branch-free code if  $m$  can be set up without branching.

## 2-21 Alternating among Two or More Values

Suppose a variable  $x$  can have only two possible values  $a$  and  $b$ , and you wish to assign to  $x$  the value other than its current one, and you wish your code to be independent of the values of  $a$  and  $b$ . For example, in a compiler  $x$  might be an opcode that is known to be either *branch true* or *branch false*, and whichever it is, you want to switch it to the other. The values of the opcodes *branch true* and *branch false* are arbitrary, probably defined by a C `#define` or `enum` declaration in a header file.

The straightforward code to do the switch is

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

or, as is often seen in C programs,

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

A far better (or at least more efficient) way to code it is either

$$\begin{aligned}x &\leftarrow a + b - x, \text{ or} \\x &\leftarrow a \oplus b \oplus x.\end{aligned}$$

If  $a$  and  $b$  are constants, these require only one or two basic RISC instructions. Of course, overflow in calculating  $a + b$  can be ignored.

This raises the question: Is there some particularly efficient way to cycle among three or more values? That is, given three arbitrary but distinct constants  $a$ ,  $b$ , and  $c$ , we seek an easy-to-evaluate function  $f$  that satisfies

$$\begin{aligned}f(a) &= b, \\f(b) &= c, \text{ and} \\f(c) &= a.\end{aligned}$$

It is perhaps interesting to note that there is always a polynomial for such a function. For the case of three constants,

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

(The idea is that if  $x = a$ , the first and last terms vanish, and the middle term simplifies to  $b$ , and so on.) This requires 14 arithmetic operations to evaluate, and for arbitrary  $a$ ,  $b$ , and  $c$ , the intermediate results exceed the computer's word size. But it is just a quadratic; if written in the usual form for a polynomial and evaluated using

Horner's rule,<sup>5</sup> it would require only five arithmetic operations (four for a quadratic with integer coefficients, plus one for a final division). Rearranging Equation (5) accordingly gives

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

This is getting too complicated to be interesting, or practical.

Another method, similar to Equation (5) in that just one of the three terms survives, is

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

This takes 11 instructions if the machine has the *equal* predicate, not counting loads of constants. Because the two addition operations are combining two 0 values with a nonzero, they can be replaced with *or* or *exclusive or* operations.

The formula can be simplified by precalculating  $a - c$  and  $b - c$ , and then using [GLS1]:

$$f(x) = ((-(x = c)) \& (a - c)) + ((-(x = a)) \& (b - c)) + c, \text{ or} \\ f(x) = ((-(x = c)) \& (a \oplus c)) \oplus ((-(x = a)) \& (b \oplus c)) \oplus c.$$

Each of these operations takes eight instructions, but on most machines these are probably no better than the straightforward C code shown below, which executes in four to six instructions for small  $a$ ,  $b$ , and  $c$ .

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

Pursuing this matter, there is an ingenious branch-free method of cycling among three values on machines that do not have comparison predicate instructions [GLS1]. It executes in eight instructions on most machines.

Because  $a$ ,  $b$ , and  $c$  are distinct, there are two bit positions,  $n_1$  and  $n_2$ , where the bits of  $a$ ,  $b$ , and  $c$  are not all the same, and where the “odd one out” (the one

---

5. Horner's rule simply factors out  $x$ . For example, it evaluates the fourth-degree polynomial  $ax^4 + bx^3 + cx^2 + dx + e$  as  $x(x(x(ax + b) + c) + d) + e$ . For a polynomial of degree  $n$  it takes  $n$  multiplications and  $n$  additions, and it is very suitable for the *multiply-add* instruction.



whose bit differs in that position from the other two) is different in positions  $n_1$  and  $n_2$ . This is illustrated below for the values 21, 31, and 20, shown in binary.

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

Without loss of generality, rename  $\mathbf{a}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  so that  $\mathbf{a}$  has the odd one out in position  $n_1$  and  $\mathbf{b}$  has the odd one out in position  $n_2$ , as shown above. Then there are two possibilities for the values of the bits at position  $n_1$ , namely  $(\mathbf{a}_{n_1}, \mathbf{b}_{n_1}, \mathbf{c}_{n_1}) = (0, 1, 1)$  or  $(1, 0, 0)$ . Similarly, there are two possibilities for the bits at position  $n_2$ , namely  $(\mathbf{a}_{n_2}, \mathbf{b}_{n_2}, \mathbf{c}_{n_2}) = (0, 1, 0)$  or  $(1, 0, 1)$ . This makes four cases in all, and formulas for each of these cases are shown below.

Case 1.  $(\mathbf{a}_{n_1}, \mathbf{b}_{n_1}, \mathbf{c}_{n_1}) = (0, 1, 1)$ ,  $(\mathbf{a}_{n_2}, \mathbf{b}_{n_2}, \mathbf{c}_{n_2}) = (0, 1, 0)$ :

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

Case 2.  $(\mathbf{a}_{n_1}, \mathbf{b}_{n_1}, \mathbf{c}_{n_1}) = (0, 1, 1)$ ,  $(\mathbf{a}_{n_2}, \mathbf{b}_{n_2}, \mathbf{c}_{n_2}) = (1, 0, 1)$ :

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

Case 3.  $(\mathbf{a}_{n_1}, \mathbf{b}_{n_1}, \mathbf{c}_{n_1}) = (1, 0, 0)$ ,  $(\mathbf{a}_{n_2}, \mathbf{b}_{n_2}, \mathbf{c}_{n_2}) = (0, 1, 0)$ :

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

Case 4.  $(\mathbf{a}_{n_1}, \mathbf{b}_{n_1}, \mathbf{c}_{n_1}) = (1, 0, 0)$ ,  $(\mathbf{a}_{n_2}, \mathbf{b}_{n_2}, \mathbf{c}_{n_2}) = (1, 0, 1)$ :

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

In these formulas, the left operand of each multiplication is a single bit. A multiplication by 0 or 1 can be converted into an *and* with a value of 0 or all 1's. Thus, the formulas can be rewritten as illustrated below for the first formula.

$$f(x) = ((x \ll (31 - n_1)) \gg 31) \& (\mathbf{a} - \mathbf{b}) + ((x \ll (31 - n_2)) \gg 31) \& (\mathbf{c} - \mathbf{a}) + \mathbf{b}$$

Because all variables except  $x$  are constants, this can be evaluated in eight instructions on the basic RISC. Here again, the additions and subtractions can be replaced with *exclusive or*.

This idea can be extended to cycling among four or more constants. The essence of the idea is to find bit positions  $n_1, n_2, \dots$ , at which the bits uniquely identify the constants. For four constants, three bit positions always suffice. Then

(for four constants) solve the following equation for  $s$ ,  $t$ ,  $u$ , and  $v$  (that is, solve the system of four linear equations in which  $f(x)$  is  $a$ ,  $b$ ,  $c$ , or  $d$ , and the coefficients  $x_{n_i}$  are 0 or 1):

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

If the four constants are uniquely identified by only two bit positions, the equation to solve is

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

## 2-22 A Boolean Decomposition Formula

In this section, we have a look at the minimum number of binary Boolean operations, or instructions, that suffice to implement any Boolean function of three, four, or five variables. By a “Boolean function” we mean a Boolean-valued function of Boolean arguments.

Our notation for Boolean algebra uses “+” for *or*, juxtaposition for *and*,  $\oplus$  for *exclusive or*, and either an overbar or a prefix  $\neg$  for *not*. These operators can be applied to single-bit operands or “bitwise” to computer words. Our main result is the following theorem:

**THEOREM.** *If  $f(x, y, z)$  is a Boolean function of three variables, then it can be decomposed into the form  $g(x, y) \oplus zh(x, y)$ , where  $g$  and  $h$  are Boolean functions of two variables.*<sup>6</sup>

*Proof* [Ditlow].  $f(x, y, z)$  can be expressed as a sum of minterms, and then  $\bar{z}$  and  $z$  can be factored out of their terms, giving

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

Because the operands to “+” cannot both be 1, the *or* can be replaced with *exclusive or*, giving

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

where we have twice used the identity  $(a \oplus b)c = ac \oplus bc$ .

---

6. Logic designers will recognize this as Reed-Muller, a.k.a positive Davio, decomposition. According to Knuth [Knu4, 7.1.1], it was known to I. I. Zhegalkin [Matematicheskii Sbornik 35 (1928), 311–369]. It is sometimes referred to as the Russian decomposition.

This is in the required form with  $g(x, y) = f_0(x, y)$  and  $h(x, y) = f_0(x, y) \oplus f_1(x, y)$ .  $f_0(x, y)$ , incidentally, is  $f(x, y, z)$  with  $z = 0$ , and  $f_1(x, y)$  is  $f(x, y, z)$  with  $z = 1$ .

*COROLLARY. If a computer's instruction set includes an instruction for each of the 16 Boolean functions of two variables, then any Boolean function of three variables can be implemented with four (or fewer) instructions.*

One instruction implements  $g(x, y)$ , another implements  $h(x, y)$ , and these are combined with *and* and *exclusive or*.

As an example, consider the Boolean function that is 1 if exactly two of  $x, y$ , and  $z$  are 1:

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

Before proceeding, the interested reader might like to try to implement  $f$  with four instructions, without using the theorem.

From the proof of the theorem,

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

which is four instructions.

Clearly, the theorem can be extended to functions of four or more variables. That is, any Boolean function  $f(x_1, x_2, \dots, x_n)$  can be decomposed into the form  $g(x_1, x_2, \dots, x_{n-1}) \oplus x_n h(x_1, x_2, \dots, x_{n-1})$ . Thus, a function of four variables can be decomposed as follows:

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

This shows that a computer that has an instruction for each of the 16 binary Boolean functions can implement any function of four variables with ten instructions. Similarly, any function of five variables can be implemented with 22 instructions.

However, it is possible to do much better. For functions of four or more variables there is probably no simple plug-in equation like the theorem gives, but exhaustive computer searches have been done. The results are that any Boolean function of four variables can be implemented with seven binary Boolean instructions, and any such function of five variables can be implemented with 12 such instructions [Knu4, 7.1.2].

In the case of five variables, only 1920 of the  $2^{25} = 4,294,967,296$  functions require 12 instructions, and these 1920 functions are all essentially the same function. The variations are obtained by permuting the arguments, replacing some arguments with their complements, or complementing the value of the function.

## 2-23 Implementing Instructions for All 16 Binary Boolean Operations

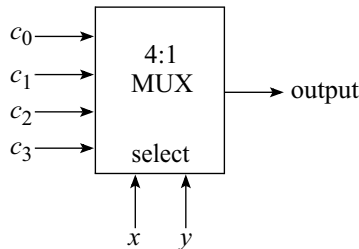
The instruction sets of some computers include all 16 binary Boolean operations. Many of the instructions are useless in that their function can be accomplished with another instruction. For example, the function  $f(x, y) = 0$  simply clears a register, and most computers have a variety of ways to do that. Nevertheless, one reason a computer designer might choose to implement all 16 is that there is a simple and quite regular circuit for doing it.

Refer to Table 2-1 on page 17, which shows all 16 binary Boolean functions. To implement these functions as instructions, choose four of the opcode bits to be the same as the function values shown in the table. Denoting these opcode bits by  $c_0$ ,  $c_1$ ,  $c_2$ , and  $c_3$ , reading from the bottom up in the table, and the input registers by  $x$  and  $y$ , the circuit for implementing all 16 binary Boolean operations is described by the logic expression

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

For example, with  $c_0 = c_1 = c_2 = c_3 = 0$ , the instruction computes the zero function,  $f(x, y) = 0$ . With  $c_0 = 1$  and the other opcode bits 0 it is the *and* instruction. With  $c_0 = c_3 = 0$  and  $c_1 = c_2 = 1$  it is *exclusive or*, and so forth.

This can be implemented with  $n$  4:1 MUXs, where  $n$  is the word size of the machine. The data bits of  $x$  and  $y$  are the select lines, and the four opcode bits are the data inputs to each MUX. The MUX is a standard building block in today's technology, and it is usually a very fast circuit. It is illustrated below.



The function of the circuit is to select  $c_0$ ,  $c_1$ ,  $c_2$ , or  $c_3$  to be the output, depending on whether  $x$  and  $y$  are 00, 01, 10, or 11, respectively. It is like a four-position rotary switch.

Elegant as this is, it is somewhat expensive in opcode points, using 16 of them. There are a number of ways to implement all 16 Boolean operations using only eight opcode points, at the expense of less regular logic. One such scheme is illustrated in Table 2-3.

TABLE 2-3. EIGHT SUFFICIENT BOOLEAN INSTRUCTIONS

Function Values	Formula	Instruction Mnemonic (Name)
0001	$xy$	<i>and</i>
0010	$x\bar{y}$	<i>andc (and with complement)</i>
0110	$x \oplus y$	<i>xor (exclusive or)</i>
0111	$x + y$	<i>or</i>
1110	$\overline{xy}$	<i>nand (negative and)</i>
1101	$\overline{x\bar{y}}$ , or $\bar{x} + y$	<i>cor (complement and or)</i>
1001	$\overline{x \oplus y}$ , or $x \equiv y$	<i>eqv (equivalence)</i>
1000	$\overline{x + y}$	<i>nor (negative or)</i>

The eight operations not shown in the table can be done with the eight instructions shown, by interchanging the inputs or by having both register fields of the instruction refer to the same register. See exercise 13.

IBM's POWER architecture uses this scheme, with the minor difference that POWER has *or with complement* rather than *complement and or*. The scheme shown in Table 2-3 allows the last four instructions to be implemented by complementing the result of the first four instructions, respectively.

### Historical Notes

The algebra of logic expounded in George Boole's *An Investigation of the Laws of Thought* (1854)<sup>7</sup> is somewhat different from what we know today as "Boolean algebra." Boole used the *integers* 1 and 0 to represent truth and falsity, respectively, and he showed how they could be manipulated with the methods of ordinary numerical algebra to formalize natural language statements involving "and," "or," and "except." He also used ordinary algebra to formalize statements in set theory involving intersection, union of disjoint sets, and complementation. He also formalized statements in probability theory, in which the variables take on real number values from 0 to 1. The work often deals with questions of philosophy, religion, and law.

Boole is regarded as a great thinker about logic because he formalized it, allowing complex statements to be manipulated mechanically and flawlessly with the familiar methods of ordinary algebra.

Skipping ahead in history, there are a few programming languages that include all 16 Boolean operations. IBM's PL/I (ca. 1966) includes a built-in function named BOOL. In  $\text{BOOL}(x, y, z)$ ,  $z$  is a bit string of length four (or converted to that

---

7. The entire 335-page work is available at [www.gutenberg.org/etext/15114](http://www.gutenberg.org/etext/15114).

if necessary), and  $x$  and  $y$  are bit strings of equal length (or converted to that if necessary). Argument  $z$  specifies the Boolean operation to be performed on  $x$  and  $y$ . Binary 0000 is the zero function, 0001 is  $xy$ , 0010 is  $x\bar{y}$ , and so forth.

Another such language is Basic for the Wang System 2200B computer (ca. 1974), which provides a version of BOOL that operates on character strings rather than on bit strings or integers [Neum].

Still another such language is MIT PDP-6 Lisp, later called MacLisp [GLS1].

## Exercises

1. David de Kloet suggests the following code for the snoob function, for  $x \neq 0$ , where the final assignment to  $y$  is the result:

```

y ← x + (x & -x)
x ← x & -y
while((x & 1) = 0) x ← x  $\overset{s}{\gg}$  1
x ← x  $\overset{s}{\gg}$  1
y ← y | x

```

This is essentially the same as Gosper's code (page 15), except the right shift is done with a *while*-loop rather than with a *divide* instruction. Because division is usually costly in time, this might be competitive with Gosper's code if the *while*-loop is not executed too many times. Let  $n$  be the length of the bit strings  $x$  and  $y$ ,  $k$  the number of 1-bits in the strings, and assume the code is executed for all values of  $x$  that have exactly  $k$  1-bits. Then for each invocation of the function, how many times, on average, will the body of the *while*-loop be executed?

2. The text mentions that a left shift by a variable amount is not right-to-left computable. Consider the function  $x \ll (x \& 1)$  [Knu8]. This is a left shift by a variable amount, but it can be computed by

$$\begin{aligned}
 &x + (x \& 1) * x, \text{ or} \\
 &x + (x \& (-(x \& 1))),
 \end{aligned}$$

which are all right-to-left computable operations. What is going on here? Can you think of another such function?

3. Derive Dietz's formula for the average of two unsigned integers,

$$(x \& y) + ((x \oplus y) \overset{u}{\gg} 1).$$

4. Give an overflow-free method for computing the average of four unsigned integers,  $\lfloor (a + b + c + d)/4 \rfloor$ .
5. Many of the comparison predicates shown on page 23 can be simplified substantially if bit 31 of either  $x$  or  $y$  is known. Show how the seven-instruction expression for  $x \leq y$  can be simplified to three basic RISC, non-comparison, instructions if  $y_{31} = 0$ .
6. Show that if two numbers, possibly distinct, are added with “end-around carry,” the addition of the carry bit cannot generate another carry out of the high-order position.
7. Show how end-around carry can be used to do addition if negative numbers are represented in one’s-complement notation. What is the maximum number of bit positions that a carry (from any bit position) might be propagated through?
8. Show that the MUX operation,  $(x \& m) \mid (y \& \sim m)$ , can be done in three instructions on the basic RISC (which does not have the *and with complement* instruction).
9. Show how to implement  $x \oplus y$  in four instructions with *and-or-not* logic.
10. Given a 32-bit word  $x$  and two integer variables  $i$  and  $j$  (in registers), show code to copy the bit of  $x$  at position  $i$  to position  $j$ . The values of  $i$  and  $j$  have no relation, but assume that  $0 \leq i, j \leq 31$ .
11. How many binary Boolean instructions are sufficient to evaluate any  $n$ -variable Boolean function if it is decomposed recursively by the method of the theorem?
12. Show that alternative decompositions of Boolean functions of three variables are
  - (a)  $f(x, y, z) = g(x, y) \oplus \bar{z}h(x, y)$  (the “negative Davio decomposition”), and
  - (b)  $f(x, y, z) = g(x, y) \oplus (z + h(x, y))$ .
13. It is mentioned in the text that all 16 binary Boolean operations can be done with the eight instructions shown in Table 2-3, by interchanging the inputs or by having both register fields of the instruction refer to the same register. Show how to do this.
14. Suppose you are not concerned about the six Boolean functions that are really constants or unary functions, namely  $f(x, y) = 0, 1, x, y, \bar{x}$ , and  $\bar{y}$ , but you want your instruction set to compute the other ten functions with one instruction. Can this be done with fewer than eight binary Boolean instruction types (opcodes)?
15. Exercise 13 shows that eight instruction types suffice to compute any of the 16 two-operand Boolean operations with one R-R (register-register) instruction. Show that six instruction types suffice in the case of R-I (register-immediate)

instructions. With R-I instructions, the input operands cannot be interchanged or equated, but the second input operand (the immediate field) can be complemented or, in fact, set to any value at no cost in execution time. Assume for simplicity that the immediate fields are the same length as the general purpose registers.

- 16.** Show that not all Boolean functions of three variables can be implemented with three binary logical instructions.



# INDEX

0-bits, leading zeros. *See* nlz function.  
0-bits, trailing zeros. *See also* ntz (*number of trailing zeros*) function.  
    counting, 107–114.  
    detecting, 324. *See also* CRC (cyclic redundancy check).  
    plots and graphs, 466  
0-bytes, finding, 117–121  
1-bits, counting. *See* Counting bits.  
3:2 compressor, 90–95  
The 16 Boolean binary operations, 53–57

## A

Absolute value  
    computing, 18  
    multibyte, 40–41  
    negative of, 23–26  
*add* instruction  
    condition codes, 36–37  
    propagating arithmetic bounds, 70–73  
Addition  
    arithmetic tables, 453  
    combined with logical operations, 16–17  
    double-length, 38–39  
    multibyte, 40–41  
    of negabinary numbers, 301–302  
    overflow detection, 28–29  
    plots and graphs, 461  
    in various number encodings, 304–305  
Advanced Encryption Standard, 164  
Alternating among values, 48–51  
Alverson's method, 237–238  
*and*  
    plots and graphs, 459  
    in three instructions, 17  
*and with complement*, 131  
Answers to exercises, by chapter  
    1: Introduction, 405–406  
    2: Basics, 407–415  
    3: Power-of-2 Boundaries, 415–416  
    4: Arithmetic Bounds, 416–417  
    5: Counting Bits, 417–418  
    6: Searching words, 418–423  
    7: Rearranging Bits and Bytes, 423–425  
    8: Multiplication, 425–428  
    9: Integer Division, 428–430  
    10: Integer Division by Constants, 431–434  
    11: Some Elementary Functions, 434–435  
    12: Unusual Bases for Number Systems, 435–439  
    13: Gray Code, 439–441  
    14: Cyclic Redundancy Check, 441–442  
    15: Error-Correcting Codes, 442–445  
    16: Hilbert's Curve, 446  
    17: Floating-Point, 446–448  
    18: Formulas for Primes, 448–452  
Arithmetic, computer vs. ordinary, 1  
Arithmetic bounds  
    checking, 67–69  
    of expressions, 70–71  
    propagating through, 70–73  
    range analysis, 70  
    searching for values in, 122  
Arithmetic tables, 4-bit machine, 453–456  
Arrays  
    checking bounds. *See* Arithmetic bounds.  
    counting 1-bits, 89–96  
    indexes, checking. *See* Arithmetic bounds.  
    indexing a sparse array, 95  
    permutation, 161–163  
    rearrangements, 165–166  
    of short integers, 40–41  
Autodin-II polynomial, 323  
Average, computing, 19, 55–56

## B

Base  $-1 + i$  number system, 306–308  
    extracting real and imaginary parts, 310  
Base  $-1 - i$  number system, 308–309

- Base -2 number system, 299–306
  - Gray code, 315
  - rounding down, 310
- Basic RISC instruction set, 5–6
- Basic, Wang System 2200B, 55
- Big-endian format, converting to little-endian, 129
- Binary decomposition, integer exponentiation, 288–290
- Binary forward error-correcting block codes (FEC), 331
- Binary search
  - counting leading 0's, 99–104
  - integer logarithm, 291–297
  - integer square root, 279–287
- Bit matrices, multiplying, 98
- Bit operations
  - compress* operation, 150–156
  - computing parity. *See* Parity.
  - counting bits. *See* Counting bits.
  - finding strings of 1-bits, 123–128
  - flipping bits, 135
  - general permutations, 161–165
  - generalized bit reversal, 135
  - generalized extract, 150–156
  - half shuffle, 141
  - inner perfect shuffle, plots and graphs, 468–469
  - inner perfect unshuffle, plots and graphs, 468
  - inner shuffle, 139–141
  - numbering schemes, 1
  - outer shuffle, 139–141, 373
  - perfect shuffle, 139–141
  - reversing bits. *See* Reversing bits and bytes.
  - on rightmost bits. *See* Rightmost bits.
  - searching words for bit strings, 107, 123–128
  - sheep and goats operation, 161–165
  - shuffling bits, 139–141, 165–166
  - transposing a bit matrix, 141–150
  - unshuffling bits, 140–141, 150, 162
- Bit reversal function, plots and graphs, 467
- Bit vectors, 1
- bitgather** instruction, 163–165
- Bits. *See specific topics.*
- bitsize function, 106–107
- Bliss, Robert D., xv
- Bonzini, Paolo, 263
- BOOL function, 54–55
- Boole, George, 54
- Boolean binary operations, all 16, 53–57
- Boolean decomposition formula, 51–53, 56–57
- Boundary crossings, powers of 2, 63–64
- Bounds, arithmetic. *See* Arithmetic bounds.
- Bounds checking. *See* Checking arithmetic bounds.
- branch on carry and register result non-zero* instruction, 63
- Bytes. *See also specific topics.*
  - definition, 1
  - finding first 0-byte, 117–121
- C**
- C language
  - arithmetic on pointers, 105, 240
  - GNU extensions, 105
  - iterative statements, 4, 10
  - referring to same location with different types, 104
  - representation of character strings, 117
  - summary of elements, 2–4
- Caches, 166–167
- Carry-save adder (CSA) circuit, 90–95
- CCITT (Le Comité Consultatif International...), 321
- Ceiling function, identities, 183–184
- Chang, Albert, 123
- Character strings, 117
- Check bits
  - Hamming code, 332
  - SEC-DED code, 334–335
- Checking arithmetic bounds, 67–69
- Chinese ring puzzle, 315
- Chipkill technology, 336
- Code, definition, 343
- Code length, 331, 343
- Code rate, 343
- Code size, 343
- Comparison predicates
  - from the carry bit, 26–27
  - definition, 23

- number of leading zeros* (nlz) function, 23–24, 107
  - signed comparisons, from unsigned, 25
  - true/false results, 23
  - using negative absolute values, 23–26
  - Comparisons
    - computer evaluation of, 27
    - floating-point comparisons using integer operations, 381–382
    - three-valued compare* function, 21–22. *See also* sign function.
  - Compress function, plots and graphs, 464–465
  - compress* operation, 119, 150–161
    - with *insert* and *extract* instructions, 155–156
  - Computability test, right-to-left, 13–14, 55
  - Computer algebra, 2–4
  - Computer arithmetic
    - definition, 1
    - plots and graphs, 461–463
  - Condition codes, 36–37
  - Constants
    - dividing by. *See* Division of integers by constants.
    - multiplying by, 175–178
  - Counting bits. *See also* ntz (*number of trailing zeros*) function; nlz (*number of leading zeros*) function; *population count* function.
  - 1-bits in
    - 7- and 8-bit quantities, 87
    - an array, 89–95
    - a word, 81–88
  - bitsize function, 106–107
  - comparing two words, 88–89
  - divide and conquer strategy, 81–82
  - leading 0's, with
    - binary search method, 99–100
    - floating-point methods, 104–106
    - population count* instruction, 101–102
  - rotate and sum method, 85–86
  - search tree method, 109
  - with table lookup, 86–87
  - trailing 0's, 107–114
  - by turning off 1-bits, 85
  - CRC (cyclic redundancy check)
    - background, 319–320
    - check bits, generating, 319–320
    - checksum, computing
      - generator polynomials, 322–323, 329
    - with hardware, 324–326
    - with software, 327–329
    - with table lookup, 328–329
    - techniques for, 320
  - code vector, 319
  - definition, 319
  - feedback shift register* circuit, 325–326
  - generator polynomial, choosing, 322–323, 329
  - parity bits, 319–320
  - practice
    - hardware checksums, 324–326
    - leading zeros, detecting, 324
    - overview, 323–324
    - residual/residue, 324
    - software checksums, 327–329
    - trailing zeros, detecting, 324
    - theory, 320–323
  - CRC codes, generator polynomials, 322, 323
  - CRC-CITT polynomial, 323
  - Cryptography
    - Advanced Encryption Standard, 164
    - bitgather* instruction, 164–165
    - DES (Data Encryption Standard), 164
    - Rijndael algorithm, 164
    - SAG method, 162–165
    - shuffling bits, 139–141, 165
    - Triple DES, 164
  - CSA (carry-save add) circuit, 90–95
  - Cube root, approximate, floating-point, 389
  - Cube root, integer, 287–288
  - Curves. *See also* Hilbert's curve.
    - Peano, 371–372
    - space-filling, 355–372
  - Cycling among values, 48–51
- ## D
- Davio decomposition, 51–53, 56–57
  - de Bruijn cycles, 111–112
  - de Kloet, David, 55

- De Morgan's laws, 12–13
- DEC PDP-10 computer, xiii, 84
- Decryption. *See* Cryptography.
- DES (Data Encryption Standard), 164
- Dietz's formula, 19, 55
- difference or zero* (doz) function, 41–45
- Distribution of leading digits, 385–387
- Divide and conquer strategy, 81–82
- Division
  - arithmetic tables, 455
  - doubleword
    - from long division, 197–202
    - signed, 201–202
    - by single word, 192–197
    - unsigned, 197–201
  - floor, 181–182, 237
  - modulus, 181–182, 237
  - multiword, 184–188
  - of negabinary numbers, 302–304
  - nonrestoring algorithm, 192–194
  - notation, 181
  - overflow detection, 34–36
  - plots and graphs, 463–464
  - restoring algorithm, 192–193
  - shift-and-subtract algorithms (hardware), 192–194
  - short, 189–192, 195–197
  - signed
    - computer, 181
    - doubleword, 201–202
    - long, 189
    - multiword, 188
    - short, 190–192
  - unsigned
    - computer, 181
    - doubleword, 197–201
    - long, 192–197
    - short from signed, 189–192
- Division of integers by constants
  - by 3, 207–209, 276–277
  - by 5 and 7, 209–210
  - exact division
    - converting to, 274–275
    - definition, 240
    - multiplicative inverse, Euclidean algorithm, 242–245
    - multiplicative inverse, Newton's method, 245–247
    - multiplicative inverse, samples, 247–248
- floor division, 237
- incorporating into a compiler, signed, 220–223
- incorporating into a compiler, unsigned, 232–234
- magic numbers
  - Alverson's method, 237–238
  - calculating, signed, 212–213, 220–223
  - calculating, unsigned, 231–234
  - definition, 211
  - sample numbers, 238–239
  - table lookup, 237
  - uniqueness, 224
- magicu algorithm, 232–234
- magicu2 algorithm, 236
- modulus division, 237
- remainder by multiplication and shifting right
  - signed, 273–274
  - unsigned, 268–272
- remainder by summing digits
  - signed, 266–268
  - unsigned, 262–266
- signed
  - by divisors  $\leq -2$ , 218–220
  - by divisors  $\geq 2$ , 210–218
  - by powers of 2, 205–206
  - incorporating into a compiler, 220–223
  - not using `mulhs` (*multiply high signed*), 259–262
  - remainder by multiplication and shifting right, 273–274
  - remainder by summing digits, 266–268
  - remainder from powers of 2, 206–207
  - test for zero remainder, 250–251
  - uniqueness, 224
- timing test, 276
- unsigned
  - best programs for, 234–235
  - by 3 and 7, 227–229
  - by divisors  $\geq 1$ , 230–232
  - by powers of 2, 227

Division of integers by constants, unsigned  
     (*continued*)  
     incorporating into a compiler,  
         232–234  
     incremental division and remainder  
         technique, 232–234  
     not using `mulhu` (*multiply high  
         unsigned*) instruction, 251–259  
     remainder by multiplication and  
         shifting right, 268–272  
     remainder by summing digits,  
         262–266  
     remainder from powers of 2, 227  
     test for zero remainder, 248–250

Double buffering, 46

Double-length addition/subtraction, 38–39

Double-length shifts, 39–40

Doubleword division  
     by single word, 192–197  
     from long division, 197–202  
     signed, 201–202  
     unsigned, 197–201

Doublewords, definition, 1

`doz` (*difference or zero*) function, 41–45

Dubé, Danny, 112

## E

ECCs (error-correcting codes)  
     check bits, 332  
     code, definition, 343  
     code length, 331, 343  
     code rate, 343  
     code size, 343  
     coding theory problem, 345–351  
     efficiency, 343

FEC (binary forward error-correcting  
     block codes), 331

Gilbert-Varshamov bound, 348–350

Hamming bound, 348, 350

Hamming code, 332–342  
     converting to SEC-DED code,  
         334–337  
     extended, 334–337  
     history of, 335–337  
     overview, 332–334  
     SEC-DED on 32 information bits,  
         337–342

Hamming distance, 95, 343–345

information bits, 332

linear codes, 348–349

overview, 331, 342–343

perfect codes, 333, 349, 352

SEC (single error-correcting) codes,  
     331

SEC-DED (single error-correcting,  
     double error-detecting) codes  
     on 32 information bits, 337–342  
     check bits, minimum required, 335  
     converting from Hamming code,  
         334–337  
     definition, 331  
     singleton bound, 352  
     sphere-packing bound, 348, 350  
     spheres, 347–351

Encryption. *See* Cryptography.

End-around-carry, 38, 56, 304–305

Error detection, digital data. *See* CRC  
     (cyclic redundancy check).

Estimating multiplication overflow, 33–34

Euclidean algorithm, 242–245

Euler, Leonhard, 392

Even parity, 96

Exact division  
     definition, 240  
     multiplicative inverse, Euclidean algo-  
         rithm, 242–245  
     multiplicative inverse, Newton's  
         method, 245–247  
     multiplicative inverse, samples,  
         247–248  
     overview, 240–242

Exchanging  
     conditionally, 47  
     corresponding register fields, 46  
     two fields in same register, 47  
     two registers, 45–46

*exclusive or*  
     plots and graphs, 460  
     propagating arithmetic bounds through,  
         77–78  
     scan operation on an array of bits, 97  
     in three instructions, 17

Execution time model, 9–10

Exercise answers. *See* Answers to exercises.

*Expand* operation, 156–157, 159–161

Exponentiation

- by binary decomposition, 288–290
- in Fortran, 290

Extended Hamming code, 334–342

- on 32 information bits, 337–342

Extract, generalized, 150–156

## F

Factoring, 178

FEC (binary forward error-correcting block codes), 331

*feedback shift register* circuit, 325–326

Fermat numbers, 391

FFT (Fast Fourier Transform), 137–139

*find leftmost 0-byte*, 117–121

*find rightmost 0-byte*, 118–121

Finding

- decimal digits, 122
- first 0-byte, 117–121
- first uppercase letter, 122
- length of character strings, 117
- next higher number, same number of 1-bits, 14–15
- the  $n$ th prime, 391–398, 403
- strings of 1-bits
  - first string of a given length, 123–125
  - longest string, 125–126
  - shortest string, 126–128
- values within arithmetic bounds, 122

Flipping bits, 135

Floating-point numbers, 375–389

- distribution of leading digits, 385–387
- formats (single/double), 375–376
- gradual underflow, 376
- IEEE arithmetic standard, 375
- IEEE format, 375–377
- NaN (not a number), 375–376
- normalized, 375–377
- subnormal numbers, 375–377
- table of miscellaneous values, 387–389
- ulp (unit in the last position), 378

Floating-point operations

- approximate cube root, 389
- approximate reciprocal square root, 383–385

approximate square root, 389

comparing using integer operations, 381–382

conversion table, 378–381

converting to/from integers, 377–381

counting leading 0's with, 104–106

simulating, 107

Floor division, 181–182, 237

Floor function, identities, 183, 202–203

Floyd, R. W., 114

Formula functions, 398–403

Formulas for primes, 391–403

Fortran

IDIM function, 44

integer exponentiation, 290

ISIGN function, 22

MOD function, 182

Fractal triangles, plots and graphs, 460

Full adders, 90

Full RISC instruction set, 7

Fundamental theorem of arithmetic, 404

## G

Gardner, Martin, 315

Gaudet, Dean, 110

Gaudet's algorithm, 110

*generalized extract* operation, 150–156

Generalized unshuffle. *See* SAG (sheep and goats) operation.

Generator polynomials, CRC codes, 321–323

Gilbert-Varshamov bound, 348–350

Golay, M. J. E., 331

Goryavsky, Julius, 103

Gosper, R. W.

- iterating through subsets, 14–15
- loop-detection, 114–116

Gradual underflow, 376

Graphics-rendering, Hilbert's curve, 372–373

Graphs. *See* Plots and graphs.

Gray, Frank, 315

Gray code

- applications, 315–317
- balanced, 317
- converting integers to, 97, 312–313
- cyclic, 312

- definition, 311
- history of, 315–317
- incrementing Gray-coded integers, 313–315
- negabinary Gray code, 315
- plots and graphs, 466
- reflected, 311–312, 315
- single track (STGC), 316–317

- Greatest common divisor* function, plots and graphs, 464

- GRP instruction, 165

## H

- Hacker, definition, xvi
- HAKMEM (hacks memo), xiii
- Half shuffle, 141
- Halfwords, 1
- Hamiltonian paths, 315
- Hamming, R. W., 331
- Hamming bound, 348, 350
- Hamming code
  - on 32 information bits, 337–342
  - converting to SEC-DED code, 334–337
  - extended, 334–337
  - history of, 335–337
  - overview, 332–334
  - perfect, 333, 352
- Hamming distance, 95, 343–345
  - triangle inequality, 352
- Hardware checksums, 324–326
- Harley, Robert, 90, 101
- Harley's algorithm, 101, 103
- Hexadecimal floating-point, 385
- High-order half of product, 173–174
- Hilbert, David, 355
- Hilbert's curve. *See also* Space-filling curves.
  - applications, 372–373
  - coordinates from distance
    - curve generator driver program, 359
  - description, 358–366
  - Lam and Shapiro method, 362–364, 368
  - parallel prefix operation, 3
    - 65–366
  - state transition table, 361, 367
  - description, 355–356

- distance from coordinates, 366–368
- generating, 356–358
- illustrations, 355, 357
- incrementing coordinates, 368–371
- non-recursive generation, 371
- ray tracing, 372
- three-dimensional analog, 373
- Horner's rule, 49

## I

### IBM

- Chipkill technology, 336
- Harvest computer, 336
- PCs, error checking, 336
- PL/I language, 54
- Stretch computer, 81, 336
- System/360 computer, 385
- System/370 computer, 63
- IDIM function, 44
- IEEE arithmetic standard, 375
- IEEE format, floating-point numbers, 375–377
- IEEE Standard for Floating-Point Arithmetic*, 375
- Image processing, Hilbert's curve, 372
- Incremental division and remainder technique, 232–234
- Inequalities, logical and arithmetic expressions, 17–18
- Information bits, 332
- Inner perfect shuffle* function, plots and graphs, 468–469
- Inner perfect unshuffle* function, plots and graphs, 468
- Inner shuffle, 139–141
- insert* instruction, 155–156
- Instruction level parallelism, 9
- Instruction set for this book, 5–8
- integer cube root* function, 287–288, 297
- Integer exponentiation, 288–290
- integer fourth root* function, 297
- integer log base 2* function, 106, 291
- integer log base 10* function, 292–297
- Integer quotient function, plots and graphs, 463
- integer remainder* function, 463
- integer square root* function, 279–287

Integers. *See also specific operations on integers.*

complex, 306–309

converting to/from floating-point, 377–381

converting to/from Gray code, 97, 312–313

reversed, incrementing, 137–139

reversing, 129–137

Inverse Gray code function

formula, 312

plots and graphs, 466

*An Investigation of the Laws of Thought*, 54

ISIGN (transfer of sign) function, 22

Iterating through subsets, 14–15

ITU-TSS (International Telecommunications Union...), 321

ITU-TSS polynomial, 323

## K

Knuth, Donald E., 132

Knuth's Algorithm D, 184–188

Knuth's Algorithm M, 171–172, 174–175

Knuth's mod operator, 181

Kronecker, Leopold, 375

## L

Lam and Shapiro method, 362–364, 368

Landry, F., 391

Leading 0's, counting, 99–106. *See also* `nlz` (*number of leading zeros*) function.

Leading 0's, detecting, 324. *See also* CRC (cyclic redundancy check).

Leading digits, distribution, 385–387

Least common multiple function, plots and graphs, 464

Linear codes, 348–349

Little-endian format, converting to/from big-endian, 129

*load word byte-reverse* (`lwbrrx`) instruction, 118

Logarithms

binary search method, 292–293

definition, 291

log base 2, 106–107, 291

log base 10, 291–297

table lookup, 292, 294–297

Logical operations

with addition and subtraction, 16–17

*and*, plots and graphs, 459

binary, table of, 17

*exclusive or*, plots and graphs, 460

*or*, plots and graphs, 459

propagating arithmetic bounds through, 74–76, 78

tight bounds, 74–78

Logical operators on integers, plots and graphs, 459–460

Long Division, definition, 189

Loop detection, 114–115

LRU (least recently used) algorithm, 166–169

`lwbrrx` (*load word byte-reverse*) instruction, 118

## M

MacLisp, 55

*magic* algorithm

incremental division and remainder technique, 232–234

signed division, 220–223

unsigned division, 232–234

Magic numbers

Alverson's method, 237–238

calculating, signed, 212–213, 220–223

calculating, unsigned, 232–234

calculating, Python code for

definition, 211

samples, 238–239

table lookup, 237

uniqueness, 224

*magicu* algorithm, 232–234

in Python, 240

*magicu2* algorithm, 236–237

`max` function, 41–45

Mills, W. H., 403

Mills's theorem, 403–404

`min` function, 41–45

MIT PDP-6 Lisp, 55

MOD function (Fortran), 182

`modu` (*unsigned modulus*) function, 98

Modulus division, 181–182, 237

Moore, Eliakim Hastings, 371–372



**mulhs** (*multiply high signed*) instruction  
 division with, 207–210, 212, 218, 222, 235  
 implementing in software, 173–174  
 not using, 259–262

**mulhu** (*multiply high unsigned*) instruction  
 division with, 228–229, 234–235, 238  
 implementing in software, 173  
 not using, 251–259

**Multibyte absolute value**, 40–41

**Multibyte addition/subtraction**, 40–41

**Multiplication**  
 arithmetic tables, 454  
 of complex numbers, 178–179  
 by constants, 175–178  
 factoring, 178  
 low-order halves independent of signs, 178  
 high-order half of 64-bit product, 173–174  
 high-order product signed from/to unsigned, 174–175  
 multiword, 171–173  
 of negabinary numbers, 302  
 overflow detection, 31–34  
 plots and graphs, 462

**Multiplicative inverse**  
 Euclidean algorithm, 242–245  
 Newton's method, 245–247, 278  
 samples, 247–248

***multiply* instruction**, condition codes, 36–37

**Multiword division**, 184–189

**Multiword multiplication**, 171–173

**MUX operation in three instructions**, 56

**mux** (*multiplex*) instruction, 406

## N

**NAK** (negative acknowledgment), 319

**NaN** (not a number), 375–376

**Negabinary number system**, 299–306  
 Gray code, 315

**Negative absolute value**, 23–26

**Negative overflow**, 30

**Newton-Raphson calculation**, 383

**Newton's method**, 457–458  
 integer cube root, 287–288

integer square root, 279–283

multiplicative inverse, 245–248

**Next higher number**, same number of 1-bits, 14–15

**Nibbles**, 1

**nlz** (*number of leading zeros*) function  
 applications, 79, 107, 128  
 bitsize function, 106–107  
 comparison predicates, 23–24, 107  
 computing, 99–106  
 for counting trailing 0's, 107  
 finding 0-bytes, 118  
 finding strings of 1-bits, 123–124  
 incrementing reversed integers, 138  
 and *integer log base 2* function, 106  
 rounding to powers of 2, 61

**Nonrestoring algorithm**, 192–194

**Normalized numbers**, 376

**Notation used in this book**, 1–4

***n*th prime**, finding  
 formula functions, 398–401  
 Willans's formulas, 393–397  
 Wormell's formula, 397–398

**ntz** (*number of trailing zeros*) function  
 applications, 114–116  
 from counting leading 0's, 107  
 loop detection, 114–115  
 ruler function, 114

**Number systems**  
 base  $-1 + i$ , 306–308  
 base  $-1 - i$ , 308–309  
 base  $-2$ , 299–306, 315  
 most efficient base, 309–310  
 negabinary, 299–306, 315

## O

**Odd parity**, 96

**1-bits**, counting. *See* Counting bits.

**or**  
 plots and graphs, 459  
 in three instructions, 17

**Ordinary arithmetic**, 1

**Ordinary rational division**, 181

***Outer perfect shuffle bits* function**, plots and graphs, 469

***Outer perfect shuffle* function**, plots and graphs, 467

*Outer perfect unshuffle* function, plots and graphs, 468

Outer shuffle, 139–141, 373

Overflow detection

definition, 28

division, 34–36

estimating multiplication overflow,  
33–34

multiplication, 31–34

negative overflow, 30

signed add/subtract, 28–30

unsigned add/subtract, 31

## P

Parallel prefix operation

definition, 97

Hilbert's curve, 364–366

inverse, 116

parity, 97

Parallel suffix operation

*compress* operation, 150–155

*expand* operation, 156–157, 159–161

generalized extract, 150–156

inverse, 116

Parity

adding to 7-bit quantities, 98

applications, 98

computing, 96–98

definition, 96

parallel prefix operation, 97

scan operation, 97

two-dimensional, 352

Parity bits, 319–320

PCs, error checking, 336

Peano, Giuseppe, 355

Peano curves, 371–372. *See also* Hilbert's curve.

Peano-Hilbert curve. *See* Hilbert's curve.

Perfect codes, 333, 349

Perfect shuffle, 139–141, 373

Permutations on bits, 161–165. *See also* Bit operations.

Planar curves, 355. *See also* Hilbert's curve.

Plots and graphs, 459–469

addition, 461

bit reversal function, 467

*compress* function, 464–465

division, 463–464

fractal triangles, 460

Gray code function, 466

*greatest common divisor* function, 464

inner perfect shuffle, 468–469

inner perfect unshuffle, 468

integer quotient function, 463

inverse Gray code function, 466

*least common multiple* function, 464

logical *and* function, 459

logical *exclusive or* function, 460

logical operators on integers, 459–460

logical *or* function, 459

multiplication, 462

number of trailing zeros, 466

outer perfect shuffle, 467–469

outer perfect unshuffle, 468

*population count* function, 467

*remainder* function, 463

*rotate left* function, 465

ruler function, 466

SAG (sheep and goats) function,  
464–465

self-similar triangles, 460

Sierpinski triangle, 460

subtraction, 461

unary functions, 466–469

unsigned product of  $x$  and  $y$ , 462

Poetry, 278, 287

*population count* function. *See also* Counting bits.

applications, 95–96

computing Hamming distance, 95

counting 1-bits, 81

counting leading 0's, 101–102

counting trailing 0's, 107–114

plots and graphs, 467

Position sensors, 315–317

Powers of 2

boundary crossings, detecting, 63–64

rounding to, 59–62, 64

signed division, 205–206

unsigned division, 227

PPERM instruction, 165

Precision, loss of, 385–386

Prime numbers

Fermat numbers, 391

- finding the  $n$ th prime
    - formula functions, 398–403
    - Willans's formulas, 393–397
    - Wormell's formula, 397–398
  - formulas for, 391–403
  - from polynomials, 392
  - Propagating arithmetic bounds
    - add* and *subtract* instructions, 70–73
    - logical operations, 73–78
    - signed numbers, 71–73
    - through *exclusive or*, 77–78
  - PSHUFB (*Shuffle Packed Bytes*) instruction, 163
  - PSHUFD (*Shuffle Packed Doublewords*) instruction, 163
  - PSHUFW (*Shuffle Packed Words*) instruction, 163
- Q**
- Quicksort, 81
- R**
- Range analysis, 70
  - Ray tracing, Hilbert's curve, 372
  - Rearrangements and index transformations, 165–166
  - Reed-Muller decomposition, 51–53, 56–57
  - Reference matrix method (LRU), 166–169
  - Reflected binary Gray code, 311–312, 315
  - Registers
    - exchanging, 45–46
    - exchanging conditionally, 47
    - exchanging fields of, 46–47
    - reversing contents of, 129–135
  - RISC computers, 5
  - Reiser, John, 113
  - Reiser's algorithm, 113–114
  - Remainder* function, plots and graphs, 463
  - Remainders
    - arithmetic tables, 456
    - of signed division
      - by multiplication and shifting right, 273–274
      - by summing digits, 266–268
    - from non-powers of 2, 207–210
    - from powers of 2, 206–207
    - test for zero, 248–251
  - of unsigned division
    - by multiplication and shifting right, 268–272
    - by summing digits, 262–266
    - and immediate* instruction, 227
    - incremental division and remainder technique, 232–234
    - test for zero, 248–250
  - remu function, 119, 135–136
  - Residual/residue, 324
  - Restoring algorithm, 192–193
  - Reversing bits and bytes, 129–137
    - 6-, 7-, 8-, and 9-bit quantities, 135–137
    - 32-bit words, 129–135
    - big-endian format, converting to little-endian, 129
  - definition, 129
  - generalized, 135
  - load word byte-reverse (lwbrx)* instruction, 118
  - rightmost 16 bits of a word, 130
  - with rotate shifts, 129–133
  - small integers, 135–137
  - table lookup, 134
  - Riemann hypothesis, 404
  - Right justify* function, 116
  - Rightmost bits, manipulating, 11–12, 15
    - De Morgan's laws, 12–13
    - right-to-left computability test, 13–14, 55
  - Rijndael algorithm, 164
  - RISC
    - basic instruction set, 5–6
    - execution time model, 9–10
    - extended mnemonics, 6, 8
    - full instruction set, 7–8
    - registers, 5–6
  - Rotate and sum method, 85–86
  - Rotate left* function, plots and graphs, 464–465
  - Rotate shifts, 37–38, 129–133
  - Rounding to powers of 2, 59–62, 64
  - Ruler function, 114, 466
  - Russian decomposition, 51–53, 56–57

## S

- SAG (sheep and goats) operation
  - description, 162–165
  - plots and graphs, 464–465
- Scan operation, 97
- Seal, David, 90, 110
- Search tree method, 109
- Searching. *See* Finding.
- SEC (single error-correcting) codes, 331
- SEC-DED (single error-correcting, double error-detecting) codes
  - on 32 information bits, 337–342
  - check bits, minimum required, 335
  - converting from Hamming code, 334–335
  - definition, 331
- Select* instruction, 406
- Self-reproducing program, xvi
- Self-similar triangles, plots and graphs, 460
- shift left double* operation, 39
- shift right double signed* operation, 39–40
- shift right double unsigned* operation, 39
- shift right extended immediate (shrxi)* instruction, 228–229
- shift right signed* instruction
  - alternative to, for sign extension, 19–20
  - division by power of 2, 205–206
  - from unsigned, 20
- Shift-and-subtract algorithm
  - hardware, 192–194
  - integer square root, 285–287
- Shifts
  - double-length, 39–40
  - rotate, 37–38
- Short division, 189–192, 195–196
- Shroepel's formula, 305–306
- shrxi* (*shift right extended immediate*) instruction, 228–229
- Shuffle Packed Bytes (PSHUFb)* instruction, 163
- Shuffle Packed Doublewords (PSHUFd)* instruction, 163
- Shuffle Packed Words (PSHUFw)* instruction, 163
- Shuffling
  - arrays, 165–166
  - bits
    - half shuffle, 141
    - inner perfect shuffle, plots and graphs, 468–469
    - inner perfect unshuffle, plots and graphs, 468
    - inner shuffle, 139–141
    - outer shuffle, 139–141, 373
    - perfect shuffle, 139–141
    - shuffling bits, 139–141, 165–166
    - unshuffling, 140–141, 150, 162, 165–166
- Sierpinski triangle, plots and graphs, 460
- Sign extension, 19–20
- sign function, 20–21. *See also three-valued compare* function.
- Signed bounds, 78
- Signed comparisons, from unsigned, 25
- Signed computer division, 181–182
- Signed division
  - arithmetic tables, 455
  - computer, 181
  - doubleword, 201–202
  - long, 189
  - multiword, 188
  - short, 190–192
- Signed division of integers by constants
  - best programs for, 225–227
  - by divisors  $\leq -2$ , 218–220
  - by divisors  $\geq 2$ , 210–218
  - by powers of 2, 205–206
  - incorporating into a compiler, 220–223
  - remainder from non-powers of 2, 207–210
  - remainder from powers of 2, 206–207
  - test for zero remainder, 250–251
  - uniqueness of magic number, 224
- Signed long division, 189
- Signed numbers, propagating arithmetic bounds, 71–73
- Signed short division, 190–192
- signum function, 20–21
- Single error-correcting, double error-detecting (SEC-DED) codes. *See* SEC-DED (single error-correcting, double error-detecting) codes.
- Single error-correcting (SEC) codes, 331
- snoob function, 14–15
- Software checksums, 327–329

Space-filling curves, 371–372. *See also* Hilbert's curve.

Sparse array indexing, 95

Sphere-packing bound, 348–350

Spheres, ECCs (error-correcting codes), 347–350

Square root, integer  
     binary search, 281–285  
     hardware algorithm, 285–287  
     Newton's method, 279–283  
     shift-and-subtract algorithm, 285–287

Square root, approximate, floating-point, 389

Square root, approximate reciprocal, floating-point, 383–385

Stibitz, George, 308

Strachey, Christopher, 130

Stretch computer, 81, 336

Strings. *See* Bit operations; Character strings.

**strlen** (*string length*) C function, 117

Subnormal numbers, 376

Subnorms, 376

*subtract* instruction  
     condition codes, 36–37  
     propagating arithmetic bounds, 70–73

Subtraction  
     arithmetic tables, 453  
     *difference or zero* (doz) function, 41–45  
     double-length, 38–39  
     combined with logical operations, 16–17  
     multibyte, 40–41  
     of negabinary numbers, 301–302  
     overflow detection, 29–31  
     plots and graphs, 461

Swap-and-complement method, 362–365

Swapping pointers, 46

System/360 computer, 385

System/370 computer, 63

## T

Table lookup, counting bits, 86–87

*three-valued compare* function, 21–22. *See also* sign function.

Tight bounds  
     *add* and *subtract* instructions, 70–73  
     logical operations, 74–79

Timing test, division of integers by constants, 276

Toggling among values, 48–51

Tower of Hanoi puzzle, 116, 315

Trailing zeros. *See also* *ntz* (*number of trailing zeros*) function.  
     counting, 107–114  
     detecting, 324. *See also* CRC (cyclic redundancy check).  
     plots and graphs, 466

Transfer of sign (ISIGN) function, 22

Transposing a bit matrix  
     8 x 8, 141–145  
     32 x 32, 145–149

Triangles  
     fractal, 460  
     plots and graphs, 460  
     self-similar, 460  
     Sierpinski, 460

Triple DES, 164

True/false comparison results, 23

Turning off 1-bits, 85

## U

Ulp (unit in the last position), 378

Unaligned load, 65

Unary functions, plots and graphs, 466–469

Uniqueness, of magic numbers, 224

Unshuffling  
     arrays, 162  
     bits, 140–141, 162, 468

Unsigned division  
     arithmetic tables, 455  
     computer, 181  
     doubleword, 197–201  
     long, 192–197  
     short from signed, 189–192

Unsigned division of integers by constants  
     best programs for, 234–235  
     by 3 and 7, 227–229  
     by divisors  $\geq 1$ , 230–232  
     by powers of 2, 227  
     incorporating into a compiler, 232–234  
     incremental division and remainder technique, 232–234  
     remainders, from powers of 2, 227  
     test for zero remainder, 248–250

*unsigned modulus* (modu) function, 84

Unsigned product of  $x$  and  $y$ , plots and graphs, 462

Uppercase letters, finding, 122

## V

Voorhies, Douglas, 373

## W

Willans, C. P., 393

Willans's formulas, 393–397

Wilson's theorem, 393, 403

Word parity. *See* Parity.

Words

counting bits, 81–87

definition, 1

division

doubleword by single word, 192–197

Knuth's Algorithm D, 184–188

multiword, 184–189

signed, multiword, 188

multiplication, multiword, 171–173

reversing, 129–134

searching for

first 0-byte, 117–121

first uppercase letter, 122

strings of 1-bits, 123–128

a value within a range, 122

word parallel operations, 13

Wormell, C. P., 397

Wormell's formula, 397–398

## Z

zbytél function, 117–121

zbyter function, 117–121

Zero means  $2n$ , 22–23