# Computer Arithmetic

## IN THIS CHAPTER

## 10-1    Introduction

Arithmetic instructions in digital computers manipulate data to produce results necessary for the solution of computational problems. These instructions perform arithmetic calculations and are responsible for the bulk of activity involved in processing data in a computer. The four basic arithmetic operations are addition, subtraction, multiplication, and division. From these four basic operations, it is possible to formulate other arithmetic functions and solve scientific problems by means of numerical analysis methods.

An arithmetic processor is the part of a processor unit that executes arithmetic operations. The data type assumed to reside in processor registers during the execution of an arithmetic instruction is specified in the definition of the instruction. An arithmetic instruction may specify binary or decimal data, and in each case the data may be in fixed-point or floating-point form. Fixed-point numbers may represent integers or fractions. Negative numbers may be in signed-magnitude or signed-complement representation. The arithmetic processor is very simple if only a binary fixed-point *add* instruction is included. It would be more complicated if it includes all four arithmetic oper-

ations for binary and decimal data in fixed-point and floating-point representation.

At an early age we are taught how to perform the basic arithmetic operations in signed-magnitude representation. This knowledge is valuable when the operations are to be implemented by hardware. However, the designer must be thoroughly familiar with the sequence of steps to be followed in order to carry out the operation and achieve a correct result. The solution to any problem that is stated by a finite number of well-defined procedural steps is called an *algorithm*. An algorithm was stated in Sec. 3-3 for the addition of two fixed-point binary numbers when negative numbers are in signed-2's complement representation. This is a simple algorithm since all it needs for its implementation is a parallel binary adder. When negative numbers are in signed-magnitude representation, the algorithm is slightly more complicated and its implementation requires circuits to add and subtract, and to compare the signs and the magnitudes of the numbers. Usually, an algorithm will contain a number of procedural steps which are dependent on results of previous steps. A convenient method for presenting algorithms is a flowchart. The computational steps are specified in the flowchart inside rectangular boxes. The decision steps are indicated inside diamond-shaped boxes from which two or more alternate paths emerge.

*algorithm*

In this chapter we develop the various arithmetic algorithms and show the procedure for implementing them with digital hardware. We consider addition, subtraction, multiplication, and division for the following types of data:

1. Fixed-point binary data in signed-magnitude representation
2. Fixed-point binary data in signed-2's complement representation
3. Floating-point binary data
4. Binary-coded decimal (BCD) data

## 10-2   Addition and Subtraction

As stated in Sec. 3-3, there are three ways of representing negative fixed-point binary numbers: signed-magnitude, signed-1's complement, or signed-2's complement. Most computers use the signed-2's complement representation when performing arithmetic operations with integers. For floating-point operations, most computers use the signed-magnitude representation for the mantissa. In this section we develop the addition and subtraction algorithms for data represented in signed-magnitude and again for data represented in signed-2's complement.

It is important to realize that the adopted representation for negative numbers refers to the representation of numbers in the registers before and

after the execution of the arithmetic operation. It does not mean that complement arithmetic may not be used in an intermediate step. For example, it is convenient to employ complement arithmetic when performing a subtraction operation with numbers in signed-magnitude representation. As long as the initial minuend and subtrahend, as well as the final difference, are in signed-magnitude form the fact that complements have been used in an intermediate step does not alter the fact that the representation is in signed-magnitude.

### Addition and Subtraction with Signed-Magnitude Data

The representation of numbers in signed-magnitude is familiar because it is used in everyday arithmetic calculations. The procedure for adding or subtracting two signed binary numbers with paper and pencil is simple and straightforward. A review of this procedure will be helpful for deriving the hardware algorithm.

We designate the magnitude of the two numbers by $A$ and $B$. When the signed numbers are added or subtracted, we find that there are eight different conditions to consider, depending on the sign of the numbers and the operation performed. These conditions are listed in the first column of Table 10-1. The other columns in the table show the actual operation to be performed with the *magnitude* of the numbers. The last column is needed to prevent a negative zero. In other words, when two equal numbers are subtracted, the result should be $+0$ not $-0$.

*magnitude*

The algorithms for addition and subtraction are derived from the table and can be stated as follows (the words inside parentheses should be used for the subtraction algorithm):

*addition (subtraction) algorithm*

*Addition (subtraction) algorithm:* when the signs of $A$ and $B$ are identical (different), add the two magnitudes and attach the sign of $A$ to the result. When the signs of $A$ and $B$ are different (identical), compare the magnitudes and

TABLE 10-1 Addition and Subtraction of Signed-Magnitude Numbers

| Operation | Add Magnitudes | Subtract Magnitudes | | |
|---|---|---|---|---|
| | | When $A > B$ | When $A < B$ | When $A = B$ |
| $(+A) + (+B)$ | $+(A + B)$ | | | |
| $(+A) + (-B)$ | | $+(A - B)$ | $-(B - A)$ | $+(A - B)$ |
| $(-A) + (+B)$ | | $-(A - B)$ | $+(B - A)$ | $+(A - B)$ |
| $(-A) + (-B)$ | $-(A + B)$ | | | |
| $(+A) - (+B)$ | | $+(A - B)$ | $-(B - A)$ | $+(A - B)$ |
| $(+A) - (-B)$ | $+(A + B)$ | | | |
| $(-A) - (+B)$ | $-(A + B)$ | | | |
| $(-A) - (-B)$ | | $-(A - B)$ | $+(B - A)$ | $+(A - B)$ |

subtract the smaller number from the larger. Choose the sign of the result to be the same as $A$ if $A > B$ or the complement of the sign of $A$ if $A < B$. If the two magnitudes are equal, subtract $B$ from $A$ and make the sign of the result positive.

The two algorithms are similar except for the sign comparison. The procedure to be followed for identical signs in the addition algorithm is the same as for different signs in the subtraction algorithm, and vice versa.

### Hardware Implementation

To implement the two arithmetic operations with hardware, it is first necessary that the two numbers be stored in registers. Let $A$ and $B$ be two registers that hold the magnitudes of the numbers, and $A_s$ and $B_s$ be two flip-flops that hold the corresponding signs. The result of the operation may be transferred to a third register: however, a saving is achieved if the result is transferred into $A$ and $A_s$. Thus $A$ and $A_s$ together form an accumulator register.

Consider now the hardware implementation of the algorithms above. First, a parallel-adder is needed to perform the microoperation $A + B$. Second, a comparator circuit is needed to establish if $A > B$, $A = B$, or $A < B$. Third, two parallel-subtractor circuits are needed to perform the microoperations $A - B$ and $B - A$. The sign relationship can be determined from an exclusive-OR gate with $A_s$ and $B_s$ as inputs.

This procedure requires a magnitude comparator, an adder, and two subtractors. However, a different procedure can be found that requires less equipment. First, we know that subtraction can be accomplished by means of complement and add. Second, the result of a comparison can be determined from the end carry after the subtraction. Careful investigation of the alternatives reveals that the use of 2's complement for subtraction and comparison is an efficient procedure that requires only an adder and a complementer.

Figure 10-1 shows a block diagram of the hardware for implementing the addition and subtraction operations. It consists of registers $A$ and $B$ and sign flip-flops $A_s$ and $B_s$. Subtraction is done by adding $A$ to the 2's complement of $B$. The output carry is transferred to flip-flop $E$, where it can be checked to determine the relative magnitudes of the two numbers. The add-overflow flip-flop $AVF$ holds the overflow bit when $A$ and $B$ are added. The $A$ register provides other microoperations that may be needed when we specify the sequence of steps in the algorithm.

The addition of $A$ plus $B$ is done through the parallel adder. The $S$ (sum) output of the adder is applied to the input of the $A$ register. The complementer provides an output of $B$ or the complement of $B$ depending on the state of the mode control $M$. The complementer consists of exclusive-OR gates and the parallel adder consists of full-adder circuits as shown in Fig. 4-7 in Chap. 4. The $M$ signal is also applied to the input carry of the adder. When $M = 0$, the output of $B$ is transferred to the adder, the input carry is 0, and the output of
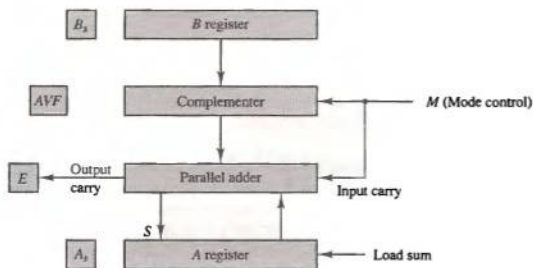
Figure 10-1  Hardware for signed-magnitude addition and subtraction.

the adder is equal to the sum $A + B$. When $M = 1$, the 1's complement of $B$ is applied to the adder, the input carry is 1, and output $S = A + \overline{B} + 1$. This is equal to $A$ plus the 2's complement of $B$, which is equivalent to the subtraction $A - B$.

## Hardware Algorithm

The flowchart for the hardware algorithm is presented in Fig. 10-2. The two signs $A_s$ and $B_s$ are compared by an exclusive-OR gate. If the output of the gate is 0, the signs are identical; if it is 1, the signs are different. For an *add* operation, identical signs dictate that the magnitudes be added. For a *subtract* operation, different signs dictate that the magnitudes be added. The magnitudes are added with a microoperation $EA \leftarrow A + B$, where $EA$ is a register that combines $E$ and $A$. The carry in $E$ after the addition constitutes an overflow if it is equal to 1. The value of $E$ is transferred into the add-overflow flip-flop $AVF$.

The two magnitudes are subtracted if the signs are different for an *add* operation or identical for a *subtract* operation. The magnitudes are subtracted by adding $A$ to the 2's complement of $B$. No overflow can occur if the numbers are subtracted so $AVF$ is cleared to 0. A 1 in $E$ indicates that $A \geq B$ and the number in $A$ is the correct result. If this number is zero, the sign $A_s$ must be made positive to avoid a negative zero. A 0 in $E$ indicates that $A < B$. For this case it is necessary to take the 2's complement of the value in $A$. This operation can be done with one microoperation $A \leftarrow \overline{A} + 1$. However, we assume that *complement and increment* the $A$ register has circuits for microoperations *complement* and *increment*, so the 2's complement is obtained from these two microoperations. In other paths of the flowchart, the sign of the result is the same as the sign of $A$, so no change in $A_s$ is required. However, when $A < B$, the sign of the result is the complement of the original sign of $A$. It is then necessary to complement $A_s$ to obtain
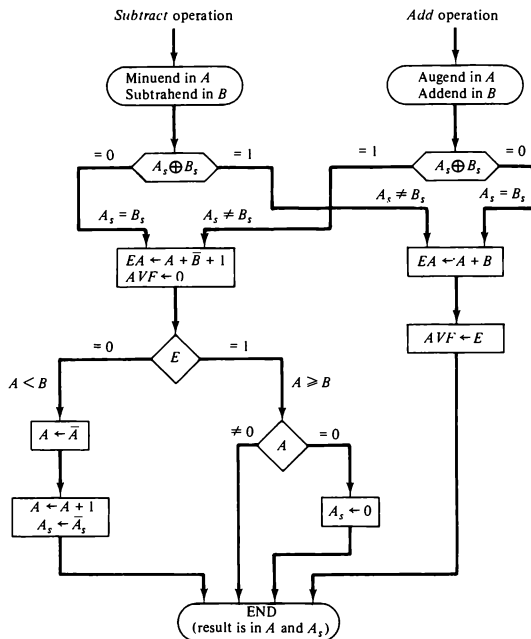
Figure 10-2 Flowchart for add and subtract operations.

the correct sign. The final result is found in register $A$ and its sign in $A_s$. The value in $AVF$ provides an overflow indication. The final value of $E$ is immaterial.

### Addition and Subtraction with Signed-2's Complement Data

The signed-2's complement representation of numbers together with arithmetic algorithms for addition and subtraction are introduced in Sec. 3-3. They are summarized here for easy reference. The leftmost bit of a binary number represents the sign bit: 0 for positive and 1 for negative. If the sign bit is 1, the entire number is represented in 2's complement form. Thus +33 is represented

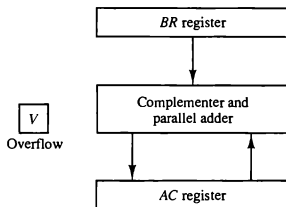as 00100001 and $-33$ as 11011111. Note that 11011111 is the 2's complement of 00100001, and vice versa.

The addition of two numbers in signed-2's complement form consists of adding the numbers with the sign bits treated the same as the other bits of the number. A carry-out of the sign-bit position is discarded. The subtraction consists of first taking the 2's complement of the subtrahend and then adding it to the minuend.

When two numbers of $n$ digits each are added and the sum occupies $n + 1$ digits, we say that an overflow occurred. The effect of an overflow on the sum of two signed-2's complement numbers is discussed in Sec. 3-3. An overflow can be detected by inspecting the last two carries out of the addition. When the two carries are applied to an exclusive-OR gate, the overflow is detected when the output of the gate is equal to 1.

The register configuration for the hardware implementation is shown in Fig. 10-3. This is the same configuration as in Fig. 10-1 except that the sign bits are not separated from the rest of the registers. We name the $A$ register $AC$ (accumulator) and the $B$ register $BR$. The leftmost bit in $AC$ and $BR$ represent the sign bits of the numbers. The two sign bits are added or subtracted together with the other bits in the complementer and parallel adder. The overflow flip-flop $V$ is set to 1 if there is an overflow. The output carry in this case is discarded.

The algorithm for adding and subtracting two binary numbers in signed-2's complement representation is shown in the flowchart of Fig. 10-4. The sum is obtained by adding the contents of $AC$ and $BR$ (including their sign bits). The overflow bit $V$ is set to 1 if the exclusive-OR of the last two carries is 1, and it is cleared to 0 otherwise. The subtraction operation is accomplished by adding the content of $AC$ to the 2's complement of $BR$. Taking the 2's complement of $BR$ has the effect of changing a positive number to negative, and vice versa. An overflow must be checked during this operation because the two numbers added could have the same sign. The programmer must realize that if an overflow occurs, there will be an erroneous result in the $AC$ register.

**Figure 10-3** Hardware for signed-2's complement addition and subtraction.
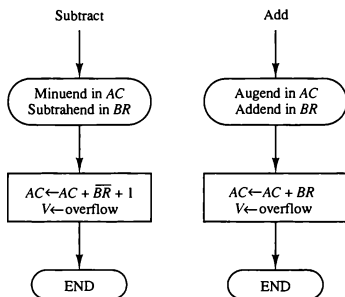
**Figure 10-4** Algorithm for adding and subtracting numbers in signed-2's complement representation.

Comparing this algorithm with its signed-magnitude counterpart, we note that it is much simpler to add and subtract numbers if negative numbers are maintained in signed-2's complement representation. For this reason most computers adopt this representation over the more familiar signed-magnitude.

## 10-3   Multiplication Algorithms

Multiplication of two fixed-point binary numbers in signed-magnitude representation is done with paper and pencil by a process of successive shift and add operations. This process is best illustrated with a numerical example.

$$
\begin{array}{rrl}
23 & 10111 & \text{Multiplicand} \\
19 & \times 10011 & \text{Multiplier} \\
\hline
 & 10111 & \\
 & 10111 & \\
 & 00000 & + \\
 & 00000 & \\
 & 10111 & \\
\hline
437 & 110110101 & \text{Product}
\end{array}
$$

The process consists of looking at successive bits of the multiplier, least significant bit first. If the multiplier bit is a 1, the multiplicand is copied down; otherwise, zeros are copied down. The numbers copied down in successive lines are shifted one position to the left from the previous number. Finally, the numbers are added and their sum forms the product.

The sign of the product is determined from the signs of the multiplicand and multiplier. If they are alike, the sign of the product is positive. If they are unlike, the sign of the product is negative.
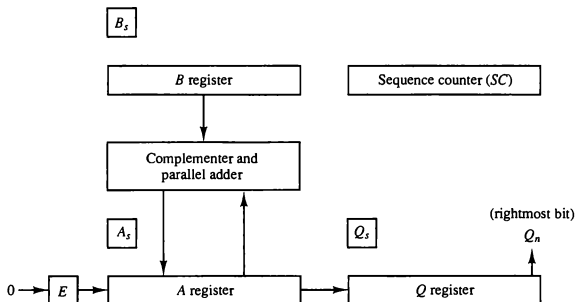
## Hardware Implementation for Signed-Magnitude Data

When multiplication is implemented in a digital computer, it is convenient to change the process slightly. First, instead of providing registers to store and add simultaneously as many binary numbers as there are bits in the multiplier, it is convenient to provide an adder for the summation of only two binary numbers and successively accumulate the partial products in a register. Second, instead of shifting the multiplicand to the left, the partial product is shifted to the right, which results in leaving the partial product and the multiplicand in the required relative positions. Third, when the corresponding bit of the multiplier is 0, there is no need to add all zeros to the partial product since it will not alter its value.

The hardware for multiplication consists of the equipment shown in Fig. 10-1 plus two more registers. These registers together with registers $A$ and $B$ are shown in Fig. 10-5. The multiplier is stored in the $Q$ register and its sign in $Q_s$. The sequence counter $SC$ is initially set to a number equal to the number of bits in the multiplier. The counter is decremented by 1 after forming each partial product. When the content of the counter reaches zero, the product is formed and the process stops.

Initially, the multiplicand is in register $B$ and the multiplier in $Q$. The sum of $A$ and $B$ forms a partial product which is transferred to the $EA$ register. Both partial product and multiplier are shifted to the right. This shift will be denoted by the statement shr $EAQ$ to designate the right shift depicted in Fig. 10-5. The

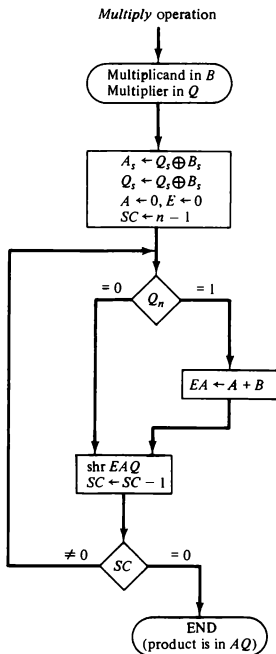Figure 10-5  Hardware for multiply operation.

least significant bit of $A$ is shifted into the most significant position of $Q$, the bit from $E$ is shifted into the most significant position of $A$, and 0 is shifted into $E$. After the shift, one bit of the partial product is shifted into $Q$, pushing the multiplier bits one position to the right. In this manner, the rightmost flip-flop in register $Q$, designated by $Q_n$, will hold the bit of the multiplier, which must be inspected next.

### Hardware Algorithm

Figure 10-6 is a flowchart of the hardware multiply algorithm. Initially, the multiplicand is in $B$ and the multiplier in $Q$. Their corresponding signs are in $B_s$ and $Q_s$, respectively. The signs are compared, and both $A$ and $Q$ are set to

**Figure 10-6**  Flowchart for multiply operation.

correspond to the sign of the product since a double-length product will be stored in registers $A$ and $Q$. Registers $A$ and $E$ are cleared and the sequence counter $SC$ is set to a number equal to the number of bits of the multiplier. We are assuming here that operands are transferred to registers from a memory unit that has words of $n$ bits. Since an operand must be stored with its sign, one bit of the word will be occupied by the sign and the magnitude will consist of $n - 1$ bits.

After the initialization, the low-order bit of the multiplier in $Q_n$ is tested. If it is a 1, the multiplicand in $B$ is added to the present partial product in $A$. If it is a 0, nothing is done. Register $EAQ$ is then shifted once to the right to form the new partial product. The sequence counter is decremented by 1 and its new value checked. If it is not equal to zero, the process is repeated and a new partial product is formed. The process stops when $SC = 0$. Note that the partial product formed in $A$ is shifted into $Q$ one bit at a time and eventually replaces the multiplier. The final product is available in both $A$ and $Q$, with $A$ holding the most significant bits and $Q$ holding the least significant bits.

The previous numerical example is repeated in Table 10-2 to clarify the hardware multiplication process. The procedure follows the steps outlined in the flowchart.

## Booth Multiplication Algorithm

Booth algorithm gives a procedure for multiplying binary integers in signed-2's complement representation. It operates on the fact that strings of 0's in the multiplier require no addition but just shifting, and a string of 1's in the multiplier from bit weight $2^k$ to weight $2^m$ can be treated as $2^{k+1} - 2^m$. For example, the binary number 001110 (+14) has a string of 1's from $2^3$ to $2^1$

TABLE 10-2 Numerical Example for Binary Multiplier

| Multiplicand $B = 10111$ | $E$ | $A$ | $Q$ | $SC$ |
|---|---|---|---|---|
| Multiplier in $Q$ | 0 | 00000 | 10011 | 101 |
| $Q_n = 1$; add $B$ | | 10111 | | |
| First partial product | 0 | 10111 | | |
| Shift right $EAQ$ | 0 | 01011 | 11001 | 100 |
| $Q_n = 1$; add $B$ | | 10111 | | |
| Second partial product | 1 | 00010 | | |
| Shift right $EAQ$ | 0 | 10001 | 01100 | 011 |
| $Q_n = 0$; shift right $EAQ$ | 0 | 01000 | 10110 | 010 |
| $Q_n = 0$; shift right $EAQ$ | 0 | 00100 | 01011 | 001 |
| $Q_n = 1$; add $B$ | | 10111 | | |
| Fifth partial product | 0 | 11011 | | |
| Shift right $EAQ$ | 0 | 01101 | 10101 | 000 |
| Final product in $AQ = 0110110101$ | | | | |

$(k = 3, m = 1)$. The number can be represented as $2^{k+1} - 2^m = 2^4 - 2^1 = 16 - 2 = 14$. Therefore, the multiplication $M \times 14$, where $M$ is the multiplicand and 14 the multiplier, can be done as $M \times 2^4 - M \times 2^1$. Thus the product can be obtained by shifting the binary multiplicand $M$ four times to the left and subtracting $M$ shifted left once.
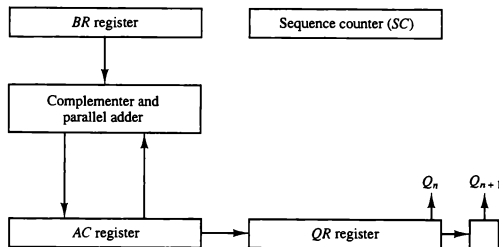
As in all multiplication schemes, Booth algorithm requires examination of the multiplier bits and shifting of the partial product. Prior to the shifting, the multiplicand may be added to the partial product, subtracted from the partial product, or left unchanged according to the following rules:

1. The multiplicand is subtracted from the partial product upon encountering the first least significant 1 in a string of 1's in the multiplier.
2. The multiplicand is added to the partial product upon encountering the first 0 (provided that there was a previous 1) in a string of 0's in the multiplier.
3. The partial product does not change when the multiplier bit is identical to the previous multiplier bit.

The algorithm works for positive or negative multipliers in 2's complement representation. This is because a negative multiplier ends with a string of 1's and the last operation will be a subtraction of the appropriate weight. For example, a multiplier equal to $-14$ is represented in 2's complement as 110010 and is treated as $-2^4 + 2^2 - 2^1 = -14$.

The hardware implementation of Booth algorithm requires the register configuration shown in Fig. 10-7. This is similar to Fig. 10-5 except that the sign bits are not separated from the rest of the registers. To show this difference, we rename registers $A$, $B$, and $Q$, as $AC$, $BR$, and $QR$, respectively. $Q_n$ designates the least significant bit of the multiplier in register $QR$. An extra flip-flop $Q_{n+1}$ is appended to $QR$ to facilitate a double bit inspection of the multiplier. The flowchart for Booth algorithm is shown in Fig. 10-8. $AC$ and the appended

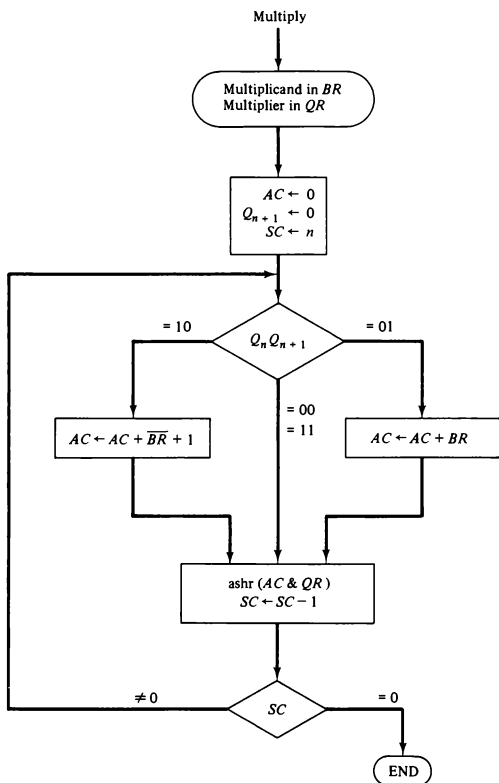Figure 10-7   Hardware for Booth algorithm.

**Figure 10-8**  Booth algorithm for multiplication of signed-2's complement numbers.

bit $Q_{n+1}$ are initially cleared to 0 and the sequence counter $SC$ is set to a number $n$ equal to the number of bits in the multiplier. The two bits of the multiplier in $Q_n$ and $Q_{n+1}$ are inspected. If the two bits are equal to 10, it means that the first 1 in a string of 1's has been encountered. This requires a subtraction of the multiplicand from the partial product in $AC$. If the two bits are equal to 01, it means that the first 0 in a string of 0's has been encountered. This requires the addition of the multiplicand to the partial product in $AC$. When the two bits are equal, the partial product does not change. An overflow cannot occur because the addition and subtraction of the multiplicand follow each other. As a consequence, the two numbers that are added always have opposite signs, a condition that excludes an overflow. The next step is to shift right the partial product and the multiplier (including bit $Q_{n+1}$). This is an arithmetic shift right (ashr) operation which shifts $AC$ and $QR$ to the right and leaves the sign bit in $AC$ unchanged (see Sec. 4-6). The sequence counter is decremented and the computational loop is repeated $n$ times.

A numerical example of Booth algorithm is shown in Table 10-3 for $n = 5$. It shows the step-by-step multiplication of $(-9) \times (-13) = +117$. Note that the multiplier in $QR$ is negative and that the multiplicand in $BR$ is also negative. The 10-bit product appears in $AC$ and $QR$ and is positive. The final value of $Q_{n+1}$ is the original sign bit of the multiplier and should not be taken as part of the product.

## Array Multiplier

Checking the bits of the multiplier one at a time and forming partial products is a sequential operation that requires a sequence of add and shift microoperations. The multiplication of two binary numbers can be done with one microoperation by means of a combinational circuit that forms the product bits all

TABLE 10-3 Example of Multiplication with Booth Algorithm

| $Q_n Q_{n+1}$ | $BR = 10111$<br>$\overline{BR} + 1 = 01001$ | $AC$ | $QR$ | $Q_{n+1}$ | $SC$ |
|---|---|---|---|---|---|
| | Initial | 00000 | 10011 | 0 | 101 |
| 1  0 | Subtract $BR$ | $\underline{01001}$ | | | |
| | | $\overline{01001}$ | | | |
| | ashr | 00100 | 11001 | 1 | 100 |
| 1  1 | ashr | 00010 | 01100 | 1 | 011 |
| 0  1 | Add $BR$ | $\underline{10111}$ | | | |
| | | $\overline{11001}$ | | | |
| | ashr | 11100 | 10110 | 0 | 010 |
| 0  0 | ashr | 11110 | 01011 | 0 | 001 |
| 1  0 | Subtract $BR$ | $\underline{01001}$ | | | |
| | | $\overline{00111}$ | | | |
| | ashr | 00011 | 10101 | 1 | 000 |

at once. This is a fast way of multiplying two numbers since all it takes is the time for the signals to propagate through the gates that form the multiplication array. However, an array multiplier requires a large number of gates, and for this reason it was not economical until the development of integrated circuits.

To see how an array multiplier can be implemented with a combinational circuit, consider the multiplication of two 2-bit numbers as shown in Fig. 10-9. The multiplicand bits are $b_1$ and $b_0$, the multiplier bits are $a_1$ and $a_0$, and the product is $c_3 c_2 c_1 c_0$. The first partial product is formed by multiplying $a_0$ by $b_1 b_0$. The multiplication of two bits such as $a_0$ and $b_0$ produces a 1 if both bits are 1; otherwise, it produces a 0. This is identical to an AND operation and can be implemented with an AND gate. As shown in the diagram, the first partial product is formed by means of two AND gates. The second partial product is formed by multiplying $a_1$ by $b_1 b_0$ and is shifted one position to the left. The two partial products are added with two half-adder (HA) circuits. Usually, there are more bits in the partial products and it will be necessary to use full-adders to produce the sum. Note that the least significant bit of the product does not have to go through an adder since it is formed by the output of the first AND gate.

A combinational circuit binary multiplier with more bits can be constructed in a similar fashion. A bit of the multiplier is ANDed with each bit of the multiplicand in as many levels as there are bits in the multiplier. The binary output in each level of AND gates is added in parallel with the partial product of the previous level to form a new partial product. The last level produces the product. For $j$ multiplier bits and $k$ multiplicand bits we need $j \times k$ AND gates and $(j - 1)$ $k$-bit adders to produce a product of $j + k$ bits.

As a second example, consider a multiplier circuit that multiplies a binary number of four bits with a number of three bits. Let the multiplicand be



Figure 10-9    2-bit by 2-bit array multiplier.

**Figure 10-10** 4-bit by 3-bit array multiplier.

represented by $b_3 b_2 b_1 b_0$ and the multiplier by $a_2 a_1 a_0$. Since $k = 4$ and $j = 3$, we need 12 AND gates and two 4-bit adders to produce a product of seven bits. The logic diagram of the multiplier is shown in Fig. 10-10.

## 10-4 Division Algorithms

Division of two fixed-point binary numbers in signed-magnitude representation is done with paper and pencil by a process of successive compare, shift, and subtract operations. Binary division is simpler than decimal division be-

cause the quotient digits are either 0 or 1 and there is no need to estimate how many times the dividend or partial remainder fits into the divisor. The division process is illustrated by a numerical example in Fig. 10-11. The divisor $B$ consists of five bits and the dividend $A$, of ten bits. The five most significant bits of the dividend are compared with the divisor. Since the 5-bit number is smaller than $B$, we try again by taking the six most significant bits of $A$ and compare this number with $B$. The 6-bit number is greater than $B$, so we place a 1 for the quotient bit in the sixth position above the dividend. The divisor is then shifted once to the right and subtracted from the dividend. The difference

*partial remainder*    is called a *partial remainder* because the division could have stopped here to obtain a quotient of 1 and a remainder equal to the partial remainder. The process is continued by comparing a partial remainder with the divisor. If the partial remainder is greater than or equal to the divisor, the quotient bit is equal to 1. The divisor is then shifted right and subtracted from the partial remainder. If the partial remainder is smaller than the divisor, the quotient bit is 0 and no subtraction is needed. The divisor is shifted once to the right in any case. Note that the result gives both a quotient and a remainder.

## Hardware Implementation for Signed-Magnitude Data

When the division is implemented in a digital computer, it is convenient to change the process slightly. Instead of shifting the divisor to the right, the dividend, or partial remainder, is shifted to the left, thus leaving the two numbers in the required relative position. Subtraction may be achieved by adding $A$ to the 2's complement of $B$. The information about the relative magnitudes is then available from the end-carry.

The hardware for implementing the division operation is identical to that required for multiplication and consists of the components shown in Fig. 10-5. Register $EAQ$ is now shifted to the left with 0 inserted into $Q_n$ and the previous value of $E$ lost. The numerical example is repeated in Fig. 10-12 to clarify the

**Figure 10-11**   Example of binary division.

| Divisor: | | 11010 | Quotient = $Q$ |
|---|---|---|---|
| $B$ = 10001 | ⟌ | 0111000000 | Dividend = $A$ |
| | | 01110 | 5 bits of $A < B$, quotient has 5 bits |
| | | 011100 | 6 bits of $A \geqslant B$ |
| | | −10001 | Shift right $B$ and subtract; enter 1 in $Q$ |
| | | −010110 | 7 bits of remainder $\geqslant B$ |
| | | −−10001 | Shift right $B$ and subtract; enter 1 in $Q$ |
| | | −−001010 | Remainder $< B$; enter 0 in $Q$; shift right $B$ |
| | | −−−010100 | Remainder $\geqslant B$ |
| | | −−−−10001 | Shift right $B$ and subtract; enter 1 in $Q$ |
| | | −−−−000110 | Remainder $< B$; enter 0 in $Q$ |
| | | −−−−−00110 | Final remainder |

Divisor $B = 10001$,          $\bar{B} + 1 = 01111$

| | $E$ | $A$ | $Q$ | $SC$ |
|---|---|---|---|---|
| Dividend: | | 01110 | 00000 | 5 |
| shl $EAQ$ | 0 | 11100 | 00000 | |
| add $\bar{B} + 1$ | | 01111 | | |
| $E = 1$ | 1 | 01011 | | |
| Set $Q_n = 1$ | 1 | 01011 | 00001 | 4 |
| shl $EAQ$ | 0 | 10110 | 00010 | |
| Add $\bar{B} + 1$ | | 01111 | | |
| $E = 1$ | 1 | 00101 | | |
| Set $Q_n = 1$ | 1 | 00101 | 00011 | 3 |
| shl $EAQ$ | 0 | 01010 | 00110 | |
| Add $\bar{B} + 1$ | | 01111 | | |
| $E = 0$; leave $Q_n = 0$ | 0 | 11001 | 00110 | |
| Add $B$ | | 10001 | | |
| Restore remainder | 1 | 01010 | | 2 |
| shl $EAQ$ | 0 | 10100 | 01100 | |
| Add $\bar{B} + 1$ | | 01111 | | |
| $E = 1$ | 1 | 00011 | | |
| Set $Q_n = 1$ | 1 | 00011 | 01101 | 1 |
| shl $EAQ$ | 0 | 00110 | 11010 | |
| Add $\bar{B} + 1$ | | 01111 | | |
| $E = 0$; leave $Q_n = 0$ | 0 | 10101 | 11010 | |
| Add $B$ | | 10001 | | |
| Restore remainder | 1 | 00110 | 11010 | 0 |
| Neglect $E$ | | | | |
| Remainder in $A$: | | 00110 | | |
| Quotient in $Q$: | | | 11010 | |

Figure 10-12 Example of binary division with digital hardware.

proposed division process. The divisor is stored in the $B$ register and the double-length dividend is stored in registers $A$ and $Q$. The dividend is shifted to the left and the divisor is subtracted by adding its 2's complement value. The information about the relative magnitude is available in $E$. If $E = 1$, it signifies that $A \geq B$. A quotient bit 1 is inserted into $Q_n$ and the partial remainder is shifted to the left to repeat the process. If $E = 0$, it signifies that $A < B$ so the quotient in $Q_n$ remains a 0 (inserted during the shift). The value of $B$ is then added to restore the partial remainder in $A$ to its previous value. The partial remainder is shifted to the left and the process is repeated again until all five quotient bits are formed. Note that while the partial remainder is shifted left, the quotient bits are shifted also and after five shifts, the quotient is in $Q$ and the final remainder is in $A$.

Before showing the algorithm in flowchart form, we have to consider the sign of the result and a possible overflow condition. The sign of the quotient is determined from the signs of the dividend and the divisor. If the two signs

are alike, the sign of the quotient is plus. If they are unalike, the sign is minus. The sign of the remainder is the same as the sign of the dividend.

## Divide Overflow

The division operation may result in a quotient with an overflow. This is not a problem when working with paper and pencil but is critical when the operation is implemented with hardware. This is because the length of registers is finite and will not hold a number that exceeds the standard length. To see this, consider a system that has 5-bit registers. We use one register to hold the divisor and two registers to hold the dividend. From the example of Fig. 10-11 we note that the quotient will consist of six bits if the five most significant bits of the dividend constitute a number greater than the divisor. The quotient is to be stored in a standard 5-bit register, so the overflow bit will require one more flip-flop for storing the sixth bit. This divide-overflow condition must be avoided in normal computer operations because the entire quotient will be too long for transfer into a memory unit that has words of standard length, that is, the same as the length of registers. Provisions to ensure that this condition is detected must be included in either the hardware or the software of the computer, or in a combination of the two.
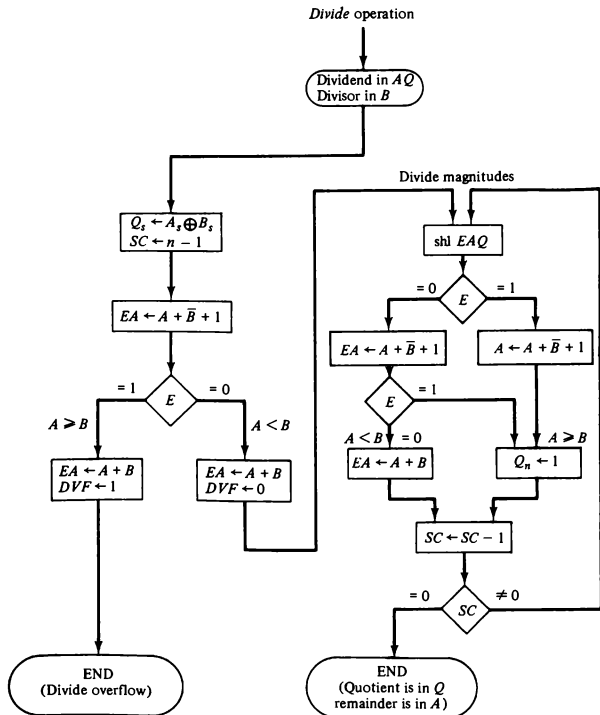
When the dividend is twice as long as the divisor, the condition for overflow can be stated as follows: A divide-overflow condition occurs if the high-order half bits of the dividend constitute a number greater than or equal to the divisor. Another problem associated with division is the fact that a division by zero must be avoided. The divide-overflow condition takes care of this condition as well. This occurs because any dividend will be greater than or equal to a divisor which is equal to zero. Overflow condition is usually detected when a special flip-flop is set. We will call it a divide-overflow flip-flop and label it $DVF$.

The occurrence of a divide overflow can be handled in a variety of ways. In some computers it is the responsibility of the programmers to check if $DVF$ is set after each divide instruction. They then can branch to a subroutine that takes a corrective measure such as rescaling the data to avoid overflow. In some older computers, the occurrence of a divide overflow stopped the computer and this condition was referred to as a *divide stop*. Stopping the operation of the computer is not recommended because it is time consuming. The procedure in most computers is to provide an interrupt request when $DVF$ is set. The interrupt causes the computer to suspend the current program and branch to a service routine to take a corrective measure. The most common corrective measure is to remove the program and type an error message explaining the reason why the program could not be completed. It is then the responsibility of the user who wrote the program to rescale the data or take any other corrective measure. The best way to avoid a divide overflow is to use floating-point data. We will see in Sec. 10-5 that a divide overflow can be handled very simply if numbers are in floating-point representation.

### Hardware Algorithm

The hardware divide algorithm is shown in the flowchart of Fig. 10-13. The dividend is in $A$ and $Q$ and the divisor in $B$. The sign of the result is transferred into $Q_s$ to be part of the quotient. A constant is set into the sequence counter $SC$ to specify the number of bits in the quotient. As in multiplication, we assume that operands are transferred to registers from a memory unit that has

**Figure 10-13** Flowchart for divide operation.

words of $n$ bits. Since an operand must be stored with its sign, one bit of the word will be occupied by the sign and the magnitude will consist of $n-1$ bits.

A divide-overflow condition is tested by subtracting the divisor in $B$ from half of the bits of the dividend stored in $A$. If $A \geq B$, the divide-overflow flip-flop $DVF$ is set and the operation is terminated prematurely. If $A < B$, no divide overflow occurs so the value of the dividend is restored by adding $B$ to $A$.

The division of the magnitudes starts by shifting the dividend in $AQ$ to the left with the high-order bit shifted into $E$. If the bit shifted into $E$ is 1, we know that $EA > B$ because $EA$ consists of a 1 followed by $n-1$ bits while $B$ consists of only $n-1$ bits. In this case, $B$ must be subtracted from $EA$ and 1 inserted into $Q_n$ for the quotient bit. Since register $A$ is missing the high-order bit of the dividend (which is in $E$), its value is $EA - 2^{n-1}$. Adding to this value the 2's complement of $B$ results in

$$(EA - 2^{n-1}) + (2^{n-1} - B) = EA - B$$

The carry from this addition is not transferred to $E$ if we want $E$ to remain a 1.

If the shift-left operation inserts a 0 into $E$, the divisor is subtracted by adding its 2's complement value and the carry is transferred into $E$. If $E = 1$, it signifies that $A \geq B$; therefore, $Q_n$ is set to 1. If $E = 0$, it signifies that $A < B$ and the original number is restored by adding $B$ to $A$. In the latter case we leave a 0 in $Q_n$ (0 was inserted during the shift).

This process is repeated again with register $A$ holding the partial remainder. After $n - 1$ times, the quotient magnitude is formed in register $Q$ and the remainder is found in register $A$. The quotient sign is in $Q_s$ and the sign of the remainder in $A_s$ is the same as the original sign of the dividend.

### Other Algorithms

*restoring method*

The hardware method just described is called the *restoring method*. The reason for this name is that the partial remainder is restored by adding the divisor to the negative difference. Two other methods are available for dividing numbers, the *comparison* method and the *nonrestoring* method. In the comparison method $A$ and $B$ are compared *prior* to the subtraction operation. Then if $A \geq B$, $B$ is subtracted from $A$. If $A < B$ nothing is done. The partial remainder is shifted left and the numbers are compared again. The comparison can be determined prior to the subtraction by inspecting the end-carry out of the parallel-adder prior to its transfer to register $E$.

*comparison and nonrestoring method*

In the nonrestoring method, $B$ is not added if the difference is negative but instead, the negative difference is shifted left and then $B$ is added. To see why this is possible consider the case when $A < B$. From the flowchart in Fig. 9-11 we note that the operations performed are $A - B + B$; that is, $B$ is sub-

tracted and then added to restore $A$. The next time around the loop, this number is shifted left (or multiplied by 2) and $B$ subtracted again. This gives $2(A - B + B) - B = 2A - B$. This result is obtained in the nonrestoring method by leaving $A - B$ as is. The next time around the loop, the number is shifted left and $B$ *added* to give $2(A - B) + B = 2A - B$, which is the same as before. Thus, in the nonrestoring method, $B$ is subtracted if the previous value of $Q_n$ was a 1, but $B$ is added if the previous value of $Q_n$ was a 0 and no restoring of the partial remainder is required. This process saves the step of adding the divisor if $A$ is less than $B$, but it requires special control logic to remember the previous result. The first time the dividend is shifted, $B$ must be subtracted. Also, if the last bit of the quotient is 0, the partial remainder must be restored to obtain the correct final remainder.

## 10-5   Floating-Point Arithmetic Operations

Many high-level programming languages have a facility for specifying floating-point numbers. The most common way is to specify them by a *real* declaration statement as opposed to fixed-point numbers, which are specified by an *integer* *integer declaration* declaration statement. Any computer that has a compiler for such high-level *statement* programming language must have a provision for handling floating-point arithmetic operations. The operations are quite often included in the internal hardware. If no hardware is available for the operations, the compiler must be designed with a package of floating-point software subroutines. Although the hardware method is more expensive, it is so much more efficient than the software method that floating-point hardware is included in most computers and is omitted only in very small ones.

### Basic Considerations

Floating-point representation of data was introduced in Sec. 3-4. A floating-point number in computer registers consists of two parts: a mantissa $m$ and an exponent $e$. The two parts represent a number obtained from multiplying $m$ times a radix $r$ raised to the value of $e$; thus

$$m \times r^e$$

The mantissa may be a fraction or an integer. The location of the radix point and the value of the radix $r$ are assumed and are not included in the registers. For example, assume a fraction representation and a radix 10. The decimal number 537.25 is represented in a register with $m = 53725$ and $e = 3$ and is interpreted to represent the floating-point number

$$.53725 \times 10^3$$

A floating-point number is normalized if the most significant digit of the mantissa is nonzero. In this way the mantissa contains the maximum possible number of significant digits. A zero cannot be normalized because it does not have a nonzero digit. It is represented in floating-point by all 0's in the mantissa and exponent.

Floating-point representation increases the range of numbers that can be accommodated in a given register. Consider a computer with 48-bit words. Since one bit must be reserved for the sign, the range of fixed-point integer numbers will be $\pm(2^{47} - 1)$, which is approximately $\pm10^{14}$. The 48 bits can be used to represent a floating-point number with 36 bits for the mantissa and 12 bits for the exponent. Assuming fraction representation for the mantissa and taking the two sign bits into consideration, the range of numbers that can be accommodated is

$$\pm(1 - 2^{-35}) \times 2^{2047}$$

This number is derived from a fraction that contains 35 1's, an exponent of 11 bits (excluding its sign), and the fact that $2^{11} - 1 = 2047$. The largest number that can be accommodated is approximately $10^{615}$, which is a very large number. The mantissa can accommodate 35 bits (excluding the sign) and if considered as an integer it can store a number as large as $(2^{35} - 1)$. This is approximately equal to $10^{10}$, which is equivalent to a decimal number of 10 digits.

Computers with shorter word lengths use two or more words to represent a floating-point number. An 8-bit microcomputer may use four words to represent one floating-point number. One word of 8 bits is reserved for the exponent and the 24 bits of the other three words are used for the mantissa.

Arithmetic operations with floating-point numbers are more complicated than with fixed-point numbers and their execution takes longer and requires more complex hardware. Adding or subtracting two numbers requires first an alignment of the radix point since the exponent parts must be made equal before adding or subtracting the mantissas. The alignment is done by shifting one mantissa while its exponent is adjusted until it is equal to the other exponent. Consider the sum of the following floating-point numbers:

$$.5372400 \times 10^2$$
$$+ \ .1580000 \times 10^{-1}$$

It is necessary that the two exponents be equal before the mantissas can be added. We can either shift the first number three positions to the left, or shift the second number three positions to the right. When the mantissas are stored in registers, shifting to the left causes a loss of most significant digits. Shifting to the right causes a loss of least significant digits. The second method is preferable because it only reduces the accuracy, while the first method may cause an error. The usual alignment procedure is to shift the mantissa that has

the smaller exponent to the right by a number of places equal to the difference between the exponents. After this is done, the mantissas can be added:

$$\begin{array}{r} .5372400 \times 10^2 \\ +.0001580 \times 10^2 \\ \hline .5373980 \times 10^2 \end{array}$$

When two normalized mantissas are added, the sum may contain an overflow digit. An overflow can be corrected easily by shifting the sum once to the right and incrementing the exponent. When two numbers are subtracted, the result may contain most significant zeros as shown in the following example:

$$\begin{array}{r} .56780 \times 10^5 \\ -.56430 \times 10^5 \\ \hline .00350 \times 10^5 \end{array}$$

A floating-point number that has a 0 in the most significant position of the mantissa is said to have an *underflow*. To normalize a number that contains an underflow, it is necessary to shift the mantissa to the left and decrement the exponent until a nonzero digit appears in the first position. In the example above, it is necessary to shift left twice to obtain $.35000 \times 10^3$. In most computers, a normalization procedure is performed after each operation to ensure that all results are in a normalized form.

Floating-point multiplication and division do not require an alignment of the mantissas. The product can be formed by multiplying the two mantissas and adding the exponents. Division is accomplished by dividing the mantissas and subtracting the exponents.

The operations performed with the mantissas are the same as in fixed-point numbers, so the two can share the same registers and circuits. The operations performed with the exponents are compare and increment (for aligning the mantissas), add and subtract (for multiplication and division), and decrement (to normalize the result). The exponent may be represented in any one of the three representations: signed-magnitude, signed-2's complement, or signed-1's complement.

A fourth representation employed in many computers is known as a *biased* exponent. In this representation, the sign bit is removed from being a separate entity. The bias is a positive number that is added to each exponent as the floating-point number is formed, so that internally all exponents are positive. The following example may clarify this type of representation. Consider an exponent that ranges from $-50$ to 49. Internally, it is represented by two digits (without a sign) by adding to it a bias of 50. The exponent register contains the number $e + 50$, where $e$ is the actual exponent. This way, the exponents are represented in registers as positive numbers in the range of 00

to 99. Positive exponents in registers have the range of numbers from 99 to 50. The subtraction of 50 gives the positive values from 49 to 0. Negative exponents are represented in registers in the range from 49 to 00. The subtraction of 50 gives the negative values in the range of −1 to −50.

The advantage of biased exponents is that they contain only positive numbers. It is then simpler to compare their relative magnitude without being concerned with their signs. As a consequence, a magnitude comparator can be used to compare their relative magnitude during the alignment of the mantissa. Another advantage is that the smallest possible biased exponent contains all zeros. The floating-point representation of zero is then a zero mantissa and the smallest possible exponent.

In the examples above, we used decimal numbers to demonstrate some of the concepts that must be understood when dealing with floating-point numbers. Obviously, the same concepts apply to binary numbers as well. The algorithms developed in this section are for binary numbers. Decimal computer arithmetic is discussed in the next section.
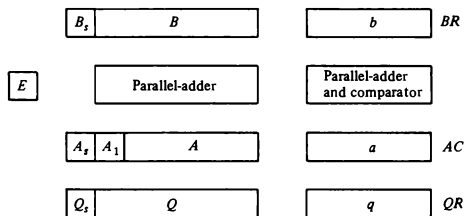
## Register Configuration

The register configuration for floating-point operations is quite similar to the layout for fixed-point operations. As a general rule, the same registers and adder used for fixed-point arithmetic are used for processing the mantissas. The difference lies in the way the exponents are handled.

The register organization for floating-point operations is shown in Fig. 10-14. There are three registers, $BR$, $AC$, and $QR$. Each register is subdivided into two parts. The mantissa part has the same uppercase letter symbols as in fixed-point representation. The exponent part uses the corresponding lower-case letter symbol.

It is assumed that each floating-point number has a mantissa in signed-magnitude representation and a biased exponent. Thus the $AC$ has a mantissa

Figure 10-14   Registers for floating-point arithmetic operations.

whose sign is in $A_s$ and a magnitude that is in $A$. The exponent is in the part of the register denoted by the lowercase letter symbol $a$. The diagram shows explicitly the most significant bit of $A$, labeled by $A_1$. The bit in this position must be a 1 for the number to be normalized. Note that the symbol $AC$ represents the entire register, that is, the concatenation of $A_s$, $A$, and $a$.

Similarly, register $BR$ is subdivided into $B_s$, $B$, and $b$, and $QR$ into $Q_s$, $Q$, and $q$. A parallel-adder adds the two mantissas and transfers the sum into $A$ and the carry into $E$. A separate parallel-adder is used for the exponents. Since the exponents are biased, they do not have a distinct sign bit but are represented as a biased positive quantity. It is assumed that the floating-point numbers are so large that the chance of an exponent overflow is very remote, and for this reason the exponent overflow will be neglected. The exponents are also connected to a magnitude comparator that provides three binary outputs to indicate their relative magnitude.

The number in the mantissa will be taken as a *fraction*, so the binary point is assumed to reside to the left of the magnitude part. Integer representation for floating-point causes certain scaling problems during multiplication and division. To avoid these problems, we adopt a fraction representation.

The numbers in the registers are assumed to be initially normalized. After each arithmetic operation, the result will be normalized. Thus all floating-point operands coming from and going to the memory unit are always normalized.

### Addition and Subtraction

During addition or subtraction, the two floating-point operands are in $AC$ and $BR$. The sum or difference is formed in the $AC$. The algorithm can be divided into four consecutive parts:

1. Check for zeros.
2. Align the mantissas.
3. Add or subtract the mantissas.
4. Normalize the result.

A floating-point number that is zero cannot be normalized. If this number is used during the computation, the result may also be zero. Instead of checking for zeros during the normalization process we check for zeros at the beginning and terminate the process if necessary. The alignment of the mantissas must be carried out prior to their operation. After the mantissas are added or subtracted, the result may be unnormalized. The normalization procedure ensures that the result is normalized prior to its transfer to memory.

The flowchart for adding or subtracting two floating-point binary numbers is shown in Fig. 10-15. If $BR$ is equal to zero, the operation is terminated, with the value in the $AC$ being the result. If $AC$ is equal to zero, we transfer

**Figure 10-15**   Addition and subtraction of floating-point numbers.

the content of $BR$ into $AC$ and also complement its sign if the numbers are to be subtracted. If neither number is equal to zero, we proceed to align the mantissas.

The magnitude comparator attached to exponents $a$ and $b$ provides three outputs that indicate their relative magnitude. If the two exponents are equal, we go to perform the arithmetic operation. If the exponents are not equal, the mantissa having the smaller exponent is shifted to the right and its exponent incremented. This process is repeated until the two exponents are equal.

The addition and subtraction of the two mantissas is identical to the fixed-point addition and subtraction algorithm presented in Fig. 10-2. The magnitude part is added or subtracted depending on the operation and the signs of the two mantissas. If an overflow occurs when the magnitudes are added, it is transferred into flip-flop $E$. If $E$ is equal to 1, the bit is transferred into $A_1$ and all other bits of $A$ are shifted right. The exponent must be incremented to maintain the correct number. No underflow may occur in this case because the original mantissa that was not shifted during the alignment was already in a normalized position.

If the magnitudes were subtracted, the result may be zero or may have an underflow. If the mantissa is zero, the entire floating-point number in the $AC$ is made zero. Otherwise, the mantissa must have at least one bit that is equal to 1. The mantissa has an underflow if the most significant bit in position $A_1$ is 0. In that case, the mantissa is shifted left and the exponent decremented. The bit in $A_1$ is checked again and the process is repeated until it is equal to 1. When $A_1 = 1$, the mantissa is normalized and the operation is completed.

### Multiplication

The multiplication of two floating-point numbers requires that we multiply the mantissas and add the exponents. No comparison of exponents or alignment of mantissas is necessary. The multiplication of the mantissas is performed in the same way as in fixed-point to provide a double-precision product. The double-precision answer is used in fixed-point numbers to increase the accuracy of the product. In floating-point, the range of a single-precision mantissa combined with the exponent is usually accurate enough so that only single-precision numbers are maintained. Thus the half most significant bits of the mantissa product and the exponent will be taken together to form a single-precision floating-point product.
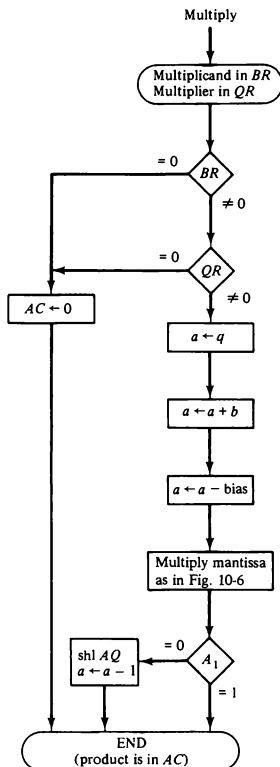
The multiplication algorithm can be subdivided into four parts:

1. Check for zeros.
2. Add the exponents.
3. Multiply the mantissas.
4. Normalize the product.

Steps 2 and 3 can be done simultaneously if separate adders are available for the mantissas and exponents.

The flowchart for floating-point multiplication is shown in Fig. 10-16. The two operands are checked to determine if they contain a zero. If either operand is equal to zero, the product in the $AC$ is set to zero and the operation is

**Figure 10-16** Multiplication of floating-point numbers.

terminated. If neither of the operands is equal to zero, the process continues with the exponent addition.

The exponent of the multiplier is in $q$ and the adder is between exponents $a$ and $b$. It is necessary to transfer the exponents from $q$ to $a$, add the two exponents, and transfer the sum into $a$. Since both exponents are biased by the addition of a constant, the exponent sum will have double this bias. The correct biased exponent for the product is obtained by subtracting the bias number from the sum.

The multiplication of the mantissas is done as in the fixed-point case with the product residing in $A$ and $Q$. Overflow cannot occur during multiplication, so there is no need to check for it.

The product may have an underflow, so the most significant bit in $A$ is checked. If it is a 1, the product is already normalized. If it is a 0, the mantissa in $AQ$ is shifted left and the exponent decremented. Note that only one normalization shift is necessary. The multiplier and multiplicand were originally normalized and contained fractions. The smallest normalized operand is 0.1, so the smallest possible product is 0.01. Therefore, only one leading zero may occur.

Although the low-order half of the mantissa is in $Q$, we do not use it for the floating-point product. Only the value in the $AC$ is taken as the product.

## Division

Floating-point division requires that the exponents be subtracted and the mantissas divided. The mantissa division is done as in fixed-point except that the dividend has a single-precision mantissa that is placed in the $AC$. Remember that the mantissa dividend is a fraction and not an integer. For integer representation, a single-precision dividend must be placed in register $Q$ and register $A$ must be cleared. The zeros in $A$ are to the left of the binary point and have no significance. In fraction representation, a single-precision dividend is placed in register $A$ and register $Q$ is cleared. The zeros in $Q$ are to the right of the binary point and have no significance.

The check for divide-overflow is the same as in fixed-point representation. However, with floating-point numbers the divide-overflow imposes no problems. If the dividend is greater than or equal to the divisor, the dividend fraction is shifted to the right and its exponent incremented by 1. For normalized operands this is a sufficient operation to ensure that no mantissa divide-*dividend alignment*    overflow will occur. The operation above is referred to as a *dividend alignment*.

The division of two normalized floating-point numbers will always result in a normalized quotient provided that a dividend alignment is carried out before the division. Therefore, unlike the other operations, the quotient obtained after the division does not require a normalization.

The division algorithm can be subdivided into five parts:

1. Check for zeros.
2. Initialize registers and evaluate the sign.

3. Align the dividend.
4. Subtract the exponents.
5. Divide the mantissas.

The flowchart for floating-point division is shown in Fig. 10-17. The two operands are checked for zero. If the divisor is zero, it indicates an attempt to divide by zero, which is an illegal operation. The operation is terminated with an error message. An alternative procedure would be to set the quotient in $QR$ to the most positive number possible (if the dividend is positive) or to the most negative possible (if the dividend is negative). If the dividend in $AC$ is zero, the quotient in $QR$ is made zero and the operation terminates.

If the operands are not zero, we proceed to determine the sign of the quotient and store it in $Q_s$. The sign of the dividend in $A_s$ is left unchanged to be the sign of the remainder. The $Q$ register is cleared and the sequence counter $SC$ is set to a number equal to the number of bits in the quotient.

The dividend alignment is similar to the divide-overflow check in the fixed-point operation. The proper alignment requires that the fraction dividend be smaller than the divisor. The two fractions are compared by a subtraction test. The carry in $E$ determines their relative magnitude. The dividend fraction is restored to its original value by adding the divisor. If $A \geq B$, it is necessary to shift $A$ once to the right and increment the dividend exponent. Since both operands are normalized, this alignment ensures that $A < B$.

Next, the divisor exponent is subtracted from the dividend exponent. Since both exponents were originally biased, the subtraction operation gives the difference without the bias. The bias is then added and the result transferred into $q$ because the quotient is formed in $QR$.

The magnitudes of the mantissas are divided as in the fixed-point case. After the operation, the mantissa quotient resides in $Q$ and the remainder in $A$. The floating-point quotient is already normalized and resides in $QR$. The exponent of the remainder should be the same as the exponent of the dividend. The binary point for the remainder mantissa lies $(n - 1)$ positions to the left of $A_1$. The remainder can be converted to a normalized fraction by subtracting $n - 1$ from the dividend exponent and by shift and decrement until the bit in $A_1$ is equal to 1. This is not shown in the flow chart and is left as an exercise.

## 10-6    Decimal Arithmetic Unit

The user of a computer prepares data with decimal numbers and receives results in decimal form. A CPU with an arithmetic logic unit can perform arithmetic microoperations with binary data. To perform arithmetic operations with decimal data, it is necessary to convert the input decimal numbers to binary, to perform all calculations with binary numbers, and to convert the results into decimal. This may be an efficient method in applications requiring a large number of calculations and a relatively smaller amount of input and
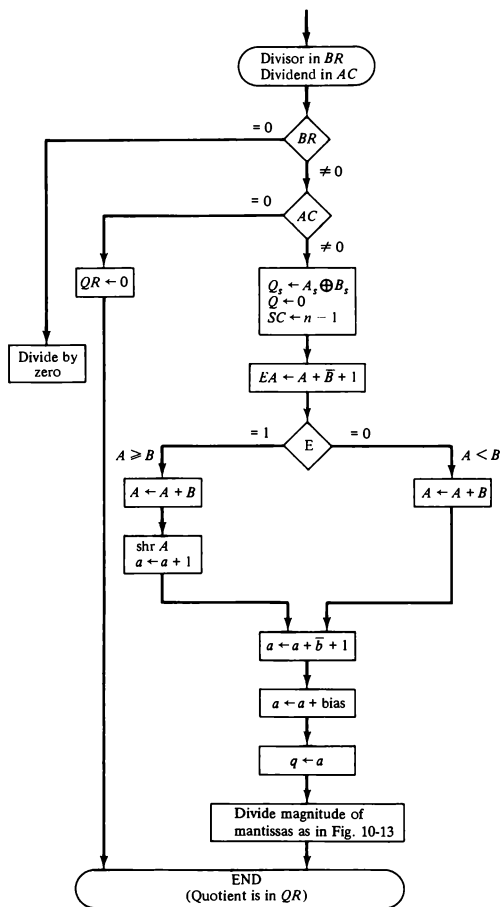
**Figure 10-17** Division of floating-point numbers.

output data. When the application calls for a large amount of input–output and a relatively smaller number of arithmetic calculations, it becomes convenient to do the internal arithmetic directly with the decimal numbers. Computers capable of performing decimal arithmetic must store the decimal data in binary-coded form. The decimal numbers are then applied to a decimal arithmetic unit capable of executing decimal arithmetic microoperations.

Electronic calculators invariably use an internal decimal arithmetic unit since inputs and outputs are frequent. There does not seem to be a reason for converting the keyboard input numbers to binary and again converting the displayed results to decimal, since this process requires special circuits and also takes a longer time to execute. Many computers have hardware for arithmetic calculations with both binary and decimal data. Users can specify by pro-grammed instructions whether they want the computer to perform calculations with binary or decimal data.

A decimal arithmetic unit is a digital function that performs decimal microoperations. It can add or subtract decimal numbers, usually by forming the 9's or 10's complement of the subtrahend. The unit accepts coded decimal numbers and generates results in the same adopted binary code. A single-stage decimal arithmetic unit consists of nine binary input variables and five binary output variables, since a minimum of four bits is required to represent each coded decimal digit. Each stage must have four inputs for the augend digit, four inputs for the addend digit, and an input-carry. The outputs include four terminals for the sum digit and one for the output-carry. Of course, there is a wide variety of possible circuit configurations dependent on the code used to represent the decimal digits.

## BCD Adder

Consider the arithmetic addition of two decimal digits in BCD, together with a possible carry from a previous stage. Since each input digit does not exceed 9, the output sum cannot be greater than $9 + 9 + 1 = 19$, the 1 in the sum being an input-carry. Suppose that we apply two BCD digits to a 4-bit binary adder. The adder will form the sum in *binary* and produce a result that may range from 0 to 19. These binary numbers are listed in Table 10-4 and are labeled by symbols $K$, $Z_8$, $Z_4$, $Z_2$, and $Z_1$. $K$ is the carry and the subscripts under the letter $Z$ represent the weights 8, 4, 2, and 1 that can be assigned to the four bits in the BCD code. The first column in the table lists the binary sums as they appear in the outputs of a 4-bit *binary* adder. The output sum of two *decimal* numbers must be represented in BCD and should appear in the form listed in the second column of the table. The problem is to find a simple rule by which the binary number in the first column can be converted to the correct BCD digit representation of the number in the second column.

In examining the contents of the table, it is apparent that when the binary sum is equal to or less than 1001, the corresponding BCD number is identical

TABLE 10-4 Derivation of BCD Adder

| Binary Sum | | | | | BCD Sum | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| $K$ | $Z_8$ | $Z_4$ | $Z_2$ | $Z_1$ | $C$ | $S_8$ | $S_4$ | $S_2$ | $S_1$ | Decimal |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 |
| 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 2 |
| 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 3 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 4 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 5 |
| 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 6 |
| 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 7 |
| 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 8 |
| 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 9 |
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 10 |
| 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 11 |
| 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 12 |
| 0 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 13 |
| 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 14 |
| 0 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 15 |
| 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 16 |
| 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 17 |
| 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 18 |
| 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 19 |

and therefore no conversion is needed. When the binary sum is greater than 1001, we obtain a nonvalid BCD representation. The addition of binary 6 (0110) to the binary sum converts it to the correct BCD representation and also produces an output-carry as required.

One method of adding decimal numbers in BCD would be to employ one 4-bit binary adder and perform the arithmetic operation one digit at a time. The low-order pair of BCD digits is first added to produce a binary sum. If the result is equal or greater than 1010, it is corrected by adding 0110 to the binary sum. This second operation will automatically produce an output-carry for the next pair of significant digits. The next higher-order pair of digits, together with the input-carry, is then added to produce their binary sum. If this result is equal to or greater than 1010, it is corrected by adding 0110. The procedure is repeated until all decimal digits are added.

The logic circuit that detects the necessary correction can be derived from the table entries. It is obvious that a correction is needed when the binary sum has an output carry $K = 1$. The other six combinations from 1010 to 1111 that need a correction have a 1 in position $Z_8$. To distinguish them from binary 1000 and 1001 which also have a 1 in position $Z_8$, we specify further that either $Z_4$
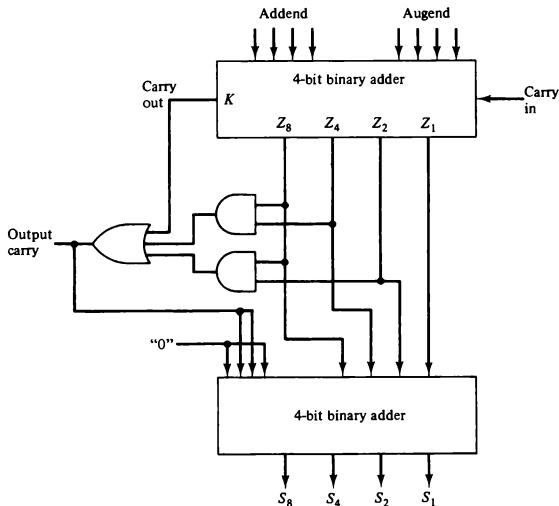
or $Z_2$ must have a 1. The condition for a correction and an output-carry can be expressed by the Boolean function

$$C = K + Z_8 Z_4 + Z_8 Z_2$$

When $C = 1$, it is necessary to add 0110 to the binary sum and provide an output-carry for the next stage.

A BCD adder is a circuit that adds two BCD digits in parallel and produces a sum digit also in BCD. A BCD adder must include the correction logic in its internal construction. To add 0110 to the binary sum, we use a second 4-bit binary adder as shown in Fig. 10-18. The two decimal digits, together with the input-carry, are first added in the top 4-bit binary adder to produce the binary sum. When the output-carry is equal to 0, nothing is added to the binary sum. When it is equal to 1, binary 0110 is added to the binary sum through the bottom 4-bit binary adder. The output-carry generated from the bottom binary adder may be ignored, since it supplies information already available in the output-carry terminal.

**Figure 10-18**    Block diagram of BCD adder.

A decimal parallel-adder that adds $n$ decimal digits needs $n$ BCD adder stages with the output-carry from one stage connected to the input-carry of the next-higher-order stage. To achieve shorter propagation delays, BCD adders include the necessary circuits for carry look-ahead. Furthermore, the adder circuit for the correction does not need all four full-adders, and this circuit can be optimized.

## BCD Subtraction

A straight subtraction of two decimal numbers will require a subtractor circuit that will be somewhat different from a BCD adder. It is more economical to perform the subtraction by taking the 9's or 10's complement of the subtrahend and adding it to the minuend. Since the BCD is not a self-complementing code, the 9's complement cannot be obtained by complementing each bit in the code. It must be formed by a circuit that subtracts each BCD digit from 9.

The 9's complement of a decimal digit represented in BCD may be obtained by complementing the bits in the coded representation of the digit provided a correction is included. There are two possible correction methods. In the first method, binary 1010 (decimal 10) is added to each complemented digit and the carry discarded after each addition. In the second method, binary 0110 (decimal 6) is added before the digit is complemented. As a numerical illustration, the 9's complement of BCD 0111 (decimal 7) is computed by first complementing each bit to obtain 1000. Adding binary 1010 and discarding the carry, we obtain 0010 (decimal 2). By the second method, we add 0110 to 0111 to obtain 1101. Complementing each bit, we obtain the required result of 0010. Complementing each bit of a 4-bit binary number $N$ is identical to the subtraction of the number from 1111 (decimal 15). Adding the binary equivalent of decimal 10 gives $15 - N + 10 = 9 - N + 16$. But 16 signifies the carry that is discarded, so the result is $9 - N$ as required. Adding the binary equivalent of decimal 6 and then complementing gives $15 - (N + 6) = 9 - N$ as required.

The 9's complement of a BCD digit can also be obtained through a combinational circuit. When this circuit is attached to a BCD adder, the result is a BCD adder/subtractor. Let the subtrahend (or addend) digit be denoted by the four binary variables $B_8$, $B_4$, $B_2$, and $B_1$. Let $M$ be a mode bit that controls the add/subtract operation. When $M = 0$, the two digits are added; when $M = 1$, the digits are subtracted. Let the binary variables $x_8$, $x_4$, $x_2$, and $x_1$ be the outputs of the 9's complementer circuit. By an examination of the truth table for the circuit, it may be observed (see Prob. 10-30) that $B_1$ should always be complemented; $B_2$ is always the same in the 9's complement as in the original digit; $x_4$ is 1 when the exclusive-OR of $B_2$ and $B_4$ is 1; and $x_8$ is 1 when $B_8 B_4 B_2 = 000$. The Boolean functions for the 9's complementer circuit are

$$x_1 = B_1 M' + B_1' M$$

$$x_2 = B_2$$

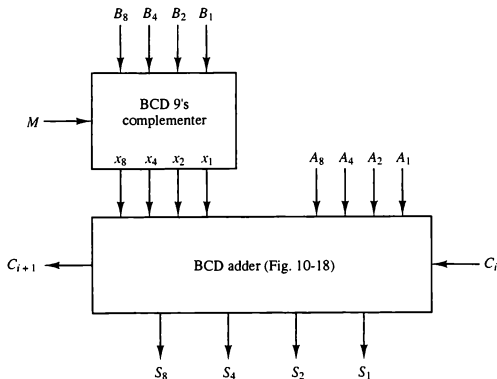$$x_4 = B_4 M' + (B_4' B_2 + B_4 B_2')M$$

$$x_8 = B_8 M' + B_8' B_4' B_2' M$$

From these equations we see that $x = B$ when $M = 0$. When $M = 1$, the $x$ outputs produce the 9's complement of $B$.

One stage of a decimal arithmetic unit that can add or subtract two BCD digits is shown in Fig. 10-19. It consists of a BCD adder and a 9's complementer. The mode $M$ controls the operation of the unit. With $M = 0$, the $S$ outputs form the sum of $A$ and $B$. With $M = 1$, the $S$ outputs form the sum of $A$ plus the 9's complement of $B$. For numbers with $n$ decimal digits we need $n$ such stages. The output carry $C_{i+1}$ from one stage must be connected to the input carry $C_i$ of the next-higher-order stage. The best way to subtract the two decimal numbers is to let $M = 1$ and apply a 1 to the input carry $C_1$ of the first stage. The outputs will form the sum of $A$ plus the 10's complement of $B$, which is equivalent to a subtraction operation if the carry-out of the last stage is discarded.

## 10-7  Decimal Arithmetic Operations

The algorithms for arithmetic operations with decimal data are similar to the algorithms for the corresponding operations with binary data. In fact, except for a slight modification in the multiplication and division algorithms, the same

**Figure 10-19**  One stage of a decimal arithmetic unit.

flowcharts can be used for both types of data provided that we interpret the microoperation symbols properly. Decimal numbers in BCD are stored in computer registers in groups of four bits. Each 4-bit group represents a decimal digit and must be taken as a unit when performing decimal microoperations.

For convenience, we will use the same symbols for binary and decimal arithmetic microoperations but give them a different interpretation. As shown in Table 10-5, a bar over the register letter symbol denotes the 9's complement of the decimal number stored in the register. Adding 1 to the 9's complement produces the 10's complement. Thus, for decimal numbers, the symbol $A \leftarrow A + \bar{B} + 1$ denotes a transfer of the decimal sum formed by adding the original content $A$ to the 10's complement of $B$. The use of identical symbols for the 9's complement and the 1's complement may be confusing if both types of data are employed in the same system. If this is the case, it may be better to adopt a different symbol for the 9's complement. If only one type of data is being considered, the symbol would apply to the type of data used.

Incrementing or decrementing a register is the same for binary and decimal except for the number of states that the register is allowed to have. A binary counter goes through 16 states, from 0000 to 1111, when incremented. A decimal counter goes through 10 states from 0000 to 1001 and back to 0000, since 9 is the last count. Similarly, a binary counter sequences from 1111 to 0000 when decremented. A decimal counter goes from 1001 to 0000.

A decimal shift right or left is preceded by the letter $d$ to indicate a shift over the four bits that hold the decimal digits. As a numerical illustration consider a register $A$ holding decimal 7860 in BCD. The bit pattern of the 12 flip-flops is

$$0111 \quad 1000 \quad 0110 \quad 0000$$

The microoperation *dshr A* shifts the decimal number one digit to the right to give 0786. This shift is over the four bits and changes the content of the register into

$$0000 \quad 0111 \quad 1000 \quad 0110$$

**TABLE 10-5** Decimal Arithmetic Microoperation Symbols

| Symbolic Designation | Description |
|---|---|
| $A \leftarrow A + B$ | Add decimal numbers and transfer sum into $A$ |
| $\bar{B}$ | 9's complement of $B$ |
| $A \leftarrow A + \bar{B} + 1$ | Content of $A$ plus 10's complement of $B$ into $A$ |
| $Q_L \leftarrow Q_L + 1$ | Increment BCD number in $Q_L$ |
| dshr $A$ | Decimal shift-right register $A$ |
| dshl $A$ | Decimal shift-left register $A$ |

## Addition and Subtraction

The algorithm for addition and subtraction of binary signed-magnitude numbers applies also to decimal signed-magnitude numbers provided that we interpret the microoperation symbols in the proper manner. Similarly, the algorithm for binary signed-2's complement numbers applies to decimal signed-10's complement numbers. The binary data must employ a binary adder and a complementer. The decimal data must employ a decimal arithmetic unit capable of adding two BCD numbers and forming the 9's complement of the subtrahend as shown in Fig. 10-19.
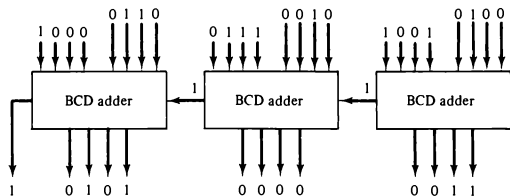
Decimal data can be added in three different ways, as shown in Fig. 10-20. The parallel method uses a decimal arithmetic unit composed of as many BCD adders as there are digits in the number. The sum is formed in parallel and requires only one microoperation. In the digit-serial bit-parallel method, the digits are applied to a single BCD adder serially, while the bits of each coded digit are transferred in parallel. The sum is formed by shifting the decimal numbers through the BCD adder one at a time. For $k$ decimal digits, this configuration requires $k$ microoperations, one for each decimal shift. In the all serial adder, the bits are shifted one at a time through a full-adder. The binary sum formed after four shifts must be corrected into a valid BCD digit. This correction, discussed in Sec. 10-6, consists of checking the binary sum. If it is greater than or equal to 1010, the binary sum is corrected by adding to it 0110 and generating a carry for the next pair of digits.

The parallel method is fast but requires a large number of adders. The digit-serial bit-parallel method requires only one BCD adder, which is shared by all the digits. It is slower than the parallel method because of the time required to shift the digits. The all serial method requires a minimum amount of equipment but is very slow.
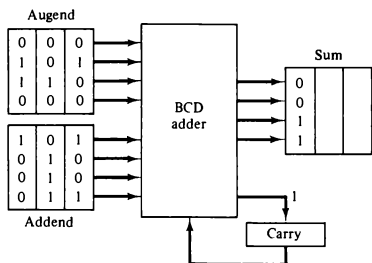
## Multiplication

The multiplication of fixed-point decimal numbers is similar to binary except for the way the partial products are formed. A decimal multiplier has digits that range in value from 0 to 9, whereas a binary multiplier has only 0 and 1 digits. In the binary case, the multiplicand is added to the partial product if the multiplier bit is 1. In the decimal case, the multiplicand must be multiplied by the digit multiplier and the result added to the partial product. This operation can be accomplished by adding the multiplicand to the partial product a number of times equal to the value of the multiplier digit.
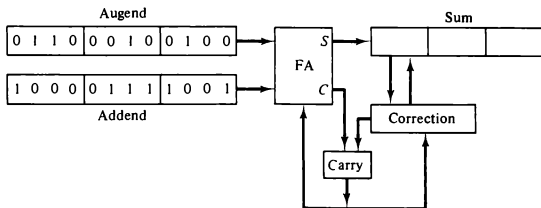
The registers organization for the decimal multiplication is shown in Fig. 10-21. We are assuming here four-digit numbers, with each digit occupying four bits, for a total of 16 bits for each number. There are three registers, $A$, $B$, and $Q$, each having a corresponding sign flip-flop $A_s$, $B_s$, and $Q_s$.

(a) Parallel decimal addition: 624 + 879 = 1503



(b) Digit-serial, bit-parallel decimal addition



(c) All serial decimal addition
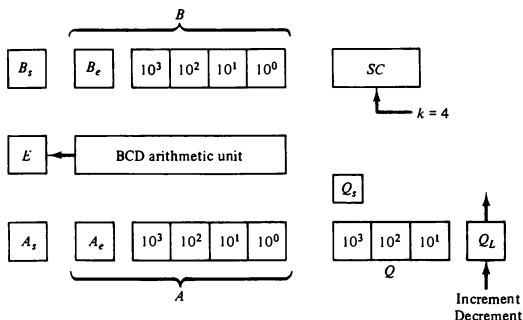
Figure 10-20 Three ways of adding decimal numbers.

Figure 10-21    Registers for decimal arithmetic multiplication and division.

Registers $A$ and $B$ have four more bits designated by $A_e$ and $B_e$ that provide an extension of one more digit to the registers. The BCD arithmetic unit adds the five digits in parallel and places the sum in the five-digit $A$ register. The end-carry goes to flip-flop $E$. The purpose of digit $A_e$ is to accommodate an overflow while adding the multiplicand to the partial product during multiplication. The purpose of digit $B_e$ is to form the 9's complement of the divisor when subtracted from the partial remainder during the division operation. The least significant digit in register $Q$ is denoted by $Q_L$. This digit can be incremented or decremented.

A decimal operand coming from memory consists of 17 bits. One bit (the sign) is transferred to $B_s$ and the magnitude of the operand is placed in the lower 16 bits of $B$. Both $B_e$ and $A_e$ are cleared initially. The result of the operation is also 17 bits long and does not use the $A_e$ part of the $A$ register.

The decimal multiplication algorithm is shown in Fig. 10-22. Initially, the entire $A$ register and $B_e$ are cleared and the sequence counter $SC$ is set to a number $k$ equal to the number of digits in the multiplier. The low-order digit of the multiplier in $Q_L$ is checked. If it is not equal to 0, the multiplicand in $B$ is added to the partial product in $A$ once and $Q_L$ is decremented. $Q_L$ is checked again and the process is repeated until it is equal to 0. In this way, the multiplicand in $B$ is added to the partial product a number of times equal to the multiplier digit. Any temporary overflow digit will reside in $A_e$ and can range in value from 0 to 9.

Next, the partial product and the multiplier are shifted once to the right. This places zero in $A_e$ and transfers the next multiplier quotient into $Q_L$. The process is then repeated $k$ times to form a double-length product in $AQ$.
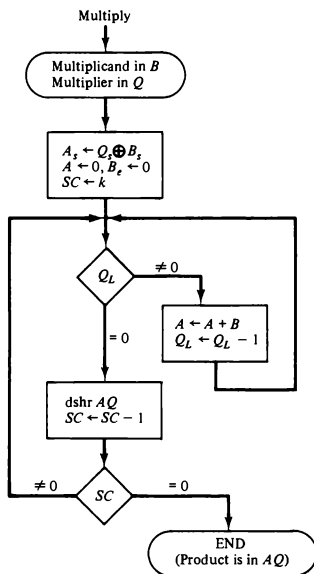
Figure 10-22 Flowchart for decimal multiplication.

## Division

Decimal division is similar to binary division except of course that the quotient digits may have any of the 10 values from 0 to 9. In the restoring division method, the divisor is subtracted from the dividend or partial remainder as many times as necessary until a negative remainder results. The correct remainder is then restored by adding the divisor. The digit in the quotient reflects the number of subtractions up to but excluding the one that caused the negative difference.

The decimal division algorithm is shown in Fig. 10-23. It is similar to the algorithm with binary data except for the way the quotient bits are formed. The dividend (or partial remainder) is shifted to the left, with its most significant digit placed in $A_e$. The divisor is then subtracted by adding its 10's complement value. Since $B_e$ is initially cleared, its complement value is 9 as required. The carry in $E$ determines the relative magnitude of $A$ and $B$. If $E = 0$, it signifies
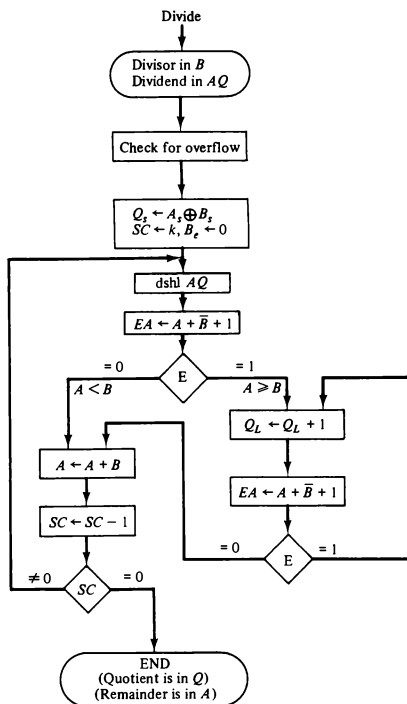
**Figure 10-23**   Flowchart for decimal division.

that $A < B$. In this case the divisor is added to restore the partial remainder and $Q_L$ stays at 0 (inserted there during the shift). If $E = 1$, it signifies that $A \geq B$. The quotient digit in $Q_L$ is incremented once and the divisor subtracted again. This process is repeated until the subtraction results in a negative difference which is recognized by $E$ being 0. When this occurs, the quotient digit is not incremented but the divisor is added to restore the positive remainder. In this way, the quotient digit is made equal to the number of times that the partial remainder "goes" into the divisor.

The partial remainder and the quotient bits are shifted once to the left and the process is repeated $k$ times to form $k$ quotient digits. The remainder is then found in register $A$ and the quotient is in register $Q$. The value of $E$ is neglected.

### Floating-Point Operations

Decimal floating-point arithmetic operations follow the same procedures as binary operations. The algorithms in Sec. 10-5 can be adopted for decimal data provided that the microoperation symbols are interpreted correctly. The multiplication and division of the mantissas must be done by the methods described above.

---

<div align="center">PROBLEMS</div>

---

**10-1.** The complementer shown in Fig. 10-1 is not needed if instead of performing $A + \overline{B} + 1$ we perform $B + \overline{A}$ ($B$ plus the 1's complement of $A$). Derive an algorithm in flowchart form for addition and subtraction of fixed-point binary numbers in signed-magnitude representation with the magnitudes subtracted by the two microoperations $A \leftarrow \overline{A}$ and $EA \leftarrow A + B$.

**10-2.** Mark each individual path in the flowchart of Fig. 10-2 by a number and then indicate the overall path that the algorithm takes when the following signed-magnitude numbers are computed. In each case give the value of $AVF$. The leftmost bit in the following numbers represents the sign bit.
   **a.** 0 101101 + 0 011111
   **b.** 1 011111 + 1 101101
   **c.** 0 101101 − 0 011111
   **d.** 0 101101 − 0 101101
   **e.** 1 011111 − 0 101101

**10-3.** Perform the arithmetic operations below with binary numbers and with negative numbers in signed-2's complement representation. Use seven bits to accommodate each number together with its sign. In each case, determine if there is an overflow by checking the carries into and out of the sign bit position.
   **a.** $(+35) + (+40)$
   **b.** $(-35) + (-40)$
   **c.** $(-35) + (+40)$

**10-4.** Consider the binary numbers when they are in signed-2's complement representation. Each number has $n$ bits: one for the sign and $k = n - 1$ for the magnitude. A negative number $-X$ is represented as $2^k + (2^k - X)$, where the first $2^k$ designates the sign bit and $(2^k - X)$ is the 2's complement of $X$. A positive number is represented as $0 + X$, where the 0 designates the sign bit, and $X$, the $k$-bit magnitude. Using these generalized symbols, prove

that the sum $(\pm X) + (\pm Y)$ can be formed by adding the numbers including their sign bits and discarding the carry-out of the sign-bit position. In other words, prove the algorithm for adding two binary numbers in signed-2's complement representation.

**10-5.** Formulate a hardware procedure for detecting an overflow by comparing the sign of the sum with the signs of the augend and addend. The numbers are in signed-2's complement representation.

**10-6.** **a.** Perform the operation $(-9) + (-6) = -15$ with binary numbers in signed-1's complement representation using only five bits to represent each number (including the sign). Show that the overflow detection procedure of checking the inequality of the last two carries fails in this case.

**b.** Suggest a modified procedure for detecting an overflow when signed-1's complement numbers are used.

**10-7.** Derive an algorithm in flowchart form for adding and subtracting two fixed-point binary numbers when negative numbers are in signed-1's complement representation.

**10-8.** Prove that the multiplication of two $n$-digit numbers in base $r$ gives a product no more than $2n$ digits in length. Show that this statement implies that no overflow can occur in the multiplication operation.

**10-9.** Show the contents of registers $E$, $A$, $Q$, and $SC$ (as in Table 10-2) during the process of multiplication of two binary numbers, 11111 (multiplicand) and 10101 (multiplier). The signs are not included.

**10-10.** Show the contents of registers $E$, $A$, $Q$, and $SC$ (as in Fig. 10-12) during the process of division of (a) 10100011 by 1011; (b) 00001111 by 0011. (Use a dividend of eight bits.)

**10-11.** Show that adding $B$ after the operation $A + \bar{B} + 1$ restores the original value of $A$. What should be done with the end carry?

**10-12.** Why should the sign of the remainder after a division be the same as the sign of the dividend?

**10-13.** Design an array multiplier that multiplies two 4-bit numbers. Use AND gates and binary adders.

**10-14.** Show the step-by-step multiplication process using Booth algorithm (as in Table 10-3) when the following binary numbers are multiplied. Assume 5-bit registers that hold signed numbers. The multiplicand in both cases is $+15$.

**a.** $(+15) \times (+13)$
**b.** $(+15) \times (-13)$

**10-15.** Derive an algorithm in flowchart form for the nonrestoring method of fixed-point binary division.

**10-16.** Derive an algorithm for evaluating the square root of a binary fixed-point number.

**10-17.** A binary floating-point number has seven bits for a biased exponent. The constant used for the bias is 64.

**a.** List the biased representation of all exponents from $-64$ to $+63$.

b. Show that a 7-bit magnitude comparator can be used to compare the relative magnitude of the two exponents.

c. Show that after the addition of two biased exponents it is necessary to subtract 64 in order to have a biased exponents sum. How would you subtract 64 by adding its 2's complement value?

d. Show that after the subtraction of two biased exponents it is necessary to add 64 in order to have a biased exponent difference.

**10-18.** Derive an algorithm in flowchart form for the comparison of two signed binary numbers when negative numbers are in signed-2's complement representation:

a. By means of a subtraction operation with the signed-2's complement numbers.

b. By scanning and comparing pairs of bits from left to right.

**10-19.** Repeat Prob. 10-18 for signed-magnitude binary numbers.

**10-20.** Let $n$ be the number of bits of the mantissa in a binary floating-point number. When the mantissas are aligned during the addition or subtraction, the exponent difference may be greater than $n - 1$. If this occurs, the mantissa with the smaller exponent is shifted entirely out of the register. Modify the mantissa alignment in Fig. 10-15 by including a sequence counter $SC$ that counts the number of shifts. If the number of shifts is greater than $n - 1$, the larger number is then used to determine the result.

**10-21.** The procedure for aligning mantissas during addition or subtraction of floating-point numbers can be stated as follows: Subtract the smaller exponent from the larger and shift right the mantissa having the smaller exponent a number of places equal to the difference between the exponents. The exponent of the sum (or difference) is equal to the larger exponent. Without using a magnitude comparator, assuming biased exponents, and taking into account that only the $AC$ can be shifted, derive an algorithm in flowchart form for aligning the mantissas and placing the larger exponent in the $AC$.

**10-22.** Show that there can be no mantissa overflow after a multiplication operation.

**10-23.** Show that the division of two normalized floating-point numbers with fractional mantissas will always result in a normalized quotient provided a dividend alignment is carried out prior to the division operation.

**10-24.** Extend the flowchart of Fig. 10-17 to provide a normalized floating-point remainder in the $AC$. The mantissa should be a fraction.

**10-25.** The algorithms for the floating-point arithmetic operations in Sec. 10-5 neglect the possibility of exponent overflow or underflow.

a. Go over the three flowcharts and find where an exponent overflow may occur.

b. Repeat (a) for exponent underflow. An exponent underflow occurs if the exponent is more negative than the smallest number that can be accommodated in the register.

c. Show how an exponent overflow or underflow can be detected by the hardware.

**10-26.** If we assume integer representation for the mantissa of floating-point numbers, we encounter certain scaling problems during multiplication and divi-

sion. Let the number of bits in the magnitude part of the mantissa be $(n - 1)$. For integer representation:

    **a.** Show that if a single-precision product is used, $(n - 1)$ must be added to the exponent product in the $AC$.

    **b.** Show that if a single-precision mantissa dividend is used, $(n - 1)$ must be subtracted from the exponent dividend when $Q$ is cleared.

**10-27.** Show the hardware to be used for the addition and subtraction of two decimal numbers in signed-magnitude representation. Indicate how an overflow is detected.

**10-28.** Show that $673 - 356$ can be computed by adding 673 to the 10's complement of 356 and discarding the end carry. Draw the block diagram of a three-stage decimal arithmetic unit and show how this operation is implemented. List all input bits and output bits of the unit.

**10-29.** Show that the lower 4-bit binary adder in Fig. 10-1 can be replaced by one full-adder and two half-adders.

**10-30.** Using combinational circuit design techniques, derive the Boolean functions for the BCD 9's complementer of Fig. 10-19. Draw the logic diagram.

**10-31.** It is necessary to design an adder for two decimal digits represented in the excess-3 code (Table 3-6). Show that the correction after adding two digits with a 4-bit binary adder is as follows:

    **a.** The output carry is equal to the uncorrected carry.

    **b.** If output carry = 1, add 0011.

    **c.** If output carry = 0, add 1101 and ignore the carry from this addition. Show that the excess-3 adder can be constructed with seven full-adders and two inverters.

**10-32.** Derive the circuit for a 9's complementer when decimal digits are represented in the excess-3 code (Table 3-6). A mode control input determines whether the digit is complemented or not. What is the advantage of using this code over BCD?

**10-33.** Show the hardware to be used for the addition and subtraction of two decimal numbers with negative numbers in signed-10's complement representation. Indicate how an overflow is detected. Derive the flowchart algorithm and try a few numbers to convince yourself that the algorithm produces correct results.

**10-34.** Show the content of registers $A$, $B$, $Q$, and $SC$ during the decimal multiplication (Fig. 10-22) of (a) $470 \times 152$ and (b) $999 \times 199$. Assume three-digit registers and take the second number as the multiplier.

**10-35.** Show the content of registers $A$, $E$, $Q$, and $SC$ during the decimal division (Fig. 10-23) of 1680/32. Assume two-digit registers.

**10-36.** Show that subregister $A_e$ in Fig. 10-21 is zero at the termination of (a) the decimal multiplication as specified in Fig. 10-22, and (b) the decimal division as specified in Fig. 10-23.

**10-37.** Change the floating-point arithmetic algorithms in Sec. 10-5 from binary to decimal data. In a table, list how each microoperation symbol should be interpreted.

## REFERENCES

1. Blaauw, G., *Digital Systems Implementation*. Englewood Cliffs, NJ: Prentice Hall, 1976.

2. Cavanagh, J. J. F., *Digital Computer Arithmetic*. New York: McGraw-Hill, 1984.

3. Hamacher, V. C., Z. G. Vranesic, and S. G. Zaky, *Computer Organization*, 3rd ed. New York: McGraw-Hill, 1990.

4. Hays, J. F., *Computer Architecture and Organization*, 2nd ed. New York: McGraw-Hill, 1988.

5. Hill, F. J., and G. R. Peterson, *Digital Systems: Hardware Organization and Design*, 3rd ed. New York: John Wiley, 1987.

6. Hwang, K., *Computer Arithmetic*. New York: John Wiley, 1979.

7. Kulisch, V. W., and W. L. Miranker, *Computer Arithmetic in Theory and Practice*. New York: Academic Press, 1980.

8. Schmid, H., *Decimal Arithmetic*. New York: John Wiley, 1979.