

(23CS304) COMPUTER ORGANIZATION & ARCHITECTURE

UNIT-III

DATA REPRESENTATION: Data types, Complements, Fixed-Point Representation, Floating-Point Representation.

COMPUTER ARITHMETIC: Addition and Subtraction, Multiplication Algorithms, Division Algorithms, Floating-point Arithmetic operations, Decimal Arithmetic unit, Decimal Arithmetic operations.

DATA REPRESENTATION

Data types: Binary information in digital computers is stored in memory or processor registers. Registers contain either data or control information. Control information is a bit or a group of bits used to specify the sequence of command signals needed for manipulation of the data in other registers. Data are numbers and other binary-coded information that are operated on, to achieve required computational results.

The data types found in the registers of digital computers may be classified as being one of the following categories:

- (1) numbers used in arithmetic computations,
- (2) letters of the alphabet used in data processing, and
- (3) other discrete symbols used for specific purposes.

All types of data, except binary numbers, are represented in computer registers in binary-coded form. This is because registers are made up of flip-flops and flip-flops are two-state devices that can store only 1's and 0's. The binary number system is the most natural system to be used in a digital computer. But sometimes it is convenient to employ different number systems, especially the decimal number system, since it is used by people to perform arithmetic computations.

A number system of base, or radix, r is a system that uses distinct symbols for r digits, numbers are represented by a string of digit symbols. The string of digits 724.5 represents the quantity

$$7 \times 10^2 + 2 \times 10^1 + 4 \times 10^0 + 5 \times 10^{-1}$$

The string of digits 101101 in the binary number system represents the quantity

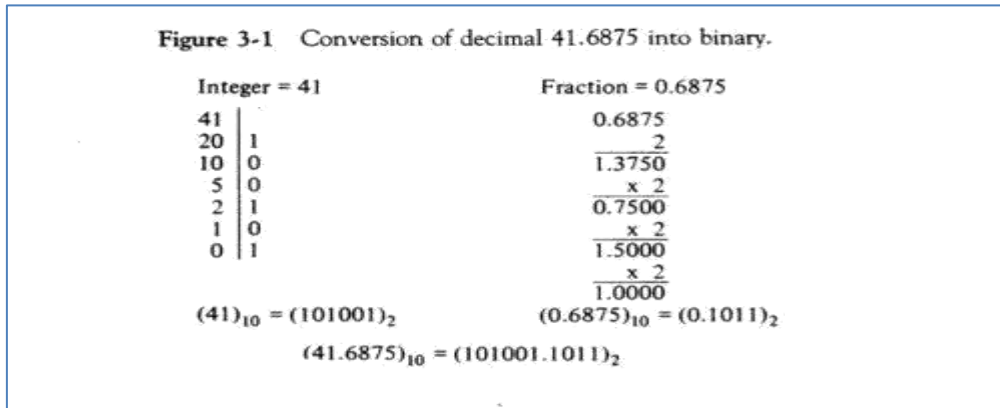
$$1 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 = 45$$
$$(101101)_2 = (45)_{10}$$

We will also use the octal (radix 8) and hexadecimal (radix 16) number systems

$$(736.4)_8 = 7 \times 8^2 + 3 \times 8^1 + 6 \times 8^0 + 4 \times 8^{-1} = (478.5)_{10}$$
$$(F3)_{16} = F \times 16^1 + 3 \times 16^0 = (243)_{10}$$

Conversion from decimal to radix r system is carried out by separating the number into its integer and fraction parts and converting each part separately.

- Divide the integer successively by r and accumulate the remainders
- Multiply the fraction successively by r until the fraction becomes zero



Each octal digit corresponds to three binary digits, each hexadecimal digit corresponds to four binary digits.

Complements:

- Complements are used in digital computers for simplifying subtraction and logical manipulation
- Two types of complements for each base r system: r 's complement and $(r - 1)$'s complement
- Given a number N in base r having n digits, the $(r - 1)$'s complement of N is defined as $(r^n - 1) - N$
- For decimal, the 9's complement of N is $(10^n - 1) - N$
- The 9's complement of 546700 is $999999 - 546700 = 453299$
- The 9's complement of 453299 is $999999 - 453299 = 546700$
- For binary, the 1's complement of N is $(2^n - 1) - N$
- The 1's complement of 1011001 is $1111111 - 1011001 = 0100110$
- The 1's complement is the true complement of the number – just toggle all bits
- The r 's complement of an n -digit number N in base r is defined as $r^n - N$
- This is the same as adding 1 to the $(r - 1)$'s complement
- The 10's complement of 2389 is $7610 + 1 = 7611$
- The 2's complement of 101100 is $010011 + 1 = 010100$
- Subtraction of unsigned n -digit numbers: $M - N$
 - Add M to the r 's complement of N – this results in $M + (r^n - N) = M - N + r^n$
 - If $M \geq N$, the sum will produce an end carry r^n which is discarded o If $M < N$, the sum does not produce an end carry and is equal to $r^n - (N - M)$, which is the r 's complement of $(N - M)$. To obtain the answer in a familiar form, take the r 's complement of the sum and place a negative sign in front.

Example: $72532 - 13250 = 59282$. The 10's complement of 13250 is 86750.

M	= 72352
10's comp. of N	= +86750
Sum	= 159282
Discard end carry	= -100000

Answer	= 59282
--------	---------

Example for $M < N$: $13250 - 72532 = -59282$

M	= 13250
10's comp. of N	= +27468
Sum	= 40718
No end carry	
Answer	= -59282 (10's comp. of 40718)

Example for $X = 1010100$ and $Y = 1000011$

X	= 1010100
2's comp. of Y	= +0111101
Sum	= 10010001
Discard end carry	= -10000000
Answer $X - Y$	= 0010001

Y	= 1000011
2's comp. of X	= +0101100
Sum	= 1101111

Fixed Point Representation:

- Positive integers and zero can be represented by unsigned numbers
- Negative numbers must be represented by signed numbers since + and – signs are not available, only 1's and 0's are
- Signed numbers have msb as 0 for positive and 1 for negative – msb is the sign bit
- Two ways to designate binary point position in a register
 - Fixed point position
 - Floating-point representation
- Fixed point position usually uses one of the two following positions
 - A binary point in the extreme left of the register to make it a fraction
 - A binary point in the extreme right of the register to make it an integer
 - In both cases, a binary point is not actually present
- The floating-point representations uses a second register to designate the position of the binary point in the first register
- When an integer is positive, the msb, or sign bit, is 0 and the remaining bits represent the magnitude
- When an integer is negative, the msb, or sign bit, is 1, but the rest of the number can be represented in one of three ways
 - Signed-magnitude representation
 - Signed-1's complement representation
 - Signed-2's complement representation
- Consider an 8-bit register and the number +14
 - The only way to represent it is 00001110
- Consider an 8-bit register and the number –14
 - Signed magnitude: 1 0001110
 - Signed 1's complement: 1 1110001
 - Signed 2's complement: 1 1110010
- Typically use signed 2's complement

- Addition of two signed-magnitude numbers follow the normal rules
 - If same signs, add the two magnitudes and use the common sign
 - Differing signs, subtract the smaller from the larger and use the sign of the larger magnitude
 - Must compare the signs and magnitudes and then either add or subtract
- Addition of two signed 2's complement numbers does not require a comparison or subtraction – only addition and complementation
 - Add the two numbers, including their sign bits
 - Discard any carry out of the sign bit position
 - All negative numbers must be in the 2's complement form
 - If the sum obtained is negative, then it is in 2's complement form

+6	00000110	-6	11111010
<u>+13</u>	<u>00001101</u>	<u>+13</u>	<u>00001101</u>
+19	00010011	+7	00000111
+6	00000110	-6	11111010
<u>-13</u>	<u>11110011</u>	<u>-13</u>	<u>11110011</u>
-7	11111001	-19	11101101

- Subtraction of two signed 2's complement numbers is as follows
 - Take the 2's complement form of the subtrahend (including sign bit)
Add it to the minuend (including the sign bit)
 - A carry out of the sign bit position is discarded
- An *overflow* occurs when two numbers of n digits each are added and the sum occupies $n + 1$ digits
- Overflows are problems since the width of a register is finite
- Therefore, a flag is set if this occurs and can be checked by the user
- Detection of an overflow depends on if the numbers are signed or unsigned
- For unsigned numbers, an overflow is detected from the end carry out of the msb
- For addition of signed numbers, an overflow cannot occur if one is positive and one is negative – both have to have the same sign
- An overflow can be detected if the carry into the sign bit position and the carry out of the sign bit position are not equal

+70	0 1000110	-70	1 0111010
<u>+80</u>	<u>0 1010000</u>	<u>-80</u>	<u>1 0110000</u>
+150	1 0010110	-150	0 1101010

- The representation of decimal numbers in registers is a function of the binary code used to represent a decimal digit
- A 4-bit decimal code requires four flip-flops for each decimal digit
- This takes much more space than the equivalent binary representation and the circuits required to perform decimal arithmetic are more complex
- Representation of signed decimal numbers in BCD is similar to the representation of signed numbers in binary

- Either signed magnitude or signed complement systems
- The sign of a number is represented with four bits
 - 0000 for +
 - 1001 for –
- To obtain the 10's complement of a BCD number, first take the 9's complement and then add one to the least significant digit
- Example: $(+375) + (-240) = +135$

□

0 375	(0000 0011 0111 1010)BCD
+9 760	(1001 0111 0110 0000)BC
0 135	(0000 0001 0011 0101) $\frac{D}{BC}$

Floating-Point Representation

- The floating-point representation of a number has two parts
- The first part represents a signed, fixed-point number – the *mantissa*
- The second part designates the position of the binary point – the *exponent*
- The mantissa may be a fraction or an integer
- Example: the decimal number +6132.789 is
 - Fraction: +0.6123789
 - Exponent: +04
 - Equivalent to $+0.6123789 \times 10^{+4}$
- A floating-point number is always interpreted to represent $m \times r^e$
- Example: the binary number +1001.11 (with 8-bit fraction and 6-bit exponent)
 - Fraction: 01001110
 - Exponent: 000100
 - Equivalent to $+(.1001110)_2 \times 2^{+4}$
- A floating-point number is said to be *normalized* if the most significant digit of the mantissa is nonzero
- The decimal number 350 is normalized, 00350 is not
- The 8-bit number 00011010 is not normalized
- Normalize it by fraction = 11010000 and exponent = -3
- Normalized numbers provide the maximum possible precision for the floating-point number

Other Binary Codes

- Digital systems can process data in discrete form only
- Continuous, or analog, information is converted into digital form by means of an analog-to-digital converter
- The reflected binary or *Gray code*, is sometimes used for the converted digital data
- The Gray code changes by only one bit as it sequences from one number to the next
- Gray code counters are sometimes used to provide the timing sequences that control the operations in a digital system

TABLE 3-6 Four Different Binary Codes for the Decimal Digit

Decimal digit	BCD 8421	2421	Excess-3	Excess-3 gray
0	0000	0000	0011	0010
1	0001	0001	0100	0110
2	0010	0010	0101	0111
3	0011	0011	0110	0101
4	0100	0100	0111	0100
5	0101	1011	1000	1100
6	0110	1100	1001	1101
7	0111	1101	1010	1111
8	1000	1110	1011	1110
9	1001	1111	1100	1010
Unused bit combinations	1010	0101	0000	0000
	1011	0110	0001	0001
	1100	0111	0010	0011
	1101	1000	1101	1000
	1110	1001	1110	1001

- Binary codes for decimal digits require a minimum of four bits
- Other codes besides BCD exist to represent decimal digits
- The 2421 code and the excess-3 code are both *self-complementing*
- The 9's complement of each digit is obtained by complementing each bit in the code
- The 2421 code is a *weighted code*
- The bits are multiplied by indicated weights and the sum gives the decimal digit
- The excess-3 code is obtained from the corresponding BCD code added to 3

Error Detection Codes

- Transmitted binary information is subject to noise that could change bits 1 to 0 and vice versa
- An *error detection code* is a binary code that detects digital errors during transmission
- The detected errors cannot be corrected, but can prompt the data to be retransmitted
- The most common error detection code used is the *parity bit*

TABLE 3-5 4-Bit Gray Code

Binary code	Decimal equivalent	Binary code	Decimal equivalent
0000	0	1100	8
0001	1	1101	9
0011	2	1111	10
0010	3	1110	11
0110	4	1010	12
0111	5	1011	13
0101	6	1001	14
0100	7	1000	15

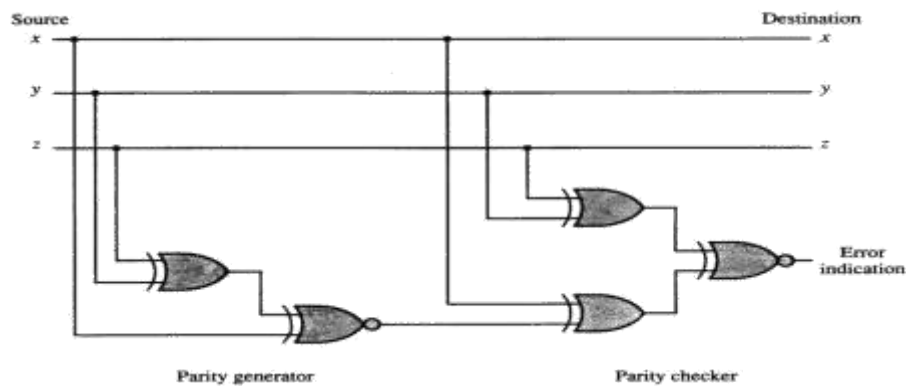
A parity bit is an extra bit included with a binary message to make the total number of 1's either odd or even

TABLE 3-7 Parity Bit Generation

Message <i>xyz</i>	<i>P</i> (odd)	<i>P</i> (even)
000	1	0
001	0	1
010	0	1
011	1	0
100	0	1
101	1	0
110	1	0
111	0	1

- The *P*(odd) bit is chosen to make the sum of 1's in all four bits odd
- The even-parity scheme has the disadvantage of having a bit combination of all 0's
- Procedure during transmission:
 - At the sending end, the message is applied to a *parity generator*
 - The message, including the parity bit, is transmitted
 - At the receiving end, all the incoming bits are applied to a *parity checker*
 - Any odd number of errors are detected
- Parity generators and checkers are constructed with XOR gates (odd function)
- An odd function generates 1 iff an odd number of input variables are 1

Figure 3-3 Error detection with odd parity bit.



COMPUTER ARITHMETIC

Introduction:

Data is manipulated by using the arithmetic instructions in digital computers. Data is manipulated to produce results necessary to give solution for the computation problems. The Addition, subtraction, multiplication and division are the four basic arithmetic operations. If we want then we can derive other operations by using these four operations.

To execute arithmetic operations there is a separate section called arithmetic processing unit in central processing unit. The arithmetic instructions are performed generally on binary or decimal data. Fixed-point numbers are used to represent integers or fractions. We can have signed or unsigned negative numbers. Fixed-point addition is the simplest arithmetic operation.

If we want to solve a problem then we use a sequence of well-defined steps. These steps are collectively called algorithm. To solve various problems we give algorithms.

In order to solve the computational problems, arithmetic instructions are used in digital computers that manipulate data. These instructions perform arithmetic calculations.

And these instructions perform a great activity in processing data in a digital computer. As we already stated that with the four basic arithmetic operations addition, subtraction, multiplication and division, it is possible to derive other arithmetic operations and solve scientific problems by means of numerical analysis methods.

A processor has an arithmetic processor(as a sub part of it) that executes arithmetic operations. The data type, assumed to reside in processor, registers during the execution of an arithmetic instruction. Negative numbers may be in a signed magnitude or signed complement representation. 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 magnitude representation for the mantissa.

Addition and Subtraction :

Addition and Subtraction with Signed –Magnitude Data

We designate the magnitude of the two numbers by A and B. Where 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 4.1. 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 present a negative zero. In other words, when two equal numbers are subtracted, the result should be +0 not -0.

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

Addition and Subtraction of Signed-Magnitude Numbers:

Eight Conditions for Signed-Magnitude Addition/Subtraction

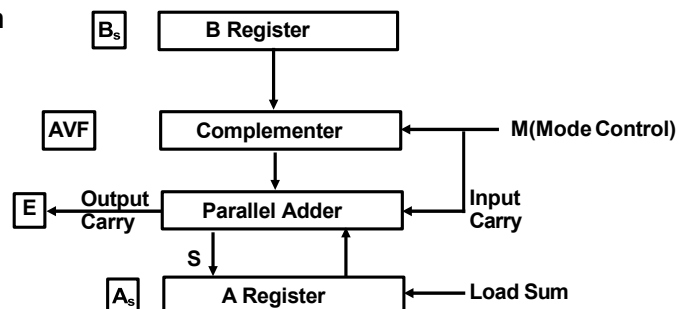
Operation	ADD Magnitudes	SUBTRACT Magnitudes		
		$A > B$	$A < B$	$A = B$
$(+A) + (+B)$	$+(A + B)$			
$(+A) + (-B)$		$+(A - B)$	$-(B - A)$	$+(A - B)$
$(-A) + (+B)$		$-(A - B)$	$+(B - A)$	$+(A - B)$
$(-A) + (-B)$	$-(A + B)$			
$(+A) - (+B)$		$+(A - B)$	$-(B - A)$	$+(A - B)$
$(+A) - (-B)$	$+(A + B)$			
$(-A) - (+B)$	$-(A + B)$			
$(-A) - (-B)$		$-(A - B)$	$+(B - A)$	$+(A - B)$

SIGNED MAGNITUDE ADDITION AND SUBTRACTION

Addition: **A + B ; A: Augend; B: Addend**

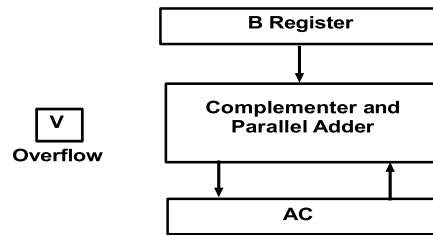
Subtraction: A - B: A: Minuend; B: Subtrahend

Hardware Implementation

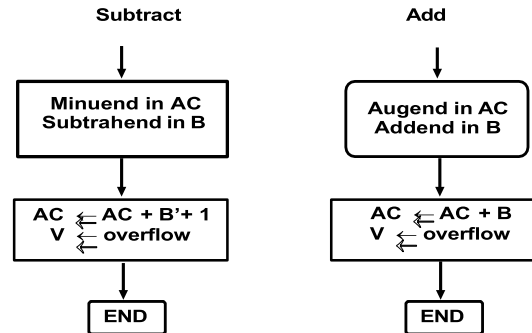


SIGNED 2'S COMPLEMENT ADDITION AND SUBTRACTION

Hardware



Algorithm



Algorithm:

- The flowchart is shown in Figure 7.1. The two signs A, and B, 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 complemented B. No overflow can occur if the numbers are subtracted so AVF is cleared to 0.
- 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 must be made positive to avoid a negative zero.
- 0 in E indicates that $A < B$. For this case it is necessary to take the 2's complement of the value in A. The operation can be done with one microoperation $A \leftarrow A' + 1$.
- However, we assume that 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 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, to obtain the correct sign.
- The final result is found in register A and its sign in As. The value in AVF provides an overflow indication. The final value of E is immaterial.
- Figure 7.2 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 As and Bs. 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 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.

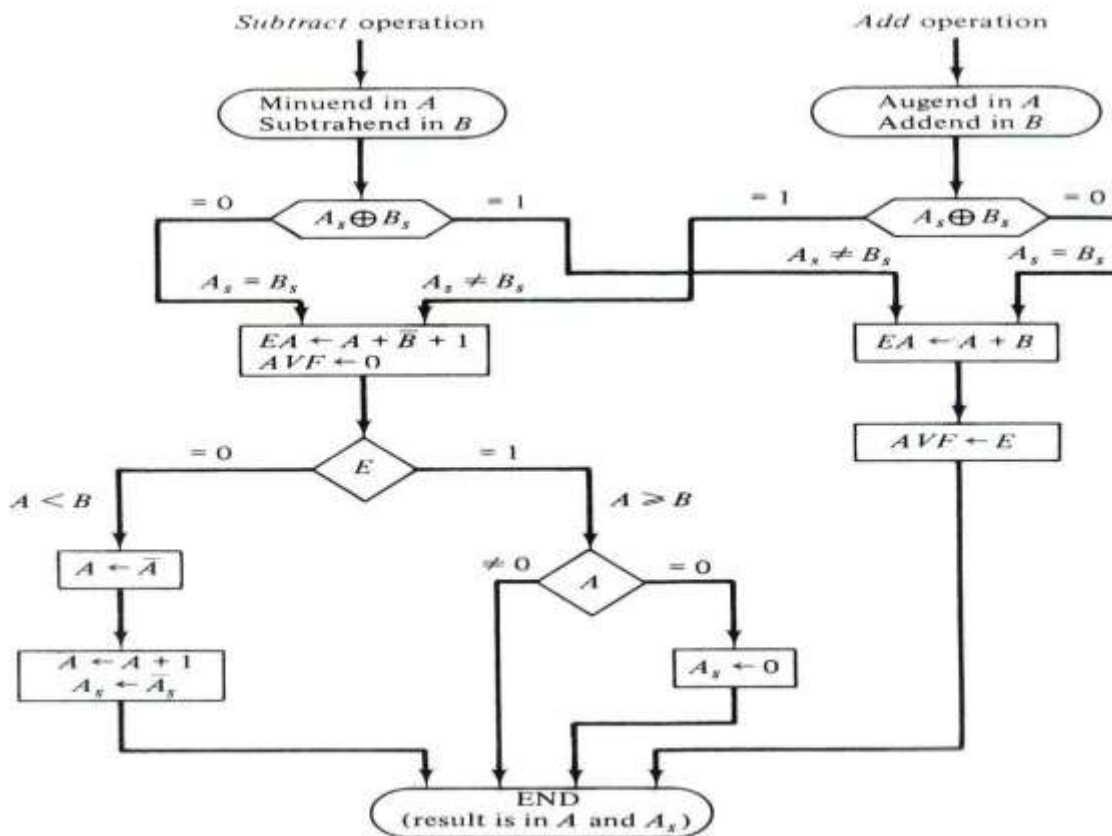


Figure 10-2 Flowchart for add and subtract operations.

Multiplication Algorithm:

In the beginning, the multiplicand is in B and the multiplier in Q. Their corresponding signs are in Bs and Qs respectively. We compare the signs of both A and Q and set to corresponding 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 the number of bits of the multiplier. 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.

Now, the low order bit of the multiplier in Qn is tested. If it is 1, the multiplicand (B) is added to present partial product (A), 0 otherwise. 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. When $SC = 0$ we stop the process.

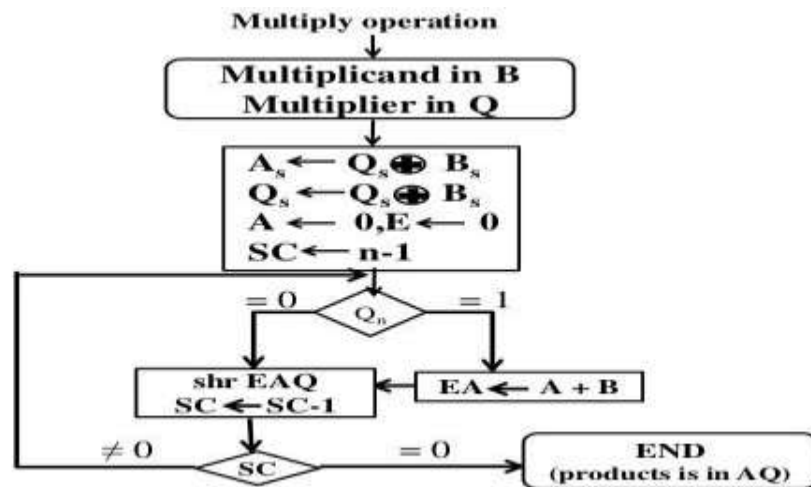
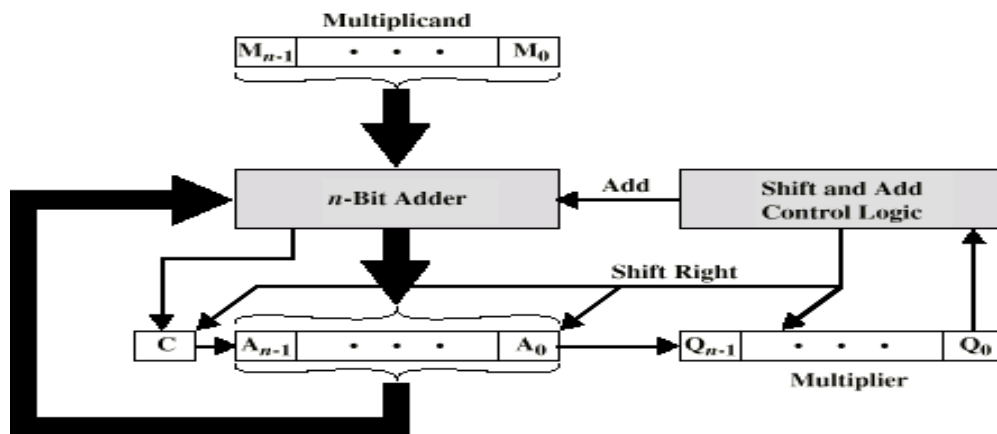


Figure: Flowchart for multiply operation.



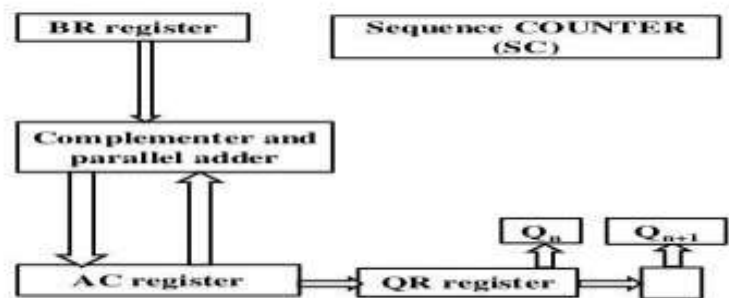
C	A	Q	M		
0	0000	1101	1011	Initial Values	
0	1011	1101	1011	Add	} First Cycle
0	0101	1110	1011	Shift	
0	0010	1111	1011	Shift	} Second Cycle
0	1101	1111	1011	Add	
0	0110	1111	1011	Shift	} Third Cycle
1	0001	1111	1011	Add	
0	1000	1111	1011	Shift	} Fourth Cycle

Booth's algorithm :

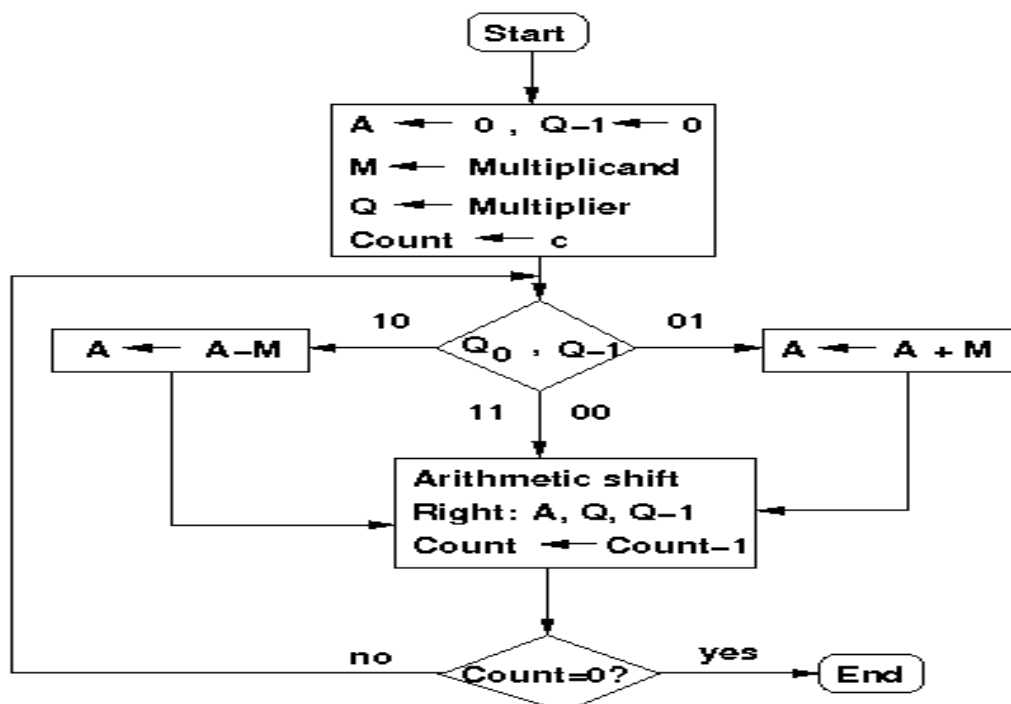
- Booth algorithm gives a procedure for multiplying binary integers in signed- 2's complement representation.
- It operates on the fact that strings of 0's in the multiplier require no addition but just shifting, and a string of 1's in the multiplier from bit weight 2^k to weight 2^m can be treated as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

Hardware for Booth Algorithm

- Sign bits are not separated from the rest of the registers
- 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
- Flip-flop Q_{n+1} is appended to QR to facilitate a double bit inspection of the multiplier



Algorithm:



- As in all multiplication schemes, booth algorithm requires examination of the multiplier bits and shifting of partial product.
- Prior to the shifting, the multiplicand may be added to the partial product, subtracted from the partial, 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 in a string of 0's in the multiplier.
 3. The partial product does not change when 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.
- 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.

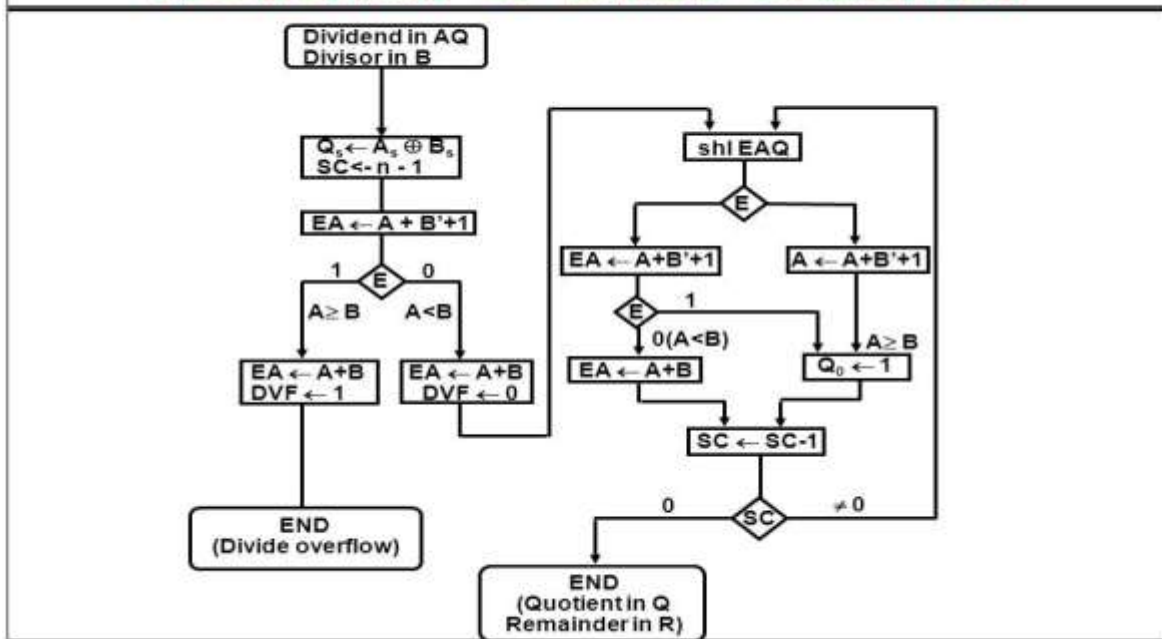
Division Algorithms

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

Divisor	1 1 0 1	$ \begin{array}{r} 000010101 \\ \overline{)100010010} \\ -1101 \\ \hline 10000 \\ -1101 \\ \hline 1110 \\ -1101 \\ \hline 1 \end{array} $	Quotient Dividend
		Remainder	

The divisor is compared with the five most significant bits of the dividend. Since the 5-bit number is

FLOWCHART OF DIVIDE OPERATION



Computer Organization

Prof. H. Yoon

Example of Binary Division with Digital Hardware

Divisor B = 10001	E	A	Q	SC
Dividend: shl EAQ add B + 1	0	01110 11100 <u>01111</u>	00000 00000	5
E = 1 Set Q _n = 1 shl EAQ Add B + 1	1	01011		
E = 1 Set Q _n = 1 shl EAQ Add B + 1	1	01011	00001	
E = 0; leave Q _n = 0 Add B	0	10110 <u>01111</u>	00010	4
E = 1 Set Q _n = 1 shl EAQ Add B + 1	1	00101		
E = 1 Set Q _n = 1 shl EAQ Add B + 1	1	00101	00011	
E = 0; leave Q _n = 0 Add B	0	01010 <u>01111</u>	00110	3
E = 1 Set Q _n = 1 shl EAQ Add B + 1	1	11001		
E = 1 Set Q _n = 1 shl EAQ Add B + 1	1	10001	00110	2
E = 0; leave Q _n = 0 Add B	0	01010 <u>01111</u>		
E = 1 Set Q _n = 1 shl EAQ Add B + 1	1	10100	01100	
E = 1 Set Q _n = 1 shl EAQ Add B + 1	1	01111		
E = 0; leave Q _n = 0 Add B	0	00011		
E = 0; leave Q _n = 0 Add B	0	00011	01101	1
E = 1 Set Q _n = 1 shl EAQ Add B + 1	1	00110	11010	
E = 0; leave Q _n = 0 Add B	0	00110		
E = 1 Set Q _n = 1 shl EAQ Add B + 1	1	10101 <u>10001</u>	11010	0
E = 0; leave Q _n = 0 Add B	0	00110		
Restore remainder Neglect E Remainder in A: Quotient in Q:	1	00110	11010	
		00110	11010	

Floating-point Arithmetic operations:

In many high-level programming languages, we have a facility for specifying floating-point numbers. The most common way is by a real declaration statement. High level programming languages must have a provision for handling floating-point arithmetic operations. The operations are generally built in the internal hardware. If no hardware is available, the compiler must be designed with a package of floating-point software subroutine. Although the hardware method is more expensive, it is much more efficient than the software method. Therefore, floating-point hardware is included in most computers and is omitted only in very small ones.

Basic Considerations:

There are two part of a floating-point number in a computer - a mantissa m and an exponent e . The two parts represent a number generated 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 position of the radix point and the value of the radix r 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 said to be normalized if the most significant digit of the mantissa is nonzero. So the mantissa contains the maximum possible number of significant digits. We cannot normalize a zero 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 for 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 $+(2^{47} - 1)$, which is approximately $+10^{14}$. 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 represented is

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

This number is derived from a fraction that contains 35 1's, an exponent of 11 bits (excluding its sign), and because $2^{11}-1 = 2047$. The largest number that can be accommodated is approximately 10^{615} .

The mantissa that can accommodated is 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 uses four words to represent one floating-point number. One word of 8 bits are reserved for the exponent and the 24 bits of the other three words are used in the mantissa.

Arithmetic operations with floating-point numbers are more complicated than with fixed-point numbers. Their execution also takes longer time 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. We do this alignment by shifting one mantissa while its exponent is adjusted until it becomes equal to the other exponent.

Consider the sum of the following floating-point numbers: $.5372400 \times 10^2$
 $+ .1580000 \times 10^{-1}$

Floating-point multiplication and division need not do an alignment of the mantissas. Multiplying the two mantissas and adding the exponents can form the product. Dividing the mantissas and subtracting the exponents perform division.

The operations done 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 compared and incremented (for aligning the mantissas), added and subtracted (for multiplication) and division), and decremented (to normalize the result). We can represent the exponent in any one of the three representations - signed-magnitude, signed 2's complement or signed 1's complement.

Register Configuration

The register configuration for floating-point operations is shown in figure 4.13. As a 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. 4.13. Three registers are there, 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 may use corresponding lower-case letter symbol.

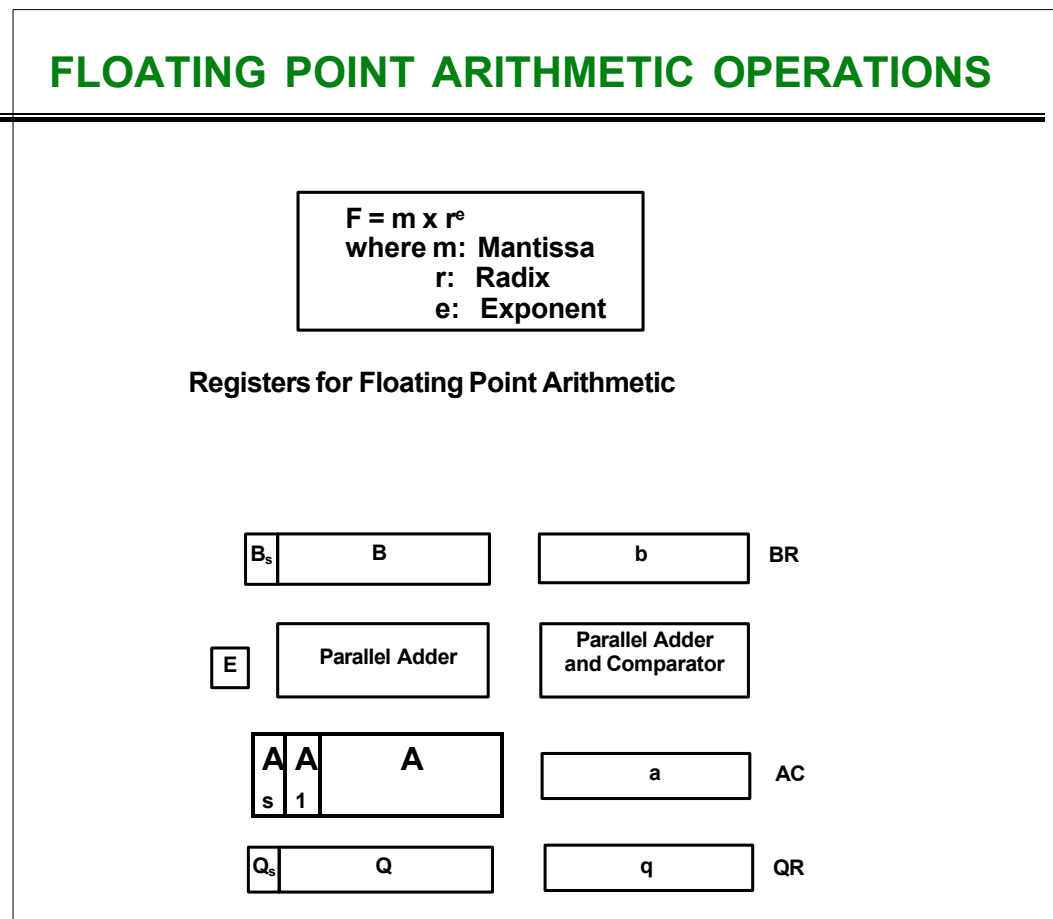


Figure 4.13: Registers for Floating Point arithmetic operations

Assuming that each floating-point number has a mantissa in signed-magnitude representation and a biased exponent. Thus the AC has a mantissa whose sign is in A_s, and a magnitude that is in A. The

diagram shows the most significant bit of A, labeled by A1. The bit in this position must be a 1 to normalize the number. Note that the symbol AC represents the entire register, that is, the concatenation of As, A and a.

In the similar way, register BR is subdivided into Bs, B, and b and QR into Qs, Q and q. A parallel-adder adds the two mantissas and loads the sum into A and the carry into E. A separate parallel adder can be used for the exponents. The exponents do not have a distinct sign bit because they are biased 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 so 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 should initially be normalized. After each arithmetic operation, the result will be normalized. Thus all floating-point operands are always normalized.

Addition and Subtraction of Floating Point Numbers

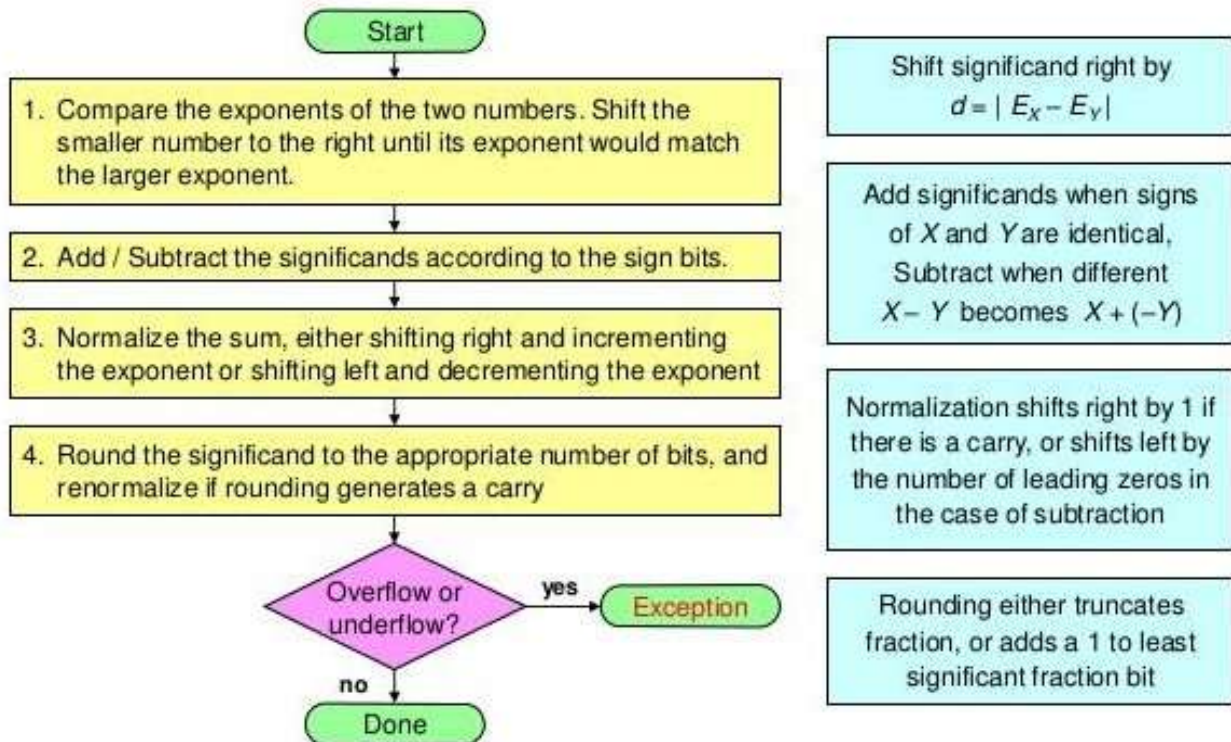
During addition or subtraction, the two floating-point operands are kept in AC and BR. The sum or difference is formed in the AC. The algorithm can be divided into four consecutive parts:

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

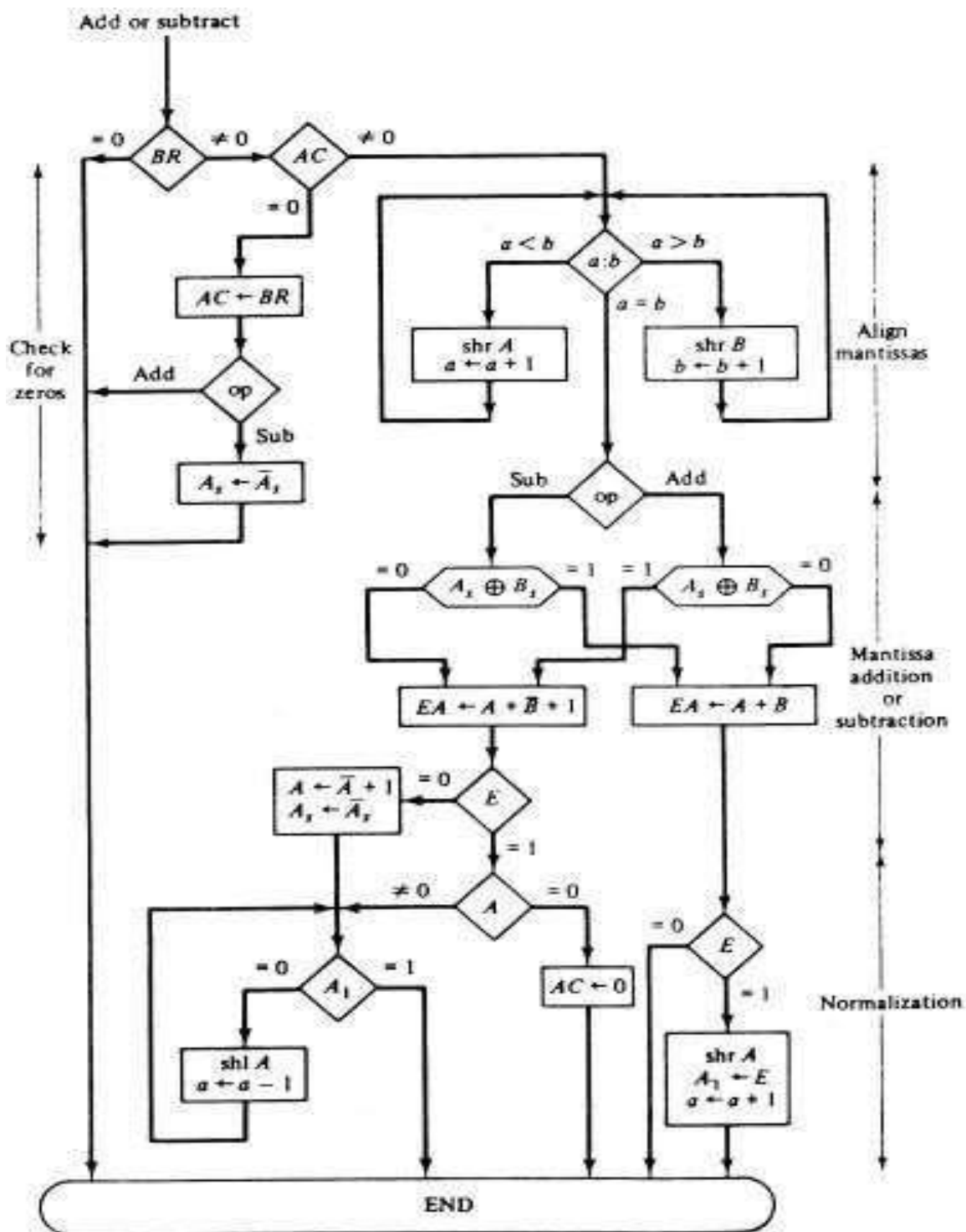
A floating-point number cannot be normalized, if it is 0. If this number is used for computation, the result may also be zero. Instead of checking for zeros during the normalization process we check for zeros at the beginning and terminate the process if necessary. The alignment of the mantissas must be carried out prior to their operation. After the mantissas are added or subtracted, the result may be un-normalized. The normalization procedure ensures that the result is normalized before it is transferred to memory.

If the magnitudes were subtracted, there may be zero or may have an underflow in the result. If the mantissa is equal to zero the entire floating-point number in the AC is cleared to 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 A1, is 0. In that case, the mantissa is shifted left and the exponent decremented. The bit in A1 is checked again and the process is repeated until $A1 = 1$. When $A1 = 1$, the mantissa is normalized and the operation is completed.

Floating Point Addition / Subtraction

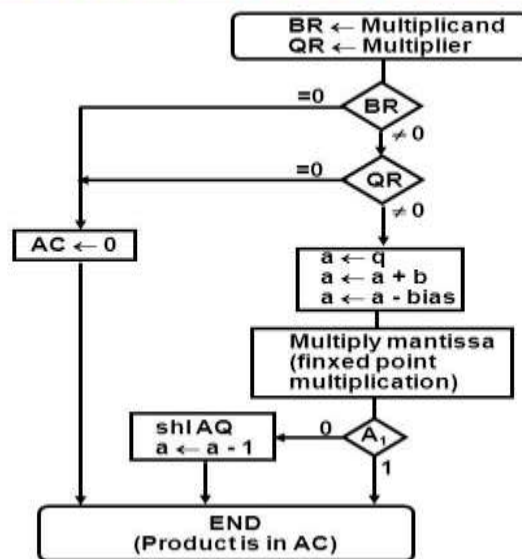


Algorithm for Floating Point Addition and Subtraction



Multiplication:

FLOATING POINT MULTIPLICATION



Division:

FLOATING POINT DIVISION

