# 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)
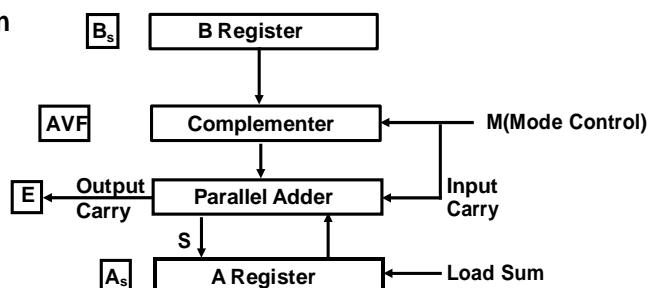
## SIGNED MAGNITUDEADDITION AND SUBTRACTION

**Addition:     A + B ; A: Augend;  B: Addend**
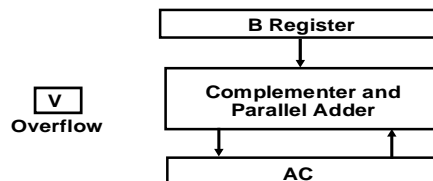**Subtraction:  A - B:   A: Minuend;  B: Subtrahend**

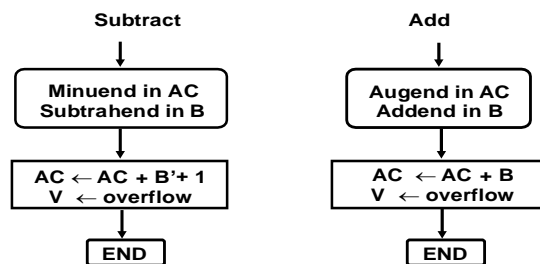| Operation | Add Magnitude | Subtract Magnitude | | |
|---|---|---|---|---|
| | | When A>B | When A<B | When A=B |
| (+A) + (+B) | +(A + B) | | | |
| (+A) + (- B) | | +(A - B) | - (B - A) | +(A - B) |
| (- A) + (+B) | | - (A - B) | +(B - A) | +(A - B) |
| (- A) + (- B) | - (A + B) | | | |
| (+A) - (+B)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

**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 >= B and the number in A is the correct result. If this numbs 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.
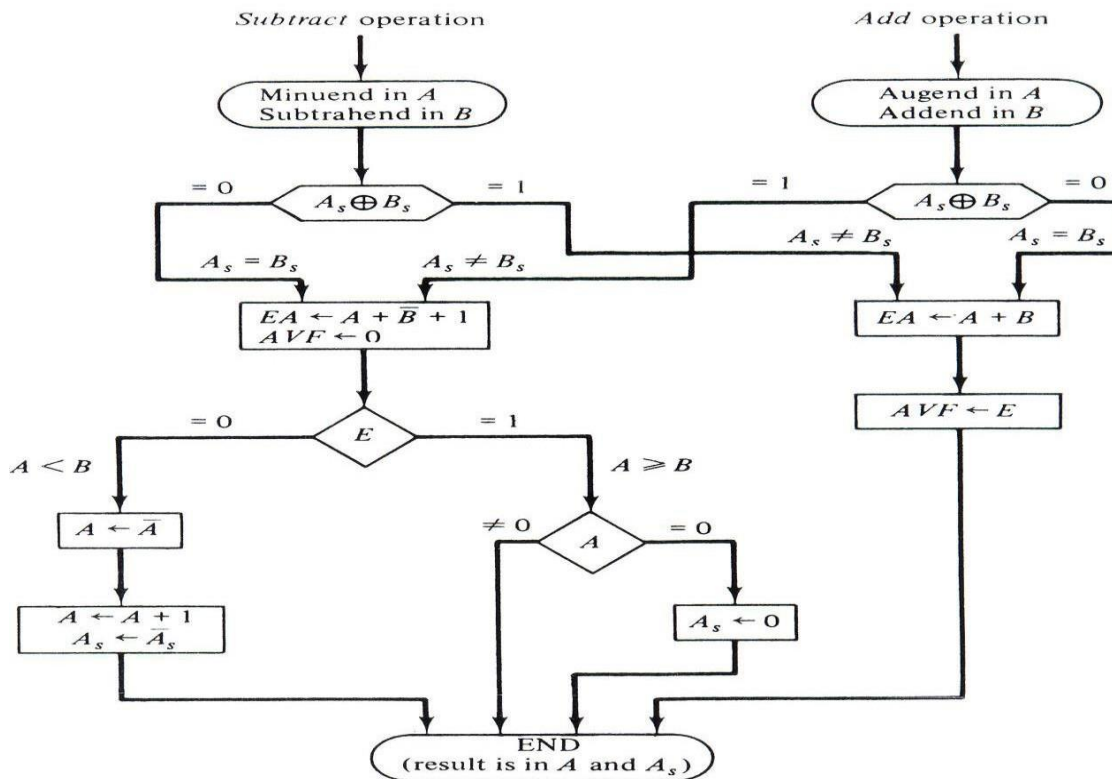
**Figure 10-2** Flowchart for add and subtract operations.

**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 stops the process.

```
C    A       Q       M
0    0000    1101    1011     Initial Values

0    1011    1101    1011     Add    ⎫  First
0    0101    1110    1011     Shift  ⎭  Cycle

0    0010    1111    1011     Shift  ⎫  Second
                                     ⎭  Cycle

0    1101    1111    1011     Add    ⎫  Third
0    0110    1111    1011     Shift  ⎭  Cycle

1    0001    1111    1011     Add    ⎫  Fourth
0    1000    1111    1011     Shift  ⎭  Cycle
```
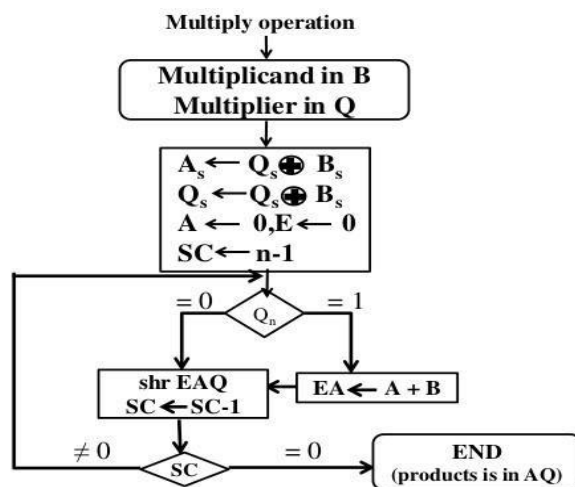


## Figure: Flowchart for multiply operation.

**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 X 14, where M is the multiplicand and 14 the multiplier, can be done as M X $2^4$ – M X $2^1$.

☐ 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 Qn+1 is appended to QR to facilitate a double bit inspection of the multiplier



☐ 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:

6

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 Qn and Qn+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

```
                000010101  Quotient
Divisor  1101)100010010  Dividend
              -1101
                10000
               -1101
                 1110
                -1101
                    1  Remainder
```

The devisor is compared with the five most significant bits of the dividend. Since the 5-bit number is smaller than B, we again repeat the same process. Now the 6-bit number is greater than B, so we place a 1 for the quotient bit in the sixth position above the dividend. Now we shift the divisor once to the right and subtract it from the dividend. The difference is known as a partial remainder because the division could have stopped here to obtain a quotient of 1 and a remainder equal to the partial

remainder. Comparing a partial remainder with the divisor continues the process. 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. Obviously the result gives both a quotient and a remainder.

Hardware Implementation for Signed-Magnitude Data

In hardware implementation for signed-magnitude data in a digital computer, it is convenient to change the process slightly. Instead of shifting the divisor to the right, two dividends, or partial remainders, are shifted to the left, thus leaving the two numbers in the required relative position. Subtraction is achieved by adding A to the 2's complement of B. End carry gives the information about the relative magnitudes.

The hardware required is identical to that of multiplication. Register EAQ is now shifted to the left with 0 inserted into Qn and the previous value of E is lost. The example is given in Figure 4.10 to clear the 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. E



Hardware Implementation for Signed-Magnitude Data

**Algorithm:**

# FLOWCHART OF DIVIDE OPERATION

```
                    Dividend in AQ
                    Divisor in B

        Qs ← As ⊕ Bs                    shl EAQ
        SC ← n - 1
                                           E
        EA ← A + B'+1
                              EA ← A+B'+1        A ← A+B'+1
         1    E    0
       A≥B        A<B              E      1
                                      0(A<B)         A≥B
    EA ← A+B    EA ← A+B      EA ← A+B         Q0 ← 1
    DVF ← 1     DVF ← 0
                                      SC ← SC-1

                                  0              ≠ 0
        END                           SC
    (Divide overflow)
                                  END
                              (Quotient in Q
                               Remainder in R)
```

Example of Binary Division with Digital Hardware

Divisor B = 10001       $\overline{B} + 1 = 01111$

| | E | A | Q | SC |
|---|---|---|---|---|
| Dividend: | | 01110 | 00000 | 5 |
| shl EAQ | 0 | 11100 | 00000 | |
| add $\overline{B} + 1$ | | 01111 | | |
| E = 1 | 1 | 01011 | | |
| Set $Q_e = 1$ | 1 | 01011 | 00001 | 4 |
| shl EAQ | 0 | 10110 | 00010 | |
| Add $\overline{B} + 1$ | | 01111 | | |
| E = 1 | 1 | 00101 | | |
| Set $Q_e = 1$ | 1 | 00101 | 00011 | 3 |
| shl EAQ | 0 | 01010 | 00110 | |
| Add $\overline{B} + 1$ | | 01111 | | |
| E = Q; leave $Q_n = 0$ | 0 | 11001 | 00110 | |
| Add B | | 10001 | | 2 |
| Restore remainder | 1 | 01010 | | |
| shl EAQ | 0 | 10100 | 01100 | |
| Add $\overline{B} + 1$ | | 01111 | | |
| E = 1 | 1 | 00011 | | |
| Set $Q_e = 1$ | 1 | 00011 | 01101 | 1 |
| Shl EAQ | 0 | 00110 | 11010 | |
| Add $\overline{B} + 1$ | | 01111 | | |
| E = 0; leave $Q_e = 0$ | 0 | 10101 | 11010 | |
| Add B | | 10001 | | |
| Restore remainder | 1 | 00110 | 11010 | 0 |
| Neglect E | | | | |
| Remainder in A: | | 00110 | | |
| Quotient in Q: | | | 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 in 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 $+ (24^7 - 1)$, which is approximately $+ 101^4$. 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.

Biased exponents have the advantage that they contain only positive numbers. Now it becomes simpler to compare their relative magnitude without bothering about their signs. 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.

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.

# FLOATING POINT ARITHMETIC OPERATIONS

$$F = m \times r^e$$
where m: Mantissa
r: Radix
e: Exponent

**Registers for Floating Point Arithmetic**

| $B_s$ | B | | b | BR |
|---|---|---|---|---|

| E | Parallel Adder | | Parallel Adder and Comparator | |
|---|---|---|---|---|

| $A_s$ | $A_1$ | A | | a | AC |
|---|---|---|---|---|---|

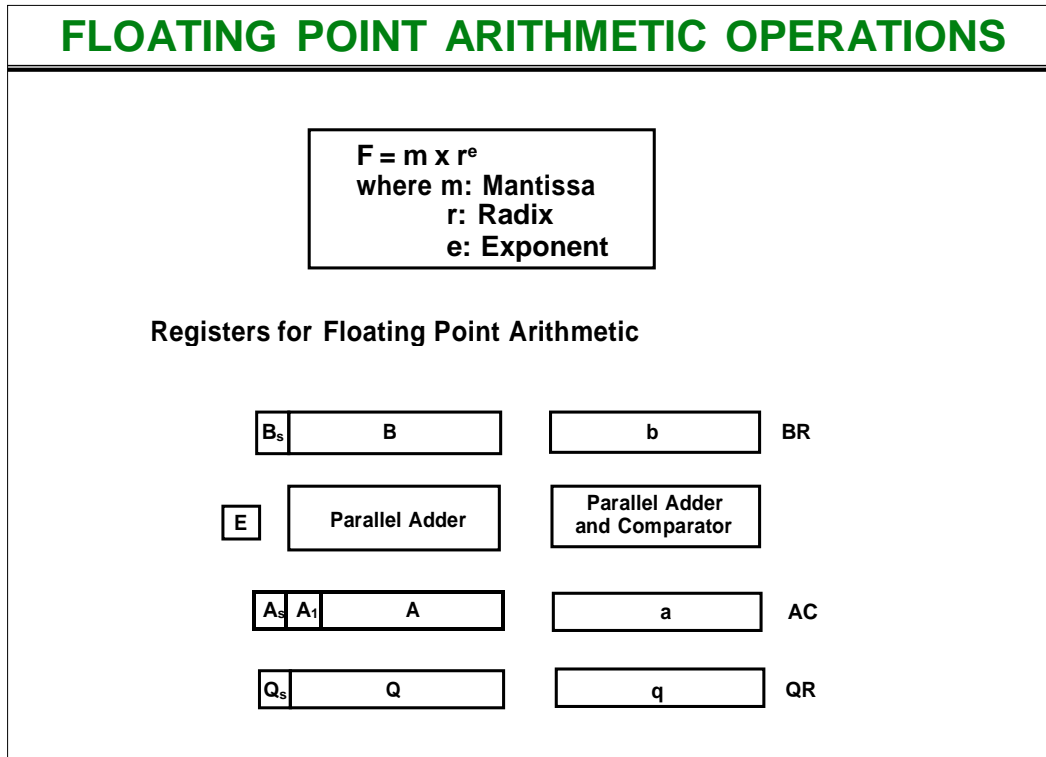| $Q_s$ | Q | | q | QR |
|---|---|---|---|---|

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 As, and a magnitude that is in A. The diagram shows the most significant bit of A, labeled by A1. The bit in his 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 district sign bit because they are biased but are represented as a biased positive quantity. It is assumed that the floating-point number 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 they 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.
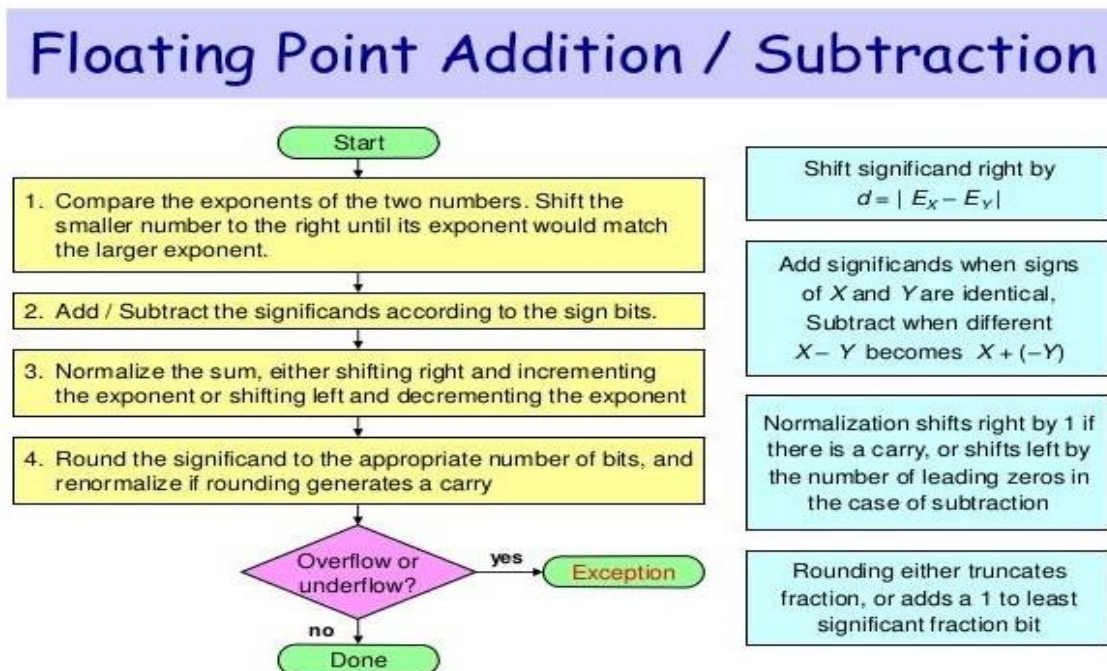
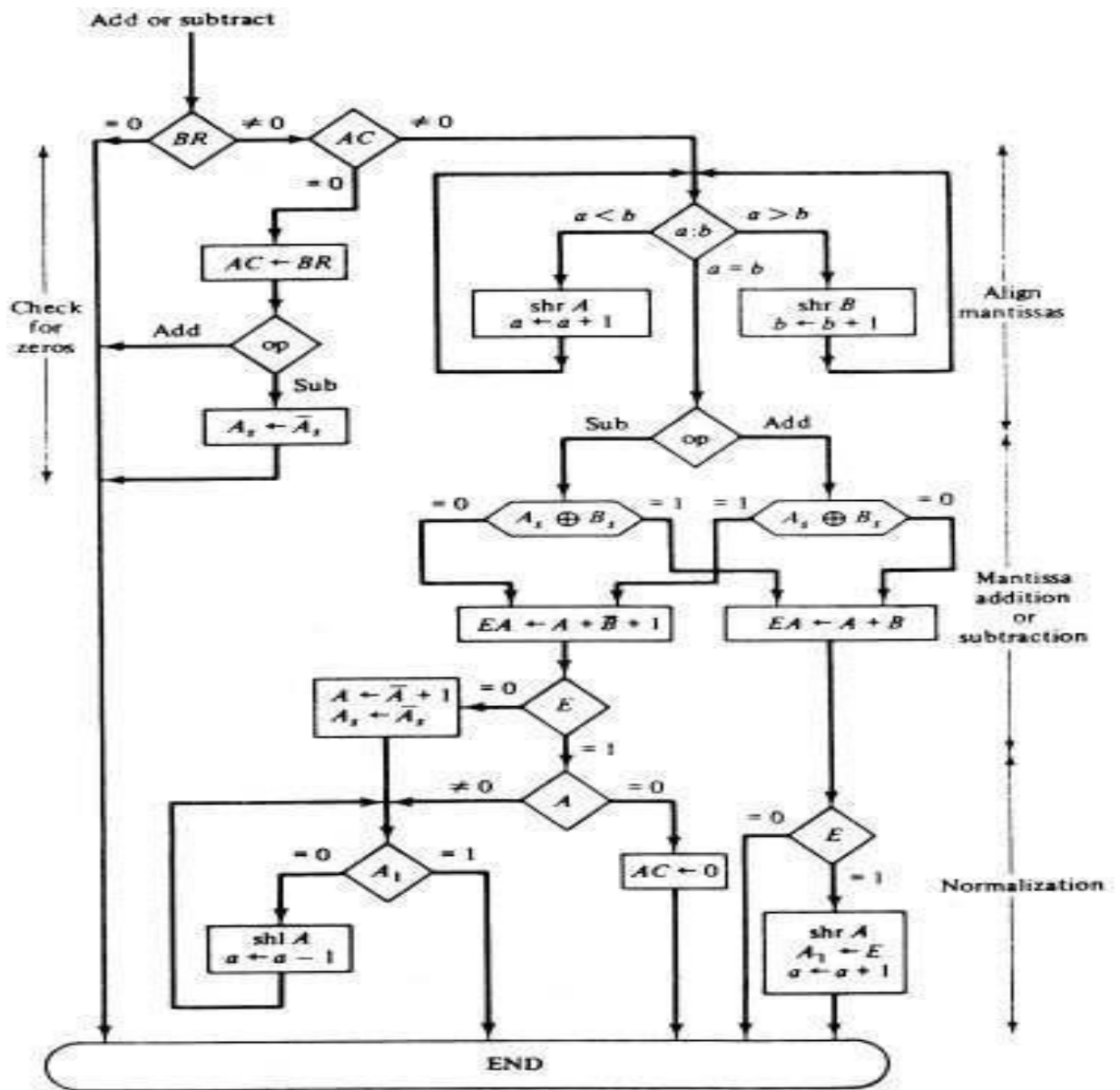**Addition and Subtraction of Floating Point Numbers**

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**

Start

1. Compare the exponents of the two numbers. Shift the smaller number to the right until its exponent would match the larger exponent.

2. Add / Subtract the significands according to the sign bits.

3. Normalize the sum, either shifting right and incrementing the exponent or shifting left and decrementing the exponent

4. Round the significand to the appropriate number of bits, and renormalize if rounding generates a carry

Overflow or underflow? — yes → Exception

no

Done

Shift significand right by
$$d = |E_x - E_y|$$

Add significands when signs of $X$ and $Y$ are identical, Subtract when different
$X - Y$ becomes $X + (-Y)$

Normalization shifts right by 1 if there is a carry, or shifts left by the number of leading zeros in the case of subtraction

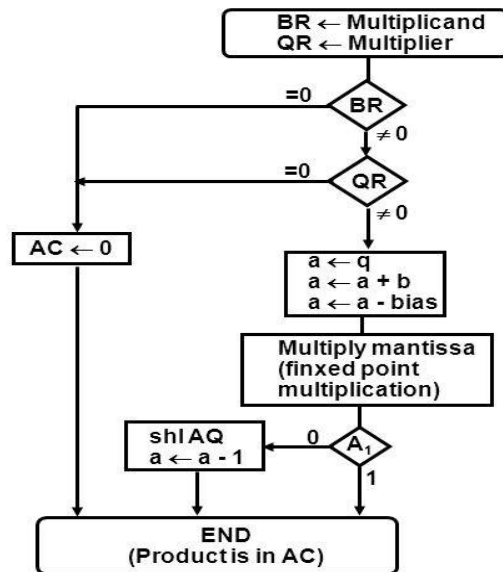Rounding either truncates fraction, or adds a 1 to least significant fraction bit

Algorithm for Floating Point Addition and Subtraction

**Multiplication:**

## FLOATING POINT MULTIPLICATION

BR ← Multiplicand
QR ← Multiplier

BR
=0
≠ 0

QR
=0
≠ 0

AC ← 0

$a \leftarrow q$
$a \leftarrow a + b$
$a \leftarrow a - bias$

Multiply mantissa
(finxed point
multiplication)

shl AQ
$a \leftarrow a - 1$

$A_1$
0
1

END
(Product is in AC)

## FLOATING POINT DIVISION

BR ← Divisor
AC ← Dividend

BR
=0
≠ 0

AC
=0
≠ 0

QR ← 0

divide
by 0

$Qs \leftarrow As + Bs$
$Q \leftarrow 0$
$SC \leftarrow n-1$

$EA \leftarrow A+B'+1$

E
1
0

A>=B
A<B

$A \leftarrow A+B$
shr A
$a \leftarrow a+1$

$A \leftarrow A+B$

$a \leftarrow a+b'+1$
$a \leftarrow a+bias$
$q \leftarrow a$

Divide Magnitude of mantissa
as in fixed point numbers