

# Chapter 9 Computer Arithmetic

## Integer Representation

### (1) Sign-Magnitude Method

$$+18 = \boxed{0}0010010$$

$$-18 = \boxed{1}0010010$$

Drawbacks -

### (2) Two's Complement Method

Ex. (-3) in 8-bit number system

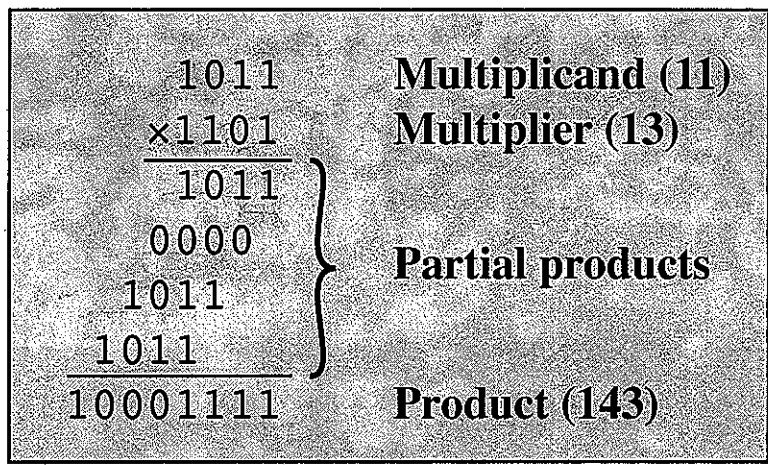
$$\begin{array}{r} +3 \quad \quad \quad 0000011 \\ \text{complement} \quad 1111100 \\ \text{add 1} \quad \quad \quad 1111101 \quad (= -3) \end{array}$$

-----  
Ex.  $5 - 3 = 5 + (-3)$

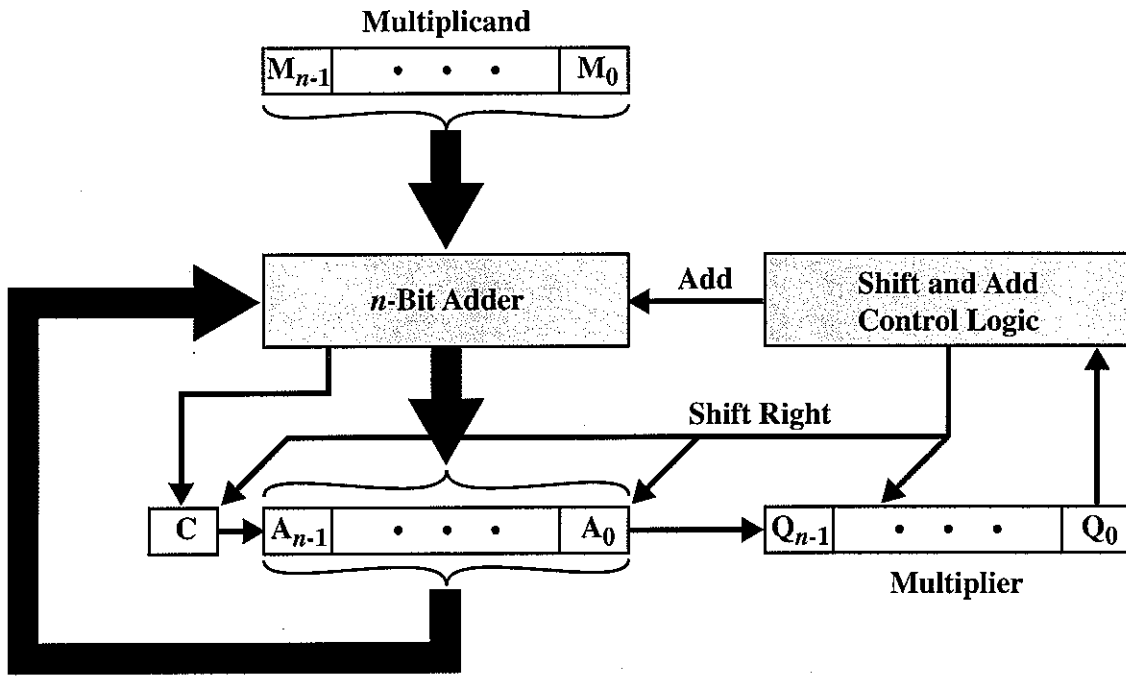
$$\begin{array}{r} 5 \quad \quad \quad 0000101 \\ -3 \quad \quad \quad 1111101 \quad + \\ \hline \boxed{1} 0000010 \quad (= 2) \\ \text{ignore} \end{array}$$

Note. No subtractor is necessary.

Note: If two numbers are added and they are both positive or both negative, then overflow occurs iff the result has the opposite sign.



**Figure 9.7 Multiplication of Unsigned Binary Integers**

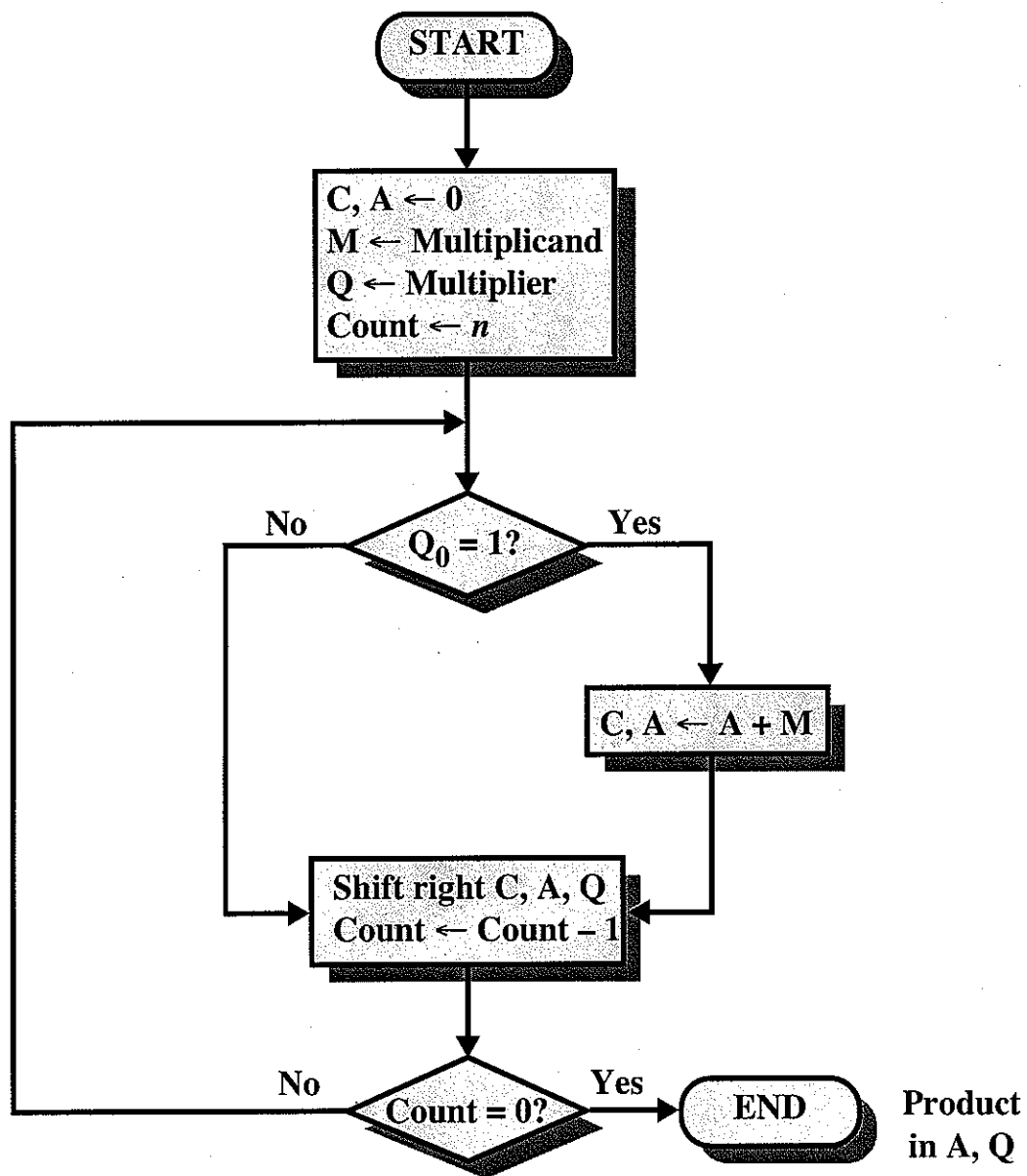


(a) Block Diagram

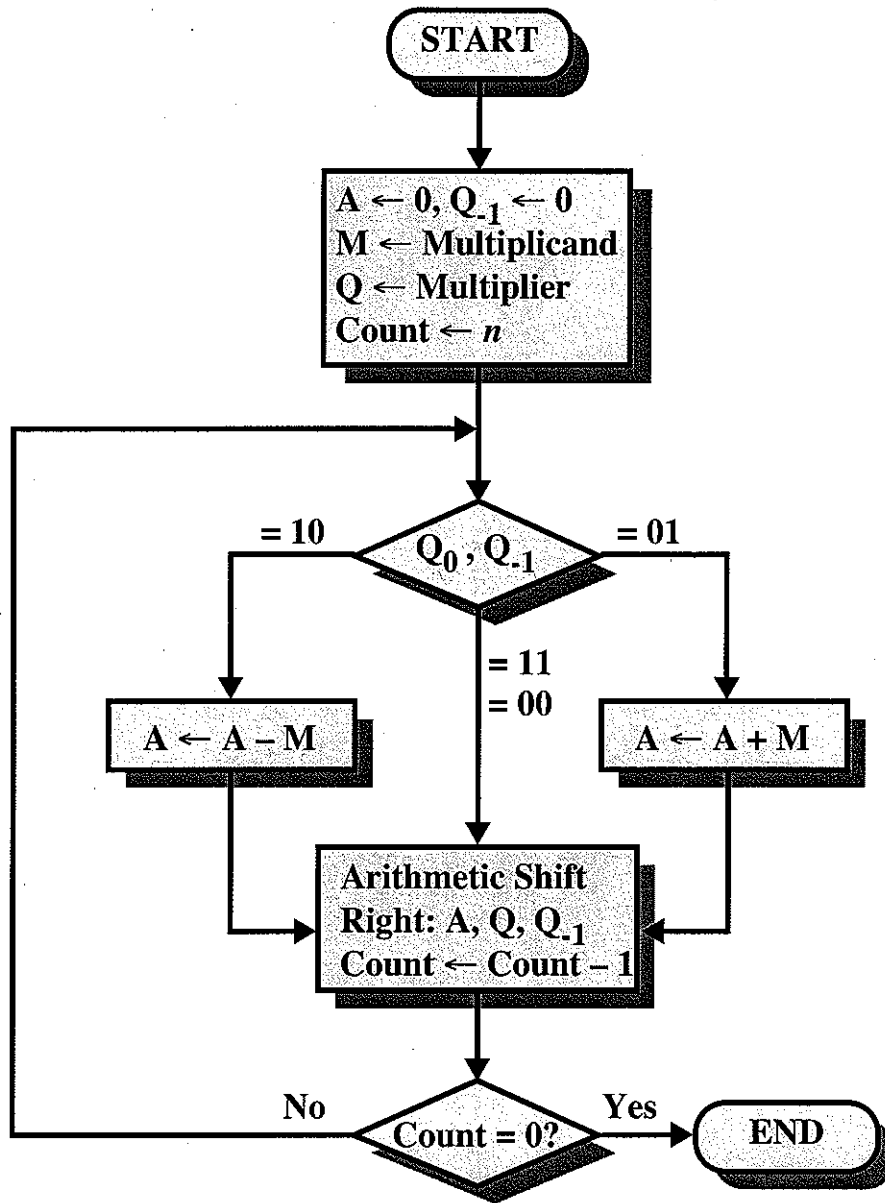
C	A	Q	M	
0	0000	1101	1011	Initial Values
0	1011	1101	1011	Add } First Cycle
0	0101	1110	1011	
0	0010	1111	1011	Shift } Second Cycle
0	1101	1111	1011	Add } Third Cycle
0	0110	1111	1011	
1	0001	1111	1011	Add } Fourth Cycle
0	1000	1111	1011	

(b) Example from Figure 9.7 (product in A, Q)

**Figure 9.8 Hardware Implementation of Unsigned Binary Multiplication**



**Figure 9.9 Flowchart for Unsigned Binary Multiplication**



**Figure 9.12 Booth's Algorithm for Two's Complement Multiplication**

A	Q	Q <sub>-1</sub>	M		
0000	0011	0	0111	Initial Values	
1001	0011	0	0111	A ← A - M Shift	} First Cycle
1100	1001	1	0111		
1110	0100	1	0111	Shift	} Second Cycle
0101	0100	1	0111	A ← A + M Shift	} Third Cycle
0010	1010	0	0111		
0001	0101	0	0111	Shift	} Fourth Cycle

**Figure 9.13 Example of Booth's Algorithm (7 × 3)**

### Booth's algorithm - continued

Ex.  $0\ 0\ \boxed{0\ 1}\ 1\ 1\ \boxed{1\ 0}$

*end of run*                      *beginning of run*

Note.  $2^4 + 2^3 + 2^2 + 2^1 \equiv 2^5 - 2^1$

Ex. (7) x (-3)

A	Q	Q <sub>-1</sub>			
0000	1101	0	<b>Initial</b>		
1001	1101	0	<b>A ← A - M</b>	}	
1100	1110	1	<b>shift</b>		First cycle
0011	1110	1	<b>A ← A + M</b>	}	
0001	1111	0	<b>shift</b>		Second cycle
1010	1111	0	<b>A ← A - M</b>	}	
1101	0111	1	<b>shift</b>		Third cycle
<span style="border: 1px solid black; padding: 2px;">1110</span>	<span style="border: 1px solid black; padding: 2px;">1011</span>	1	<b>shift</b>	}	Fourth cycle



# Floating-point Arithmetic

## Addition and Subtraction

1. Check for zero
2. Align the significands // smaller number is shifted right
3. Add or subtract the significands
4. Normalize the result // shift left until most-significant digit is nonzero.

Example. (decimal, 3 significands)

$$123 + 4.56$$

Normalized:  $1.23 \times 10^2 + 4.56 \times 10^0$

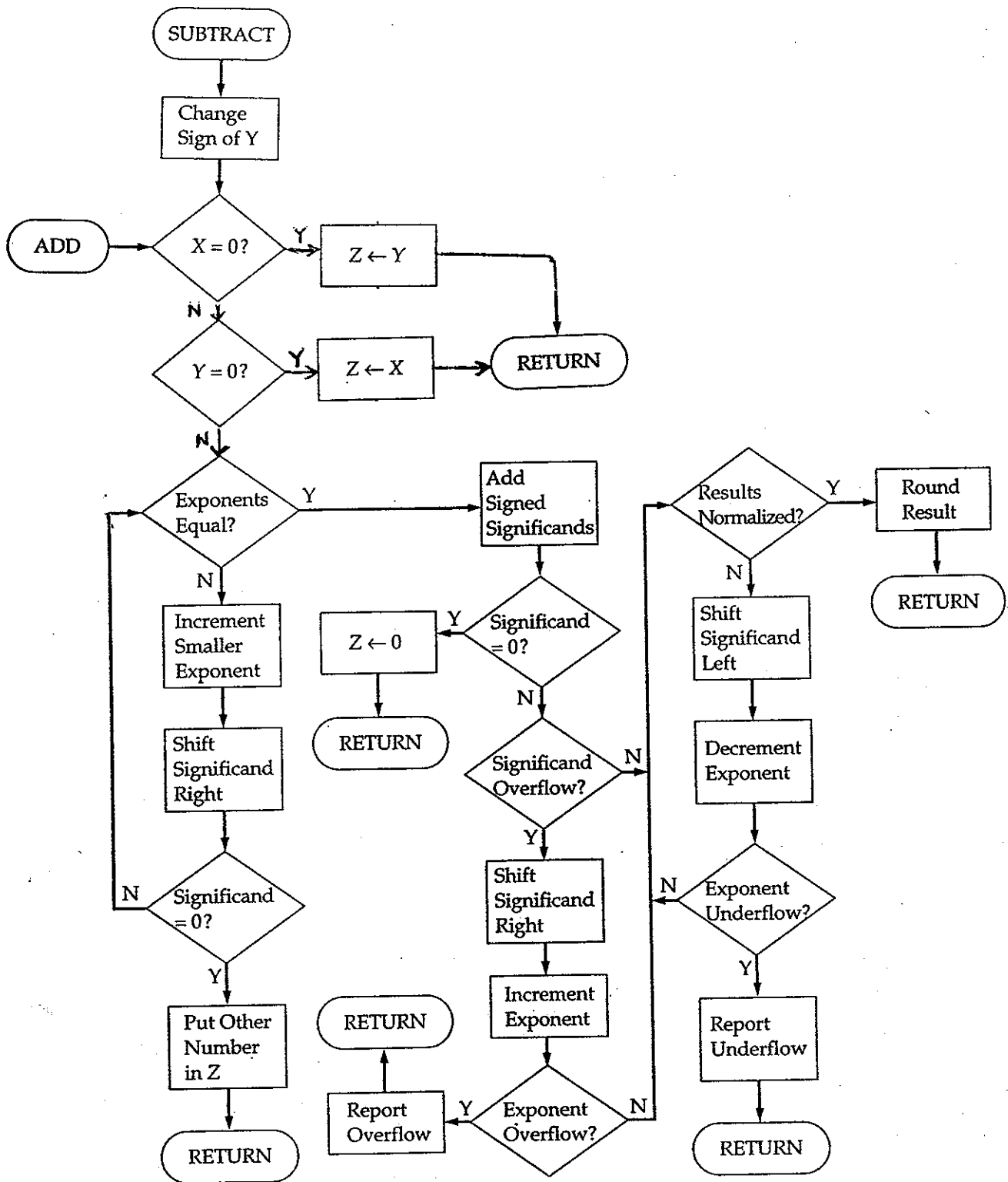
Align

$$\begin{array}{r}
 1.23 \times 10^2 \\
 +) 0.0456 \times 10^2 \\
 \hline
 1.2756 \times 10^2
 \end{array}$$

truncate:  $1.27 \times 10^2$

round:  $1.28 \times 10^2$

Floating-point addition and subtraction ( $z \leftarrow x \pm y$ )



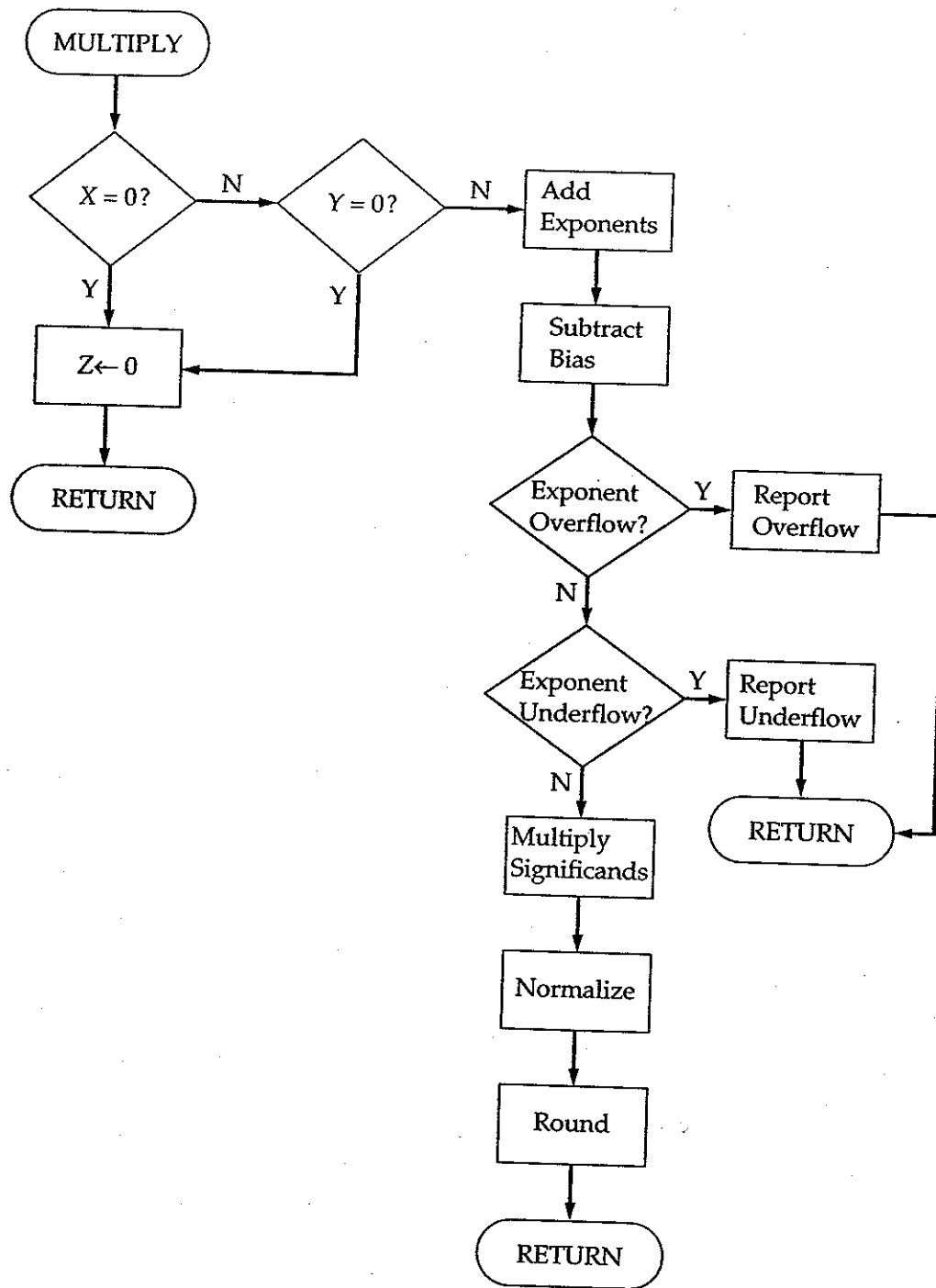


FIGURE 9.23. Floating-point multiplication ( $z \leftarrow x \times y$ )

# Precision Consideration

## 1. Guard Bits

Motivation.

$$X = 1.00\dots00 \times 2$$

$$Y = 1.11\dots11 \times 2$$

For  $X - Y$ ,  $Y$  should be shifted one bit to the right (to align).

→  $Y$  will lose one bit of significance.

$x = 1.000\dots00 \times 2^1$ $-y = 0.111\dots11 \times 2^1$ $z = 0.000\dots01 \times 2^1$ $= 1.000\dots00 \times 2^{-22}$	$x = .100000 \times 16^1$ $-y = .0FFFFFF \times 16^1$ $z = .000001 \times 16^1$ $= .100000 \times 16^{-4}$
---	---

(a) Binary example, without guard bits

(c) Hexadecimal example, without guard bits

$x = 1.000\dots00\ 0000 \times 2^1$ $-y = 0.111\dots11\ 1000 \times 2^1$ $z = 0.000\dots00\ 1000 \times 2^1$ $= 1.000\dots00\ 0000 \times 2^{-23}$	$x = .100000\ 00 \times 16^1$ $-y = .0FFFFFF\ F0 \times 16^1$ $z = .000000\ 10 \times 16^1$ $= .100000\ 00 \times 16^{-5}$
---	---

(b) Binary example, with guard bits

(d) Hexadecimal example, with guard bits

## 2. Rounding Policy

- Round to nearest number (default)

Example.

- |     |   |          |
|-----|---|----------|
| (1) | ....xxx   <sup>23</sup> 1 0 0 0 1 0 ... | roundup  |
| (2) | ....xxx   0 1 1 0 1 0 ...               | truncate |

Case of tie: .. xxx | 1 0 0 0 0 ...

Solution.

- |                       |          |
|-----------------------|----------|
| ..xx1   1 0 0 0 0 ... | roundup  |
| ..xx0   1 0 0 0 0 ... | truncate |