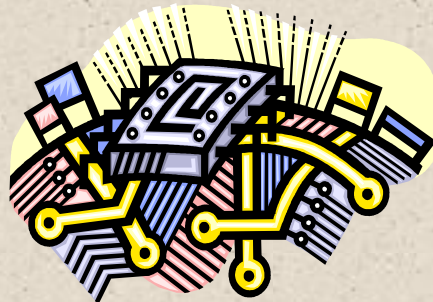# Arithmetic & Logic Unit (ALU)

- Part of the computer that actually performs arithmetic and logical operations on data

- All of the other elements of the computer system are there mainly to bring data into the ALU for it to process and then to take the results back out

- Based on the use of simple digital logic devices that can store binary digits and perform simple Boolean logic operations
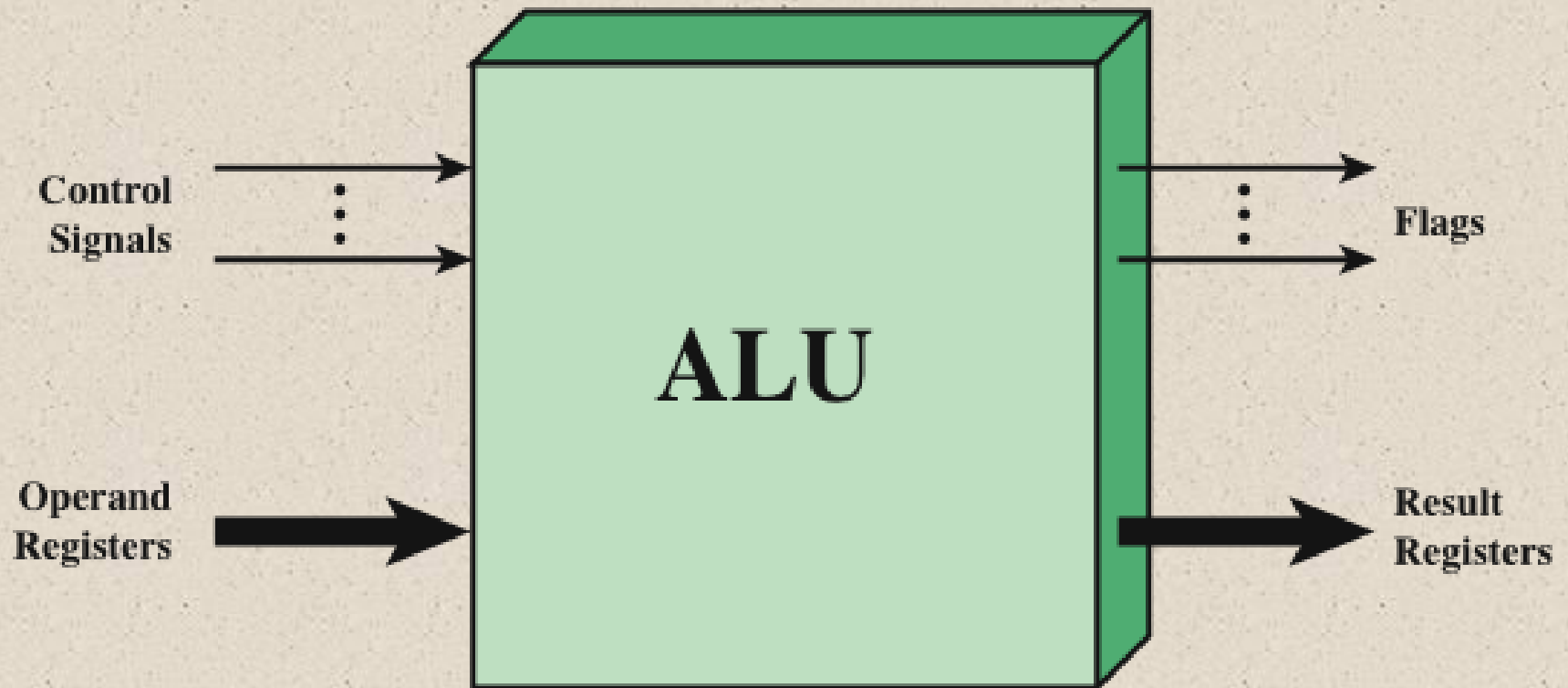
# ALU Inputs and Outputs

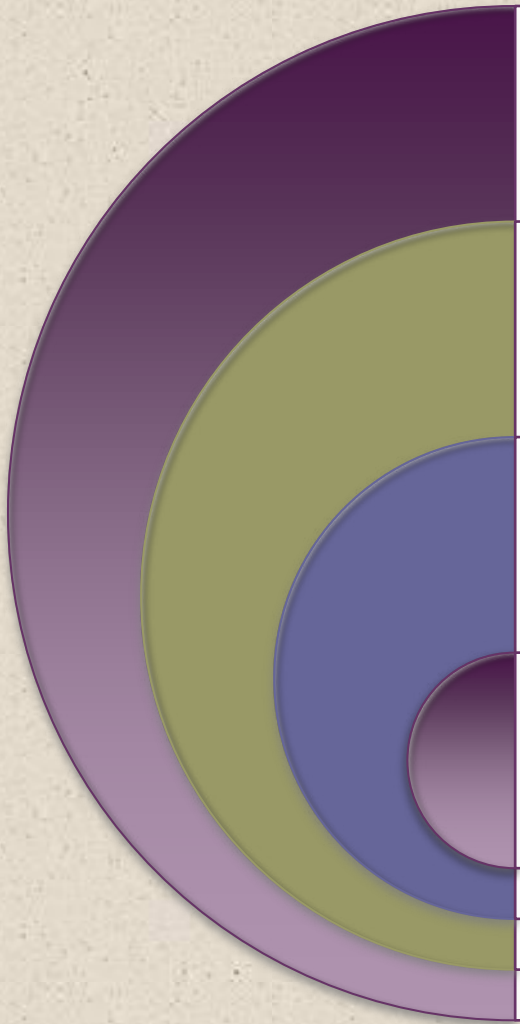

Figure 10.1  ALU Inputs and Outputs

# Integer Representation

- In the binary number system arbitrary numbers can be represented with:
  - The digits zero and one
  - The minus sign (for negative numbers)
  - The period, or *radix point* (for numbers with a fractional component)

- For purposes of computer storage and processing we do not have the benefit of special symbols for the minus sign and radix point

- Only binary digits (0,1) may be used to represent numbers

# Sign-Magnitude Representation

There are several alternative conventions used to represent negative as well as positive integers

- All of these alternatives involve treating the most significant (leftmost) bit in the word as a sign bit
- If the sign bit is 0 the number is positive
- If the sign bit is 1 the number is negative

Sign-magnitude representation is the simplest form that employs a sign bit

Drawbacks:

- Addition and subtraction require a consideration of both the signs of the numbers and their relative magnitudes to carry out the required operation
- There are two representations of 0

Because of these drawbacks, sign-magnitude representation is rarely used in implementing the integer portion of the ALU

# + Twos Complement Representation

- Uses the most significant bit as a sign bit

- Differs from sign-magnitude representation in the way that the other bits are interpreted

| Range | $-2_{n-1}$ through $2_{n-1} - 1$ |
|---|---|
| Number of Representations of Zero | One |
| Negation | Take the Boolean complement of each bit of the corresponding positive number, then add 1 to the resulting bit pattern viewed as an unsigned integer. |
| Expansion of Bit Length | Add additional bit positions to the left and fill in with the value of the original sign bit. |
| Overflow Rule | If two numbers with the same sign (both positive or both negative) are added, then overflow occurs if and only if the result has the opposite sign. |
| Subtraction Rule | To subtract $B$ from $A$, take the twos complement of $B$ and add it to $A$. |

Table 10.1  Characteristics of Twos Complement Representation and Arithmetic

# Alternative Representations for 4-Bit Integers

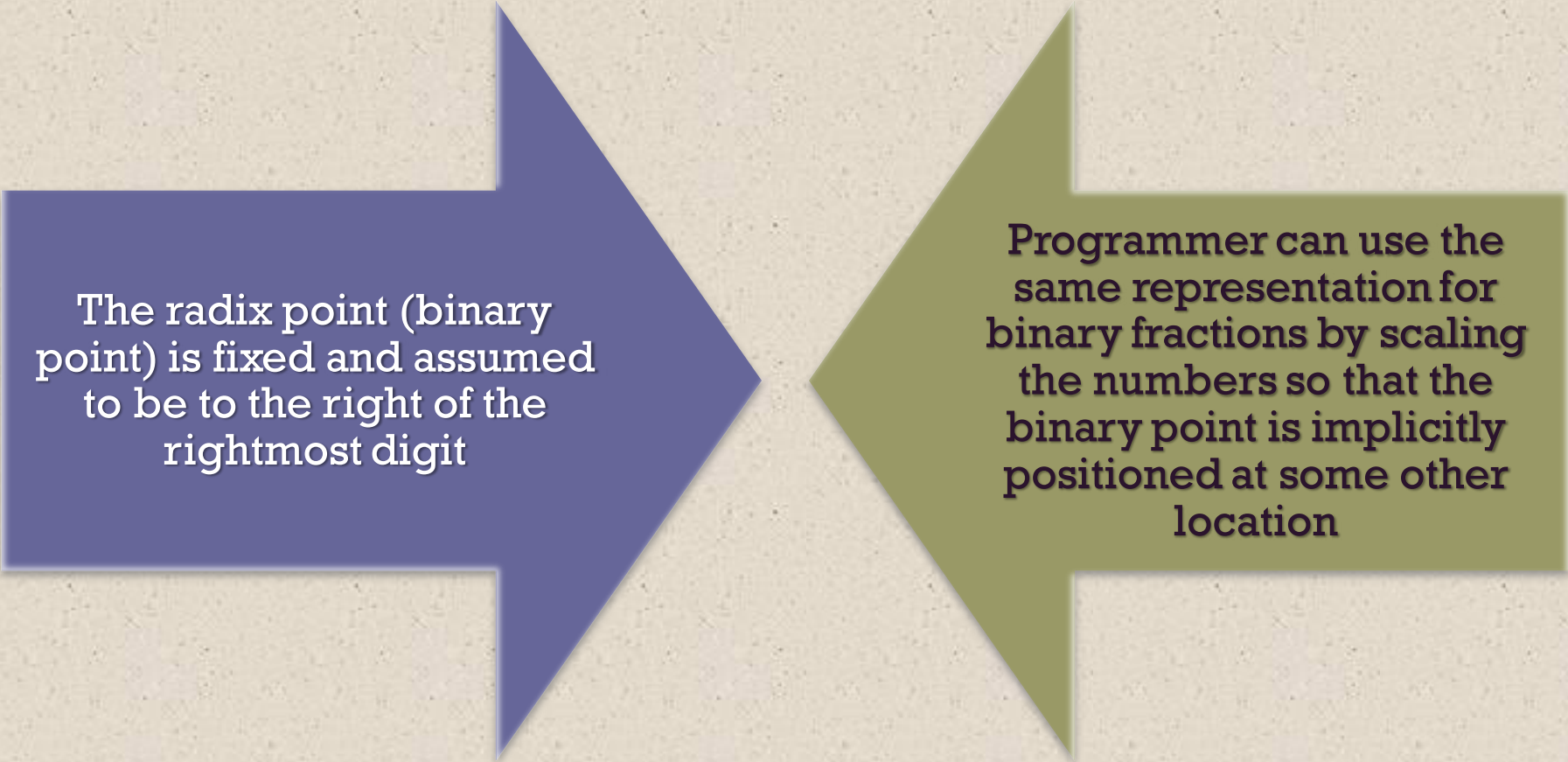| Decimal Representation | Sign-Magnitude Representation | Twos Complement Representation | Biased Representation |
|:---:|:---:|:---:|:---:|
| +8 | — | — | 1111 |
| +7 | 0111 | 0111 | 1110 |
| +6 | 0110 | 0110 | 1101 |
| +5 | 0101 | 0101 | 1100 |
| +4 | 0100 | 0100 | 1011 |
| +3 | 0011 | 0011 | 1010 |
| +2 | 0010 | 0010 | 1001 |
| +1 | 0001 | 0001 | 1000 |
| +0 | 0000 | 0000 | 0111 |
| −0 | 1000 | — | — |
| −1 | 1001 | 1111 | 0110 |
| −2 | 1010 | 1110 | 0101 |
| −3 | 1011 | 1101 | 0100 |
| −4 | 1100 | 1100 | 0011 |
| −5 | 1101 | 1011 | 0010 |
| −6 | 1110 | 1010 | 0001 |
| −7 | 1111 | 1001 | 0000 |
| −8 | — | 1000 | — |

**+**

# Range Extension

- Range of numbers that can be expressed is extended by increasing the bit length

- In sign-magnitude notation this is accomplished by moving the sign bit to the new leftmost position and fill in with zeros

- This procedure will not work for twos complement negative integers
  - Rule is to move the sign bit to the new leftmost position and fill in with copies of the sign bit
  - For positive numbers, fill in with zeros, and for negative numbers, fill in with ones
  - This is called *sign extension*

# Fixed-Point Representation

The radix point (binary point) is fixed and assumed to be to the right of the rightmost digit

Programmer can use the same representation for binary fractions by scaling the numbers so that the binary point is implicitly positioned at some other location

# Negation

- Twos complement operation
  - Take the Boolean complement of each bit of the integer (including the sign bit)
  - Treating the result as an unsigned binary integer, add 1

$$
\begin{array}{rl}
+18 = & 00010010 \text{ (twos complement)} \\
\text{bitwise complement} = & 11101101 \\
+ & \phantom{0000000}1 \\
\hline
& 11101110 = -18
\end{array}
$$

- The negative of the negative of that number is itself:

$$
\begin{array}{rl}
-18 = & 11101110 \text{ (twos complement)} \\
\text{bitwise complement} = & 00010001 \\
+ & \phantom{000000}1 \\
\hline
& 00010010 = +18
\end{array}
$$

# + Negation Special Case 1

|                        |   |           |                    |
|------------------------|---|-----------|--------------------|
| 0                      | = | 00000000  | (twos complement)  |
| Bitwise complement     | = | 11111111  |                    |
| Add 1 to LSB           |   | +        1 |                    |
| Result                 |   | 100000000 |                    |

Overflow is ignored, so:

- 0 = 0

# + Negation Special Case 2

-128     =          10000000    (twos complement)

Bitwise complement   =          01111111

Add 1 to LSB                    +              1

Result                          10000000

So:

-(-128) = -128  X

Monitor MSB (sign bit)

It should change during negation

# Addition

| | |
|---|---|
| ```<br>  1001 = −7<br>+0101 =   5<br>  1110 = −2<br>``` | ```<br>  1100 = −4<br>+0100 =   4<br>10000 =   0<br>``` |
| (a) (−7) + (+5) | (b) (−4) + (+4) |
| ```<br>  0011 = 3<br>+0100 = 4<br>  0111 = 7<br>``` | ```<br>  1100 = −4<br>+1111 = −1<br>11011 = −5<br>``` |
| (c) (+3) + (+4) | (d) (−4) + (−1) |
| ```<br>  0101 = 5<br>+0100 = 4<br>  1001 = Overflow<br>``` | ```<br>  1001 = −7<br>+1010 = −6<br>10011 = Overflow<br>``` |
| (e) (+5) + (+4) | (f) (−7) + (−6) |

**Figure 10.3  Addition of Numbers in Twos Complement Representation**

# OVERFLOW RULE:

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

Overflow

Rule

## SUBTRACTION RULE:

To subtract one number (subtrahend) from another (minuend), take the twos complement (negation) of the subtrahend and add it to the minuend.
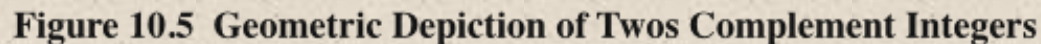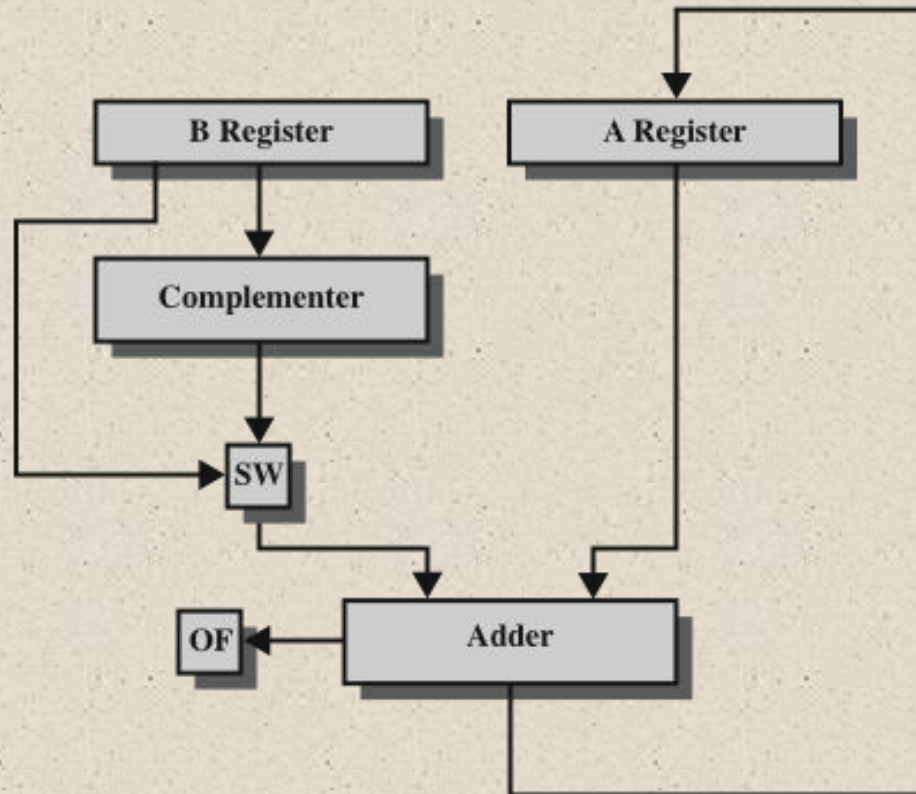
Subtraction

Rule

# Subtraction

```
          0010 =   2                      0101 =   5
         +1001 = −7                      +1110 = −2
          1011 = −5                      10011 =   3

(a) M = 2 = 0010                 (b) M = 5 = 0101
    S = 7 = 0111                     S = 2 = 0010
   −S =      1001                    −S =      1110


          1011 = −5                      0101 = 5
         +1110 = −2                      +0010 = 2
          11001 = −7                      0111 = 7

(c) M =−5 = 1011                 (d) M = 5 = 0101
    S = 2 = 0010                     S =−2 = 1110
   −S =      1110                    −S =      0010


          0111 = 7                       1010 = −6
         +0111 = 7                      +1100 = −4
          1110 = Overflow                10110 = Overflow

(e) M =  7 = 0111                (f) M = −6 = 1010
    S = −7 = 1001                    S =  4 = 0100
   −S =      0111                    −S =      1100
```

**Figure 10.4 Subtraction of Numbers in Twos Complement Representation (M − S)**

# Geometric Depiction of Twos Complement Integers



Figure 10.5  Geometric Depiction of Twos Complement Integers

# Hardware for Addition and Subtraction



OF = overflow bit
SW = Switch (select addition or subtraction)

**Figure 10.6   Block Diagram of Hardware for Addition and Subtraction**

Subtraction using r's complement:

To find M-N in base r, we add M + r's complement of N

Result is M + (rn – N)

1) If M > N then result is M – N + rn (rn is an end carry and can be neglected.

2) If M < N then result is rn –(N-M) which is the r's complement of the result.

# Signed Number Representation

- **Signed Magnitude Method**
  - $N = \pm (a_{n-1} \dots a_0.a_{-1} \dots a_{-m})_r$ is represented as
    $N = (sa_{n-1} \dots a_0.a_{-1} \dots a_{-m})_{rsm}$,                                    (1.6)
    where $s = 0$ if $N$ is positive and $s = r - 1$ otherwise.
  - $N = -(15)_{10}$
  - In binary: $N = -(15)_{10}$    $= -(1111)_2 = (1, 1111)_{2sm}$
  - In decimal: $N = -(15)_{10} = (9, 15)_{10sm}$

- **Complementary Number Systems**
  - **Radix complements** (*r*'s complements)
    $[N]_r = r^n - (N)_r$                                    (1.7)
    where $n$ is the number of digits in $(N)_r$.
  - *Positive full scale*: $r^{n-1} - 1$
  - *Negative full scale*: $-r^n - 1$
  - **Diminished radix complements** (*r-1*'s complements)

    $[N]_{r-1} = r_n - (N)_r - 1$

# + Radix Complement Number Systems (1)

- Two's complement of $(N)_2 = (101001)_2$

  $[N]_2 = 2^6 - (101001)_2 = (1000000)_2 - (101001)_2 = (010111)_2$

- $(N)_2 + [N]_2 = (101001)_2 + (010111)_2 = (1000000)_2$

  If we discard the carry, $(N)_2 + [N]_2 = 0$.

  Hence, $[N]_2$ can be used to represent $-(N)_2$.

- $[\ [N]_2\ ]_2 = [(010111)_2]_2 = (1000000)_2 - (010111)_2 = (101001)_2 = (N)_2$.

- Two's complement of $(N)_2 = (1010)_2$ for $n = 6$

  $[N]_2 = (1000000)_2 - (1010)_2 = (110110)_2$.

# + Radix Complement Number Systems (2)

- **_Algorithm 1.4_  Find $[N]_r$ given $(N)_r$ .**
  - Copy the digits of $N$, beginning with the LSD and proceeding toward the MSD until the first nonzero digit, $a_i$, has been reached
  - Replace $a_i$ with $r - a_i$ .
  - Replace each remaining digit $a_j$ , of $N$ by $(r - 1) - a_j$ until the MSD has been replaced.

- **_Example_**: 10's complement of $(56700)_{10}$ is $(43300)_{10}$

- **_Example_**: 2's complement of $(10100)_2$ is $(01100)_2$.

- **_Example_**: 2's complement of $N = (10110)_2$ for $n = 8$.
  - Put three zeros in the MSB position and apply algorithm 1.4
  - $N = 00010110$
  - $[N]_2 = (11101010)_2$

- The same rule applies to the case when $N$ contains a radix point.

# + Radix Complement Number Systems (3)

- **Algorithm 1.5 Find $[N]_r$ given $(N)_r$.**
  - First replace each digit, $a_k$, of $(N)_r$ by $(r - 1) - a_k$ and then add 1 to the resultant.

- For binary numbers ($r = 2$), complement each digit and add 1 to the result.

- **Example**: Find 2's complement of $N = (01100101)_2$.

      $N = 01100101$
          10011010  Complement the bits
                +1  Add 1
      $[N]_2 = (10011011)_{10}$

- **Example**: Find 10's complement of $N = (40960)_{10}$

      $N = 40960$
          59039  Complement the bits
              +1  Add 1
      $[N]_2 = (59040)_{10}$

# + Radix Complement Number Systems (4)

- **Two's complement number system** :
  - Positive number : $0 \leq N \leq 2^{n-1} - 1$
    - $N = +(a_{n-2}, ..., a_0)_2 = (0, a_{n-2}, ..., a_0)_{2cns}$,

      where                    .
  - Negative number: $-1 \geq N \geq -2^{n-1}$
    - $N = (a_{n-1}, a_{n-2}, ..., a_0)_2$
    - $-N = [a_{n-1}, a_{n-2}, ..., a_0]_2$  (two's complement of $N$),

      where                    .
  - **Example**: Two's complement number system representation of $\pm (N)_2$
    when $(N)_2 = (1011001)_2$ for $n = 8$:
    - $+(N)_2 = (0, 1011001)_{2cns}$
    - $-(N)_2 = [+(N)_2]_2 = [0, 1011001]_2 = (1, 0100111)_{2cns}$

# + Radix Complement Number Systems (5)

■ **Example**: Two's complement number system representation of -$(18)_{10}$ , $n$ = 8:

   ■ +$(18)_{10}$ = $(0, 0010010)_{2cns}$

   ■ -$(18)_{10}$ = $[0, 0010010]_2$ = $(1, 1101110)_{2cns}$

■ **Example**: Decimal representation of $N$ = $(1, 1101000)_{2cns}$

   ■ $N$ = $(1, 1101000)_{2cns}$ = -$[1, 1101000]_2$ = -$(0, 0011000)_{2cns}$ = -$(24)_2$ .

# Radix Complement Arithmetic (1)

- Radix complement number systems are used to convert subtraction to addition, which reduces hardware requirements (only adders are needed).

- $A - B = A + (-B)$ (add $r$'s complement of $B$ to $A$)

- Range of numbers in two's complement number system:
  $$-2^{n-1} \leq N \leq 2^{n-1} - 1$$
  , where $n$ is the number of bits.

- $2^{n-1} - 1 = (0, 11 \ldots 1)_{2cns}$ and $-2^{n-1} = (1, 00 \ldots 0)_{2cns}$

- If the result of an operation falls outside the range, an ***overflow condition*** is said to occur and the result is not valid.

- Consider three cases:
  - $A = B + C$,
  - $A = B - C$,
  - $A = -B - C$,
  
  (where $B \geq 0$ and $C \geq 0$.)

# + Radix Complement Arithmetic (2)

- **_Case 1_**: $A = B + C$
  - $(A)_2 = (B)_2 + (C)_2$
  - If $A > 2^{n-1} - 1$ (**_overflow_**), it is detected by the $n$th bit, which is set to 1.

  - **_Example_**: $(7)_{10} + (4)_{10} = ?$ using 5-bit two's complement arithmetic.
    - $+ (7)_{10} = +(0111)_2 = (0, 0111)_{2cns}$
    - $+ (4)_{10} = +(0100)_2 = (0, 0100)_{2cns}$
    - $(0, 0111)_{2cns} + (0, 0100)_{2cns} = (0, 1011)_{2cns} = +(1011)_2 = +(11)_{10}$
    - No overflow.
  - **_Example_**: $(9)_{10} + (8)_{10} = ?$
    - $+ (9)_{10} = +(1001)_2 = (0, 1001)_{2cns}$
    - $+ (8)_{10} = +(1000)_2 = (0, 1000)_{2cns}$
    - $(0, 1001)_{2cns} + (0, 1000)_{2cns} = (\boldsymbol{1}, 0001)_{2cns}$ (**_overflow_**)

# + Radix Complement Arithmetic (3)

- **Case 2**: $A = B - C$
  - $A = (B)_2 + (-(C)_2) = (B)_2 + [C]_2 = (B)_2 + 2^n - (C)_2 = 2^n + (B - C)_2$
  - If $B \geq C$, then $A \geq 2^n$ and the carry is discarded.
  - So, $(A)_2 = (B)_2 + [C]|_{\text{carry discarded}}$
  - If $B < C$, then $A = 2^n - (C - B)_2 = [C - B]_2$ or $A = -(C - B)_2$ (no carry in this case).
  - No overflow for Case 2.

  - **Example**: $(14)_{10} - (9)_{10} = ?$
    - Perform $(14)_{10} + (-(9)_{10})$
    - $(14)_{10} = +(1110)_2 = (0, 1110)_{2cns}$
    - $-(9)_{10} = -(1001)_2 = (1, 0111)_{2cns}$
    - $(14)_{10} - (9)_{10} = (0, 1110)_{2cns} + (1, 0111)_{2cns} = (0, 0101)_{2cns} + carry$
      $= +(0101)_2 = +(5)_{10}$

# + Radix Complement Arithmetic (4)

- **Example**: $(9)_{10} - (14)_{10} = ?$
  - Perform $(9)_{10} + (-(14)_{10})$
  - $(9)_{10} = +(1001)_2 = (0, 1001)_{2cns}$
  - $-(14)_{10} = -(1110)_2 = (1, 0010)_{2cns}$
  - $(9)_{10} - (14)_{10} = (0, 1001)_{2cns} + (1, 0010)_{2cns} = (1, 1011)_{2cns}$
    $= -(0101)_2 = -(5)_{10}$

<br>

- **Example**: $(0, 0100)_{2cns} - (1, 0110)_{2cns} = ?$
  - Perform $(0, 0100)_{2cns} + (-(1, 0110)_{2cns})$
  - $-(1, 0110)_{2cns} =$ two's complement of $(1, 0110)_{2cns}$
    $= (0, 1010)_{2cns}$
  - $(0, 0100)_{2cns} - (1, 0110)_{2cns} = (0, 0100)_{2cns} + (0, 1010)_{2cns}$
    $= (0, 1110)_{2cns} = +(1110)_2 = +(14)_{10}$
  - $+(4)_{10} - (-(10)_{10}) = +(14)_{10}$

# + Radix Complement Arithmetic (5)

- **Case 3**: $A = -B - C$
  - $A = [B]_2 + [C]_2 = 2^n - (B)_2 + 2^n - (C)_2 = 2^n + 2^n - (B + C)_2 = 2^n + [B + C]_2$
  - The carry bit ($2^n$) is discarded.
  - An overflow can occur, in which case the sign bit is 0.

  - **Example**: $-(7)_{10} - (8)_{10} = ?$
    - Perform $(-(7)_{10}) + (-(8)_{10})$
    - $-(7)_{10} = -(0111)_2 = (1, 1001)_{2cns}$ , $-(8)_{10} = -(1000)_2 = (1, 1000)_{2cns}$
    - $-(7)_{10} - (8)_{10} = (1, 1001)_{2cns} + (1, 1000)_{2cns} = (1, 0001)_{2cns} + carry$
          $= -(1111)_2 = -(15)_{10}$
  - **Example**: $-(12)_{10} - (5)_{10} = ?$
    - Perform $(-(12)_{10}) + (-(5)_{10})$
    - $-(12)_{10} = -(1100)_2 = (1, 0100)_{2cns}$ , $-(5)_{10} = -(0101)_2 = (1, 1011)_{2cns}$
    - $-(7)_{10} - (8)_{10} = (1, 0100)_{2cns} + (1, 1011)_{2cns} = (0, 1111)_{2cns} + carry$
    - **Overflow**, because the sign bit is 0.

**+**

# Radix Complement Arithmetic (6)

- ***Example***: $A = (25)_{10}$ and $B = -(46)_{10}$
  - $A = +(25)_{10} = (0,0011001)_{2cns}$ , $-A = (1,1100111)_{2cns}$
  - $B = -(46)_{10} = -(0,0101110)_2 = (1,1010010)_{2cns}$ , $-B = (0,0101110)_{2cns}$

  - $A + B = (0,0011001)_{2cns} + (1,1010010)_{2cns} = (1,1101011)_{2cns} = -(21)_{10}$
  - $A - B = A + (-B) = (0,0011001)_{2cns} + (0,0101110)_{2cns}$
    $$= (0,1000111)_{2cns} = +(71)_{10}$$
  - $B - A = B + (-A) = (1,1010010)_{2cns} + (1,1100111)_{2cns}$
    $$= (1,0111001)_{2cns} + carry = -(0,1000111)_{2cns} = -(71)_{10}$$
  - $-A - B = (-A) + (-B) = (1,1100111)_{2cns} + (0,0101110)_{2cns}$
    $$= (0,0010101)_{2cns} + carry = +(21)_{10}$$
  - Note: Carry bit is discarded.

# + Radix Complement Arithmetic (7)

■ Summary

| Case | Carry | Sign Bit | Condition | Overflow ? |
|------|-------|----------|-----------|------------|
| B + C | 0 | 0 | $B + C \leq 2^{n-1} - 1$ | No |
|       | 0 | 1 | $B + C > 2^{n-1} - 1$ | Yes |
| B - C | 1 | 0 | $B \leq C$ | No |
|       | 0 | 1 | $B > C$ | No |
| -B - C | 1 | 1 | $-(B + C) \geq -2^{n-1}$ | No |
|        | 1 | 0 | $-(B + C) < -2^{n-1}$ | Yes |

■ When numbers are represented using two's complement number system:

- ■ Addition: Add two numbers.
- ■ Subtraction: Add two's complement of the subtrahend to the minuend.
- ■ Carry bit is discarded, and overflow is detected as shown above.

- ■ Radix complement arithmetic can be used for any radix.

# + Diminished Radix Complement Number systems (1)

- **Diminished radix complement** $[N]_{r-1}$ of a number $(N)_r$ is:

  $[N]_{r-1} = r^n - (N)_r - 1$                                              (1.10)

- **One's complement** $(r = 2)$:

  $[N]_{2-1} = 2^n - (N)_2 - 1$                                       (1.11)

- **Example**: One's complement of $(01100101)_2$

  $[N]_{2-1} = 2^8 - (01100101)_2 - 1$

         $= (100000000)_2 - (01100101)_2 - (00000001)_2$

         $= (10011011)_2 - (00000001)_2$

         $= (10011010)_2$

# **+** Diminished Radix Complement Number systems (2)

- ***Example***: Nine's complement of (40960)

  $[N]_{2-1} = 10^5 - (40960)_{10} - 1$

  $\quad = (100000)_{10} - (40960)_{10} - (00001)_{10}$

  $\quad = (59040)_{10} - (00001)_{10}$

  $\quad = (59039)_{10}$

- ***Algorithm 1.6*** **Find $[N]_{r-1}$ given $(N)_r$ .**

  Replace each digit $a_i$ of $(N)_r$ by $r - 1 - a$. Note that when $r = 2$, this simplifies to complementing each individual bit of $(N)_r$ .

- Radix complement and diminished radix complement of a number ($N$):

  $[N]_r = [N]_{r-1} + 1$ $\qquad\qquad\qquad\qquad\qquad\qquad$ (1.12)

# + Diminished Radix Complement Arithmetic (1)

- Operands are represented using diminished radix complement number system.

- The carry, if any, is added to the result (***end-around carry***).

- ***Example***: Add  $+(1001)_2$ and $-(0100)_2$ .
  One's complement of +(1001) = 01001
  One's complement of -(0100) = 11011
  01001 + 11011 = 100100 (carry)
  Add the carry to the result: correct result is 00101.

- ***Example***: Add  $+(1001)_2$ and $-(1111)_2$ .
  One's complement of +(1001) = 01001
  One's complement of -(1111) = 10000
  01001 + 10000 = 11001 (no carry, so this is the correct result).

# + Diminished Radix Complement Arithmetic (2)

■ **Example**: Add -$(1001)_2$ and -$(0011)_2$ .
One's complement of the operands are: 10110 and 11100
$10110 + 11100 = 110010$ (carry)
Correct result is $10010 + 1 = 10011$.

■ **Example**: Add +$(75)_{10}$ and -$(21)_{10}$ .
Nine's complements of the operands are: 075 and 978
$075 + 978 = 1053$ (carry)
Correct result is $053 + 1 = 054$

■ **Example**: Add +$(21)_{10}$ and -$(75)_{10}$ .
Nine's complements of the operands are: 021 and 924
$021 + 924 = 945$ (no carry, so this is the correct result).

**Example (3):**

Using 2's complement, subtract 1010100 – 1000011

$$X \quad - \quad Y$$

| | | | |
|---|---|---|---|
| $X$ | = | | 1010100 |
| 2's complement of $Y$ | = + | | 0111101 |
| Sum | = | | 10010001 |
| Discard end carry $2^7$ | = - | | 10000000 |
| Answer: $X$ - $Y$ | = | | 0010001 |

**Example (4):**

Using 2's complement, subtract 1000011 – 1010100

$$Y \quad - \quad X$$

| | | | |
|---|---|---|---|
| $Y$ | = | | 1000011 |
| 2's complement of $X$ | = + | | 0101100 |
| Sum | = | | 1101111 |
| No end carry. | | | |
| Answer: $Y – X$ - (2's complement of 1101111) | = | | -0010001 |

**Example (5):** Using 1's complement, subtract X – Y = 1010100 – 1000011

| | | |
|---|---|---|
| $X$ | = | 1010100 |
| 1's complement of $Y$ | = + | 0111100 (+1 End-around carry) |
| Sum = | | 10010000 |
| | | + 1 |
| Answer: $X$ - $Y$ | = | 0010001 |

**Example (6):** Using 1's complement, subtract $Y – X$ 1000011 – 1010100

| | | |
|---|---|---|
| $Y$ | = | 1000011 |
| 1's complement of $X$ | = + | 0101011 |
| Sum | = | 1101110 |
| No end carry. | | |
| Answer: $Y – X$ - (1's complement of 1101110) | = | -0010001 |

# Multiplication



```
  1011        Multiplicand (11)
× 1101        Multiplier (13)
 ─────
  1011    ⎫
 0000     ⎬   Partial products
 1011     ⎪
1011      ⎭
─────────
10001111      Product (143)
```

**Figure 10.7  Multiplication of Unsigned Binary Integers**

# Hardware Implementation of Unsigned Binary Multiplication



(a) Block Diagram

| 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 | Third Cycle |
| 0 | 0110 | 1111 | 1011 | Shift | |
| 1 | 0001 | 1111 | 1011 | Add | Fourth Cycle |
| 0 | 1000 | 1111 | 1011 | Shift | |

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

**Figure 10.8 Hardware Implementation of Unsigned Binary Multiplication**

# Flowchart for Unsigned Binary Multiplication



Figure 10.9 Flowchart for Unsigned Binary Multiplication

# Twos Complement Multiplication

$$
\begin{array}{ll}
\phantom{\times}1011 & \\
\underline{\times 1101} & \\
00001011 & 1011 \times 1 \times 2^0 \\
00000000 & 1011 \times 0 \times 2^1 \\
00101100 & 1011 \times 1 \times 2^2 \\
\underline{01011000} & 1011 \times 1 \times 2^3 \\
10001111 & \\
\end{array}
$$

**Figure 10.10 Multiplication of Two Unsigned 4-Bit Integers Yielding an 8-Bit Result**

# Comparison

```
   1001  (9)                        1001  (−7)
  ×0011  (3)                       ×0011  (3)
00001001  1001 × 2⁰       11111001  (−7) × 2⁰ = (−7)
00010010  1001 × 2¹       11110010  (−7) × 2¹ = (−14)
00011011  (27)                   11101011  (−21)
```

(a) Unsigned integers          (b) Twos complement integers

**Figure 10.11  Comparison of Multiplication of Unsigned and Twos Complement Integers**

**Booth's**

**Algorithm**

**Figure 10.12  Booth's Algorithm for Twos Complement Multiplication**

# Example of Booth's Algorithm

| A | Q | $Q_{-1}$ | M | | |
|---|---|---|---|---|---|
| 0000 | 0011 | 0 | 0111 | Initial Values | |
| 1001 | 0011 | 0 | 0111 | A ← A – M | First |
| 1100 | 1001 | 1 | 0111 | Shift | Cycle |
| 1110 | 0100 | 1 | 0111 | Shift | Second Cycle |
| 0101 | 0100 | 1 | 0111 | A ← A + M | Third |
| 0010 | 1010 | 0 | 0111 | Shift | Cycle |
| 0001 | 0101 | 0 | 0111 | Shift | Fourth Cycle |

**Figure 10.13  Example of Booth's Algorithm (7✕ 3)**

# Examples Using Booth's Algorithm

```
    0111                              0111
   ×0011        (0)                  ×1101        (0)
11111001        1—0              11111001        1—0
0000000         1—1              0000111         0—1
000111          0—1              111001          1—0
00010101        (21)             11101011        (−21)
```

(a) (7) × (3) = (21)                (b) (7) × (−3) = (−21)

```
    1001                              1001
   ×0011        (0)                  ×1101        (0)
00000111        1—0              00000111        1—0
0000000         1—1              1111001         0—1
111001          0—1              000111          1—0
11101011        (−21)           00010101        (21)
```

(c) (−7) × (3) = (−21)              (d) (−7) × (−3) = (21)

**Figure 10.14  Examples Using Booth's Algorithm**

# Division



**Figure 10.15  Example of Division of Unsigned Binary Integers**

# Flowchart for Unsigned Binary Division



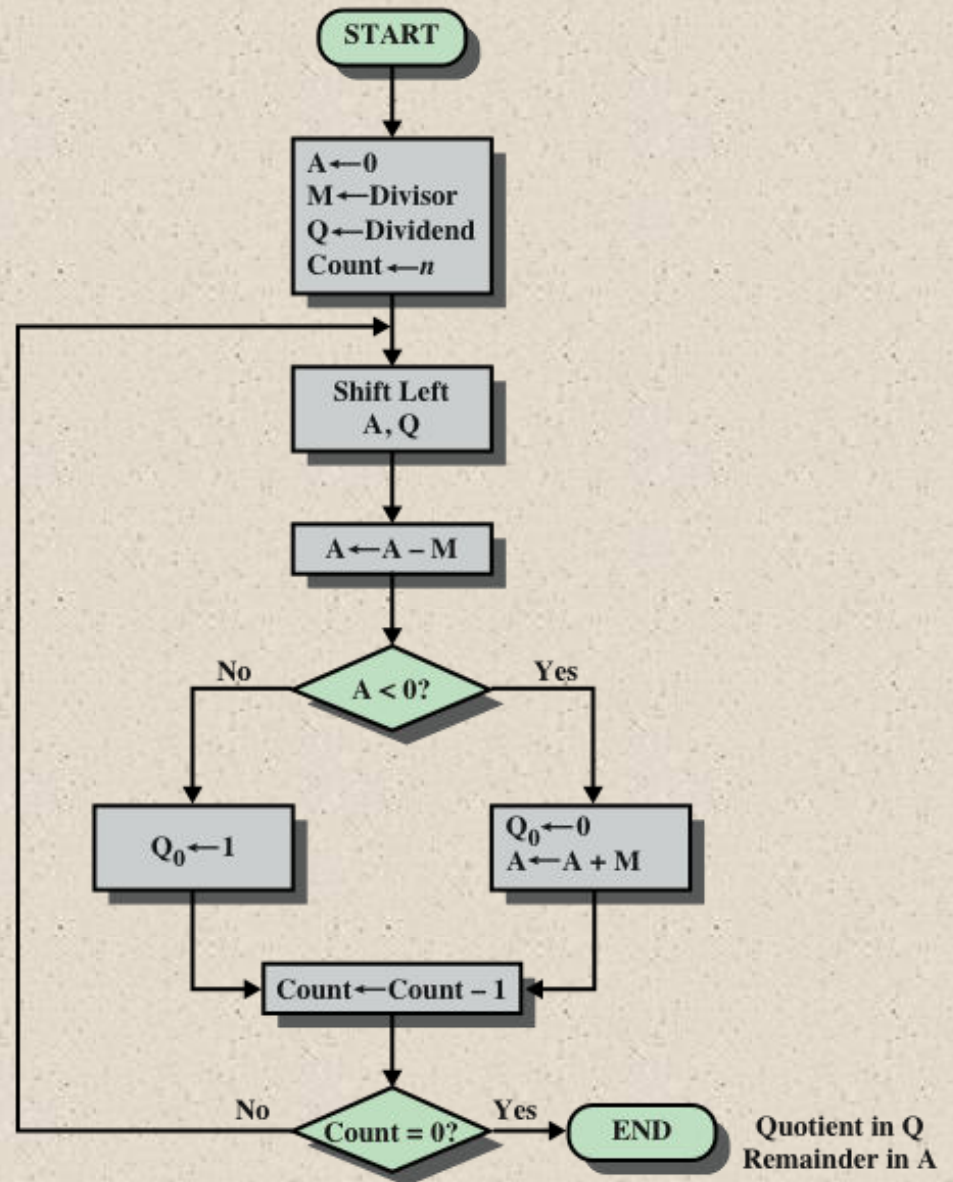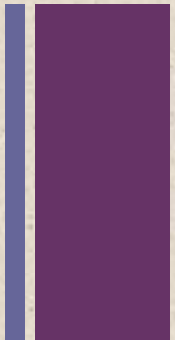**Figure 10.16  Flowchart for Unsigned Binary Division**

# Example of Restoring Twos Complement Division

| A | Q | |
|---|---|---|
| 0000 | 0111 | Initial value |
| 0000 | 1110 | Shift |
| 1101 | | Use twos complement of 0011 for subtraction |
| 1101 | | Subtract |
| 0000 | 1110 | Restore, set $Q_0 = 0$ |
| 0001 | 1100 | Shift |
| 1101 | | |
| 1110 | | Subtract |
| 0001 | 1100 | Restore, set $Q_0 = 0$ |
| 0011 | 1000 | Shift |
| 1101 | | |
| 0000 | 1001 | Subtract, set $Q_0 = 1$ |
| 0001 | 0010 | Shift |
| 1101 | | |
| 1110 | | Subtract |
| 0001 | 0010 | Restore, set $Q_0 = 0$ |

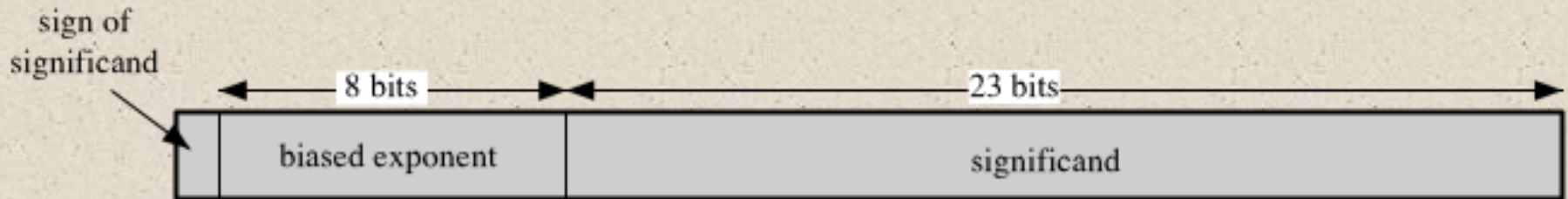**Figure 10.17  Example of Restoring Twos Complement Division (7/3)**

# + Floating-Point Representation

## Principles

- With a fixed-point notation it is possible to represent a range of positive and negative integers centered on or near 0

- By assuming a fixed binary or radix point, this format allows the representation of numbers with a fractional component as well

- Limitations:
  - Very large numbers cannot be represented nor can very small fractions
  - The fractional part of the quotient in a division of two large numbers could be lost

# Typical 32-Bit Floating-Point Format

sign of significand

← 8 bits → ← 23 bits →

| biased exponent | significand |

(a) Format

(b) Examples

$$1.1010001 \times 2^{10100} = 0\ 10010011\ 10100010000000000000000 = 1.6328125 \times 2^{20}$$
$$-1.1010001 \times 2^{10100} = 1\ 10010011\ 10100010000000000000000 = -1.6328125 \times 2^{20}$$
$$1.1010001 \times 2^{-10100} = 0\ 01101011\ 10100010000000000000000 = 1.6328125 \times 2^{-20}$$
$$-1.1010001 \times 2^{-10100} = 1\ 01101011\ 10100010000000000000000 = -1.6328125 \times 2^{-20}$$

**Figure 10.18   Typical 32-Bit Floating-Point Format**

The closest binary number to $\underline{Y}$ that can be stored by computer in 32 bits is:

$$(-1)^s \; \frac{J \,(\text{in binary})}{2^{23}} \times 2^P$$

where

$$s = 0 \text{ if } \underline{y} \le 0 \text{ and } s = 1 \text{ if } y > 1$$

$$\rightarrow Z = |y|$$

$$\rightarrow P = \text{Floor}\left(\log_2 Z\right) = \text{Floor}\left(\frac{\log Z}{\log 2}\right)$$

$$J = \text{Round}\left(Z \times 2^{23-P}\right)$$

a) Convert the given IEEE 754 formatted 32 –bit floating point number in to decimal

    1    10111011      10110000000000000000000

b) Define Normalization. Give two examples.

Give the flow chart for addition and subtraction of two floating-point numbers.

Show the IEEE 754 binary representation of the number $(-0.4375)_{ten}$ in single precision.

**+**

# Floating-Point

## Significand

- The final portion of the word

- Any floating-point number can be expressed in many ways

---

The following are equivalent, where the significand is expressed in binary form:

$$0.110 * 2^5$$
$$110 * 2^2$$
$$0.0110 * 2^6$$

---

- *Normal number*
  - The most significant digit of the significand is nonzero

# + Expressible Numbers

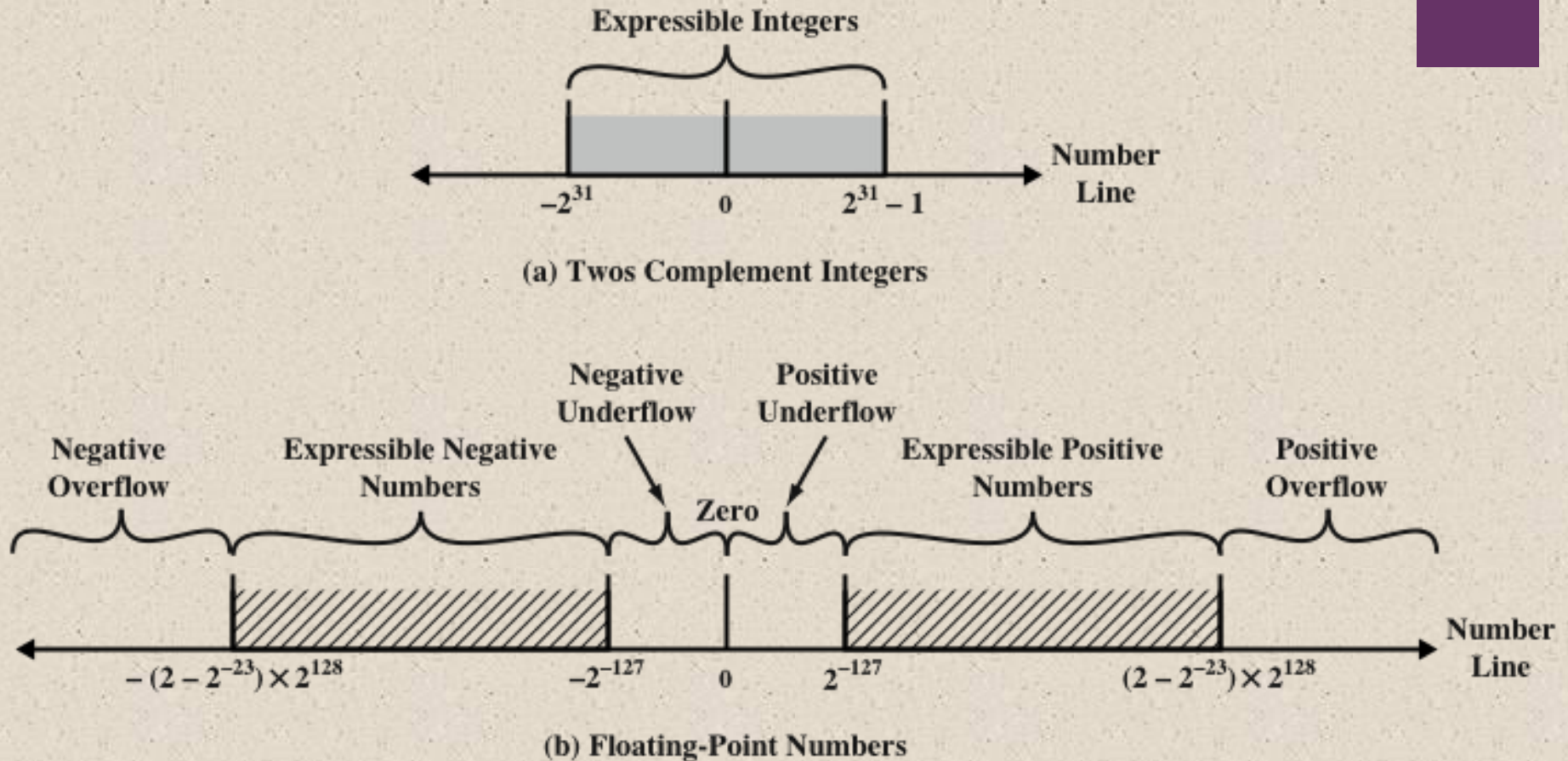**Expressible Integers**

$-2^{31}$    $0$    $2^{31} - 1$    Number Line

**(a) Twos Complement Integers**

Negative Overflow    Expressible Negative Numbers    Negative Underflow    Positive Underflow    Expressible Positive Numbers    Positive Overflow

Zero

$-(2 - 2^{-23}) \times 2^{128}$    $-2^{-127}$    $0$    $2^{-127}$    $(2 - 2^{-23}) \times 2^{128}$    Number Line

**(b) Floating-Point Numbers**

Figure 10.19  Expressible Numbers in Typical 32-Bit Formats

# Density of Floating-Point Numbers



Figure 10.20    Density of Floating-Point Numbers

# IEEE Standard 754

Most important floating-point representation is defined

Standard was developed to facilitate the portability of programs from one processor to another and to encourage the development of sophisticated, numerically oriented programs

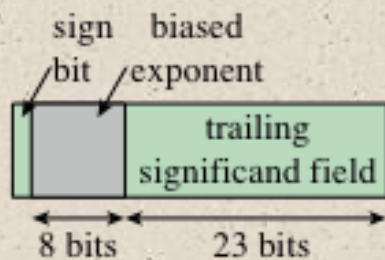Standard has been widely adopted and is used on virtually all contemporary processors and arithmetic coprocessors

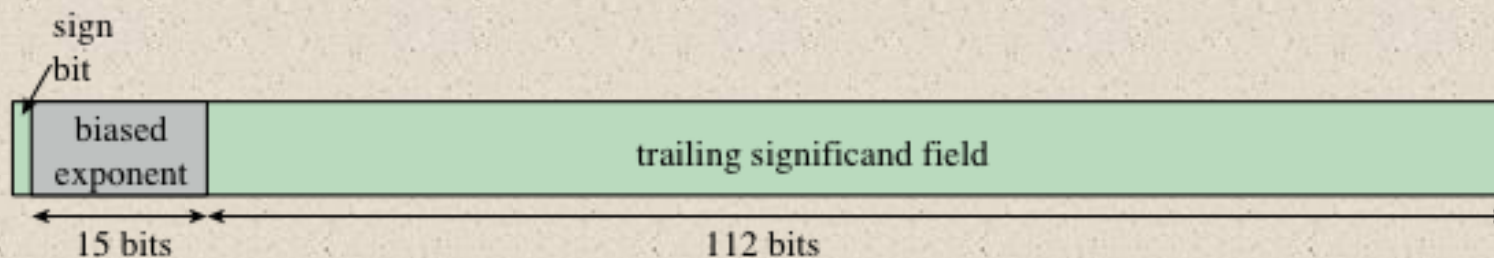IEEE 754-2008 covers both binary and decimal floating-point representations

# IEEE 754-2008

- Defines the following different types of floating-point formats:

- Arithmetic format
  - All the mandatory operations defined by the standard are supported by the format. The format may be used to represent floating-point operands or results for the operations described in the standard.

- Basic format
  - This format covers five floating-point representations, three binary and two decimal, whose encodings are specified by the standard, and which can be used for arithmetic. At least one of the basic formats is implemented in any conforming implementation.

- Interchange format
  - A fully specified, fixed-length binary encoding that allows data interchange between different platforms and that can be used for storage.

# IEEE 754 Formats



(a) binary32 format

sign bit, biased exponent — 8 bits, trailing significand field — 23 bits

(b) binary64 format

sign bit, biased exponent — 11 bits, trailing significand field — 52 bits

(c) binary128 format

sign bit, biased exponent — 15 bits, trailing significand field — 112 bits

Figure 10.21 IEEE 754 Formats

Table 10.3

IEEE 754

Format
Parameters

| Parameter | Format | | |
|---|---|---|---|
| | binary32 | binary64 | binary128 |
| Storage width (bits) | 32 | 64 | 128 |
| Exponent width (bits) | 8 | 11 | 15 |
| Exponent bias | 127 | 1023 | 16383 |
| Maximum exponent | 127 | 1023 | 16383 |
| Minimum exponent | $-126$ | $-1022$ | $-16382$ |
| Approx normal number range (base 10) | $10_{-38}, 10_{+38}$ | $10_{-308}, 10_{+308}$ | $10_{-4932}, 10_{+4932}$ |
| Trailing significand width (bits)* | 23 | 52 | 112 |
| Number of exponents | 254 | 2046 | 32766 |
| Number of fractions | $2_{23}$ | $2_{52}$ | $2_{112}$ |
| Number of values | $1.98 \times 2_{31}$ | $1.99 \times 2_{63}$ | $1.99 \times 2_{128}$ |
| Smallest positive normal number | $2_{-126}$ | $2_{-1022}$ | $2_{-16362}$ |
| Largest positive normal number | $2_{128} - 2_{104}$ | $2_{1024} - 2_{971}$ | $2_{16384} - 2_{16271}$ |
| Smallest subnormal magnitude | $2_{-149}$ | $2_{-1074}$ | $2_{-16494}$ |

* not including implied bit and not including sign bit

# + Additional Formats

## Extended Precision Formats

- Provide additional bits in the exponent (extended range) and in the significand (extended precision)

- Lessens the chance of a final result that has been contaminated by excessive roundoff error

- Lessens the chance of an intermediate overflow aborting a computation whose final result would have been representable in a basic format

- Affords some of the benefits of a larger basic format without incurring the time penalty usually associated with higher precision

## Extendable Precision Format

- Precision and range are defined under user control

- May be used for intermediate calculations but the standard places no constraint or format or length

Check it out ~2

# Table 10.4
# IEEE Formats

| Format | Format Type | | |
| --- | --- | --- | --- |
| | Arithmetic Format | Basic Format | Interchange Format |
| binary16 | | | X |
| binary32 | X | X | X |
| binary64 | X | X | X |
| binary128 | X | X | X |
| binary{k} (k = n × 32 for n > 4) | X | | X |
| decimal64 | X | X | X |
| decimal128 | X | X | X |
| decimal{k} (k = n × 32 for n > 4) | X | | X |
| extended precision | X | | |
| extendable precision | X | | |

Table 10.4   IEEE Formats

# Interpretation of IEEE 754 Floating-Point Numbers

## (a) binary 32 format

| | Sign | Biased exponent | Fraction | Value |
|---|---|---|---|---|
| positive zero | 0 | 0 | 0 | 0 |
| negative zero | 1 | 0 | 0 | $-0$ |
| plus infinity | 0 | all 1s | 0 | $\infty$ |
| Minus infinity | 1 | all 1s | 0 | $-\infty$ |
| quiet NaN | 0 or 1 | all 1s | $\neq 0$; first bit $= 1$ | qNaN |
| signaling NaN | 0 or 1 | all 1s | $\neq 0$; first bit $= 0$ | sNaN |
| positive normal nonzero | 0 | $0 < e < 255$ | f | $2_{e-127}(1.f)$ |
| negative normal nonzero | 1 | $0 < e < 255$ | f | $-2_{e-127}(1.f)$ |
| positive subnormal | 0 | 0 | $f \neq 0$ | $2_{e-126}(0.f)$ |
| negative subnormal | 1 | 0 | $f \neq 0$ | $-2_{e-126}(0.f)$ |

Table 10.5   Interpretation of IEEE 754 Floating-Point Numbers (page 1 of 3)

# Interpretation of IEEE 754 Floating-Point Numbers

## (b) binary 64 format

| | Sign | Biased exponent | Fraction | Value |
|---|---|---|---|---|
| positive zero | 0 | 0 | 0 | 0 |
| negative zero | 1 | 0 | 0 | $-0$ |
| plus infinity | 0 | all 1s | 0 | $\infty$ |
| Minus infinity | 1 | all 1s | 0 | $-\infty$ |
| quiet NaN | 0 or 1 | all 1s | $\neq 0$; first bit $= 1$ | qNaN |
| signaling NaN | 0 or 1 | all 1s | $\neq 0$; first bit $= 0$ | sNaN |
| positive normal nonzero | 0 | $0 < e < 2047$ | f | $2_{e-1023}(1.f)$ |
| negative normal nonzero | 1 | $0 < e < 2047$ | f | $-2_{e-1023}(1.f)$ |
| positive subnormal | 0 | 0 | $f \neq 0$ | $2_{e-1022}(0.f)$ |
| negative subnormal | 1 | 0 | $f \neq 0$ | $-2_{e-1022}(0.f)$ |

Table 10.5  Interpretation of IEEE 754 Floating-Point Numbers (page 2 of 3)

# Interpretation of IEEE 754 Floating-Point Numbers

# (c) binary 128 format

| | Sign | Biased exponent | Fraction | Value |
|---|---|---|---|---|
| positive zero | 0 | 0 | 0 | 0 |
| negative zero | 1 | 0 | 0 | $-0$ |
| plus infinity | 0 | all 1s | 0 | $\infty$ |
| minus infinity | 1 | all 1s | 0 | $-\infty$ |
| quiet NaN | 0 or 1 | all 1s | $\neq 0$; first bit = 1 | qNaN |
| signaling NaN | 0 or 1 | all 1s | $\neq 0$; first bit = 0 | sNaN |
| positive normal nonzero | 0 | all 1s | f | $2_{e-16383}(1.f)$ |
| negative normal nonzero | 1 | all 1s | f | $-2_{e-16383}(1.f)$ |
| positive subnormal | 0 | 0 | $f \neq 0$ | $2_{e-16383}(0.f)$ |
| negative subnormal | 1 | 0 | $f \neq 0$ | $-2_{e-16383}(0.f)$ |

Table 10.5  Interpretation of IEEE 754 Floating-Point Numbers (page 3 of 3)

# Table 10.6  Floating-Point Numbers and Arithmetic Operations

| Floating Point Numbers | Arithmetic Operations |
|---|---|
| $X = X_S \times B^{X_E}$ <br><br> $Y = Y_S \times B^{Y_E}$ | $\left. \begin{aligned} X + Y &= \left( X_S \times B^{X_E - Y_E} + Y_S \right) \times B^{Y_E} \\ X - Y &= \left( X_S \times B^{X_E - Y_E} - Y_S \right) \times B^{Y_E} \end{aligned} \right\} X_E \leq Y_E$ <br><br> $X \times Y = \left( X_S \times Y_S \right) \times B^{X_E + Y_E}$ <br><br> $\dfrac{X}{Y} = \left( \dfrac{X_S}{Y_S} \right) \times B^{X_E - Y_E}$ |

Examples:

$X = 0.3 \times 10^2 = 30$
$Y = 0.2 \times 10^3 = 200$

$X + Y = (0.3 \times 10^{2-3} + 0.2) \times 10^3 = 0.23 \times 10^3 = 230$
$X - Y = (0.3 \times 10^{2-3} - 0.2) \times 10^3 = (-0.17) \times 10^3 = -170$
$X \times Y = (0.3 \times 0.2) \times 10^{2+3} = 0.06 \times 10^5 = 6000$
$X \div Y = (0.3 \div 0.2) \times 10^{2-3} = 1.5 \times 10^{-1} = 0.15$
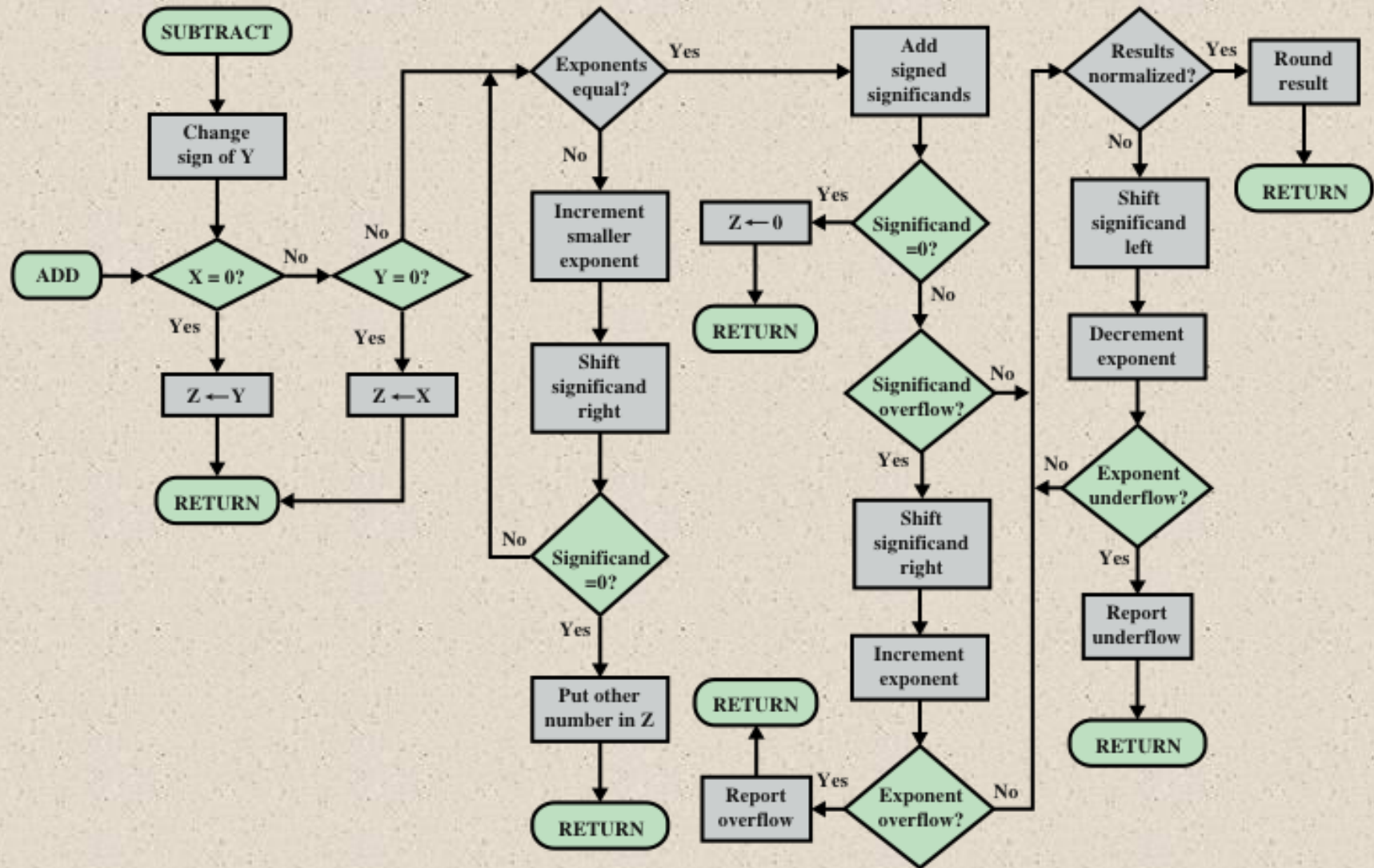
# Floating-Point Addition and Subtraction



Figure 10.22 Floating-Point Addition and Subtraction (Z← X ± Y)

# Floating-Point Multiplication



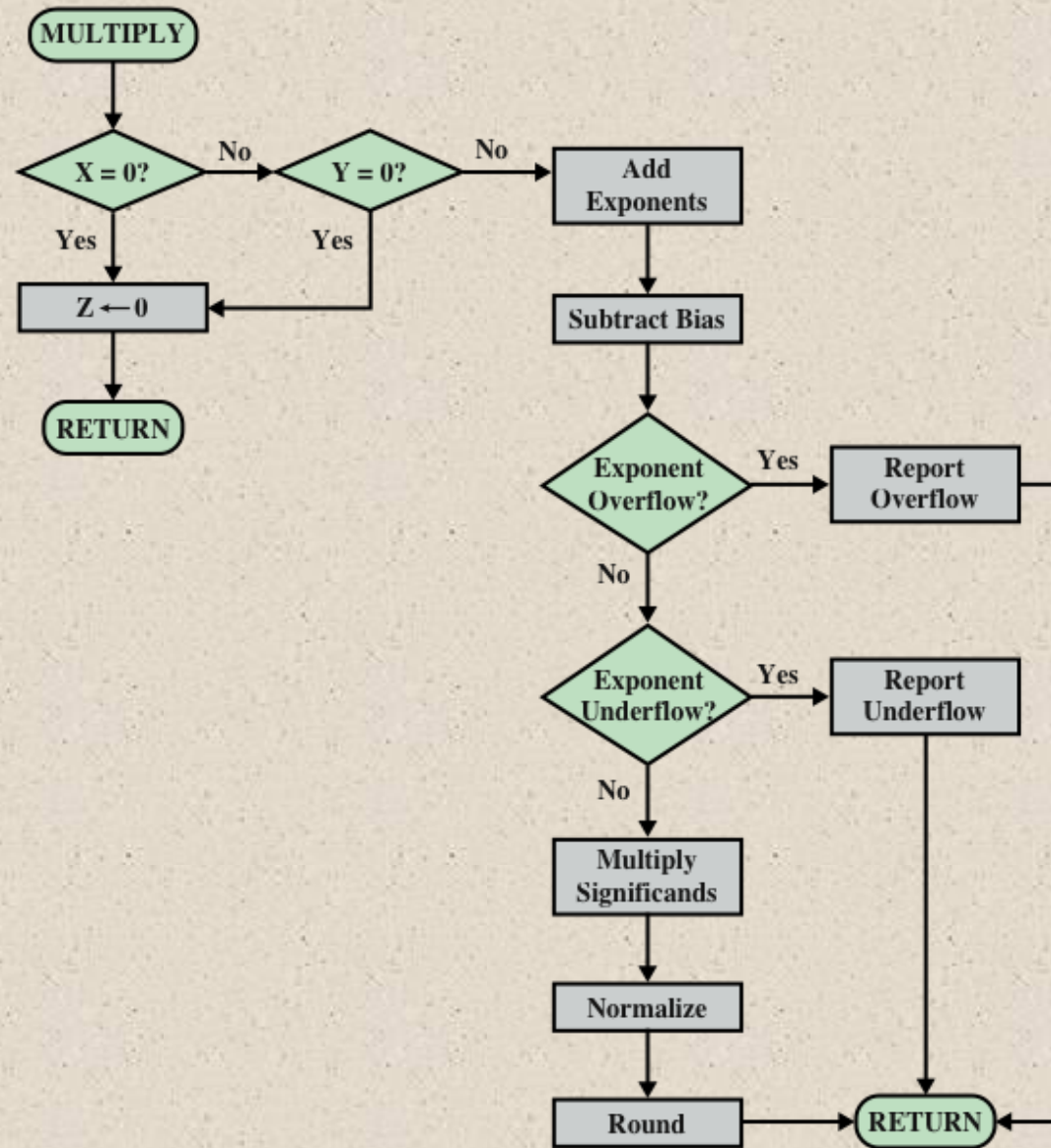**Figure 10.23  Floating-Point Multiplication (Z← X× Y)**

# Floating-Point Division



Figure 10.24 Floating-Point Division (Z← X/Y)

# Precision Considerations

## Guard Bits

$$x = 1.000.....00 \times 2^1$$
$$-y = \underline{0.111.....11} \times 2^1$$
$$z = 0.000.....01 \times 2^1$$
$$= 1.000.....00 \times 2^{-22}$$

(a) Binary example, without guard bits

$$x = 1.000.....00\ 0000 \times 2^1$$
$$-y = \underline{0.111.....11\ 1000} \times 2^1$$
$$z = 0.000.....00\ 1000 \times 2^1$$
$$= 1.000.....00\ 0000 \times 2^{-23}$$

(b) Binary example, with guard bits

$$x = .100000 \times 16^1$$
$$-y = \underline{.0FFFFF} \times 16^1$$
$$z = .000001 \times 16^1$$
$$= .100000 \times 16^{-4}$$

(c) Hexadecimal example, without guard bits

$$x = .100000\ 00 \times 16^1$$
$$-y = \underline{.0FFFFF\ F0} \times 16^1$$
$$z = .000000\ 10 \times 16^1$$
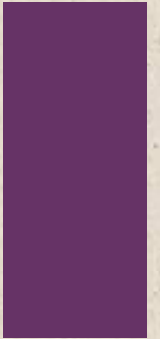$$= .100000\ 00 \times 16^{-5}$$

(d) Hexadecimal example, with guard bits

**Figure 10.25  The Use of Guard Bits**

# **Precision Considerations**

## Rounding

- IEEE standard approaches:

  - Round to nearest:
    - The result is rounded to the nearest representable number.
  - Round toward +∞ :
    - The result is rounded up toward plus infinity.
  - Round toward -∞:
    - The result is rounded down toward negative infinity.
  - Round toward 0:
    - The result is rounded toward zero.

# Interval Arithmetic

- Provides an efficient method for monitoring and controlling errors in floating-point computations by producing two values for each result

- The two values correspond to the lower and upper endpoints of an interval that contains the true result

- The width of the interval indicates the accuracy of the result

- If the endpoints are not representable then the interval endpoints are rounded down and up respectively

- If the range between the upper and lower bounds is sufficiently narrow then a sufficiently accurate result has been obtained

- *Minus infinity* and *rounding to plus* are useful in implementing interval arithmetic

# Truncation

- *Round toward zero*

- Extra bits are ignored

- Simplest technique

- A consistent bias toward zero in the operation
  - Serious bias because it affects every operation for which there are nonzero extra bits

# IEEE Standard for Binary Floating-Point Arithmetic

## Infinity

Is treated as the limiting case of real arithmetic, with the infinity values given the following interpretation:

$$- \infty < \text{(every finite number)} < + \infty$$

For example:

| | |
|---|---|
| $5 + (+ \infty ) = + \infty$ | $5 \div (+ \infty ) = +0$ |
| $5 - (+ \infty ) = - \infty$ | $(+ \infty ) + (+ \infty ) = + \infty$ |
| $5 + (- \infty ) = - \infty$ | $(- \infty ) + (- \infty) = - \infty$ |
| $5 - (- \infty ) = + \infty$ | $(- \infty ) - (+ \infty ) = - \infty$ |
| $5 * (+ \infty ) = + \infty$ | $(+ \infty ) - (- \infty ) = + \infty$ |

# IEEE Standard for Binary Floating-Point Arithmetic
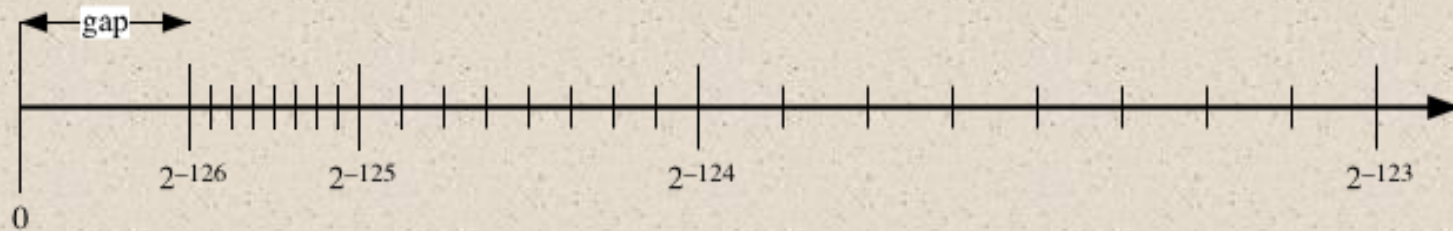
## Quiet and Signaling NaNs

- Signaling NaN signals an invalid operation exception whenever it appears as an operand

- Quiet NaN propagates through almost every arithmetic operation without signaling an exception

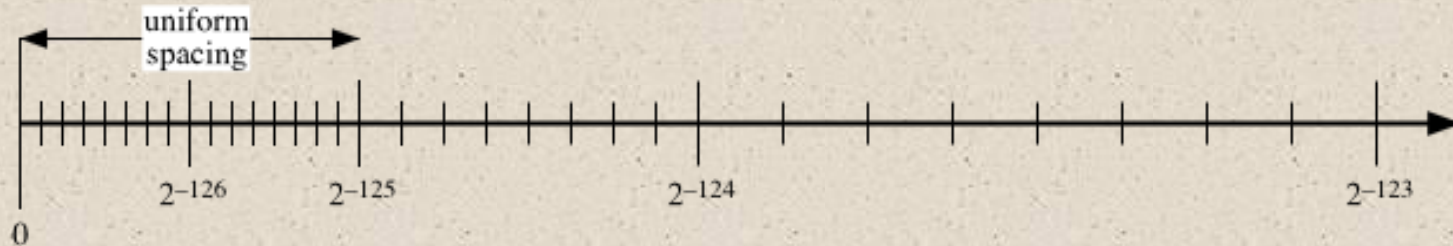| Operation | Quiet NaN Produced by |
|---|---|
| Any | Any operation on a signaling NaN |
| Add or subtract | Magnitude subtraction of infinities: $(+\infty) + (-\infty)$ $(-\infty) + (+\infty)$ $(+\infty) - (+\infty)$ $(-\infty) - (-\infty)$ |
| Multiply | $0 \times \infty$ |
| Division | $\dfrac{0}{0}$ or $\dfrac{\infty}{\infty}$ |
| Remainder | $x$ REM $0$ or $\infty$ REM $y$ |
| Square root | $\sqrt{x}$ where $x < 0$ |

Operations that
Produce a
Quiet NaN

# IEEE Standard for Binary Floating-Point Arithmetic

## Subnormal Numbers



(a) 32-bit format without subnormal numbers

(b) 32-bit format with subnormal numbers

**Figure 10.26 The Effect of IEEE 754 Subnormal Numbers**