

A number system is in general defined by the set of values and a set of rules that gives the mapping between the sequence of digits and their numerical values. A conventional number system is a non-redundant, weighted and positional. A conventional number system is fixed-radix system. When numbers are represented in machines, they have a maximum and minimum value which can be represented. If the number does not lie in this range, an overflow might occur and will produce an incorrect result.

There are two types of number representations: Fixed Point, Floating Point
The radix point “.” is not stored in the register but is understood to be in a fixed position between the k most significant digits and the m least significant digits. For this reason we call such representations fixed-point representations.

To avoid the need to tell explicitly where the radix point is, we introduce the notion of a **unit in the last position (ulp)**, which is the weight of the least significant digit.

Representation of Negative Numbers: **Sign-magnitude method**; **Biased method**; **Radix complement representation method**; **Diminished-radix complement representation method**

Sign Magnitude Method - The sign and magnitude are represented separately, where the sign is represented with the first digit while the remaining n – 1 digits represent the magnitude.

| | Minimal | Maximal | Range |
|----------|-----------------------------|-----------------------|--------------------------------|
| Positive | 0 0 ... 0 | 0 (r – 1) ... (r – 1) | [0, r ^{k-1} – ulp] |
| Negative | (r – 1) (r – 1) ... (r – 1) | (r – 1) 0 ... 0 | [-(r ^{k-1} – ulp), 0] |

Disadvantage of this method: The operation to be performed may depend on the signs of the operands

Biased Method - the basic idea of this method is to let [0, Max] to represent [-Bias, Max – Bias], where the Bias is a fixed positive integer. This method sometimes is also called “excess-Bias” method.

One major disadvantage can be implied from the following:

x + y + Bias = (x + Bias) + (y + Bias) – Bias

Complement Representation - There are two alternatives for complement methods:

Radix complement (also called 2’s complement in the binary system)

Diminished-radix complement (called 1’s complement in the binary system)

During arithmetic operations with complement numbers,

IF X and Y have opposite signs, no overflow can occur regardless whether there is a carry-out or not.

IF X and Y have the same sign and the sign of the result is different from that of the two operands, then an overflow occurs

Example 3 Add X and Y, where X = –7₁₀ = (11001)₂ and Y = –10₁₀ = (10110)₂.

$$\begin{array}{rcl} X + Y & : & \\ + & & \\ \hline & & 1 \quad 1 \quad 0 \quad 0 \quad 1 \quad X = -7 \\ & & 1 \quad 0 \quad 1 \quad 1 \quad 0 \quad Y = -10 \\ \hline & & 1 \quad 0 \quad 1 \quad 1 \quad 1 \quad 15 \end{array}$$

There is a carry-out but the result is incorrect (overflow).

Example 4 Add X and Y, where X = 7₁₀ = (00111)₂ and Y = 10₁₀ = (01010)₂.

$$\begin{array}{rcl} X + Y & : & \\ + & & \\ \hline & & 0 \quad 0 \quad 1 \quad 1 \quad 1 \quad X = 7 \\ & & 0 \quad 1 \quad 0 \quad 1 \quad 0 \quad Y = 10 \\ \hline & & 1 \quad 0 \quad 0 \quad 0 \quad 1 \quad -15 \end{array}$$

we do not have any carry-out but the result is incorrect (overflow).

Overflow can be corrected by the addition of an extra bit.

Convert C into nega-decimal number system (Method 2)

| Integer: Decimal-to-Nega-decimal | | |
|-----------------------------------|--|----------------------------------|
| Dividing-by-(-10) | Quotient | Remainder (keep it non-negative) |
| 44/(-10) | -4 | 4 = x ₀ |
| -4/(-10) = (-10 + 4)/(-10) | 1 | 6 = x ₁ |
| 1/(-10) | 0 | 1 = x ₂ |
| Fraction: Decimal-to-Nega-decimal | | |
| Multiplying-by-(-10) | Fractional part (keep it within the range) | Integral part |
| -0.3125 × (-10) | -0.875 | 4 = x ₋₁ |
| -0.875 × (-10) | -0.25 | 9 = x ₋₂ |
| -0.25 × (-10) | -0.5 | 3 = x ₋₃ |
| -0.5 × (-10) | 0 | 5 = x ₋₄ |

In the second table the range is referred to an interval that a fractional number in nega-decimal form can represent. For nega-decimal system, the range is given by [-10/11, +1/11] = [-0.909, +0.091].

For Nega-Decimal conversions, make sure that the fractional part is within the range. If not, add one to the integer part and get the fraction within range.

Range of NegaBinary Numbers:

$$\left\{ \begin{array}{l} F_{\max}^- = -(2^{-1} + 2^{-3} + 2^{-5} + 2^{-7} + \dots) = -0.667_{10} \\ F_{\max}^+ = (2^{-2} + 2^{-4} + 2^{-6} + 2^{-8} + \dots) = 0.333_{10} \end{array} \right.$$

| Integer: Decimal-to-Nega-decimal | | |
|-----------------------------------|--|----------------------------------|
| Dividing-by-(-10) | Quotient | Remainder (keep it non-negative) |
| 44/(-10) | -4 | 4 = x ₀ |
| -4/(-10) = (-10 + 4)/(-10) | 1 | 6 = x ₁ |
| 1/(-10) | 0 | 1 = x ₂ |
| Fraction: Decimal-to-Nega-decimal | | |
| Multiplying-by-(-10) | Fractional part (keep it within the range) | Integral part |
| -0.3125 × (-10) | -0.875 | 4 = x ₋₁ |
| -0.875 × (-10) | -0.25 | 9 = x ₋₂ |
| -0.25 × (-10) | -0.5 | 3 = x ₋₃ |
| -0.5 × (-10) | 0 | 5 = x ₋₄ |

Multiplication of a +ve and -ve number

| | | | | | | |
|--|---|---|---|---|---|-------------------------|
| A | 1 | 0 | 1 | 1 | | -5 |
| X | × | 0 | 0 | 1 | 1 | 3 |
| P ⁽⁰⁾ = 0 | | 0 | 0 | 0 | 0 | |
| x ₀ = 1 ⇒ Add A | + | 1 | 0 | 1 | 1 | |
| | | 1 | 0 | 1 | 1 | |
| Shift to get P ⁽¹⁾ | | 1 | 1 | 0 | 1 | 1 |
| x ₁ = 1 ⇒ Add A | + | 1 | 0 | 1 | 1 | |
| | | 1 | 0 | 0 | 1 | |
| Shift to get P ⁽²⁾ | | 1 | 1 | 0 | 0 | 1 |
| x ₂ = 0 ⇒ Shift to get P ⁽³⁾ | | 1 | 1 | 1 | 0 | 0 |
| x ₃ = 0 ⇒ output P ⁽³⁾ | | 1 | 1 | 1 | 0 | 0 |
| | | 0 | 0 | 1 | 1 | 2 ⁻³ × (-15) |

Multiplication of two -ve numbers

Example 3 The multiplier is negative and the operands are in the two’s complement form. Let X = –3 and A = –5, in n-bit 2’s complement form, where n = 4.

| | | | | | | |
|--|---|---|---|---|---|----------------------|
| A | 1 | 0 | 1 | 1 | | -5 |
| X | × | 1 | 1 | 0 | 1 | -3 |
| P ⁽⁰⁾ = 0 | | 0 | 0 | 0 | 0 | |
| x ₀ = 1 ⇒ Add A | | 1 | 0 | 1 | 1 | |
| | | 1 | 0 | 1 | 1 | |
| Shift to get P ⁽¹⁾ | | 1 | 1 | 0 | 1 | 1 |
| x ₁ = 0 ⇒ Shift to get P ⁽²⁾ | | 1 | 1 | 1 | 0 | 1 |
| x ₂ = 1 ⇒ Add A | + | 1 | 0 | 1 | 1 | |
| | | 1 | 0 | 0 | 1 | 1 |
| Shift to get P ⁽³⁾ | | 1 | 1 | 0 | 0 | 1 |
| x ₃ = 1 ⇒ Correction (Add – A) | + | 0 | 1 | 0 | 1 | |
| Get P ⁽³⁾ | | 0 | 0 | 0 | 1 | 1 |
| | | 1 | 1 | 1 | 1 | 2 ⁻³ × 15 |

Fractional Division (X < D)

Example 4 (restoring) Let X = (0.10000)₂ = 1/2 and D = (0.110)₂ = 3/4. Since X < have

$$\begin{array}{rcl} r_0 = X & & 0 \quad .1 \quad 0 \quad 0 \quad | \quad 0 \quad 0 \quad 0 \\ 2r_0 & & 0 \quad 1 \quad .0 \quad 0 \quad 0 \quad | \quad 0 \quad 0 \\ \text{Add} - D & + & 1 \quad 1 \quad .1 \quad 0 \quad 1 \quad 0 \\ r_1 = 2r_0 - D \geq 0 & & 0 \quad 0 \quad .0 \quad 1 \quad 0 \quad | \quad 0 \quad 0 \quad \text{set } q_1 = 1 \\ 2r_1 & & 0 \quad 0 \quad .1 \quad 0 \quad 0 \quad 0 \\ \text{Add} - D & + & 1 \quad 1 \quad .0 \quad 1 \quad 0 \quad 0 \\ r_1 - D < 0 & & 1 \quad 1 \quad .1 \quad 1 \quad 0 \quad 0 \quad | \quad 0 \quad \text{set } q_2 = 0 \\ r_2 = 2r_1 & & 0 \quad 0 \quad .1 \quad 0 \quad 0 \quad 0 \\ 2r_2 & & 0 \quad 1 \quad .0 \quad 0 \quad 0 \quad 0 \\ \text{Add} - D & + & 1 \quad 1 \quad .0 \quad 1 \quad 0 \quad 0 \\ r_3 = 2r_2 - D \geq 0 & & 0 \quad 0 \quad .0 \quad 1 \quad 0 \quad 0 \quad \text{set } q_3 = 1 \end{array}$$

The final results are Q = (0.101)₂ = 5/8 and r₃ = 1/4 ⇒ R = r_m2^{-m} = r₃2⁻³ = 1/32.

Integer Division (32/6)

Condition – X < 2^(a-2) * D

have

$$\begin{array}{rcl} r_0 = X & & 0 \quad 1 \quad 0 \quad 0 \quad | \quad 0 \quad 0 \quad 0 \\ 2r_0 & & 0 \quad 1 \quad 0 \quad 0 \quad 0 \quad | \quad 0 \quad 0 \\ \text{Add} - D & + & 1 \quad 1 \quad 0 \quad 1 \quad 0 \quad 0 \\ r_1 = 2r_0 - D > 0 & & 0 \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad | \quad 0 \quad 0 \quad \text{set } q_1 = 1 \\ 2r_1 & & 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 0 \\ \text{Add} - D & + & 1 \quad 1 \quad 0 \quad 1 \quad 0 \quad 0 \\ 2r_1 - D < 0 & & 1 \quad 1 \quad 1 \quad 1 \quad 0 \quad 0 \quad | \quad 0 \quad \text{set } q_2 = 0 \\ r_2 = 2r_1 & & 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 0 \\ 2r_2 & & 0 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0 \\ \text{Add} - D & + & 1 \quad 1 \quad 0 \quad 1 \quad 0 \quad 0 \\ r_3 = 2r_2 - D > 0 & & 0 \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad \text{set } q_3 = 1 \end{array}$$

The final results are Q = (101)₂ = 5 and R = r₃ = 2.

IEEE Single Precision format (32 Bits)

An IEEE single-precision floating-point number F has three parts, s, e and f, as shown below.

| s : sign | e : exponent | f : significand or mantissa |
|----------|--------------|-----------------------------|
| 1 bit | 8 bits | 23 bits |
| 32 bits | | |

IEEE floating-point number F can be evaluated by

$$F = (-1)^s \times 1.f \times 2^{e-127}.$$

| Value or meaning of an IEEE single-precision floating-point representation | | |
|--|-------|--|
| e = 0 | f = 0 | F = ±0 |
| | f ≠ 0 | F are subnormal numbers (= ±0.f × 2 ⁻¹²⁶) |
| e = 255 | f = 0 | F = ±∞ |
| | f ≠ 0 | F is NaN* (Not a Number) |
| 1 ≤ e ≤ 254 | – | F is an ordinary number and F = (–1) ^s × 1.f × 2 ^{e-127} |

* There are two types of NaN, namely signaling NaN (f₋₁ = 0) and quiet NaN (f₋₁ = 1), defined in IEEE 754-2019, where f₋₁ is the leading bit in f field.

Example 1 Convert 46.5₁₀ into IEEE single-precision standard.

Solution: First we convert the given number 46.5₁₀ into binary of the form of (1):

$$46.5_{10} = 101110.1_2 = 1.011101 \times 2^5 = (-1)^0 \times 1.011101 \times 2^{132-127}.$$

Then it is easy to obtain

$$\left\{ \begin{array}{l} s = 0 \\ e = 1000 \, 0100 \\ f = 0111 \, 0100 \, 0000 \, 0000 \, 0000 \end{array} \right.$$

IEEE Double Precision format (64 Bit)

An IEEE double-precision floating-point number F has also three parts, s, e and f, as shown below.

| s : sign | e : exponent | f : significand or mantissa |
|----------|--------------|-----------------------------|
| 1 bit | 11 bits | 52 bits |
| 64 bits | | |

IEEE floating-point number F can be evaluated by

$$F = (-1)^s \times 1.f \times 2^{e-1023}. \quad (2)$$

| Value or meaning of an IEEE double-precision floating-point representation | | |
|--|-------|---|
| e = 0 | f = 0 | F = ±0 |
| | f ≠ 0 | F are subnormal numbers (= ±0.f × 2 ⁻¹⁰²²) |
| e = 2047 | f = 0 | F = ±∞ |
| | f ≠ 0 | F is NaN (Not a Number) |
| 1 ≤ e ≤ 2046 | – | F is an ordinary number and F = (–1) ^s × 1.f × 2 ^{e-1023} |

Example 2 Convert 46.5₁₀ into IEEE double-precision standard.

Solution: First we convert the given number 46.5₁₀ into binary of the form of (2):

$$46.5_{10} = 101110.1_2 = 1.011101 \times 2^5 = (-1)^0 \times 1.011101 \times 2^{1028-1023}.$$

Then it is easy to obtain

$$\left\{ \begin{array}{l} s = 0 \\ e = 1000 \, 0000 \, 100 \\ f = 0111 \, 0100 \, \underbrace{00 \dots 00}_{51 \text{ zeros}} \end{array} \right.$$

Floating Point Arithmetic Operations:

Example 3 Convert F₁ = 4 and F₂ = 3 into IEEE floating point single precision format. Then perform floating point operations F₁ × F₂ and F₁ + F₂.

Solution:

$$F_1 = 4 = (-1)^{S_1} \times 1.f_1 \times 2^{E_1-127} = (-1)^0 \times 1.0 \times 2^{129-127},$$

$$F_2 = 3 = (-1)^{S_2} \times 1.f_2 \times 2^{E_2-127} = (-1)^0 \times 1.1 \times 2^{128-127}.$$

Then for floating point multiplication we have

$$F_3 = F_1 \times F_2 = (-1)^{S_3} \times (1.0 \times 1.1) \times 2^{(129-127+128-127)} = (-1)^0 \times 1.1 \times 2^{130-127}.$$

For floating point addition, since F₁ > F₂ we first rewrite F₂ such that its 2’s power part is the same as that of F₁.

$$F_2 = (-1)^0 \times 1.1 \times 2^{128-127} = (-1)^0 \times 0.11 \times 2^{129-127}.$$

Then it follows

$$F_4 = F_1 + F_2 = (-1)^0 \times (1.0 + 0.11) \times 2^{129-127} = (-1)^0 \times 1.11 \times 2^{129-127}.$$

For FP product (F3 = F1 x F2) - S3 = resulatant sign, F3 = (Multiplication of f1 and f2) and e3 = (e1 + e2 -127)

For addition/subtraction, f3 = (f1 +- f2). Make sure that the value of e1 and e2 are the same.

Rounding Schemes:

TRUNCATION – Remove anything after decimal point.

| Chopping scheme with d = 2 | | |
|----------------------------|---------|-------------|
| Input: | Output: | Error: |
| x | chop(x) | chop(x) – x |
| ×.00 | × | 0 |
| ×.01 | × | –1/4 |
| ×.10 | × | –1/2 |
| ×.11 | × | –3/4 |

Implementation: Its implementation is cost free

Time delay: There is no delay incurred with this rounding scheme.

ROUNDING TO NEAREST INTEGER – 0.5 is attributed to the nearest integer

| Round-to-nearest scheme with d = 2 | | |
|------------------------------------|----------|--------------|
| Input: | Output: | Error: |
| x | round(x) | round(x) – x |
| ×.00 | × | 0 |
| ×.01 | × | –1/4 |
| ×.10 | ×. + 1 | +1/2 |
| ×.11 | ×. + 1 | +1/4 |

Implementation: An implementation requires an adder and a few logic gates.

Time delay: It is equal to that of an adder of the size of the output.

ROUND TO NEAREST EVEN INTEGER – 0.5 is attributed to the nearest even integer

| Round-to-nearest-even scheme with d = 2 | | |
|---|---------|-------------|
| Input: | Output: | Error: |
| x | rtne(x) | rtne(x) – x |
| ×0.00 | ×0. | 0 |
| ×0.01 | ×0. | –1/4 |
| ×0.10 | ×0. | –1/2 |
| ×0.11 | ×1. | +1/4 |
| ×1.00 | ×1. | 0 |
| ×1.01 | ×1. | –1/4 |
| ×1.10 | ×1. + 1 | +1/2 |
| ×1.11 | ×1. + 1 | +1/4 |

Implementation and time delay: slightly higher than previous method

ROM ROUNDING – For this method, L = total number of bits, d = bits after decimal point

| ROM scheme with ℓ = 3 input bits | | |
|----------------------------------|---------|------------|
| Input: | Output: | Error: |
| x | ROM(x) | ROM(x) – x |
| ×00.0 | ×00. | 0 |
| ×00.1 | ×01. | +1/2 |
| ×01.0 | ×01. | 0 |
| ×01.1 | ×10. | +1/2 |
| ×10.0 | ×10. | 0 |
| ×10.1 | ×11. | +1/2 |
| ×11.0 | ×11. | 0 |
| ×11.1 | ×11. | –1/2 |

Implementation cost: Size of the ROM used, Time delay: Decided by the ROM reading speed

In case of ROM rounding, the final result is truncated instead of incrementing to avoid a carry chain. This needs to be done, else the system operation will be delayed. Here, we have rounded 11.1 to 11 instead of 100.