# Software Systems-1

Floating Point Arithmetic

Dr. Azam E. Al-Rawachy

COLLEGE OF COMPUTER SCIENCE AND MATHEMATICS

# 1 Introduction

Floating-point representation is an alternative technique based on scientific notation. Scientific notation is a way of expressing numbers that are too large or too small to be conveniently written in decimal form. It may be referred to as scientific form or standard index form. This base ten notation is commonly used by scientists, mathematicians, and engineers, in part because it can simplify certain arithmetic operations. In scientific notation, all numbers are written in the form of

$$m \times 10^n$$

where the exponent n is an integer, and the coefficient *m* is any real number. The real number m is called the **significand** or **mantissa**. Any given real number can be written in the form of $m \times 10^n$ in many ways:

For example, 350 can be written as $3.5 \times 10^2$ or $35 \times 10^1$ or $350 \times 10^0$.

*n* normalized scientific notation, the exponent n is chosen so that the absolute value of *m* remains at least one but less than ten ($1 \leq |m| < 10$). Thus 350 is written as $3.5 \times 10^2$. Few steps are taken to write a scientific notation as follow:

❖ Decide where the decimal must end up so that the one number is to its left.

❖ Count how many places you bounce the decimal point.

❖ Re-write in the form of $m \times 10^n$.

**Ex: Write 0.000811 in scientific notation.**

1. Move the decimal to get a number between 1 and 10 (8.11).
2. Set up scientific notation   $8.11 \times 10^?$
3. The **decimal** needs to move **Right** to change 000811 to 8.11, so the exponent will be **negative**.
4. The decimal needs to move **4** places.
5. So 0.000811 written in scientific notation is $8.11 \times 10^{-4}$.

## 2   Floating-point basics

When dealing with very large or very small numbers, a more flexible technique, known as **floating point representation**, can be employed. This allows the way that the bits are allocated to vary so that both very large and very small numbers can be represented.

**Ex:** Convert $5.5_{(10)}$ to 8-bit representation of

> **Answer**: $5.5_{(10)} = 101.1_{(2)}$

- Converts a to binary scientific notation:

$$101.1_{(2)} = 1.011 \times 2^2 \longrightarrow \text{Exponent}$$

**Mantissa**

- Applying technique for mapping that into a set of bits. In this example we'll use **8 bits** to store such a number.

| 1bit | 4bit | 3bit |
|------|------|------|
| Sign | Exponent+7 | Mantissa |

We use the first bit to represent the sign (1 for negative, 0 for positive).

- For the exponent bits

Instead of using a separate sign bit for the exponent, the standard uses a biased representation. The value of the bias is 7. Thus an exponent 0 means that $-7$ is stored in the exponent field. $(2^{n-1}-1)$ = Biased value, n=no of bits. $(2^{4-1}-1) = 7$.

$0 \rightarrow 15$ (biased)
$-7 \rightarrow 8$ (unbiased)

$7 + 2$ to $= 9_{(10)} = 1001_{(2)}$

| Sign | Exponent+7 | Mantissa |
|------|-----------|----------|
| 0 | 1001 | 011 |

Note that we omit the integer part of the mantissa: Since the mantissa must have exactly one nonzero bit to the left of its decimal point, and the only nonzero bit is 1, we know that the bit to the left of the decimal point must be a 1. There's no point in wasting space in inserting this 1 into our bit pattern, so we include only the bits of the mantissa to the right of the decimal point.
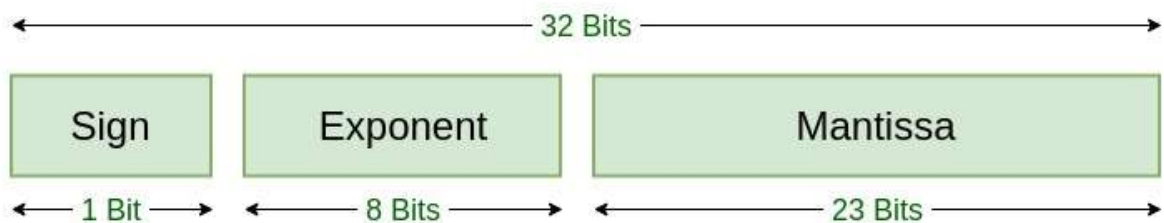
## 2.1 IEEE 754 Floating binary point

Floating point numbers are an important data type in computation which is used extensively. Yet, many users do not know the standard which is used in almost all computer hardware to store and process these. The Institute of Electrical and Electronic Engineers IEEE in 1985 and augmented in 2008 was evolved to represent floating point numbers and process them. This standard is now used by all computer manufacturers while designing floating point arithmetic units so that programs are portable among computers.
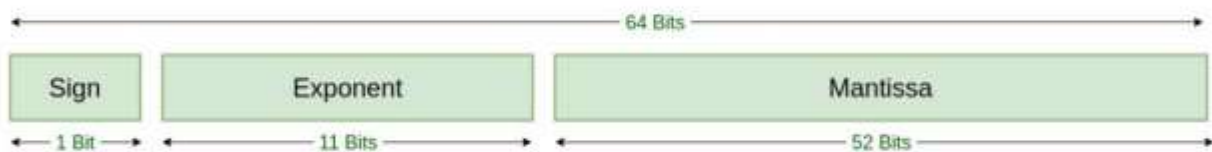
Normalized floating point representation is used for real numbers in computers. In this representation, 32 bits are divided into two parts: a part called the mantissa with its sign and the other called the exponent with its sign.

The number of bits required to represent *seven decimal digits is approximately 23*. The remaining 8 bits are allocated for the exponent, the exponent will range from *–127 to +127*.



Single precision floating point representation (32-bit)

1. Sign, for which 1 bit is allocated.

2. Mantissa (called *significand* in the standard) is allocated 23 bits.

3. Exponent is allocated 8 bits. As both positive and negative numbers are required for



Double precision floating point representation (64-bit)

## 2.2 Decimal to IEEE 754 Floating binary point

### Ex1: Convert the value 19.59375 into IEEE 754 standard 32-bit Floating binary point

- ❖ **Step1:** Determine the sign bit (0 if positive, 1 if negative).

- ❖ **Step2:** Convert to pure binary

$19.59375_{10} =$

| 16 | 8 | 4 | 2 | 1 | 0.5 | 0.25 | 0.125 | 0.0625 | 0.03125 |
|----|---|---|---|---|-----|------|-------|--------|---------|
| 1  | 0 | 0 | 1 | 1 . | 1 | 0 | 0 | 1 | 1 |

| | |
|---|---|
| $19/2 = 9$ Remainder 1 | $0.59375 * 2 = 1.1875$    1 |
| $9/2 = 4$ Remainder 1 | $0.1875 * 2 = 0.375$    0 |
| $4/2 = 2$ Remainder 0 | $0.375 * 2 = 0.75$    0 |
| $2/2 = 1$ Remainder 0 | $0.75 * 2 = 1.5$    1 |
| $1/2 = 0$ Remainder 1 | $0.5 * 2 = 1$    1 |

- ❖ **Step3:** Normalize to determine the mantissa and the ***unbiased exponent*** (place the binary point after leftmost 1)

  $1 0 0 1 1 . 1 0 0 1 1 = 1.001110011 \times 2^4$

- ❖ **Step4:** Determine the ***biased exponent.***

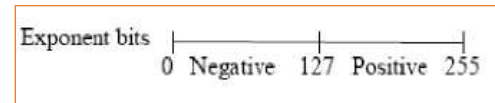  $$4+127 = 131_{10} = 10000011_2$$

- ❖ **Step5:** Removing the Leading 1 from the mantissa (Remove the leftmost1)

  1. 001 1100 1100 0000 0000 0000 = 001 1100 1100 0000 0000 0000

The value of the bias is $(2^{8-1}-1) = 127$. Thus an exponent 0 means that $-127$ is stored in the exponent field.

The range 0 to 255 of an 8-bit number is divided into two parts 0 to 127 and 128 to 255.

Thus the range of exponents is $-127$ to $+128$. Given an exponent string *exp.*, the value of the exponent will be taken as (exp $\pm127$).

Exponent bits

0 Negative   127 Positive   255

| S | Exponent | | | | | | | | Mantissa | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

**Ex2: Convert the value -123.3 into IEEE 754 standard 32-bit Floating binary point**

❖ **Step1:** Sign bit =1.

❖ **Step2:** Convert to pure binary

   **$123.3_{10}$ = 1111011.0100100011001....$_2$**

| | |
|---|---|
| $123/2 = 61$ Remainder 1 | $0.3 * 2 = 0.6$    0 |
| $61/2 = 30$ Remainder 1 | $0.6 * 2 = 1.2$    1 |
| $30/2 = 15$ Remainder 0 | $0.2 * 2 = 0.4$    0 |
| $15/2 = 7$ Remainder 1 | $0.4 * 2 = 0.8$    0 |
| $15/2 = 7$ Remainder 1 | $0.8 * 2 = 1.6$    1 |
| $3/2 = 1$ Remainder 1 | $0.6 * 2 = 1.2$    1 |
| $1/2 = 0$ Remainder 1 | $0.2 * 2 = 0.4$    0 |

❖ **Step3:** Normalize to determine the mantissa and the ***unbiased exponent*** (place the binary point after leftmost 1)

1 1 1 1 0 1 1**.** 0 1 0 0 1 1 0 0 1 <u>1 0 0 1</u> = 1**.** 1 1 1 0 1 1 0 1 0 0 1 1 0 0 1 1 0 0 1 x $2^6$

❖ **Step4:** Determine the ***biased exponent.***
   $6+127 = 133_{10} = 10000101_2$

❖ **Step5:** Removing the Leading 1 from the mantissa (Remove the leftmost1)
   1**.** 1 1 1 0 1 1 0 1 0 0 1 1 0 0 1 1 0 0 1 = 1 1 1 0 1 1 0 1 0 0 1 1 0 0 1 1 0 0 1

| S | Exponent | | | | | | | Mantissa | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 |

## 2.3    IEEE 754 Floating binary single precession point to Decimal

### Ex3: Convert 0**10000100**110100000000000000000000

- ❖ **Step1:** Sign bit = 0 (Positive number).

- ❖ Determine the exponent in Decimal: **10000100 = 132₁₀** → $\mathbf{10000100} = \mathbf{132_{10}}$

- ❖ Remove the exponent bias, 132-127=5

- ❖ Convert the mantissa to binary = 0.5+0.25+0.0625 = $0.8125_{10}$

| 0.5 | 0.25 | 0.125 | 0.0625 | 0.03125 | 0.015625 |
|-----|------|-------|--------|---------|----------|
| 1 | 1 | 0 | 1 | 0 | 0 |

$$= \frac{1}{2} + \frac{1}{4} + \frac{1}{16} = \frac{8+4+1}{16} = \frac{13}{16}$$

- ❖ Step5: add 1 to mantissa and include the sign = $1.8125_{10}$

- ❖ Step6: Calculate the final result:  $1.8125_{10} \times 2^5 = 58$

In other words, the final equation that was used in the above example is:

$$(-1)^S \times (1+ \text{Mantissa})_{10} \times 2^{exp-127}$$

### Ex4: Convert 1**10001001**000110000000000000000000

- ❖ **Step1:** Sign bit = 1 (Negative number).

- ❖ Determine the exponent in Decimal: **10001001 = 137₁₀** → $\mathbf{10001001} = \mathbf{137_{10}}$

- ❖ Remove the exponent bias, 137-127= 10

- ❖ Convert the mantissa to binary = 0.0625 + 0.03125 =$0.09375_{10}$

| 0.5 | 0.25 | 0.125 | 0.0625 | 0.03125 | 0.015625 |
|-----|------|-------|--------|---------|----------|
| 0 | 0 | 0 | 1 | 1 | 0 |

- ❖ Step5: add 1 to mantissa and include the sign = $-1.09375_{10}$

- ❖ Step6: Calculate the final result:  $-1.09375_{10} \times 2^{10} = -1120$

# 3  CHARACTER REPRESENTATION

The computer convert data to character which is a form that is meaningful to humans. The American Standard Code for Information Interchange (*ASCII*) uses 8-bits to represent 256 characters include Latin alphabet, control characters, special characters, numerical characters. ASCII codes are restricted in their abilities to provide data representation for the non-Latin alphabets used by the majority of the world's population.