

School of Engineering and Computer Science

# Numerical Applied Mathematics

## (Floating Point Numbers)

Kamel ATTAR  
[attar.kamel@gmail.com](mailto:attar.kamel@gmail.com)

Week #12 ♦ 12/JAN/2024 ♦

## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Fraction      Number

Floating point representation

Conversions

Special Values in IEEE 754-1985 Standard (32 Bits)



## Largest and Smallest Positive Floating Point Numbers: Subnormal Numbers

## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

## 1 Fraction Number

### Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

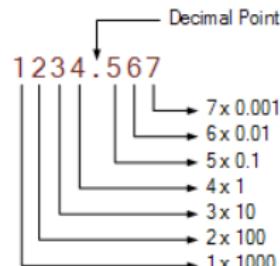
Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

- ▶ The decimal numbering system uses the concept of positional or relative weighting values producing a positional notation, where each digit represents a different weighted value depending on the position occupied either side of the decimal point.
- ▶ Then the value of any decimal number will be equal to the sum of its digits multiplied by their respective weights.
- ▶ Here in this decimal number example, the digit immediately to the right of the decimal point (number 5) is worth one tenth ( $1/10$  or  $0.1$ ) of the digit immediately to the left of the decimal point (number 4) which as a multiplication value of one ( $1$ ).
- ▶ Thus as we move through the number from left-to-right, each subsequent digit will be one tenth the value of the digit immediately to its left position, and so on.
- ▶ For our example:  $N = 1234.567_{10}$  in the weighted decimal format this will be equal to:



$$\begin{aligned}
 1234.567_{10} &= (1 \times 10^3) + (2 \times 10^2) + (3 \times 10^1) + (4 \times 10^0) + (5 \times 10^{-1}) + (6 \times 10^{-2}) + (7 \times 10^{-3}) \\
 &= (1 \times 1000) + (2 \cdot 100) + (3 \cdot 10) + (4 \cdot 1) + (5 \cdot 0.1) + (6 \cdot 0.01) + (7 \times 0.001) \\
 &= 1000 + 200 + 30 + 4 + 0.5 + 0.06 + 0.007 .
 \end{aligned}$$

## 1 Fraction Number

Fractional Decimal Number

**Fractional Binary Number**

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

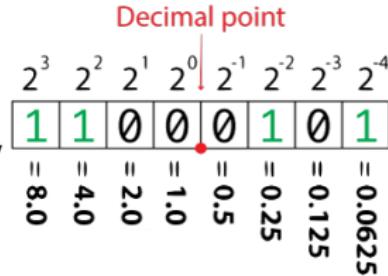
Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

- ▶ We can use the same idea of positional notation where each digit represents a different weighted value depending upon the position it occupies in the binary numbering system.
- ▶ The difference this time is that the binary number system is a positional system, where the different weighted positions of the digits are to the power of **2** (base-**2**) instead of **10**.
- ▶ Thus all the fractional digits to the right of the binary point have respective weightings which are negative powers of two, creating binary fractions.
- ▶ So for the fractional binary numbers to the right of the binary point, the weight of each digit becomes more negative giving:  $2^{-1}$ ,  $2^{-2}$ ,  $2^{-3}$ ,  $2^{-4}$ , and so on as shown.



## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

**Converting Binary Fraction to a Decimal**

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

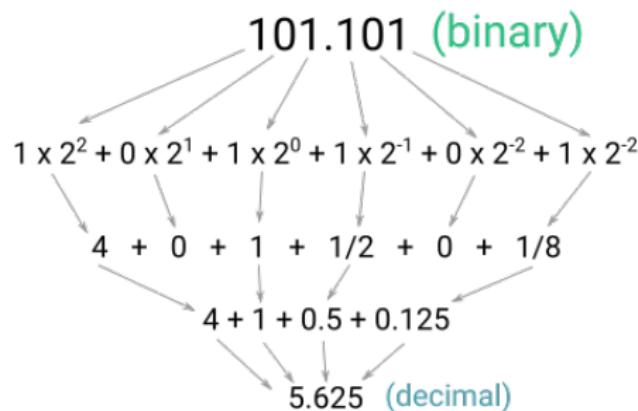
Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers



**Exercise 1.** Convert the following binary number to decimal number

- a) 1101.0111   b) 0.11   c) 11.001   d) 1011.111

## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

- ▶ The conversion of decimal fractions to binary fractions is achieved using a method similar to that we used for integers. However, this time multiplication is used instead of division with the integers instead of remainders used with the carry digit being the binary equivalent of the fractional part of the decimal number.
- ▶ When converting from decimal to binary, the integer (positive sequence right-to-left) part and the fractional (negative sequence from left-to-right) part of the decimal number are calculated separately.
- ▶ For the integer part of the number, the binary equivalent is found by successively dividing (known as successive division) the integer part of the decimal number repeatedly by 2 ( $\div 2$ ),

noting the remainders in reverse order from the least significant bit (**LSB**) to the most significant bit (**MSB**), until the value becomes “0” producing the binary equivalent.

The binary equivalent of $118_{10}$ is $1110110_2$			
118 (divide by 2)	=	59	plus remainder 0 ( <b>LSB</b> )
59 (divide by 2)	=	29	plus remainder 1 ( $\uparrow$ )
29 (divide by 2)	=	14	plus remainder 1 ( $\uparrow$ )
14 (divide by 2)	=	7	plus remainder 0 ( $\uparrow$ )
7 (divide by 2)	=	3	plus remainder 1 ( $\uparrow$ )
3 (divide by 2)	=	1	plus remainder 1 ( $\uparrow$ )
1 (divide by 2)	=	0	plus remainder 1 ( <b>MSB</b> )

- ▶ The fractional part of the number is found by successively multiplying (known as successive multiplication) the given fractional part of the decimal number repeatedly by **2** ( $\times 2$ ), noting the carries in forward order, until the value becomes “**0**” producing the binary equivalent.
- ▶ So if the multiplication process produces a product greater than **1**, the carry is a “**1**” and if the multiplication process produces a product less than “**1**”, the carry is a “**0**”.

The binary fraction equivalent of the decimal fraction: **0.8125<sub>10</sub>** is **0.1101<sub>2</sub>**

0.8125 (multiply by 2) =	1.625 =	0.625	carry 1 (MSB)
0.625 (multiply by 2) =	1.25 =	0.25	carry 1 (↓)
0.25 (multiply by 2) =	0.50 =	0.5	carry 0 (↓)
0.5 (multiply by 2) =	1.00 =	0.0	carry 1 (LSB)

- ▶ Note also that if the successive multiplication processes does not seem to be heading towards a final zero, the fractional number will have an infinite length or until the equivalent number of bits have been obtained, for example **8**-bits. or **16**-bits, etc. depending on the degree of accuracy required.
- ▶ So the binary fraction equivalent of the decimal number **118.8125<sub>10</sub>** si **1110110.1101<sub>2</sub>**

 **Exercise 2.** Convert the following decimal numbers to the indicated bases:

- a) **7562, 45** to octal    b) **1938, 257** to hexadecimal    c) **175, 175** to binary

 **Exercise 3.** Fill in the missing information in the following table:

Fractional value	Binary representation	Decimal representation
$\frac{1}{8}$	<b>0, 001</b>	<b>0, 125</b>
$\frac{3}{4}$		
$\frac{25}{16}$		
	<b>10, 1011</b>	
	<b>1, 001</b>	
		<b>5, 875</b>
		<b>3, 1875</b>

## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

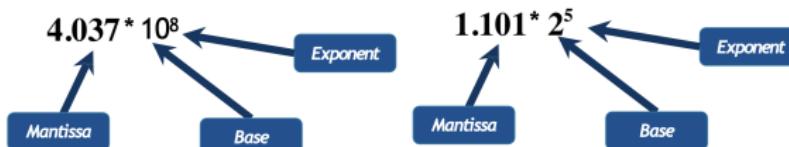
Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

- ▶ So far we know how to store **integers whole numbers**. But what if we want to store **real numbers** with **decimal fractions**. Even 27.5 needs another way to represent it.

This method is called floating point representation

- ▶ Floating point representation is used when the number contains a decimal fraction.
- ▶ Normalized decimal scientific notation has single non-zero digit to the left of the decimal (binary) point.
- ▶ Normalized binary scientific notation has single digit to the left of the radix point and there are no leading zeroes. The structure is as follows:



Binary Value	Normalized	Exp
1101.101	1.101101	3
.00101	1.01	-3
1.0001	1.0001	0
10000011.0	1.0000011	7

- ▶ Only the mantissa and the exponent are stored. The base is implied (known already) as it is not stored this will save memory capacity

## Example

To represent **6.75** into binary floating-point, first we convert the decimal value to binary numbers:

$$\begin{array}{ll} 6 = 110 \quad \text{and} & \begin{array}{rcl} 0.75 \times 2 & = 1.50 & = 1 \\ 0.50 \times 2 & = 1.00 & = 1 \\ 0.00 \times 2 & = 0 & = 0 \end{array} \implies 6.75 = 110.11 \\ & \underbrace{\phantom{0.00 \times 2 = 0}}_{0.75=11} \end{array}$$

Then the normalized floating point value of **6.75** is **110.11** = **+1.1011 × 2<sup>2</sup>**

## Example

To represent **6.75** into binary floating-point, first we convert the decimal value to binary numbers:

$$6 = 110 \quad \text{and} \quad \begin{array}{rcl} 0.75 \times 2 & = 1.50 & = 1 \\ 0.50 \times 2 & = 1.00 & = 1 \\ 0.00 \times 2 & = 0 & = 0 \end{array} \underbrace{\qquad\qquad\qquad}_{0.75=11} \implies 6.75 = 110.11$$

Then the normalized floating point value of **6.75** is  $110.11 = +1.1011 \times 2^2$

## Example

Find the normalized value of **9.5**.

## Example

To represent **6.75** into binary floating-point, first we convert the decimal value to binary numbers:

$$6 = 110 \quad \text{and} \quad \begin{array}{rcl} 0.75 \times 2 & = 1.50 & = 1 \\ 0.50 \times 2 & = 1.00 & = 1 \\ 0.00 \times 2 & = 0 & = 0 \end{array} \underbrace{\qquad\qquad\qquad}_{0.75=11} \implies 6.75 = 110.11$$

Then the normalized floating point value of **6.75** is  $110.11 = +1.1011 \times 2^2$

## Example

Find the normalized value of **9.5**.

## Solution

The normalized floating point value of **9.5** is  $1001.1 = +1.0011 \times 2^3$

Fraction Number

Floating point representation

Conversions

Special Values in IEEE 754-1985 Standard (32 Bits)

Normalized binary scientific notation

IEEE Floating Point Standard



18

## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

- ▶ The IEEE (Institute of Electrical and Electronic Engineers) is an international organization that has designed specific binary formats for storing floating point numbers.
- ▶ The IEEE defines two different formats with different precisions: single and double precision. Single precision is used by float variables in C and double precision is used by double variables.
- ▶ Intel's math coprocessor also uses a third, higher precision called extended precision. In fact, all data in the coprocessor itself is in this precision. When it is stored in memory from the coprocessor it is converted to either single or double precision automatically.
- ▶ There are **4** main different sizes of floating point numbers **16, 32, 64** and **128** bit. Sometimes referred to as Half, Single, Double and Quadruple precision

Floating point format	Total bits	Sign bits	Biased Exponent bits	Fraction bits
Half precision	<b>16</b>	<b>1</b>	<b>5</b>	<b>10</b>
Single precision	<b>32</b>	<b>1</b>	<b>8</b>	<b>23</b>
Double precision	<b>64</b>	<b>1</b>	<b>11</b>	<b>52</b>

Single precision		
Sign	Biased Exponent	Mantissa
S	E	F
1 bit	8 bits	23 bits

$$(-1)^S \times (1 + 0.F) \times 2^{(E-127)}$$

- S:** The sign of a binary floating-point number is represented by a single bit. A **1** bit indicates a negative number, and a **0** bit indicates a positive number.
- E:** The binary exponent is not stored directly. Instead, the sum of the exponent and **7F** is stored as a **8-bits**. Notice that the binary exponent is unsigned, so it cannot be negative.

We will pad the left of the exponent with 0's up to 8 bits

exp	E=e + 127	Binary Rep.
5	+132	10000100
-10	+117	01110101
128	+255	11111111
-1	+126	01111110

The largest possible exponent is **128**, because when added to **127**, produces **255**, the largest unsigned value represented by **8** bits. The approximate range exponents is

from **-127** to **128**.

The main advantage of this format is a unique representation for exponent zero.

Single precision		
Sign	Biased Exponent	Mantissa
S	E	F
1 bit	8 bits	23 bits

$$(-1)^S \times (1 + 0.F) \times 2^{(E-127)}$$

**F:** It's the first **23**-bits after the **1.** in the significant. The fraction part assumes a normalized significand (in the form **1.b<sub>10</sub>b<sub>11</sub>...b<sub>32</sub>**). Since the first bit is always an one, the leading one is not stored! This allows the storage of an additional bit at the end and so increases the precision slightly. This idea is known as the hidden one representation.

We will pad the right of the mantissa with 0's up to 23 bits.

## Example

Consider the number  $x = -1.11011 \times 2^{-48}$ . Determine the floating point representation in IEEE single precision (**32** bits).

## Solution

- 1 We immediately see that  $x$  is a negative number and so the sign is  $\sigma = +$ . Therefore the first bit in our floating point representation of this number will be  $b_1 = S = 1$ .
- 2 Now we also see that the exponent  $e = -48$  IEEE floating point single precision (32 bits) stores the number  $E = e + 127$  instead though, and hence  $E = -48 + 127 = 79 = (01001111)_2$ . Therefore  $b2b3 \dots b9 = 01001111$ .
- 3 lastly we will determine the last twenty-three digits which represent the fractional part. We note that  $F = 0.11011$ .

Sign	Exponent	Mantissa
1	01001111	1101100000000000000000000

Thus the IEEE single precision representation of  $-1.11011 \times 2^{-48}$  is

**10100111110110000000000000000000** $_2$

 **Exercise 4.** Represent the following binary numbers as single precision float. Show:

- a scientific notation representation
- the significand, exponent, and sign bits (32 bits total)
- a hexadecimal representation of these 32 bits

- a) **-1.11**   b) **+1101.101**   c) **-0.00101**   d) **+100111.0**   e) **+0.0000001101011**

Double precision		
Sign	Exponent	Mantissa
S	E	F
1 bit	11 bits	52 bits

$$(-1)^S \times (1 + 0.F) \times 2^{(E-1023)}$$

- The largest possible exponent is **1024**, because when added to **1023**, produces **2047**, the largest unsigned value represented by **11** bits.
- The approximate range is from **-1023** to **1024**.



## Floating point representation



Comparison Chart: Single Precision vs. Double Precision

	<b>Single Precision</b>	<b>Double Precision</b>
<b>Overview</b>	Uses <b>32</b> bits of memory to represent a numerical value, with one of the bits representing the sign of mantissa	Uses <b>64</b> bits of memory to represent a numerical value, with one of the bits representing the sign of mantissa
<b>Biased exponent</b>	<b>8</b> bits used for exponent	<b>11</b> bits used for exponent
<b>Mantissa</b>	Uses <b>23</b> bits for mantissa (to represent fractional part)	Uses <b>52</b> bits for mantissa (to represent fractional part)
<b>Real-world Application</b>	Often used for games or any program that requires wider representation without a high level of precision	Often used for scientific calculations and complex programs that require a high level of precision



## Floating point representation

Comparison Chart: Single Precision vs. Double Precision

### Example

Take Euler's number ( $e$ ), for example. Here are the first 50 decimal digits of  $e$ :

2.7182818284590452353602874713526624977572470936999

Here's Euler's number in binary, converted to single precision:

0 10000000 0101101111100001010100  
8 bits                    23 bits

Here's Euler's number in binary, converted to double precision:

0 100000000000 010110111110000101010001011000101000101011101101001  
11 bits                    52 bits

- ▶ If you increase the amount of bits allocated to the Mantissa you increase the accuracy/precision
- ▶ If you increase the amount of bits allocated to the exponent you increase the range of the number



## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers



## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

### Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers



The rules for converting a decimal number into floating point are as follows:

- 1 Convert the absolute value of the number to binary, perhaps with a fractional part after the binary point. This can be done by converting the integral and fractional parts separately. The integral part is converted by using successive division. The fractional part can be converted by multiplication.
- 2 Append  $\times 2^0$  to the end of the binary number (which does not change its value).
- 3 Normalize the number. Move the binary point so that it is one bit from the left. Adjust the exponent of two so that the value does not change.
- 4 Place the mantissa into the mantissa field of the number. Omit the leading one, and fill with zeros on the right.
- 5 Add the bias to the exponent of two, and place it in the exponent field. The bias is  $2^{k-1} - 1$ , where  $k$  is the number of bits in the exponent field. For the eight-bit format,  $k = 3$ , so the bias is  $2^3 - 1 = 3$ . For IEEE 32-bit,  $k = 8$ , so the bias is  $2^8 - 1 = 127$ .
- 6 Set the sign bit, **1** for negative, **0** for positive, according to the sign of the original number.

## Example (How would **23.85** be stored)

- ① We convert the integer part into binary  $23 = 10111$
- ② We convert the decimal part into binary floating point

$$\begin{array}{rcl} 0.85 \times 2 & = 1.7 & = 1 \\ 0.7 \times 2 & = 1.4 & = 1 \\ 0.4 \times 2 & = 0.8 & = 0 \\ 0.8 \times 2 & = 1.6 & = 1 \\ 0.6 \times 2 & = 1.2 & = 1 \\ 0.2 \times 2 & = 0.4 & = 0 \\ 0.4 \times 2 & = 0.8 & = 0 \\ 0.8 \times 2 & = 1.6 & = 1 \end{array} \underbrace{\qquad\qquad\qquad}_{0.85=0.110110} \Rightarrow 23.85 = 10111.11011001100110\dots$$

- ③ We normalize the binary floating number  $23.85 = 1.011111011001100110 \times 2^4$

## Example

### 4 IEEE representation:

- The sign bit is **0**, because it is positive.
- The true exponent is **4**, so the biased exponent is  $4 + 127 = 131 = 1000\ 0011$
- The fraction is **011111 0110 0110 0110 0**

Single precision		
Sign	Exponent	Mantissa
0	1000 0011	011111 0110 0110 0110 0110 0

1 bit                    8 bits                    23 bits

Thus the IEEE single precision representation of **23.85** is **0100000110111110110011001100110011002**

Exercise 5. Represent the following decimal numbers as single precision float. Show:

- a scientific notation representation
- the significand, exponent, and sign bits (**32** bits total)
- a hexadecimal representation of these **32** bits

- a) **15.0625**   b) **-35.6875**   c) **86.5625**



## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

The rules for converting a floating point number into decimal are simply the reverse of the decimal to floating point conversion:

- 1 If the original number is in hex, convert it to binary.
- 2 Separate into the sign, exponent, and mantissa fields.
- 3 Extract the mantissa from the mantissa field, and restore the leading one. You may also omit the trailing zeros.
- 4 Extract the exponent from the exponent field, and subtract the bias to recover the actual exponent of two. As before, the bias is  $2^{k-1} - 1$ , where  $k$  is the number of bits in the exponent field, giving 3 for the 8-bit format, 127 for the 32-bit and 1023 for the 64-bit.
- 5 De-normalize the number: move the binary point so the exponent is 0, and the value of the number remains unchanged.
- 6 Convert the binary value to decimal. This is done just as with binary integers, but the place values right of the binary point are fractions.
- 7 Set the sign of the decimal number according to the sign bit of the original floating point number: make it negative for 1; leave positive for 0.

If the binary exponent is very large or small, you can convert the mantissa directly to decimal without de-normalizing. Then use a calculator to raise two to the exponent, and perform the multiplication. This will give an approximate answer, but is sufficient in most cases.

## Example

Convert the **32-bit floating point number 44361000** (in hex) to decimal

## Solution

- ① Convert:  $44361000_{16} = 0100\ 0100\ 0011\ 0110\ 0001\ 0000\ 0000_2$ .

Single precision		
Sign	Exponent	Mantissa
0=+	$10001000 = 136$	011011000010000000000000

- ② Separate:
- Actual exponent of 2 is  $e = E - 127 = 136 - 127 = 9$
  - De-normalize:  $1.01101100001_2 \times 2^9 = 1011011000.01_2$

- ③ Convert  $1011011000.01 = 512 + 128 + 64 + 16 + 8 + 0.25 = 728.25$ .

- ④ Sign: positive

Result  $44361000_{16} = 728.25$

## Example

Consider the following floating point number presented in IEEE single precision (32 bits) as **01101011101101010000000000000000** (Determine the sign  $\sigma$ , exponent  $e$ , and significand/mantissa  $F$  and determine the value of  $x = \sigma \cdot 1 \cdot F \times 2^e$ )

## Solution

- ① We note that the first bit of the number given above is  $S = b_1 = 0$  It immediately follows that we have that the sign of  $x$  is  $\sigma = (-1)^0 = +$
- ② Now the next eight bits  $b_2 b_3 \dots b$  care **11010111** and represent  $E = e + 127$ . We want to find what decimal number represents the binary number  $E = (11010111)_2$ . We have that:  
$$E = 1 + 2 + 4 + 0 + 16 + 0 + 64 + 128 = 215 \implies e = E - 127 = 215 - 127 = 88$$
- ③ Lastly, recall that the twenty-three bits  $b_{10} b_{11} \dots b_{32}$  represent the fractional part, and that  $F = 0 \cdot b_{10} b_{11} \dots b_{32}$  and so:

$$F = 0.0110101$$

So the decimal representation of this number is  $x = \sigma \cdot 1 \cdot F \times 2^e = +1.0110101 \times 2^{88}$



## Example

Consider the following number presented in IEEE single precision **32 bits**

**1100110010111110001000000000000** Determine the sign  $\sigma = (-1)^S$ , exponent  $e$ , and significand/mantissa  $F$  and determine the value of  $x = \sigma \cdot 1 \cdot F \times 2^e$

## Solution

- ① Once again we immediately have that since  $S = b_1 = 1$  then the sign of  $x$  is  $\sigma = (-1)^1 = -$
- ② Now next eight bits are **10011001** These bits represent  $E = e + 127$ . Thus we have that:  
$$E = (10011001)_2 = 1+0+0+8+16+0+0+128 = 153 \implies e = E - 127 = 153 - 127 = 26$$
- ③ Lastly find the mantissa by using the last twenty-three bits of the given number. We have that:

$$F = 0.0111110001$$

So the decimal representation of this number is  $x = \sigma \cdot 1 \cdot F \times 2^e = -1.0111110001 \times 2^{26}$



☞ **Exercise 6.** Convert the following  $n$ -bit floating point number to decimal.

- a)  $D3_{16}$  with  $n = 8$
- b)  $BE580000_{16}$  with  $n = 32$
- c)  $A3358000_{16}$  with  $n = 32$
- d)  $76650000_{16}$  with  $n = 32$

☞ **Exercise 7.** Convert the decimal number **0.125** into IEEE single precision floating point.

Give your answer in hexadecimal.

Fraction Number

Floating point representation

Conversions

Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers

Subnormal Numbers



## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

Fraction Number

Floating point representation

Conversions

Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers

Subnormal Numbers



## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers



- ▶ As the significand is assumed to have a hidden 1 as the most significant bit, all 0s in the significand part of the number will be taken as **1, 00 . . . 0**
- ▶ Thus zero is represented in the IEEE Standard by all 0s for the exponent and all 0s for the significand.
- ▶ All 0s for the exponent is not allowed to be used for any other number.
- ▶ If the sign bit is 0 and all the other bits 0, the number is **+0**. If the sign bit is 1 and all the other bits 0, it is **-0**.
- ▶ Even though **+0** and **-0** have distinct representations they are assumed equal. Thus the representations of **+0** and **-0** are:

Representations of +0		
Sign	Exponent	Mantissa
0	00000000	00000000000000000000000000000000

Representations of -0		
Sign	Exponent	Mantissa
1	00000000	00000000000000000000000000000000

Fraction Number

Floating point representation

Conversions

Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers

Subnormal Numbers



## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

**Representation of Infinity**

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers



- ▶ All **1s** in the exponent field is assumed to represent infinity ( $\infty$ )
- ▶ A sign bit **0** represents  $+\infty$  and a sign bit **1** represents  $-\infty$ . Thus the representations of  $+\infty$  and  $-\infty$  are:

Representations of $+\infty$		
Sign	Exponent	Mantissa
0	11111111	00000000000000000000000000000000

Representations of $-\infty$		
Sign	Exponent	Mantissa
1	11111111	00000000000000000000000000000000

Fraction Number

Floating point representation

Conversions

Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers

Subnormal Numbers



## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

- ▶ When an operation is performed by a computer on a pair of operands, the result may not be mathematically defined.
- ▶ the IEEE Standard defines two types of NaN.
  - When the result of an operation is not defined (i.e., indeterminate) it is called a **Quiet NaN** (QNaN). For example:  $\frac{0}{0}$ ,  $+\infty - \infty$ ,  $\sqrt{-1}$
  - The other type of NaN is called a **Signalling NaN** (SNaN). This is used to give an error message. When an operation leads to a floating point underflow, i.e., the result of a computation is smaller than the smallest number that can be stored as a floating point number, or the result is an overflow, i.e., it is larger than the largest number that can be stored, SNaN is used.
- ▶ QNaN is represented by **0** or **1** as the sign bit, all **1s** as exponent, and a **0** as the left-most bit of the significand and at least one **1** in the rest of the significand.
- ▶ SNaN is represented by **0** or **1** as the sign bit, all **1s** as exponent, and a **1** as the left-most bit of the significand and any string of bits for the remaining **22** bits.
- ▶ We give below the representations of QNaN and SNaN.

QNaN		
Sign	Exponent	Mantissa
0 or 1	11111111	00010000000000000000000000000000
The most significant bit of the significand is 0. There is at least one 1 in the rest of the significand.		

SNaN		
Sign	Exponent	Mantissa
0 or 1	11111111	10000000000000000000000000000000
The most significant bit of the significand is 1. Any combination of bits is allowed for the other bits of the significand.		

Fraction Number

Floating point representation

Conversions

Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers

Subnormal Numbers

44

EPITA

## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

## Subnormal Numbers

Largest Positive Number		
Sign	Exponent	Mantissa
0	11111110	111111111111111111111111

- ▶ Significand:  $1, 1111 \dots 1 = 2^0 + 2^{-1} + 2^{-2} + \dots + 2^{-23} = 1 + 1 - 2^{-23} = 2 - 2^{-23}$ .
- ▶ Exponent:  $(254 - 127) = 127$ .
- ▶ Largest Number =
 
$$+1, 1 \dots 1 \times 2^{127} = (1 + 1 - 2^{-23}) \times 2^{127} = (2 - 2^{-23}) \times 2^{127} \approx 3.403 \times 10^{38}.$$

If the result of a computation exceeds the largest number that can be stored in the computer, then it is called an **overflow**.

Smallest Positive Number		
Sign	Exponent	Mantissa
0	00000001	00000000000000000000000000

- ▶ Significand:  $1, 0 \dots 0 = 1$
- ▶ Exponent:  $(1 - 127) = -126$ .
- ▶ The smallest normalized number is  $= 1 \times 2^{-126} \approx 1.1755 \times 10^{-38}$ .

If the result of a computation is less than the smallest number that can be stored in the computer, then it is called an **underflow**.

Fraction Number

Floating point representation

Conversions

Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers

46

EPITA

## 1 Fraction Number

Fractional Decimal Number

Fractional Binary Number

Converting Binary Fraction to a Decimal

Converting Decimal to a Binary Fraction

## 2 Floating point representation

Normalized binary scientific notation

IEEE Floating Point Standard

IEEE single precision

IEEE double precision

## 3 Conversions

Decimal to Floating-Point Conversions

Float to Decimal Conversion

## 4 Special Values in IEEE 754-1985 Standard (32 Bits)

Representation of Zero

Representation of Infinity

Representation of Non Numbers

Largest and Smallest Positive Floating Point Numbers:

Subnormal Numbers

## ★ Special Values in IEEE 754-1985 Standard (32 Bits) ★

### Subnormal Numbers

- When all the exponent bits are **0** and the leading hidden bit of the significand is **0**, then the floating point number is called a **subnormal number**.
- Thus, one logical representation of a subnormal number is

$$(-1)^S \times 0.F \times 2^{(-127)}$$

where **F** has at least one **1** (otherwise the number will be taken as **0**)

However, the standard uses  $-126$ , i.e., bias +1 for the exponent rather than  $-127$  which is the bias for some not so obvious reason, possibly because by using  $-126$  instead of  $-127$ , the gap between the largest subnormal number and the smallest normalized number is smaller.

Largest subnormal number			Smallest Positive Subnormal Number		
Sign	Exponent	Mantissa	Sign	Exponent	Mantissa
0	00000000	1111111111111111111111111111	0	00000000	00000000000000000000000000000001
$(1 - 2^{-23})2^{-126} = 0.999999988 \times 2^{-126} \approx 2^{-126}$			$2^{-23} \times 2^{-126} = 2^{-149} \approx 2^{-126}$		

## ★ Special Values in IEEE 754-1985 Standard (32 Bits) ★

Summary

Value	Sign	Exponent (8 bits)	Significand (23 bits)
+0	0	00000000	00 . . . 00
-0	1	00000000	00 . . . 00
$+1.F \times 2^{E-127}$	0	00000001 to 11111110	$aa \dots aa$ ( $a = 0$ or $1$ )
$-1.F \times 2^{E-127}$	1	00000001 to 11111110	$aa \dots aa$ ( $a = 0$ or $1$ )
$+\infty$	0	11111111	00 . . . 00
$-\infty$	1	11111111	00 . . . 00
SNaN	0 or 1	11111111	00 . . . 01 to 011 . . . 11 leading bit 0 at least one 1 in the rest
QNaN	0 or 1	11111111	1000 . . . 10 leading bit 1
L.P.F.P.N	0	11111110	11111111111111111111111111111111
S.P.F.P.N	0	00000001	00000000000000000000000000000000
Positive subnormal $0.F \times 2^{-126}$	0	00000000	00000001 to 11111111 at least one 1 in the rest

A photograph of a beach scene. In the foreground, several thatched umbrellas are set up on a sandy area. Some small tables are visible under the umbrellas. To the right, a red flag flies from a pole. The background shows the ocean with waves and a cloudy sky.

Thank you! Questions?