



We will start with a definition of floating point and double using IEEE format definitions. They are

- 32-bit: 1 sign bit, 8 exp bits, 23 frac bits
- 64-bit: 1 sign bit, 11 exp bits, 52 frac bits
- 8-bit: 1 sign bit, 4 exp bits, 3 frac bits

1. Provide the representation of the value 3.14159265 in 32-bit IEEE format.
2. Provide the representation of the value 3.14159265 in 64-bit IEEE.
3. Provide the representation of the value 3.14159265 in the 8 bit format

$$3.14159265 \Rightarrow 11.001001000011111010100111001$$

$$(3)_{10} = (11)_2$$

$$32\text{-bit bias: } 2^{(8-1)} - 1 = 127$$

$$64\text{-bit bias: } 2^{(11-1)} - 1 = 1023$$

$$8\text{-bit bias: } 2^{(4-1)} - 1 = 7$$

$$0.14159265$$

X	2
0	'2 8 3 1 '8 5 3 0
X	0 .5 '6 6 '3 2 0 6 0
X	1 .1 2 '2 2 4 1 2 0
X	0 .2 '6 5 '4 8 2 4 0
X	0 .8 3 0 9 6 4 8 0
X	1 0 6 1 9 2 9 6 0
X	0 .1 2 3 8 5 9 2 0
X	0 .2 4 5 4 4 2 6 0
X	0 .4 9 5 4 4 2 6 0
X	0 .9 9 0 8 7 3 0 0
X	1 .0 8 1 7 4 2 0 2
X	1 .9 6 3 1 4 9 4 0
X	1 .9 2 6 9 6 8 8 0
X	1 .8 5 3 1 7 2 6 0
X	1 .2 0 7 8 5 2 0 0
X	1 .4 1 5 9 9 1 0 4 0
X	0 .6 3 1 8 2 0 6 0
X	1 .6 6 3 6 4 1 6 0
X	1 .3 2 7 2 8 3 2 0
X	0 .6 5 1 5 6 6 4 0
X	1 .3 0 9 1 3 2 8 0
X	0 .6 1 8 2 6 4 6 0
X	1 .2 8 6 5 3 1 2 0 0
X	0 .4 2 3 0 6 2 0 0 0
X	0 .8 4 6 1 2 4 8 0
X	1 .8 9 2 4 9 6 0

0	1.8 9 > 2 14 9 6 0
X	1.7 6 > 1 9 4 9 2 1
X	1.5 0 > 1 8 9 1 8 4 0
X	1.1 3 > 1 9 9 6 8 0
X	0.2 > 1 5 9 1 2 6 0
X	0.5 5 > 1 0 6 2 4 0
X	0.1 3 > 0 9 8 4 6 0
X	0.2 0 > 0 9 4 8 8 0
X	0.4 1 5 8 9 2 6 0
X	0.8 3 1 3 9 5 2 0
X	1.0 6 3 5 9 9 9 0
X	1.2 3 1 1 8 6 0 0
X	0.6 0 9 7 9 6 0 0
X	1.7 0 5 2 7 2 2 0
X	0.6 1 2 4 6 1 0
X	1.2 3 1 8 9 2 9 0
X	0.4 6 4 9 6 5 6 0
X	0.8 3 1 9 5 2 1 0
X	1.8 > 9 14 2 4 0
X	1.7 5 6 2 6 4 8 0
X	1.9 1 6 5 9 4 6 0
X	1.0 3 3 1 9 2 0
X	0.0 6 0 2 8 4 0
X	0.1 3 2 5 6 4 9 0
X	0.2 6 5 1 7 0 0
X	0.5 7 0 2 9 2 0
X	1.0 6 0 4 5 4 8 0

$$\begin{aligned} &1.100 \dots \times 10^1 \\ &\epsilon: 1 \\ &E = 1 + 1023 = 1024 \\ &64 \quad (1024)_{10} = (10000000000)_2 = E \end{aligned}$$

$$32 \quad E = 1 + 127 = 128$$

$$(128)_{10} = (10000000)_2 = E$$

$$8: \quad E = 1 + 7 = 8$$

$$(8)_{10} = (1000)_2 = E$$

We will start with a definition of floating point and double using IEEE format definitions. They are

- 32-bit: 1 sign bit, 8 exp bits, 23 frac bits
 - 64-bit: 1 sign bit, 11 exp bits, 52 frac bits
 - 8-bit: 1 sign bit, 4 exp bits, 3 frac bits

1. Provide the representation of the value 3.14159265 in 32-bit IEEE format.

2. Provide the representation of the value 3.14159265 in 64-bit IEEE.

3. Provide the representation of the value 3.14159265 in the 8 bit format

| 0100 0000 | 0 0 | 00 | 00 | 0 0 0 0 | 1 1 1 | 1 B | 0 1 1 |
 4 0 4 9 0 F D B ↑
 0x40490FDB B/C Rounding
 ↓
 1 0 0 | 0 0 | 0 0 0 0 | 1 1 1 | 1 B | 1 0 1 0 | 0 0 1 1 | 0 0 | 1 0 1 0 | 1 0 0 1 1 1 | 0 0 0 1
 | 0100 0000 0000 | 1 0 0 | 0 0 | 1 0 | 0 0 0 0 | 1 1 1 | 1 B | 1 0 1 0 | 0 0 1 1 | 0 0 | 0 0 0 1 1 1 | 0 0 0 1
 4 0 0 9 2 1 F B 5 3 C 8 D 4 F 1
 0x400921FB53C8D4F1
 | 0100 0000 | 1 0 1
 8 5
 0x45 ← B/C Rounding

Converting Back:

64 Bit value

S: 0

$$E: 1000000000 \Rightarrow 1024 \therefore 1024 - 1023 = 1 = E_p$$

M: 10000000001111010100111001010011110001

$$\mu_{\text{normalized}} = 1,1000001000001111010100111001010011110001$$

$$1^S \times M_{\text{normalized}} \times 2^{E_M}$$

$$1 \times 1,1000001000001111010100111001010011110001 \times 2$$

$$2^k \quad 1 \quad 0 \quad 1 \\ 1,1000001000001111010100111001010011110001$$

↓

$$2^{15} + 2^{14} + 2^{13} + 2^{12} + 2^{11} + 2^{10} + 2^9 + 2^8 + 2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^1 + 2^0$$

$$= 3.14159265$$

32 bit value

S: 0

$$E: 1000000 \rightarrow 128 : 128 - 127 = 1$$

M: 00100000111101011

M-Normalized: 1.0010000111101011

$$1^s \times M_{\text{normalized}} \times 2^{E_{\text{exp}}}$$

$$\times 1.0010000111101011 \times 2^{\text{exp}}$$

$$2^{\text{exp}} \quad | \quad 0 \cdot 1 \\ | \quad 1.0010000111101011$$

C

$$2^{-2} + 2^{-3} + 2^{-6} + 2^{-11} + 2^{-12} + 2^{-13} + 2^{-14} + 2^{-15} + 2^{-16} + 2^{-17} + 2^{-18} + 2^{-19} + 2^{-20} + 2^{-21} + 2^{-22}$$
$$= 3.141592651$$

8 bit value

S: 0

F: 1000 \Rightarrow 8 i.e. 8-7 = 1 = E_b

M: 101

M-Normalized: 1.101

$$1 \times M-\text{normalized} \times 2^{\text{exp}}$$

$$(\times 1.101 \times 2^{\text{exp}})$$

$$\begin{array}{r} 2^1 \\ 2^0 \\ 1.01 \end{array}$$

$$2^1 + 2^0 + 2^{-2}$$

$$= 3.125$$

Precision Loss:

$$64 : 3.14159265$$

$$32 : 3.141592741$$

$$6 : 3.125$$

$$3.14159265 - 3.141592741 = -9.1 \times 10^{-8}$$

$$3.14159265 - 3.125 = 0.01659265$$

0.333

32bit

50

F. 2+12>123 : 0 111101

$$2^{-2} + 2^{-4} + 2^{-6} + 2^{-8} + 2^{-10} + 2^{-12} + 2^{-14} + 2^{-16} + 2^{-18} + 2^{-20} + 2^{-22} + 2^{-24}$$

C. left

5:0

$$\overline{z} = -2 + 2i = (0, 0, 1)_2$$

M. CIO

$$1.01 \times 10^{-2}$$

$$2^{-2} + 3^{-4} = 0.3125$$

$$0.\overline{3}333333 - 0.\overline{3}\overline{125} = 0.\overline{0208333135}$$