

Reserved Exponent Values

Exponent Value	Mantissa	Represents
✓ 11111111	All zeros	Infinity (∞)
✓ 11111111	Not all zeros	Not a number (NaN)
00000000	All zeros	Zero
<u>00000000</u>	Not all zeros	Subnormal (very small)

Practice problems

Convert the following decimal numbers to 32-bit IEEE 754 format by hand:

- a. 1.0
- b. -0.1
- c. 2016.0
- d. 0.00390625
- e. -3125.3125
- f. 0.33
- g. -0.67
- h. 3.14

[▼ Answer](#)

- a. 3f800000
- b. bdc00000
- c. 44fc0000
- d. 3b800000
- e. c5435500
- f. 3ea8f5c3
- g. bf2b851f
- h. 4048f5c3

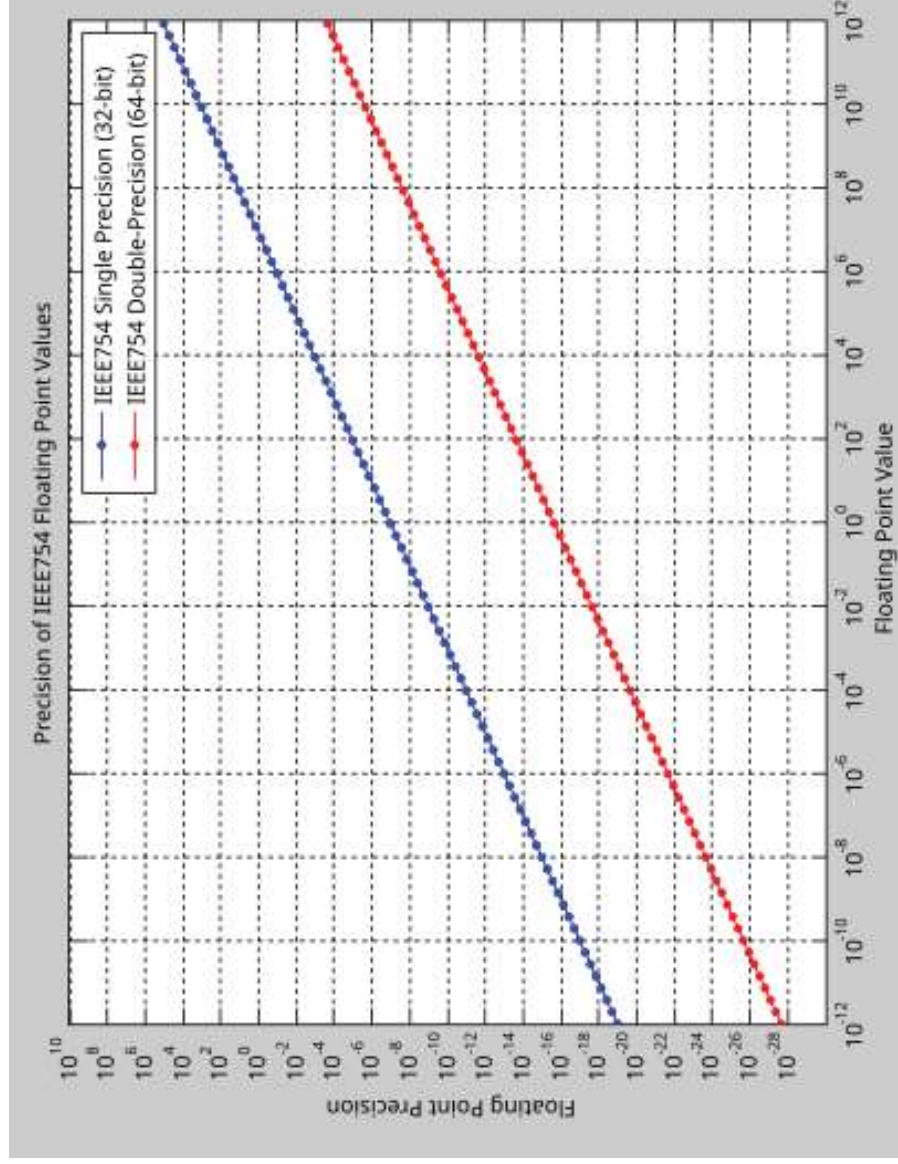
Convert the following hexadecimal numbers to decimal by hand using the 32-bit IEEE 754 format:

- a. 40000000
- b. bf800000
- c. 3d800000
- d. c1804000
- e. 42c81000
- f. 3f99999a
- g. 42f6e666
- h. c25948b4

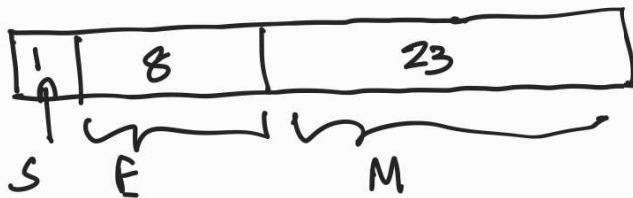
[▼ Answer](#)

- a. +2.0
- b. -1.0
- c. +0.0625
- d. -16.03125
- e. 100.03125
- f. 1.2
- g. 123.449997
- h. -54.320999

Precision and range in IEEE 754



https://en.wikipedia.org/wiki/IEEE_754



IEEE 754

Range: More is better
Resolution: Less is better.

Biased Exponent: 1 to 254

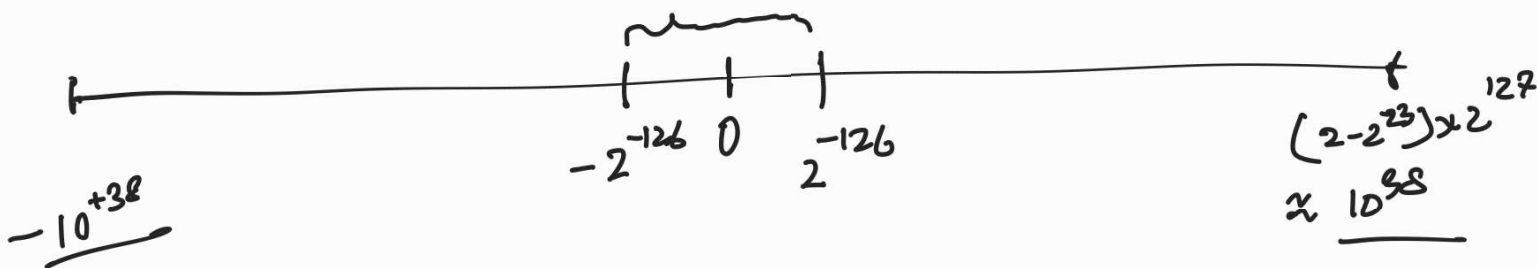
True exponent: -126 to 127

$$\begin{aligned} \text{Smallest value} &: 1.0000 \dots 0 \times 2^{-126} = 2^{-126} \\ \text{Largest value} &: 1.1111 \dots 1 \times 2^{127} = (2 - 2^{-23}) \times 2^{127} \end{aligned}$$

$$\begin{array}{r} 10.0000 \dots 0 \\ -0.0000 \dots 1 \end{array} \quad \begin{array}{l} \nearrow 2 \\ \searrow 2^{-23} \end{array}$$

2

Subnormal numbers



Range of signed numbers.

$$\begin{array}{c} -2^{-31} \text{ to } 2^{31} - 1 \\ \approx 10^{10} \end{array}$$

Resolution.

$$\text{Smallest} \left\{ \begin{array}{l} -1.0000 \dots 0 \times 2^{-126} \\ -1.0000 \dots 1 \times 2^{-126} \\ 1.0000 \dots 1 \times 2^{-126} \end{array} \right\}$$

Signif. res.

Resolution = 1 ✓

$$2^{-23} \times 2^{-126} = 2^{-149} \approx 10^{-45}$$

$$\begin{array}{r} 1.000 \dots 0 \times 2^0 \\ 1.000 \dots 1 \times 2^0 \\ \hline 1.000 \dots 1 \times 2^0 \end{array}$$

$$\rightarrow 2^{-23} \times 2^0 = 2^{-23}$$

$$\begin{array}{r}
 1.00000 \dots 0 \times 2^{25} \\
 1.00000 \dots 1 \times 2^{25} \\
 \hline
 2^{-23} \times 2^{25} = 2^2 = \boxed{4}
 \end{array}$$

Worst case:

$$\begin{array}{r}
 1.000 \dots 0 \times 2^{127} \\
 1.000 \dots 1 \times 2^{127} \\
 \hline
 2^{-23} \times 2^{127} = 2^{105} \approx \boxed{10^{31}} \checkmark
 \end{array}$$

Floating point addition.

$$\begin{array}{r}
 \rightarrow 9.876 \times 10^4 \\
 \rightarrow + 5.678 \times 10^3 \rightarrow
 \end{array}$$

① Exponent Comparison.

$$4 > 3$$

② Mantissa alignment:

$$\begin{array}{r}
 9.876 \times 10^4 \\
 0.5678 \times 10^4
 \end{array}$$

③ Perform operation on mantissa.

$$\begin{array}{r}
 9.876 \\
 0.5678 \\
 \hline
 10.4438
 \end{array}$$

④ Apply the common exponent to mantissa.

$$\begin{array}{r}
 10.4438 \times 10^4 \\
 \hline
 1.04438 \times 10^5
 \end{array}$$

Subtraction:

$$\begin{array}{r}
 9.876 \times 10^4 \\
 - 5.678 \times 10^3 \\
 \vdots \\
 9.876 \\
 - 0.5678 \\
 \hline
 \end{array}$$

Convention:

the larger exponent value is chosen for alignment.

$$\begin{array}{r}
 \downarrow \\
 9.3082 \times 10^4 \\
 \downarrow
 \end{array}$$

Binary Addition

α/β			
Saw	Fire	T	F
	T	✓	X
	F	X	✓

α

β

(more serious,
fire but alarm
does not sound).

$$\begin{array}{r}
 0.8 \\
 \hline
 11001100. \\
 100000 \\
 \hline
 11001111 \\
 + 1111 \\
 \hline
 11010
 \end{array}$$