Edwin Huang
11/30/2021

1. Convert following decimal #s to binary #s.

a) 5.75

$\rightarrow$ 5.00 + 0.75

$5_{10} = 0101$ $\longrightarrow$ $\boxed{101.11}$

$0.75 = \frac{3}{4} = .11$

b) $\frac{63}{64}$

$64_{10} = 1\ 000\ 000$ 

$\phantom{64_{10} =}\ -\ \phantom{000\ 0}1$ 

$\begin{array}{c} 64 \\ -1 \\ \hline 63 \end{array}$

$\phantom{64_{10} =}\ 0\ 111\ \ \ \ 111 \rightarrow \boxed{0.111\ 111}$

c) 9.8125

$\rightarrow 9_{10} = 8_{10} + 1_{10} = \boxed{1001}$

$.8125 = \frac{13}{16} = \boxed{.1101}$

$.8125 \cdot 2 = \textcircled{1}.625$

$0.625 \cdot 2 = \textcircled{1}.25$ $\quad \rightarrow 1101$

$0.25 \cdot 2 = 0.5$

$0.5 \cdot 2 = \textcircled{1}$

$= \boxed{1001.1101}$

2. Convert 34.890625 into the IEEE 754 floating-point rep.

1) Sign: $\oplus$, Positive

2) Exponent: $34 \rightarrow 100010$

$\phantom{2) Exponent:}\ .890625 \rightarrow 111001$

$34.890625 = 100010.111001 = 1.001011001 \times 2^5$

3) Mantissa: 001011001

| Sign: | 0 |
|---|---|
| Exponent: | 0000 0101 |
| Mantissa: | 001011001 |

3. Convert $0, 01111011, 00000...00$ to decimal.
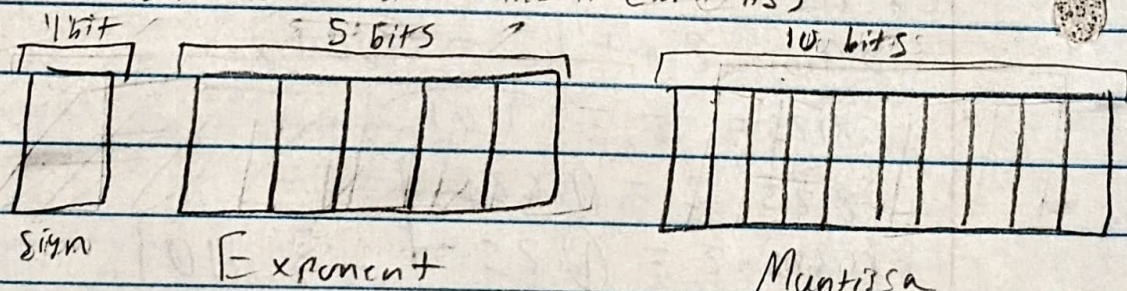
→ Sign: 0 (+)

   Exponent: $0111 1011_2$

     $= 64 + 32 + 16 + 8 + 2 + 1 = 123_2$

→ $123 - 127 = -4$ (since the bias in 32-bits is 127)
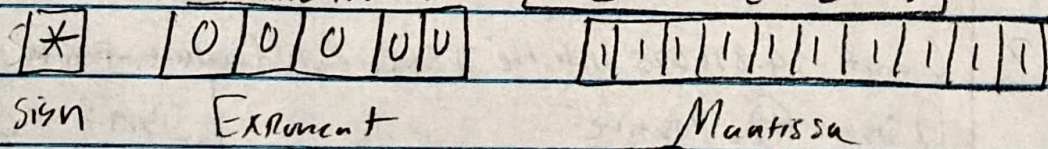
   $= 1.000\ 0000\ 0000\ 0000\ 0000\ 0000 \times 2^{-4}$

Mantissa: $0.001 = \frac{1}{2^{-4}} = \frac{1}{16} = \boxed{0.0625}$

4. Explain the definition of denormalized number and show the largest denormalized # and smallest normalized # (for ① #s)



| 1 bit | 5 bits | 10 bits |
| --- | --- | --- |
| Sign | Exponent | Mantissa |

A denormalized number means that it is a number that is less than 0 and is used to act as a offset in floating point arithmetic.

Largest Denormalized Number: $\boxed{\pm 2^{-14} \cdot (1 - 2^{-10})}$

| * | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Sign | Exponent | | | | | Mantissa | | | | | | | | | |

Smallest Normalized Number: $\boxed{\pm 2^{-127}}$

| * | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Sign | Exponent | | | | | Mantissa | | | | | | | | | |