# Lab 4 Pre-lab: Floating Point Conversion

Justin Gou (jyg2qhc)

February 6, 2020

## 1 Float to Hex

**Your magic (32 bit) floating point number is 70.25 This is the number that needs to be converted to (little endian) binary, and expressed in hexadecimal.**

### 1.1 Solution

We know that 70.25 in binary scientific notation uses base $2^6$, as 64 is the largest power of 2 in 70.25.
By dividing $70\frac{1}{4}$ by 64, we get $\frac{281}{256} = 1.09766$, which after subtracting 1 is just $\frac{25}{256}$.
To produce the fraction using negative powers of two, this fraction can be written as

$$\frac{25}{256} = \frac{16}{256} + \frac{8}{256} + \frac{1}{256}$$
$$= \frac{1}{16} + \frac{1}{32} + \frac{1}{256}$$
$$= (\frac{1}{2})^4 + (\frac{1}{2})^5 + (\frac{1}{2})^8$$

Based on this, we know the mantissa is 000 1100 1000 0000 0000 0000
Earlier, we said the exponent would be 6 but it is represented by adding 127, so we convert 133 to binary, which is 1000 0101.
Finally, we know the sign of the float is positive, meaning it can be represented as a 0.
By combining the three pieces, we find the binary string

0100 0010 1000 1100 1000 0000 0000 0000

Write each group of four digits as a hex digit to get 42 8c 80 00
The prompt requests the answer in little endian, which is simply

$$\boxed{0x00808c42}$$

## 2 Hex to Float

**Your other magic floating point number is, in hex, 0x00009ec3 This is the number that needs to be converted to a (32 bit) floating point number. Note that the hexadecimal printed above is in little-endian format!**

### 2.1 Solution

First, let's convert the number to big-endian: **c3 9e 00 00**
Convert the number from hex to binary: **1100 0011 1001 1110 0000 0000 0000 0000**
Split this binary number into its sign bit, exponent, and mantissa

Sign bit : 1

Exponent : 1000 0111

Mantissa : 001 1110 0000 0000 0000 0000

Sign bit is 1, meaning the number is negative.
Exponent binary converted to decimal is $128 + 4 + 2 + 1 = 135$; subtract 127 from this number to get the exponent : 8 (this tells us the number will be multiplied by $2^8$
Mantissa converted to a decimal is $(\frac{1}{2})^3 + (\frac{1}{2})^4 + (\frac{1}{2})^5 + (\frac{1}{2})^6 = 0.234375$ — add one to this number to get 1.234375

Putting the three pieces together, we get $-1.234375 \cdot 2^8 = -1.234375 \cdot 256 = \boxed{-316}$