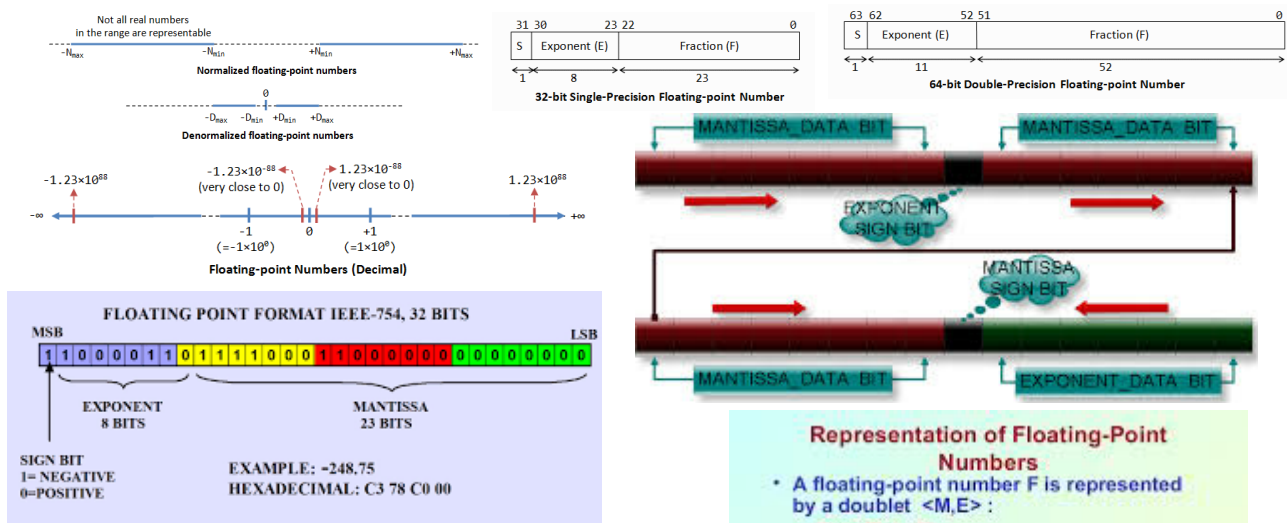# FLOATING POINT VARIABLES

Integers and floats are two different kinds of numerical data. An integer is a number without a decimal point where as a float is a floating point number, i.e. it is a number with a decimal point.Floating point numbers are stored in this format: M x b^e, where m is the mantissa (an integer number), b is base and e is exponent. C supports two floating types: float and double. The float and double are represented using 32-bit single precision and 64-bit double precision. For single precision floating point we have: 1 sign bit, 8 exponent bits, 23 mantissa bits. For double precision floating points we have: 1 sign bit, 11 exponent bits, 52 mantissa bits.Following figure illustrate how floating point number is stored in memory:

Five important rules to be followed:

Rule 1: To find the mantissa and exponent, we convert data into scientific form.

Rule 2: Before the storing of exponent, 127 is added to the exponent.

Rule 3: Exponent is stored in memory in first byte from right to left side.

Rule 4: If exponent is negative number it will be stored in 2's complement form.

Rule 5: Mantissa is stored in the memory in second byte from right to left side.

Example: Memory representation of float a= -10.3f