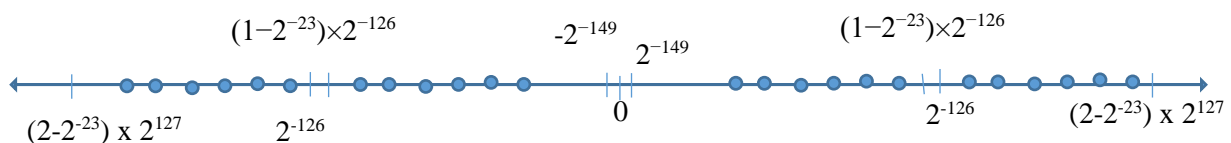# ARM Architecture

Q.1 Does any of the above three components play a role in the defining the Precession of the number? If so which are the component or Components which play the role in defining precession and how ? Explain this with example in your own words.

Ans: Here The sign bit does not affect the precision. The size of mantissa and the signed exponent combined gives details of precession. The digits in mantissa will tell up to what extent the number in close to the result. The exponent represents the range in which the number is in(Very small or very big.)

Q.2 What is Normal and Subnormal Values as per IEEE 754 standards? Explain this with the help of number line.

Ans: For single precession, the numbers in range $\pm 2^{-126}$ to $(2-2^{-23})$ x $2^{127}$ are called normal numbers. While the numbers less than $2^{-126}$ are subnormal numbers. $(\pm 2^{-149}$ to $(1-2^{-23})\times 2^{-126})$. Drawing on number line….



Q.3 IEEE 754vv defines standards for rounding floating point numbers to a representable value. There are five methods defines by IEEE for this – Take time and understand what these five methods and explain it in your words using diagrams, illustrations of your own.

Ans:

i.   Rounding to Nearest, ties to even: Here the number is rounded off to nearest value. If number falls in between, the least significant bit written as zero.
   a.   Eg. 15.1 to 16.9 will be rounded to 16
   b.   Here in binary LSB will be zero.
ii.  Rounding to Nearest, ties away from 0: Here the number is rounded off to nearest value. If number falls in between, the LSB is will be the LSB of nearest larger integer(+ve Numbers).
   a.   Like 1.56 = 2,
   b.   4.6 = 5
   Here for negative numbers, the lesser integer is accounted for.
   c.   -1.56 = -2
   d.   -4.6 = -5

iii. Round to 0: Here the nearest integer less than given number for +ve numbers and greater for -ve will be considered.
   a.   1.5 = 1
   b.   -1.5 = -1

iv.    Round to $+\infty$: Here for all numbers immediate greater integer is considered.

      a.  5.3 = 6
      b.  -5.3 = 5

v.    Round to $-\infty$: Here for all numbers immediate lesser integer is considered.

      a.  5.3 = 5
      b.  -5.3 = -6