

215: fixed point  
1010: number

### Representation of floating point numbers

Floating point Number: Contains two parts:

- (i) Mantissa: Contains signed fixed point number
- (ii) Exponent: specify the position of decimal (binary) point.

$$\boxed{\pm m \times r^e}$$

↓      ↓      ↓  
Sign mantissa   radix  
(Base)

Examples:

$$+215.37 = +0.\underline{21537} \times 10^3$$

↓      ↓      ↓  
Sign   Mantissa   Base  
(radix)

$$+1000.110 = +0.\underline{1000110} \times 2^4$$

↓      ↓      ↓  
Sign   Mantissa   Base  
(radix)

$$-10101 = -0.\underline{10101} \times 2^5$$

↓      ↓      ↓  
Sign   Mantissa   Base

These numbers are called floating point numbers because the position of the decimal (binary) point is fluctuating (floating).

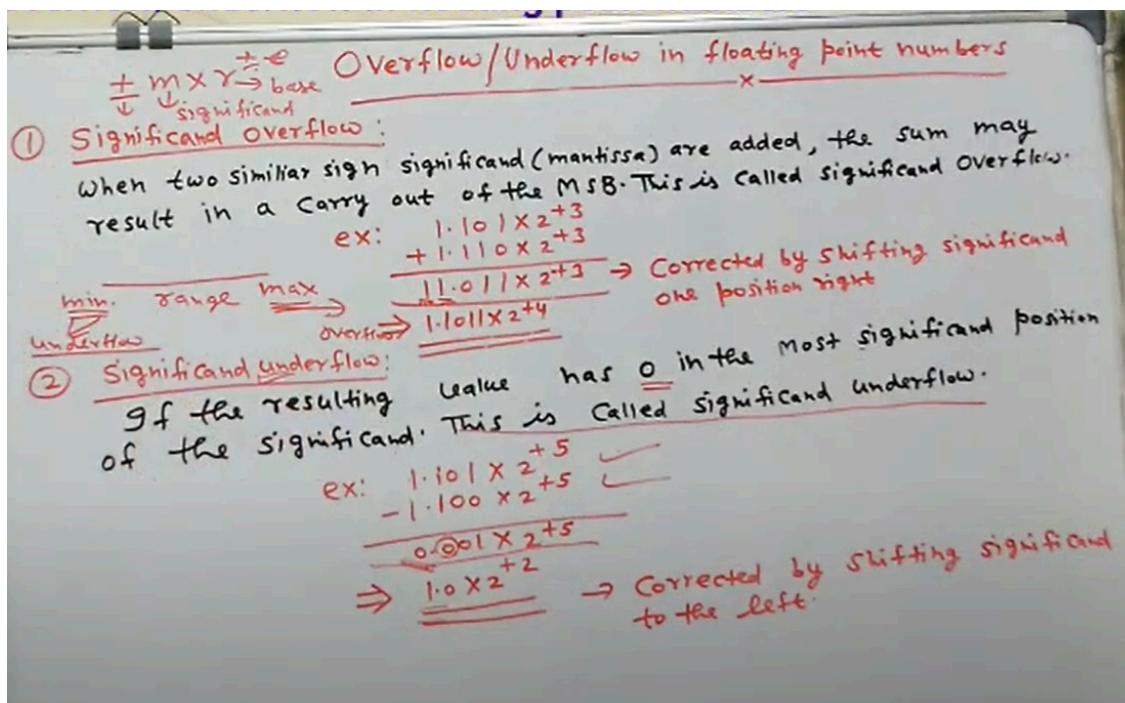
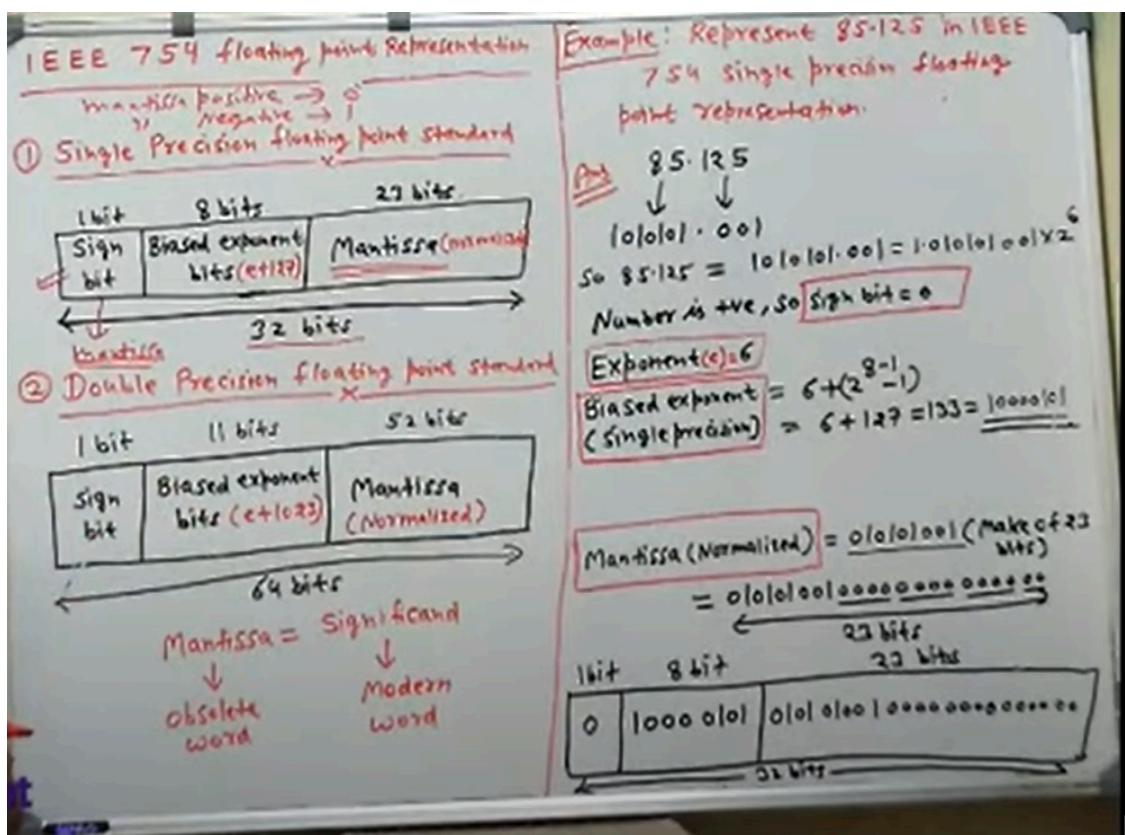
### Normalization:

A floating point number is said to be normalized, when we force the integer part of its mantissa to be 1, and allow its fractional part varying.

ex:  $13.25 = 1101.01_2 = 1.10101 \times 2^4$  (Normalized)

↓      ↓  
integer   fractional  
part   part

Normalized floating point number =  $1.\underline{fffff} \times 2^{exponent}$   
not stored in representation



### ③ Exponent Overflow:

An exponent overflow occurs, when a resulting positive exponent exceeds the maximum possible exponent value.

### ④ Exponent Underflow:

An exponent underflow occurs, when a resulting negative exponent is less than the minimum possible value.

Floating point Addition/Subtraction

Algorithm for floating point addition/subtraction

- ① Check for zeros: If either no. is zero, result will be other numbers with appropriate sign.
- ② Align the Mantissas: If exponents are equal, perform the arithmetic operation; else shift the mantissa with smaller exponent to the right until its exponent equal to the larger exponent.
- ③ Perform the addition/subtraction depending on the operation and sign of the two mantissas.

Let  $X = X_A * 2^{E_A}$   
 $Y = X_B * 2^{E_B}$  }  $Z = ?$

① During addition: If  $Y=0$ ,  $Z = X = X_A * 2^{E_A}$   
If  $X=0$ ,  $Z = Y = X_B * 2^{E_B}$

During subtraction: If  $Y=0$ ,  $Z = X = X_A * 2^{E_A}$   
If  $X=0$ ,  $Z = -Y = -X_B * 2^{E_B}$

② Let  $X = 1.10 * 2^3$  (Larger exponent)  
 $Y = 1.01 * 2^1$  (Smaller exponent)

↓  
mantissa

Shift the mantissa of  $Y$  to the right, until its exponent is equal to the exponent of  $X$ .

$Y = 0.10 * 2^2$  (One time shift)  
 $Y = 0.01 * 2^3$  (Two time shift)

NOW  $X = 1.10 * 2^3$   
 $Y = 0.01 * 2^3$

Go to step 3

\* During addition, if two significand are equal with opposite sign, the result will be zero.

\* There is also a possibility of significand (mantissa) overflow by 1 digit, then the mantissa of the result is shifted right & exponent is incremented.

Significant overflow occurs, when the addition of two significant (mantissa) of same sign may result in a carry at MSB of mantissa.

as we increment exponent, there is also a possibility of exponent overflow. Then result will show "error."

(4) Normalize the result: If result is not normalized, then we shift mantissa left and decrement the exponent until value "Left of binary point is 1." As we are decreasing the exponent, so exponent can also be underflow. Then again error is reported.

$$X = 1.01 * 2^3$$

$$Y = -1.01 * 2^3$$

$$\text{then } X+Y = \frac{1.01 * 2^3}{-1.01 * 2^3}$$

$$Z = X+Y = \frac{0 * 2^3}{0 * 2^3} = 0$$

$$\text{If } X = 1.01 * 2^3$$

$$Y = 1.10 * 2^3$$

$$Z = X+Y = \underline{\underline{1.011 * 2^3}}$$

↓  
Carry at MSB

$$\text{Correction} = \underline{\underline{1.011 * 2^4}}$$

$$\text{If } X = 1.01 * 2^3$$

$$Y = -1.00 * 2^3$$

$$Z = X+Y = \frac{0.01 * 2^3}{1.0 * 2^1}$$

(Not normalized)  
(Normalized)

Next "flow chart for floating point add / sub"

