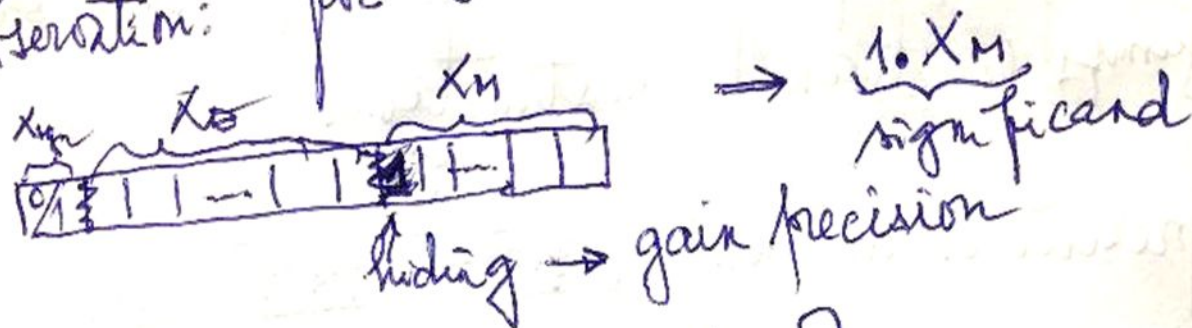
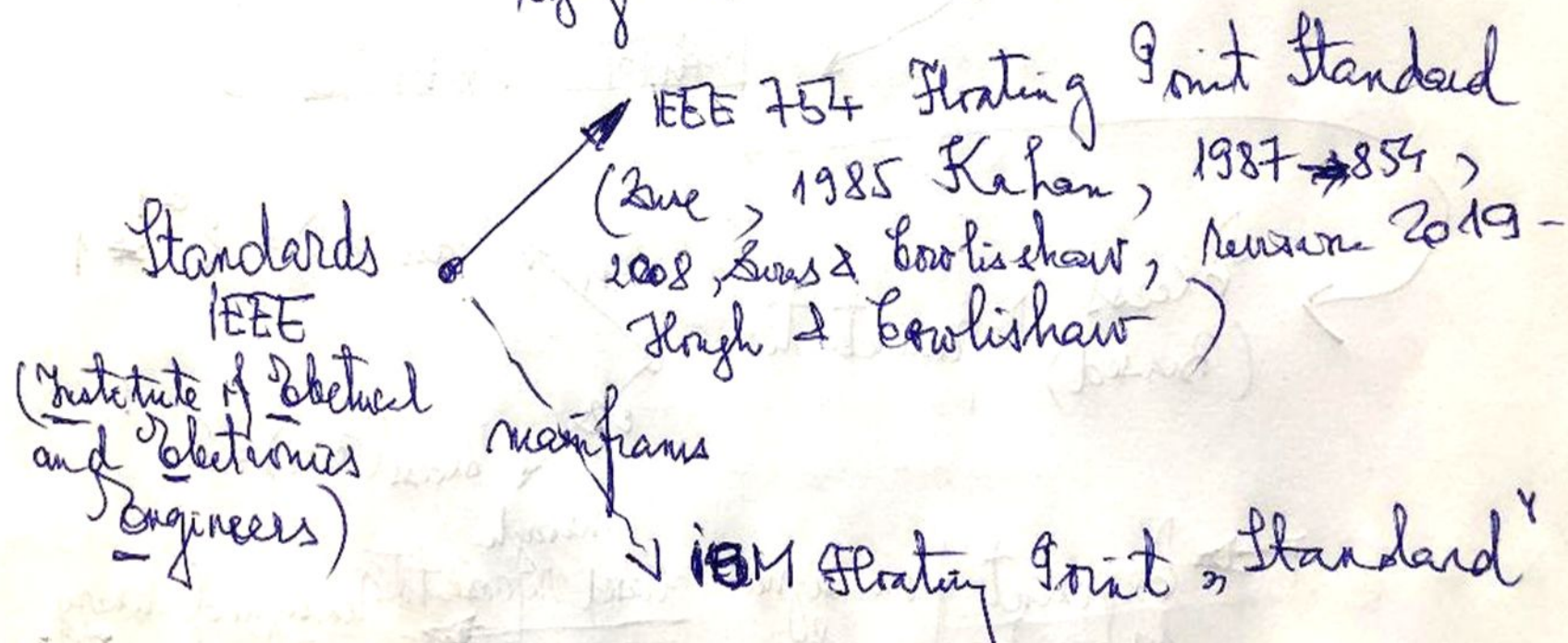


Exponent bit pattern	Unsigned value	Signed value	
		Bias = 127	Bias = 128
11111111	255	+128	+127
11111110	254	+127	+126
⋮	⋮	⋮	⋮
10000001	129	+2	+1
10000000	128	+1	0
01111111	127	0	-1
01111110	126	-1	-2
⋮	⋮	⋮	⋮
00000001	1	-126	-127
00000000	0	-127	-128

* Observation: for SM mantissas

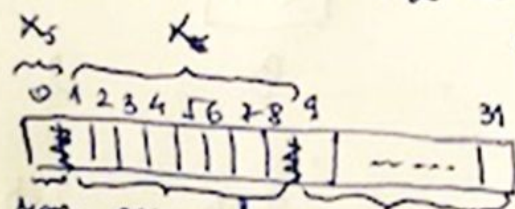


$$1 \leq \underbrace{1.X_M}_{\text{significand}} < 2$$



IEEE 754 Floating Point Standard

32 bits



sign
exponent
excess 127
sign-magnitude
binary integer

fraction part
of sign-magnitude
binary significand
with hidden integer bit

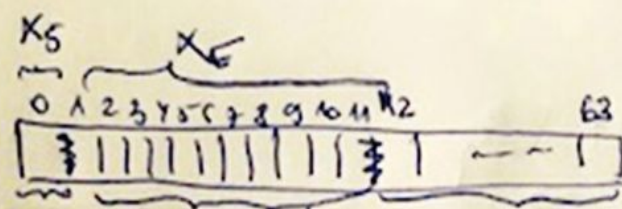
$$X = (-1)^{X_s} \cdot 2^{X_E - 127} \cdot (1.X_f)$$

where $0 < X_E < 255$

hidden bit
significand

IEEE 754

64 bits



sign
exponent

$$X = (-1)^{X_s} \cdot 2^{X_E - 1023} \cdot (1.X_f)$$

where $0 < X_E < 2047$

negative
overflow

X_1

negative
underflow

X_2

0

positive
underflow

X_3

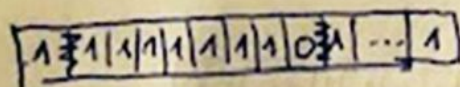
$X_3 = |X_2|$

valid
numbers

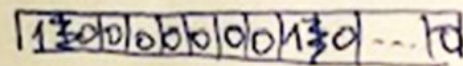
positive
overflow

X_4

$X_4 = |X_1|$



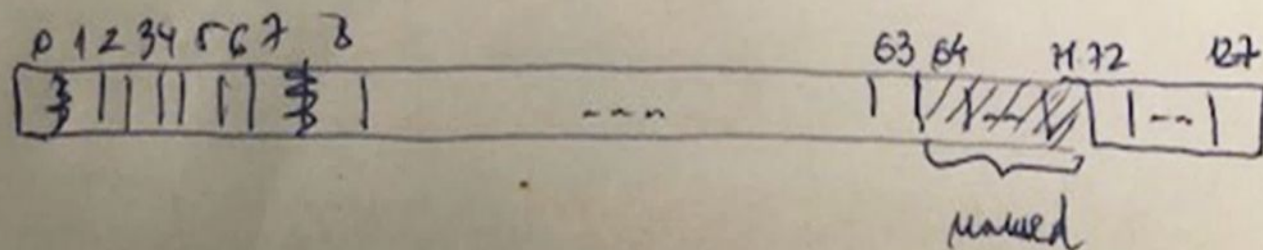
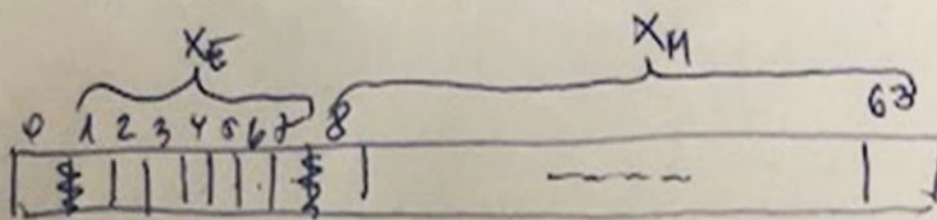
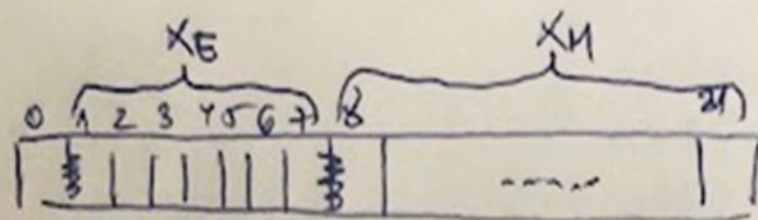
$$X_1 = (-1)^1 \cdot 2^{255-127} \cdot (1 + \frac{1}{2} + \dots + \frac{1}{2^{23}}) = -2^{127} (2 - 2^{-23})$$



$$X_2 = (-1)^0 \cdot 2^{1-127} \cdot (1 + 0 + \dots + 0) = -2^{-126}$$

IBM Floating Point Standard

$$X = (-1)^{X_S} \cdot 16^{X_E - 64} \cdot (0.X_M)$$



IBM

32 bits

64 bits

128 bits

$$X = -794, 08984375_{10} = -1100011010, 00010111_2$$

$$X = (-1)^{x_s} \cdot 2^{x_E - 127} \cdot (1.x_n)$$

794 | 2
 0 | 397 | 2
 1 | 198 | 2
 0 | 99 | 2
 1 | 49 | 2
 1 | 24 | 2
 0 | 12 | 2
 0 | 6 | 2
 0 | 3 | 2
 1 | 1

$$0,08984375 \times 2$$

$$0,17968750 \times 2$$

$$0,35937500 \times 2$$

$$0,71875000 \times 2$$

$$1,43750000 \times 2$$

$$0,87500000 \times 2$$

$$1,75000000 \times 2$$

$$1,50000000 \times 2$$

$$1,00000000$$

$$X = -1,10001101000010111 \times 2^9$$

$$X_E = 127 = 9 \Rightarrow X_E = 136$$

1	1	0	0	0	1	0	0	0	1	0	0	0	1	1	0	1	0	0	0	0	1	0	1	1	1	0	0	0	0		
C				4				4				6				8				5				C				0			

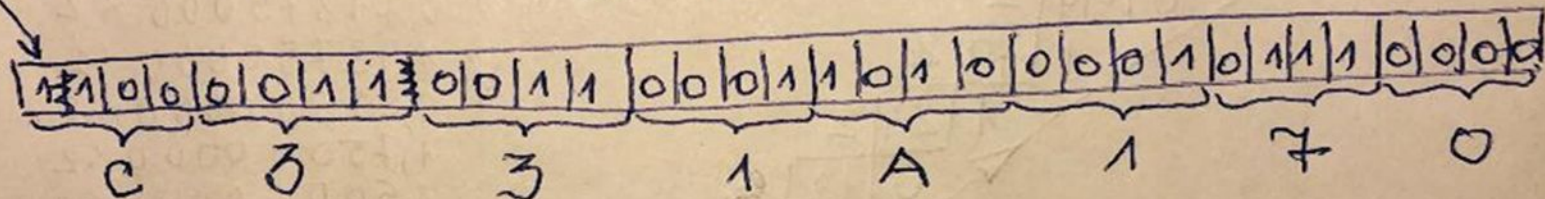
$$X_{16} = C44685C0_{16}$$

$$X = -794,08984375_{10} = -1100011010,00010111_2 =$$

$$= -0,00110001101000010111 \times 16^3$$

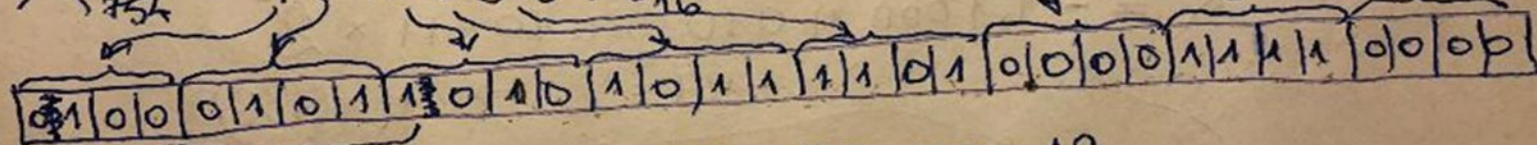
$$X = (-1)^{X_S} \cdot 16^{X_E-64} (0 X_M)$$

$$X_E - 64 = 3 \Rightarrow X_E = 67$$



$$X_{IBM} = C331A170_{16}$$

$$X_{754} = 45ABD0F0_{16}$$



$$139 = X_E \Rightarrow X_E - 127 = 139 - 127 = 12$$

$$X_{754} = (-1)^{X_S} \cdot 2^{X_E-127} \cdot 1,0101011110100001111 =$$

hidden bit

$$= 1010101111010,0001111_2 =$$

$$= + (4096 + 1024 + 256 + 64 + 32 + 16 + 8 + 2 + \frac{15}{16}) =$$

$$= + 5498,1171875_{10}$$