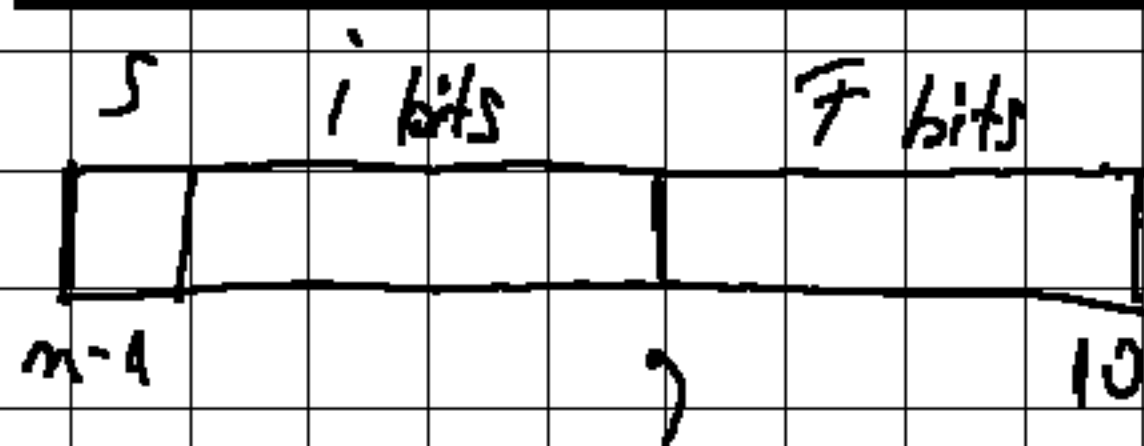


Fixed point repres. of real numbers

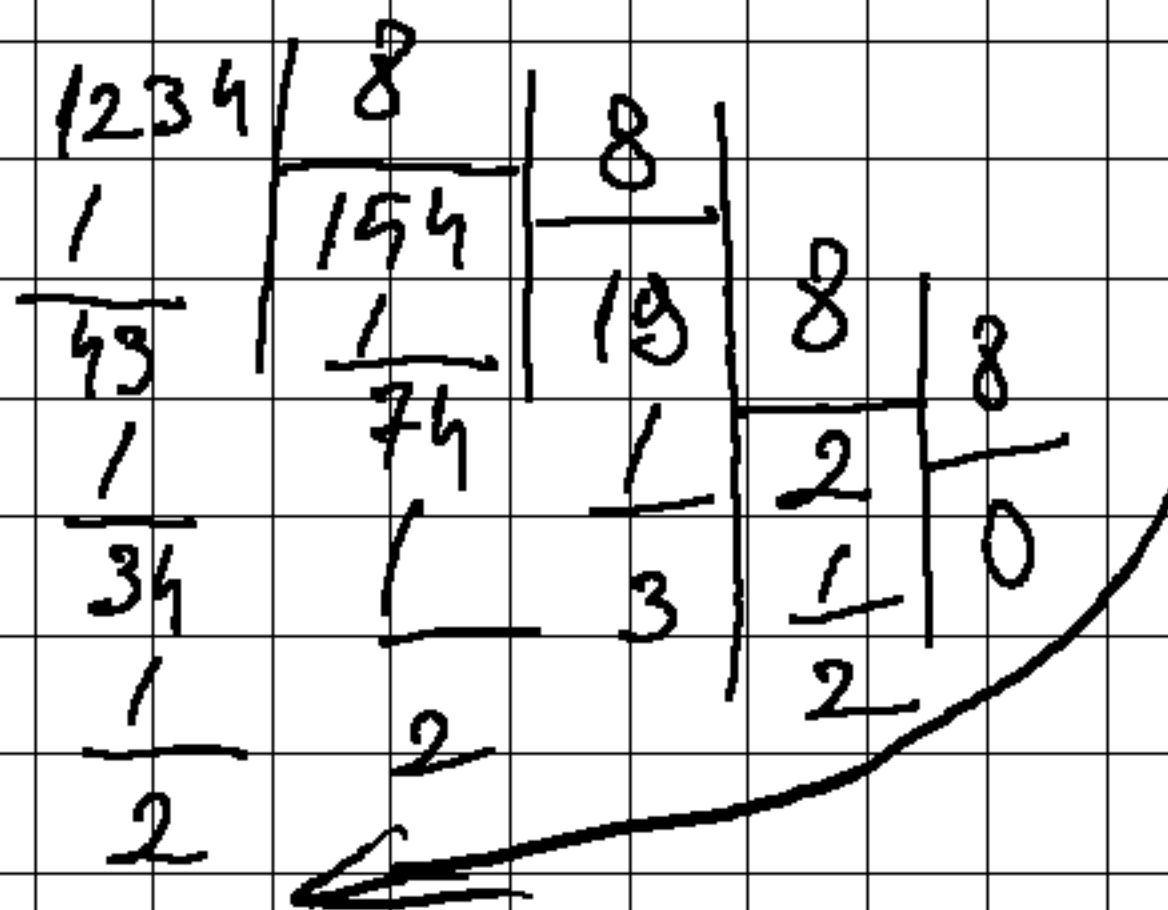


$$n = 1 + i + F \text{ bits}$$

$$n = 32 \text{ bits}, \quad i = 15 \text{ bits}, \quad F = 16 \text{ bits}$$

$$X = 1234, 56$$

$$1234 = 2322_{(8)} = 010011010010_{(2)}$$



$$0,56 \cdot 8 = (4)48$$

$$0,48 \cdot 8 = (3)84$$

$$0,84 \cdot 8 = (6)72$$

$$0,72 \cdot 8 = (5)76$$

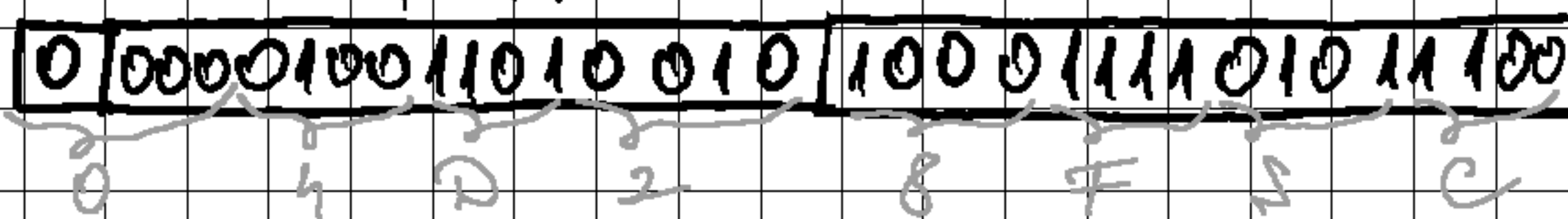
$$0,76 \cdot 8 = (6)08$$

$$0,08 \cdot 8 = (0)64$$

$$0,56 = 0,436560_{(8)} =$$

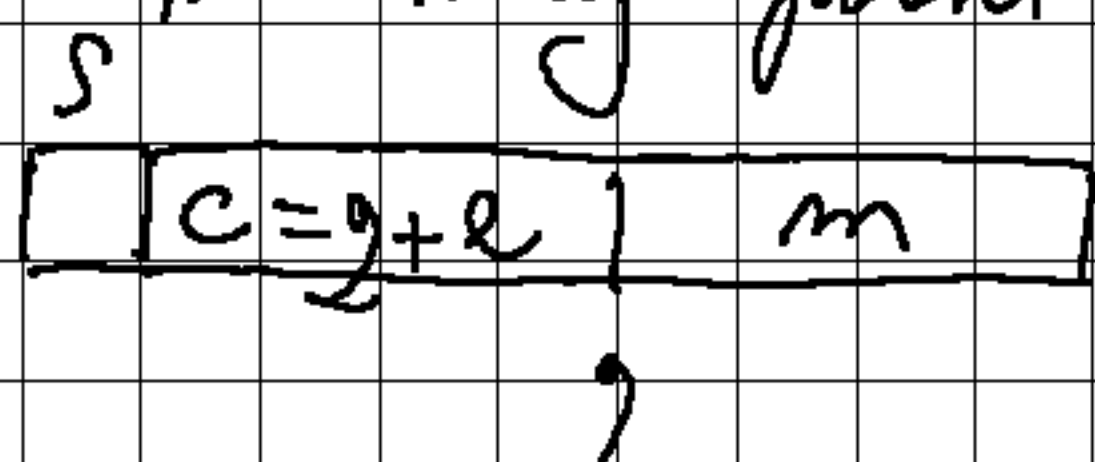
$$= 0,100011110101110000_{(2)}$$

$$S \quad i = 15 \text{ bits} \quad F = 16 \text{ bits}$$



Homework: 10.5.

Ex2: Floating-point representation



$$1234,56 = 0,123456 \cdot 10^4 \rightarrow \text{exp}$$

$$0,013 = 0,13 \cdot 10^{-1}$$

Single precision: 32 bits, $2=127$, C on 8 bits, m on 23 bits

double precision: 64 bits, $2=1023$, C on 11 bits, m on 52 bits

$$X = -465,78, \quad sp, \quad m > 1$$

465	8	8	8
1	58	7	0
65	1	1	
1	2	7	
1			

$$456 = 721_{(8)}$$

$$0,78 = 0,61727_{(8)}$$

$$0,78 \cdot 8 = 6,24$$

$$0,24 \cdot 8 = 1,92$$

$$0,92 \cdot 8 = 7,36$$

$$0,36 \cdot 8 = 2,88$$

$$0,88 \cdot 8 = 7,04$$

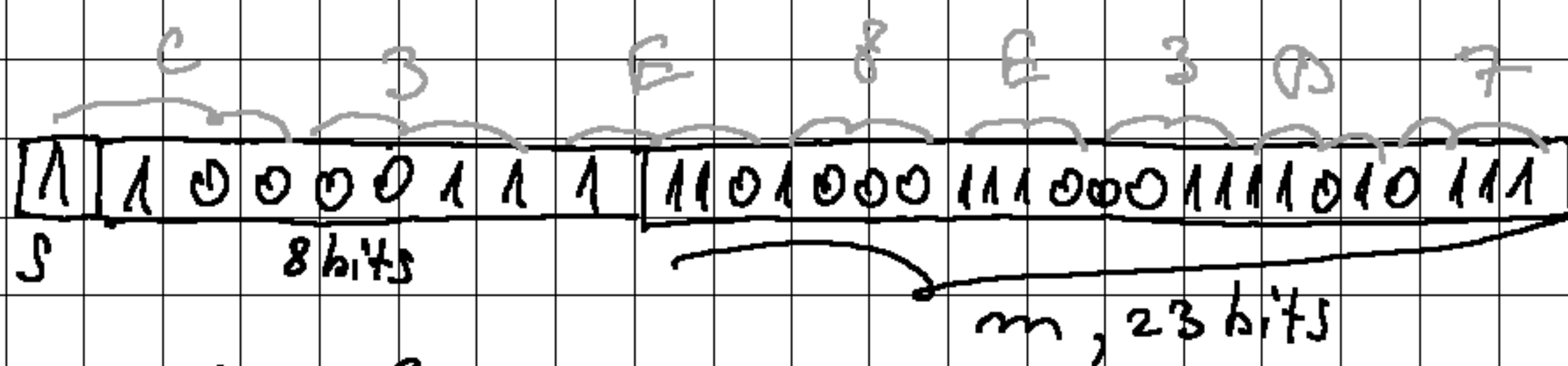
$$X = -465,78 = -721,61727$$

$$= -111010001,110001111010111$$

$$= -1,11010001110001111010111$$

hidden bit

m



$$C = 127 + 8 = 135$$

$$135 = 128 + 4 + 2 + 1 = 2^7 + 2^2 + 2^1 + 2^0 = 10000111_2$$

Ex 3: $X = 0,17$, Sp , $m < 1$

$$0,17 \cdot 8 = 1,36$$

$$0,36 \cdot 8 = 2,88$$

$$0,88 \cdot 8 = 7,64$$

$$0,64 \cdot 8 = 5,12$$

$$0,12 \cdot 8 = 0,96$$

$$0,96 \cdot 8 = 7,68$$

$$0,68 \cdot 8 = 5,44$$

$$0,44 \cdot 8 = 3,52$$

$$0,52 \cdot 8 = 4,16$$

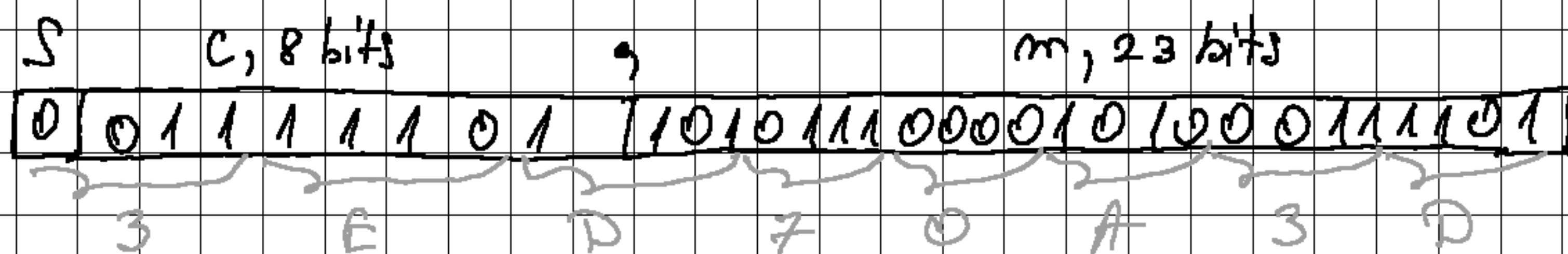
$$0,17 = 0,127024369_{(8)}$$

$$= 0,001010111000010100011110101_{(2)}$$

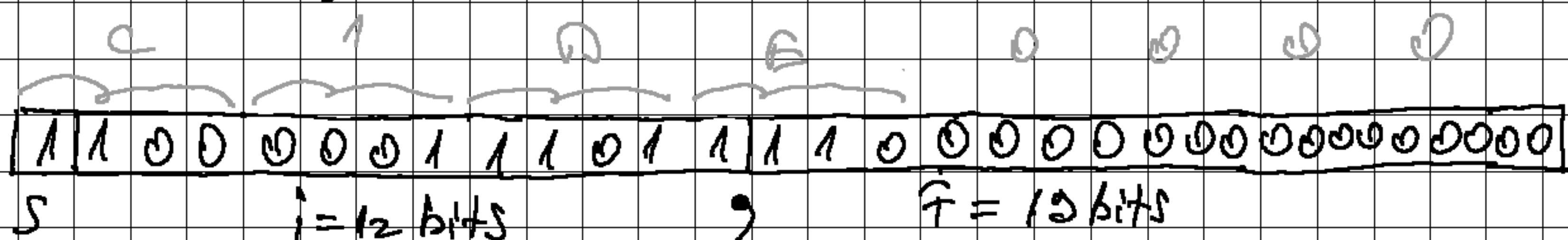
$$= 0,1010111000010100011110101_{(2)} \cdot 10^{-2} \rightarrow \text{exp}$$

$$C = 2 + 2 = -2 + 127 = 125 = 64 + 32 + 16 + 8 + 4 + 1 =$$

$$= 2^6 + 2^5 + 2^4 + 2^3 + 2^2 + 2^0 = 01111101_{(2)}$$



Ex 4: Find the real number x having $C1DE0000_{(16)}$ as it's fixed-point repes on 32 bits with $i = 12$ bits and $F = 19$ bits.

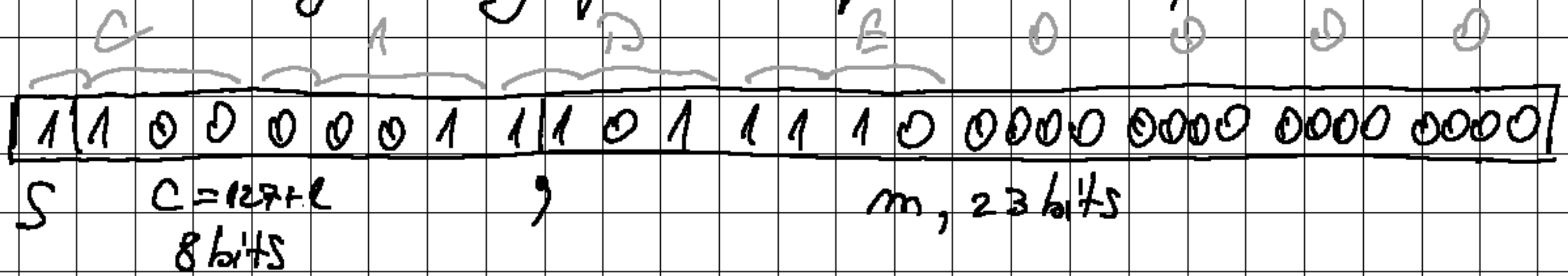


$$X = -1000000111011, 11_{(2)} =$$

$$= - \left(2^{11} + 2^5 + 2^4 + 2^3 + 2^1 + 2^0 + 2^{-1} + 2^{-2} \right) =$$

$$= - \left(2048 + 32 + 16 + 8 + 2 + 1 + \frac{1}{2} + \frac{1}{4} \right) = -2107,75$$

Ex: Find the real number x having $C1DE0000_{(16)}$ as it's floating point repres. $SP, m > 1$



$$C = 10000001_{(2)} = 2^7 + 2 + 1 = 128 + 3 = 131$$

$$e = 131 - 127 = 4$$

$$X = -\underset{\substack{\text{hidden} \\ \text{bit}}}{1}, 101111 \cdot 10^4 = -11011, 11_{(2)}$$

$$-11011, 11_{(2)} = -\left(2^4 + 2^3 + 2 + 1 + 0,5 + 0,25\right) =$$

$$= -27,75_{(10)}$$