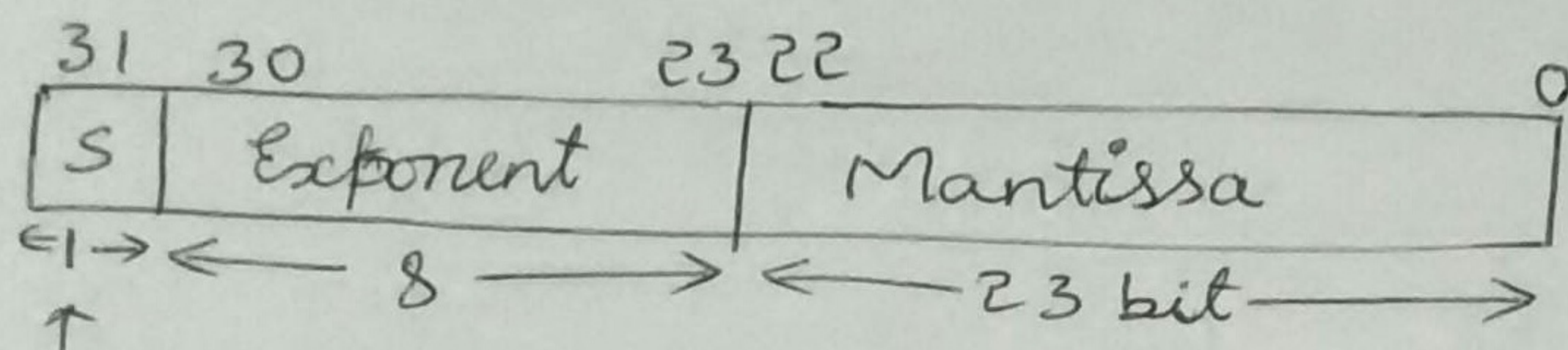


Floating point data type

- Float (32 bit)
- Double (64 bit)

Float 32 bit Format (Single Precision)



Sign bit

1 → -ve number

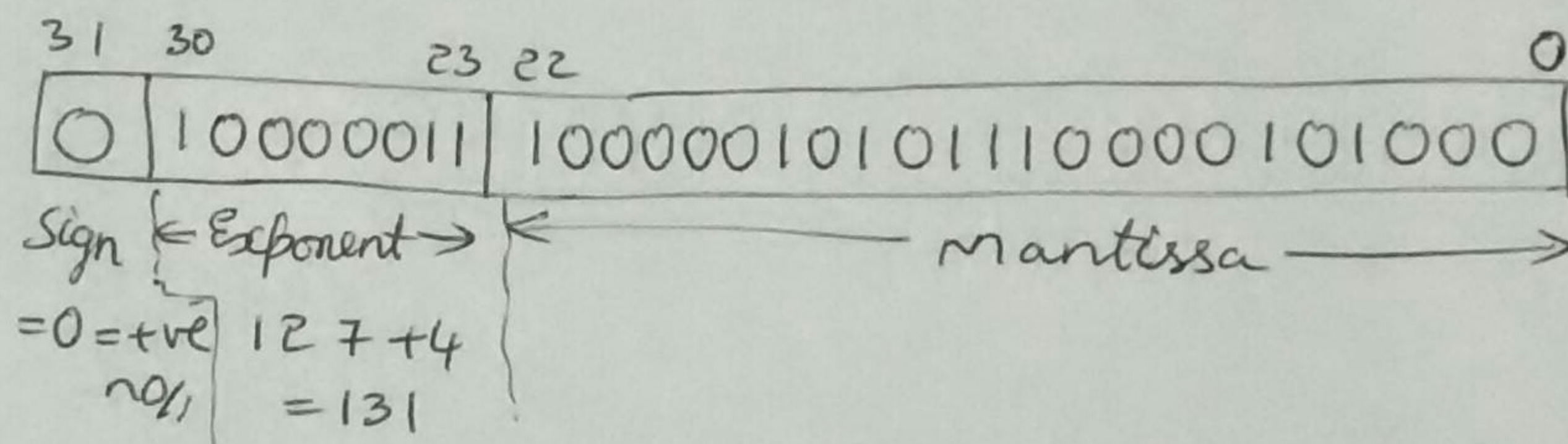
0 → +ve number

Ex: float a;

a = 24.17;

+ 24.17 =

For -ve no.,
make S=1



24.17

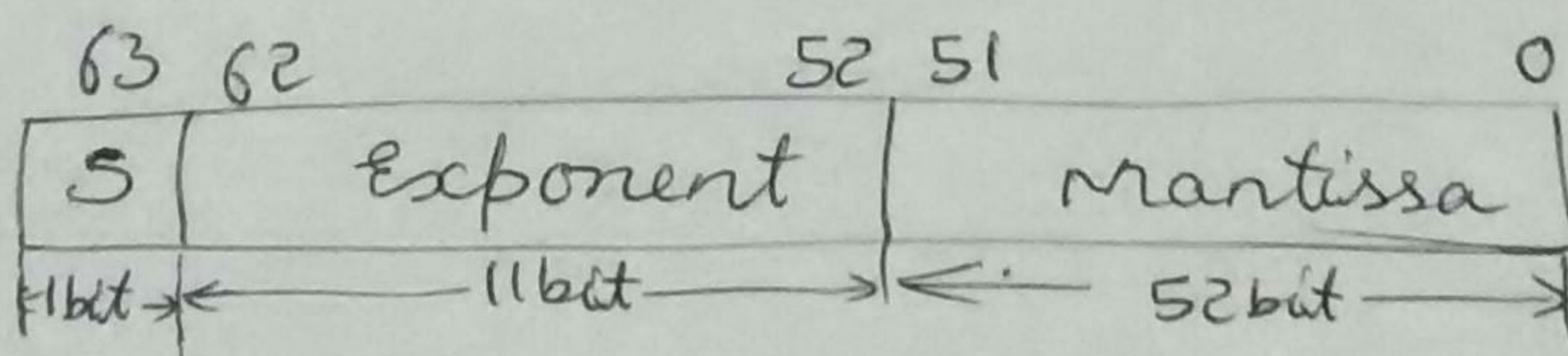
Excess = 127

11000.0010101110000101000
 1.10000010101110000101000, $\times 2^4$

Mantissa

$\Rightarrow \left(\begin{matrix} 127 \\ +4 \end{matrix} \right) = \text{Exponent} = 131$

Double 64 bit Format (Double Precision)



S=0 → +ve no.,

S=1 → -ve no.,

Exponent = (1023 + power of 2)

Excess = 1023

Mantissa is the fraction part