

Name	p	k	bias	Standard
float8	4	4	7	No
float16	11	5	15	IEEE754
bfloat16	8	8	127	No
float32	24	8	127	IEEE754
float64	53	11	1023	IEEE754
float128	113	15	16383	IEEE754
float256	237	19	262143	IEEE754

Table 1: IEEE 754 specification for bits allocation.  
*p*: bits of fraction part (included implicit one),  
*k*: bits of exponent part,  
*bias*: exponent bias constant.

15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
S	Exponent								Fraction						
$+-$	$2^7$	$2^6$	$2^5$	$2^4$	$2^3$	$2^2$	$2^1$	$2^0$	$2^{-1}$	$2^{-2}$	$2^{-3}$	$2^{-4}$	$2^{-5}$	$2^{-6}$	$2^{-7}$
	128	64	32	16	8	4	2	1	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{32}$	$\frac{1}{64}$	$\frac{1}{128}$

Table 2: bfloat16 bits subdivision, powers and actual factors.

7	6	5	4	3	2	1	0
S	Exponent				Fraction		
$+-$	$2^3$	$2^2$	$2^1$	$2^0$	$2^{-1}$	$2^{-2}$	$2^{-3}$
	8	4	2	1	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$

Table 3: Fictional Float8 bits subdivision, powers and actual factors.

$$\begin{aligned}
sign &= bit_{(p+k-1)} \\
S &= (-1)^{sign} \\
E &= \sum_{n=0}^{k-1} bit_{(p-1+n)} \cdot 2^n \\
F &= \sum_{n=1}^{p-1} bit_{(p-1-n)} \cdot 2^{-n} \\
value &= S \cdot 2^{E-bias} \cdot (1 + F)
\end{aligned}$$

Equation 1: General formula.

$$\begin{aligned}
sign &= bit_{31} \\
S &= (-1)^{sign} \\
E &= \sum_{n=0}^7 bit_{(23+n)} \cdot 2^n \\
F &= \sum_{n=1}^{23} bit_{(23-n)} \cdot 2^{-n} \\
value &= S \cdot 2^{E-127} \cdot (1 + F)
\end{aligned}$$

Equation 2: `float32` formula.

$$\begin{aligned}
S &= (-1)^{bit_0} \\
E &= bit_1 \cdot 2^3 + bit_2 \cdot 2^2 + bit_3 \cdot 2^1 + bit_4 \cdot 2^0 \\
F &= bit_5 \cdot 2^{-1} + bit_6 \cdot 2^{-2} + bit_7 \cdot 2^{-3} \\
V &= S \cdot 2^{E-7} \cdot (1 + F)
\end{aligned}$$

Equation 3: Expanded formula for an imaginary 8-bit float type ( $p = 4, k = 4$ ).

$$\begin{aligned}
sign &= bit_{31} \\
S &= (-1)^{sign} \\
number &= \sum_{n=1}^{p+k} bit_{n-1} \cdot 2^{n-p} \\
E &= \lfloor number \rfloor \\
F &= number \bmod 1 \\
value &= S \cdot 2^{E-bias} \cdot (1 + F)
\end{aligned}$$