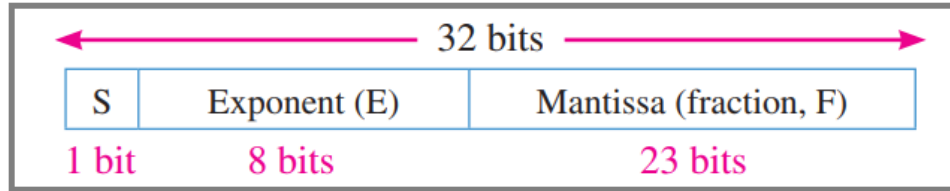


Single-Precision Floating-Point Binary Numbers

In the standard format for a single-precision binary number, the **sign bit (S)** is the left-most bit, the **exponent (E)** includes the next eight bits, and the mantissa or **fractional part (F)** includes the remaining 23 bits, as shown next.



In the mantissa or fractional part, the binary point is understood to be to the left of the 23 bits. Effectively, there are 24 bits in the mantissa because in any binary number the left-most (most significant) bit is always a 1. Therefore, this 1 is understood to be there although it does not occupy an actual bit position.

The eight bits in the exponent represent a biased exponent, which is obtained by adding 127 to the actual exponent. The purpose of the bias is to allow very large or very small numbers without requiring a separate sign bit for the exponents. The biased exponent allows a range of actual exponent values from -126 to +128.

To illustrate, let's consider the following floating-point binary number:

S	E	F
1	10010001	10001110001000000000000

The sign bit is 1. The biased exponent is $10010001 = 145$. Applying the formula, we get

$$\begin{aligned}\text{Number} &= (-1)^1 (1.10001110001)(2^{145-127}) \\ &= (-1)(1.10001110001)(2^{18}) = -11000111000100000000\end{aligned}$$

This floating-point binary number is equivalent to -407,688 in decimal. Since the exponent can be any number between -126 and +128, extremely large and small numbers can be expressed. A 32-bit floating-point number can replace a binary integer number having 129 bits. Because the exponent determines the position of the binary point, numbers containing both integer and fractional parts can be represented.

There are two exceptions to the format for floating-point numbers: The number 0.0 is represented by all 0s, and infinity is represented by all 1s in the exponent and all 0s in the mantissa.

EXAMPLE 2-18

Convert the decimal number 3.248×10^4 to a single-precision floating-point binary number.

Solution

Convert the decimal number to binary.

$$3.248 \times 10^4 = 32480 = 11111011100000_2 = 1.1111011100000 \times 2^{14}$$

The MSB will not occupy a bit position because it is always a 1. Therefore, the mantissa is the fractional 23-bit binary number 11110111000000000000000 and the biased exponent is

$$14 + 127 = 141 = 10001101_2$$

The complete floating-point number is

0	10001101	11110111000000000000000
---	----------	-------------------------