

bits. As a result, a single-precision number has a maximum value of approximately 3.40×10^{38} , with a precision of about 6 decimal digits.

The IEEE standard also describes two other formats, single extended precision and double extended precision. The standard doesn't specify the number of bits in these formats, although it requires that the single extended type occupy at least 43 bits and the double extended type at least 79 bits. For more information about the IEEE standard and floating-point arithmetic in general, see "What every computer scientist should know about floating-point arithmetic" by David Goldberg (*ACM Computing Surveys*, vol. 23, no. 1 (March 1991): 5–48).

subnormal numbers ➤ 23.4

Table 7.4 shows the characteristics of the floating types when implemented according to the IEEE standard. (The table shows the smallest positive *normalized* values. Subnormal numbers can be smaller.) The long double type isn't shown in the table, since its length varies from one machine to another, with 80 bits and 128 bits being the most common sizes.

Table 7.4
Floating Type
Characteristics
(IEEE Standard)

Type	Smallest Positive Value	Largest Value	Precision
float	1.17549×10^{-38}	3.40282×10^{38}	6 digits
double	2.22507×10^{-308}	1.79769×10^{308}	15 digits

On computers that don't follow the IEEE standard, Table 7.4 won't be valid. In fact, on some machines, float may have the same set of values as double, or double may have the same values as long double. Macros that define the characteristics of the floating types can be found in the <float.h> header.

<float.h> header ➤ 23.1

C99

In C99, the floating types are divided into two categories. The float, double, and long double types fall into one category, called the *real floating types*. Floating types also include the *complex types* (float _Complex, double _Complex, and long double _Complex), which are new in C99.

complex types ➤ 27.3

Floating Constants

Floating constants can be written in a variety of ways. The following constants, for example, are all valid ways of writing the number 57.0:

57.0 57. 57.0e0 57E0 5.7e1 5.7e+1 .57e2 570.e-1

A floating constant must contain a decimal point and/or an exponent; the exponent indicates the power of 10 by which the number is to be scaled. If an exponent is present, it must be preceded by the letter E (or e). An optional + or - sign may appear after the E (or e).

By default, floating constants are stored as double-precision numbers. In other words, when a C compiler finds the constant 57.0 in a program, it arranges for the number to be stored in memory in the same format as a double variable. This rule generally causes no problems, since double values are converted automatically to float when necessary.

Q&A