

Summary
10/08/2022 13:35:26

Differences exist between documents.

New Document:

[Pages from n3047_C2X-draft](#)

5 pages (309 KB)

10/08/2022 13:35:18

Used to display results.

Old Document:

[Pages from n2478_C3-Feb2020-draft](#)

4 pages (221 KB)

10/08/2022 13:35:18

[Get started: first change is on page 1.](#)

No pages were deleted

How to read this report

Highlight indicates a change.

Deleted indicates deleted content.

 indicates pages were changed.

 indicates pages were moved.

- 29 **EXAMPLE 2** The following describes floating-point representations that also meet the requirements for single-precision and double-precision numbers in IEC 60559,²⁹⁾ and the appropriate values in a `<float.h>` header for types `float` and `double`:

$$x_f = s2^e \sum_{k=1}^{24} f_k 2^{-k}, \quad -125 \leq e \leq +128$$

$$x_d = s2^e \sum_{k=1}^{53} f_k 2^{-k}, \quad -1021 \leq e \leq +1024$$

<code>FLT_IS_IEC_60559</code>	2		
<code>FLT_RADIX</code>	2		
<code>FLT_MANT_DIG</code>	24		
<code>FLT_EPSILON</code>	1.19209290E-07F	// decimal constant	
<code>FLT_EPSILON</code>	0X1P-23F	// hex constant	
<code>FLT_DECIMAL_DIG</code>	9		
<code>FLT_DIG</code>	6		
<code>FLT_MIN_EXP</code>	-125		
<code>FLT_MIN</code>	1.17549435E-38F	// decimal constant	
<code>FLT_MIN</code>	0X1P-126F	// hex constant	
<code>FLT_TRUE_MIN</code>	1.40129846E-45F	// decimal constant	
<code>FLT_TRUE_MIN</code>	0X1P-149F	// hex constant	
<code>FLT_HAS_SUBNORM</code>	1		
<code>FLT_MIN_10_EXP</code>	-37		
<code>FLT_MAX_EXP</code>	+128		
<code>FLT_MAX</code>	3.40282347E+38F	// decimal constant	
<code>FLT_MAX</code>	0X1.fffffeP127F	// hex constant	
<code>FLT_MAX_10_EXP</code>	+38		
<code>DBL_MANT_DIG</code>	53		
<code>DBL_IS_IEC_60559</code>	2		
<code>DBL_EPSILON</code>	2.2204460492503131E-16	// decimal constant	
<code>DBL_EPSILON</code>	0X1P-52	// hex constant	
<code>DBL_DECIMAL_DIG</code>	17		
<code>DBL_DIG</code>	15		
<code>DBL_MIN_EXP</code>	-1021		
<code>DBL_MIN</code>	2.2250738585072014E-308	// decimal constant	
<code>DBL_MIN</code>	0X1P-1022	// hex constant	
<code>DBL_TRUE_MIN</code>	4.9406564584124654E-324	// decimal constant	
<code>DBL_TRUE_MIN</code>	0X1P-1074	// hex constant	
<code>DBL_HAS_SUBNORM</code>	1		
<code>DBL_MIN_10_EXP</code>	-307		
<code>DBL_MAX_EXP</code>	+1024		
<code>DBL_MAX</code>	1.7976931348623157E+308	// decimal constant	
<code>DBL_MAX</code>	0X1.fffffffffffffP1023	// hex constant	
<code>DBL_MAX_10_EXP</code>	+308		

Forward references: conditional inclusion (6.10.1), predefined macro names (6.10.9), complex arithmetic `<complex.h>` (7.3), extended multibyte and wide character utilities `<wchar.h>` (7.31), floating-point environment `<fenv.h>` (7.6), general utilities `<stdlib.h>` (7.24), input/output `<stdio.h>` (7.23), mathematics `<math.h>` (7.12), IEC 60559 floating-point arithmetic (Annex F), IEC 60559-compatible complex arithmetic (Annex G).

5.2.4.2.3 Characteristics of decimal floating types in `<float.h>`

- 1 This subclause specifies macros in `<float.h>` that provide characteristics of decimal floating types (an optional feature) in terms of the model presented in 5.2.4.2.2. An implementation that does not support decimal floating types shall not provide these macros. The prefixes `DEC32_`, `DEC64_`, and `DEC128_` denote the types `_Decimal32`, `_Decimal64`, and `_Decimal128` respectively.
- 2 `DEC_EVAL_METHOD` is the decimal floating-point analog of `FLT_EVAL_METHOD` (5.2.4.2.2). Its implementation-defined value characterizes the use of evaluation formats for decimal floating

²⁹⁾The floating-point model in that standard sums powers of b from zero, so the values of the exponent limits are one less than shown here.

types:

- 1 indeterminable;
 - 0 evaluate all operations and constants just to the range and precision of the type;
 - 1 evaluate operations and constants of type `_Decimal32` and `_Decimal64` to the range and precision of the `_Decimal64` type, evaluate `_Decimal128` operations and constants to the range and precision of the `_Decimal128` type;
 - 2 evaluate all operations and constants to the range and precision of the `_Decimal128` type.
- 3 Each of the decimal signaling NaN macros

```
DEC32_SNAN
DEC64_SNAN
DEC128_SNAN
```

expands to a constant expression of the respective decimal floating type representing a signaling NaN. If an optional unary + or - operator followed by a signaling NaN macro is used for initializing an object of the same type that has static or thread storage duration, the object is initialized with a signaling NaN value.

- 4 The macro

```
DEC_INFINITY
```

expands to a constant expression of type `_Decimal32` representing positive infinity.

- 5 The macro

```
DEC_NAN
```

expands to a constant expression of type `_Decimal32` representing a quiet NaN.

- 6 The integer values given in the following lists shall be replaced by constant expressions suitable for use in `#if` preprocessing directives:

- radix of exponent representation, $b(=10)$

For the standard floating types, this value is implementation-defined and is specified by the macro `FLT_RADIX`. For the decimal floating types there is no corresponding macro, since the value 10 is an inherent property of the types. Wherever `FLT_RADIX` appears in a description of a function that has versions that operate on decimal floating types, it is noted that for the decimal floating-point versions the value used is implicitly 10, rather than `FLT_RADIX`.

- number of digits in the coefficient

```
DEC32_MANT_DIG    7
DEC64_MANT_DIG    16
DEC128_MANT_DIG   34
```

- minimum exponent

```
DEC32_MIN_EXP     -94
DEC64_MIN_EXP     -382
DEC128_MIN_EXP    -6142
```

- maximum exponent

decimal floating-point arithmetic, representations that have the same numerical value but different quantum exponents, e.g., $(+1, 10, -1)$ representing 1.0 and $(+1, 100, -2)$ representing 1.00, are distinguishable. To facilitate exact fixed-point calculation, operation results that are of decimal floating type have a *preferred quantum exponent*, as specified in IEC 60559, which is determined by the quantum exponents of the operands if they have decimal floating types (or by specific rules for conversions from other types). The table below gives rules for determining preferred quantum exponents for results of IEC 60559 operations, and for other operations specified in this document. When exact, these operations produce a result with their preferred quantum exponent, or as close to it as possible within the limitations of the type. When inexact, these operations produce a result with the least possible quantum exponent. For example, the preferred quantum exponent for addition is the minimum of the quantum exponents of the operands. Hence $(+1, 123, -2) + (+1, 4000, -3) = (+1, 5230, -3)$ or $1.23 + 4.000 = 5.230$.

- 10 The following table shows, for each operation delivering a result in decimal floating-point format, how the preferred quantum exponents of the operands, $Q(x)$, $Q(y)$, etc., determine the preferred quantum exponent of the operation result, provided the table formula is defined for the arguments. For the cases where the formula is undefined and the function result is $\pm\infty$, the preferred quantum exponent is immaterial because the quantum exponent of $\pm\infty$ is defined to be infinity. For the other cases where the formula is undefined and the function result is finite, the preferred quantum exponent is unspecified³⁰⁾.

Preferred quantum exponents

Operation	Preferred quantum exponent of result
roundeven, round, trunc, ceil, floor, rint, nearbyint	$\max(Q(x), 0)$
nextup, nextdown, nextafter, nexttoward	least possible
remainder	$\min(Q(x), Q(y))$
fmin, fmax, fminimum, fmaximum, fminimum_mag, fmaximum_mag, fminimum_num, fmaximum_num, fminimum_mag_num, fmaximum_mag_num	$Q(x)$ if x gives the result, $Q(y)$ if y gives the result
scalbn, scalbln	$Q(x) + n$
ldexp	$Q(x) + p$
logb	0
+, d32add, d64add	$\min(Q(x), Q(y))$
-, d32sub, d64sub	$\min(Q(x), Q(y))$
*, d32mul, d64mul	$Q(x) + Q(y)$
/, d32div, d64div	$Q(x) - Q(y)$
sqrt, d32sqrt, d64sqrt	$\lfloor Q(x)/2 \rfloor$
fma, d32fma, d64fma	$\min(Q(x) + Q(y), Q(z))$
conversion from integer type	0
exact conversion from non-decimal floating type	0
inexact conversion from non-decimal floating type	least possible
conversion between decimal floating types	$Q(x)$
*cx returned by canonicalize	$Q(*x)$
strt0, wcst0, scanf , floating constants of decimal floating type	see 7.24.1.6
-(x), +(x)	$Q(x)$
fabs	$Q(x)$
copysign	$Q(x)$
quantize	$Q(y)$

³⁰⁾Although unspecified in IEC 60559, a preferred quantum exponent of 0 for these cases would be a reasonable implementation choice.

quantum	$Q(x)$
* encptr returned by encodedec , encodebin	$Q(*xptr)$
* xptr returned by decodedec , decodebin	$Q(*encptr)$
fmod	$\min(Q(x), Q(y))$
fdim	$\min(Q(x), Q(y))$ if $x > y$, 0 if $x \leq y$
cbrt	$\lfloor Q(x)/3 \rfloor$
hypot	$\min(Q(x), Q(y))$ †
pow	$\lfloor y \times Q(x) \rfloor$
modf	$Q(\text{value})$
* iptr returned by modf	$\max(Q(\text{value}), 0)$
frexp	$Q(\text{value})$ if $\text{value} = 0$, $-(\text{length of coefficient of value})$ otherwise
* res returned by setpayload , setpayloadsig	0 if pl does not represent a valid payload, not applicable otherwise (NaN returned)
getpayload	0 if * x is a NaN, unspecified otherwise
compoundn	$\lfloor n \times \min(0, Q(x)) \rfloor$
pown	$\lfloor n \times Q(x) \rfloor$
powr	$\lfloor y \times Q(x) \rfloor$
rootn	$\lfloor Q(x)/n \rfloor$
rsqrt	$-\lfloor Q(x)/2 \rfloor$
transcendental functions	0

A function family listed in the table above indicates the functions for all decimal floating types, where the function family is represented by the name of the functions without a suffix. For example, **ceil** indicates the functions **ceild32**, **ceild64**, and **ceild128**.

Forward references: extended multibyte and wide character utilities `<wchar.h>` (7.31), floating-point environment `<fenv.h>` (7.6), general utilities `<stdlib.h>` (7.24), input/output `<stdio.h>` (7.23), mathematics `<math.h>` (7.12), † type-generic mathematics `<tgmath.h>` (7.27), IEC 60559 floating-point arithmetic (Annex F).