

## Annex X

(normative)

# IEC 60559 interchange and extended types

5

### X.1 Introduction

[1] This annex extends programming language C to include types with the arithmetic interchange and extended floating-point formats specified in ISO/IEC/IEEE 60559:2011, and to include functions that support the non-arithmetic interchange formats in that standard. This annex was adapted from ISO/IEC TS 18661-3:2015, *Floating-point extensions for C — Interchange and extended types*.

[2] An implementation that defines `__STDC_IEC_60559_TYPES__` to `20yyymmL` shall conform to the specifications in this annex. An implementation may define `__STDC_IEC_60559_TYPES__` only if it defines `__STDC_IEC_60559_BFP__`, indicating support for IEC 60559 binary floating-point arithmetic, or defines `__STDC_IEC_60559_DFP__`, indicating support for IEC 60559 decimal floating-point arithmetic (or defines both). Where a binding between the C language and IEC 60559 is indicated, the IEC 60559-specified behavior is adopted by reference, unless stated otherwise.

=====

#### Change to C2X working draft all-20190708:

In 6.10.8.3#1, add:

`__STDC_IEC_60559_TYPES__` The integer constant `20yyymmL`, intended to indicate conformance to the specification in Annex X (IEC 60559 interchange and extended types).

=====

### X.2 Types

[1] This clause specifies types that support IEC 60559 arithmetic interchange and extended formats. The encoding conversion functions (X.11.3) and numeric conversion functions for encodings (X.12.3, X.12.4) support the non-arithmetic interchange formats specified in IEC 60559.

#### X.2.1 Interchange floating types

[1] IEC 60559 specifies interchange formats, identified by their width, which can be used for the exchange of floating-point data between implementations. The two tables below give parameters for the IEC 60559 interchange formats.

### Binary interchange format parameters

Parameter	binary16	binary32	binary64	binary128	binary $N$ ( $N \geq 128$ )
$N$ , storage width in bits	16	32	64	128	multiple of 32
$p$ , precision in bits	11	24	53	113	$N - \text{round}(4 \times \log_2(N)) + 13$
$emax$ , maximum exponent $e$	15	127	1023	16383	$2^{(N-p-1)} - 1$
<i>Encoding parameters</i>					
$bias, E-e$	15	127	1023	16383	$emax$
sign bit	1	1	1	1	1
$w$ , exponent field width in bits	5	8	11	15	$\text{round}(4 \times \log_2(N)) - 13$
$t$ , trailing significand field width in bits	10	23	52	112	$N - w - 1$
$N$ , storage width in bits	16	32	64	128	$1 + w + t$

The function  $\text{round}()$  in the table above rounds to the nearest integer. For example, binary256 would have  $p = 237$  and  $emax = 262143$ .

5

### Decimal interchange format parameters

Parameter	decimal32	decimal64	decimal128	decimal $N$ ( $N \geq 32$ )
$N$ , storage width in bits	32	64	128	multiple of 32
$p$ , precision in digits	7	16	34	$9 \times N/32 - 2$
$emax$ , maximum exponent $e$	96	384	6144	$3 \times 2^{(N/16 + 3)}$
<i>Encoding parameters</i>				
$bias, E-e$	101	398	6176	$emax + p - 2$
sign bit	1	1	1	1
$W+5$ , combination field width in bits	11	13	17	$N/16 + 9$
$t$ , trailing significand field width in bits	20	50	110	$15 \times N/16 - 10$
$N$ , storage width in bits	32	64	128	$1 + 5 + w + t$

For example, decimal256 would have  $p = 70$  and  $emax = 1572864$ .

[2] Types designated

Float $N$ , where  $N$  is 16, 32, 64, or  $\geq 128$  and a multiple of 32

10 and types designated

Decimal $N$ , where  $N \geq 32$  and a multiple of 32

are collectively called the *interchange floating types*. Each interchange floating type has the IEC 60559 interchange format corresponding to its width ( $N$ ) and radix (2 for `_Float $N$` , 10 for `_Decimal $N$` ). Each interchange floating type is not compatible with any other type.

[3] An implementation that defines `__STDC_IEC_60559_BFP__` and `__STDC_IEC_60559_TYPES__` shall provide `_Float32` and `_Float64` as interchange floating types with the same representation and alignment requirements as `float` and `double`, respectively. If the implementation's `long double` type supports an IEC 60559 interchange format of width  $N > 64$ , then the implementation shall also provide the type `_Float $N$`  as an interchange floating type with the same representation and alignment requirements as `long double`. The implementation may provide other binary interchange floating types; the set of such types supported is implementation-defined.

[4] An implementation that defines `__STDC_IEC_60559_DFP__` shall provide the types `_Decimal32`, `_Decimal64`, and `_Decimal128`. If the implementation also defines `__STDC_IEC_60559_TYPES__`, it may provide other decimal interchange floating types; the set of such types supported is implementation-defined.

## X.2.2 Non-arithmetic interchange formats

[1] Providing an interchange floating type entails supporting it as an IEC 60559 arithmetic format. An implementation supports IEC 60559 non-arithmetic interchange formats by providing the associated encoding-to-encoding conversion functions (X.11.3.2), string-from-encoding functions (X.12.3), and string-to-encoding functions (X.12.4).

[2] An implementation that defines `__STDC_IEC_60559_BFP__` and `__STDC_IEC_60559_TYPES__` shall support the IEC 60559 binary16 format, at least as a non-arithmetic interchange format. Otherwise, the set of non-arithmetic interchange formats supported is implementation-defined.

## X.2.3 Extended floating types

[1] For each of its basic formats, IEC 60559 specifies an extended format whose maximum exponent and precision exceed those of the basic format it is associated with. The table below gives the minimum values of these parameters:

**Extended format parameters for floating-point numbers**

	Extended formats associated with:				
Parameter	binary32	binary64	binary128	decimal64	decimal128
$p \text{ digits} \geq$	32	64	128	22	40
$emax \geq$	1023	16383	65535	6144	24576

[2] Types designated `_Float32x`, `_Float64x`, `_Float128x`, `_Decimal64x`, and `_Decimal128x` support the corresponding IEC 60559 extended formats and are collectively called the *extended floating types*. Each extended floating type is not compatible with any other type. An implementation that defines `__STDC_IEC_60559_BFP__` and `__STDC_IEC_60559_TYPES__` shall provide `_Float32x`, which may have the same set of values as `double`, and may provide any of the other two binary extended floating types. An implementation that defines `__STDC_IEC_60559_DFP__` and `__STDC_IEC_60559_TYPES__` shall provide: `_Decimal64x`, which may have the same set of values as `_Decimal128`, and may provide `_Decimal128x`. Which (if any) of the optional extended floating types are provided is implementation-defined.

NOTE IEC 60559 does not specify an extended format associated with the decimal32 format, nor does this annex specify an extended type associated with the `_Decimal32` type.

## X.2.4 Classification of real floating types

[1] 6.2.5 defines standard floating types as a collective name for the types **float**, **double**, and **long double** and it defines decimal floating types as a collective name for the types **\_Decimal32**, **\_Decimal64**, and **\_Decimal128**. This subclause broadens the definition of decimal floating types and defines *binary floating types* to be collective names for types for all the appropriate IEC 60559 arithmetic formats. Note that standard floating types (which have an implementation-defined radix) are not included in either decimal floating types (which all have radix 10) or binary floating types (which all have radix 2).

[2] The *real floating types* are broadened to include all interchange floating types and extended floating types, as well as standard floating types.

[3] The interchange floating types designated **\_FloatN** and the extended floating types designated **\_FloatNx** are collectively called the *binary floating types*. The interchange floating types designated **\_DecimalN** and the extended floating types designated **\_DecimalNx** are collectively called the *decimal floating types*. Note that binary floating types and the decimal floating types are real floating types.

[4] Thus real floating types are classified as follows:

standard floating types:

**float**  
**double**  
**long double**

binary floating types:

**\_FloatN**  
**\_FloatNx**

decimal floating types:

**\_DecimalN**  
**\_DecimalNx**

## X.2.5 Complex types

[1] This subclause broadens the C complex types (6.2.5) to also include similar types whose corresponding real parts are binary floating types. For the interchange floating types **\_FloatN**, and the extended floating types **\_FloatNx**, there are complex types designated respectively as **\_FloatN\_Complex** and **\_FloatNx\_Complex**. (Complex types are a conditional feature that implementations need not support; see 6.10.8.3.)

## X.2.6 Imaginary types

[1] This subclause broadens the C imaginary types (G.2) to also include similar types whose corresponding real parts are binary floating types. For the interchange floating types **\_FloatN** and the extended floating types **\_FloatNx**, there are imaginary types designated respectively as **\_FloatN\_Imaginary** and **\_FloatNx\_Imaginary**. The imaginary types (along with the real floating and complex types) are floating types. (Annex G, including imaginary types, is a conditional feature that implementations need not support; see 6.10.8.3.)

## X.3 Characteristics in <float.h>

[1] This subclause specifies new <float.h> macros, analogous to the macros for standard floating types, that characterize this annex's binary floating types in terms of the model presented in 5.2.4.2.2.

It generalizes the specification of characteristics in 5.2.4.2.3 of the three decimal floating types to include this annex's decimal floating types. The prefix **FLT $N$ \_** indicates a binary interchange floating type or a non-arithmetic binary interchange format of width  $N$ . The prefix **FLT $N$ X\_** indicates a binary extended floating type that extends a basic format of width  $N$ . The prefix **DEC $N$ \_** indicates a decimal interchange floating type of width  $N$ . The prefix **DEC $N$ X\_** indicates a decimal extended floating type that extends a basic format of width  $N$ . The type parameters  $p$ ,  $e_{max}$ , and  $e_{min}$  for extended floating types are for the extended floating type itself, not for the basic format that it extends. For each interchange or extended floating type that the implementation provides, **<float.h>** shall define the associated macros in the following lists. Conversely, for each such type that the implementation does not provide, **<float.h>** shall not define the associated macros in the following list, except, the implementation shall define the macros **FLT $N$ \_DECIMAL\_DIG** and **FLT $N$ \_DIG** if it supports IEC 60559 non-arithmetic binary interchange formats of width  $N$  by providing the encoding-to-encoding conversion functions (X.11.3.2) in **<math.h>** and the string-from-encoding (X.12.3) and string-to-encoding (X.12.4) functions in **<stdlib.h>**.

[2] If **FLT\_RADIX** is 2, the value of the macro **FLT\_EVAL\_METHOD** (5.2.4.2.2) characterizes the use of evaluation formats for standard floating types and for binary interchange and extended floating types:

- 1 indeterminate;
  - 0 evaluate all operations and constants, whose semantic type has at most the range and precision of **float**, to the range and precision of **float**; evaluate all other operations and constants to the range and precision of the semantic type;
  - 1 evaluate operations and constants, whose semantic type has at most the range and precision of **double**, to the range and precision of **double**; evaluate all other operations and constants to the range and precision of the semantic type;
  - 2 evaluate operations and constants, whose semantic type has at most the range and precision of **long double**, to the range and precision of **long double**; evaluate all other operations and constants to the range and precision of the semantic type;
- $N$ , where **\_Float $N$**  is a supported interchange floating type  
evaluate operations and constants, whose semantic type has at most the range and precision of the **\_Float $N$**  type, to the range and precision of the **\_Float $N$**  type;  
evaluate all other operations and constants to the range and precision of the semantic type;
- $N + 1$ , where **\_Float $N$ x** is a supported extended floating type  
evaluate operations and constants, whose semantic type has at most the range and precision of the **\_Float $N$ x** type, to the range and precision of the **\_Float $N$ x** type;  
evaluate all other operations and constants to the range and precision of the semantic type.

If **FLT\_RADIX** is not 2, the use of evaluation formats for operations and constants of binary interchange and extended floating types is implementation-defined.

[3] The implementation-defined value of the macro **DEC\_EVAL\_METHOD** (5.2.4.2.3) characterizes the use of evaluation formats for decimal interchange and extended floating types:

- 1 indeterminate;
- 0 evaluate all operations and constants just to the range and precision of the type;
- 1 evaluate operations and constants, whose semantic type has at most the range and precision of the **\_Decimal64** type, to the range and precision of the **\_Decimal64**

type; evaluate all other operations and constants to the range and precision of the semantic type;

- 2 evaluate operations and constants, whose semantic type has at most the range and precision of the `_Decimal128` type, to the range and precision of the `_Decimal128` type; evaluate all other operations and constants to the range and precision of the semantic type;

$N$ , where `_Decimal $N$`  is a supported interchange floating type evaluate operations and constants, whose semantic type has at most the range and precision of the `_Decimal $N$`  type, to the range and precision of the `_Decimal $N$`  type; evaluate all other operations and constants to the range and precision of the semantic type;

$N + 1$ , where `_Decimal $N$ x` is a supported extended floating type evaluate operations and constants, whose semantic type has at most the range and precision of the `_Decimal $N$ x` type, to the range and precision of the `_Decimal $N$ x` type; evaluate all other operations and constants to the range and precision of the semantic type;

[4] 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,  $b$  (= 2 for binary, 10 for decimal)

For the standard floating types, this value is implementation-defined and is specified by the macro `FLT_RADIX`. For the interchange and extended floating types there is no corresponding macro; the radix is an inherent property of the types.

— number of bits in the floating-point significand,  $p$

`FLT $N$ _MANT_DIG`  
`FLT $N$ X_MANT_DIG`

— number of digits in the coefficient,  $p$

`DEC $N$ _MANT_DIG`  
`DEC $N$ X_MANT_DIG`

— number of decimal digits,  $n$ , such that any floating-point number with  $p$  bits can be rounded to a floating-point number with  $n$  decimal digits and back again without change to the value,  $\lceil 1 + p \log_{10} 2 \rceil$

`FLT $N$ _DECIMAL_DIG`  
`FLT $N$ X_DECIMAL_DIG`

— number of decimal digits,  $q$ , such that any floating-point number with  $q$  decimal digits can be rounded into a floating-point number with  $p$  bits and back again without change to the  $q$  decimal digits,  $\lfloor (p - 1) \log_{10} 2 \rfloor$

`FLT $N$ _DIG`  
`FLT $N$ X_DIG`

— minimum negative integer such that the radix raised to one less than that power is a normalized floating-point number,  $e_{min}$

5  
**FLT<sub>N</sub>\_MIN\_EXP**  
**FLT<sub>X</sub>\_MIN\_EXP**  
**DEC<sub>N</sub>\_MIN\_EXP**  
**DEC<sub>X</sub>\_MIN\_EXP**

— minimum negative integer such that 10 raised to that power is in the range of normalized floating-point numbers,  $\lceil \log_{10} 2^{e_{min}-1} \rceil$

10  
**FLT<sub>N</sub>\_MIN\_10\_EXP**  
**FLT<sub>X</sub>\_MIN\_10\_EXP**

— maximum integer such that the radix raised to one less than that power is a representable finite floating-point number,  $e_{max}$

15  
**FLT<sub>N</sub>\_MAX\_EXP**  
**FLT<sub>X</sub>\_MAX\_EXP**  
**DEC<sub>N</sub>\_MAX\_EXP**  
**DEC<sub>X</sub>\_MAX\_EXP**

20  
— maximum integer such that 10 raised to that power is in the range of representable finite floating-point numbers,  $\lfloor \log_{10}((1 - 2^{-p})2^{e_{max}}) \rfloor$

**FLT<sub>N</sub>\_MAX\_10\_EXP**  
**FLT<sub>X</sub>\_MAX\_10\_EXP**

25  
— maximum representable finite floating-point number,  $(1 - b^{-p})b^{e_{max}}$

**FLT<sub>N</sub>\_MAX**  
**FLT<sub>X</sub>\_MAX**  
**DEC<sub>N</sub>\_MAX**  
**DEC<sub>X</sub>\_MAX**

30  
— the difference between 1 and the least value greater than 1 that is representable in the given floating-point type,  $b^{1-p}$

35  
**FLT<sub>N</sub>\_EPSILON**  
**FLT<sub>X</sub>\_EPSILON**  
**DEC<sub>N</sub>\_EPSILON**  
**DEC<sub>X</sub>\_EPSILON**

— minimum normalized positive floating-point number,  $b^{e_{min}-1}$

40  
**FLT<sub>N</sub>\_MIN**  
**FLT<sub>X</sub>\_MIN**  
**DEC<sub>N</sub>\_MIN**  
**DEC<sub>X</sub>\_MIN**

— minimum positive subnormal floating-point number,  $b^{e_{min}-p}$

45  
**FLT<sub>N</sub>\_TRUE\_MIN**  
**FLT<sub>X</sub>\_TRUE\_MIN**  
**DEC<sub>N</sub>\_TRUE\_MIN**  
**DEC<sub>X</sub>\_TRUE\_MIN**

## X.4 Conversions

[1] This subclause enhances the usual arithmetic conversions (6.3.1.8) to handle interchange and extended floating types. It supports the IEC 60559 recommendation against allowing implicit conversions of operands to obtain a common type where the conversion is between types where neither is a subset of (or equivalent to) the other.

[2] This subclause also broadens the operation binding in F.3 for the IEC 60559 `convertFormat` operation to apply to IEC 60559 arithmetic and non-arithmetic formats

### X.4.1 Real floating and integer

[1] When a finite value of interchange or extended floating type is converted to an integer type other than `_Bool`, the fractional part is discarded (i.e., the value is truncated toward zero). If the value of the integral part cannot be represented by the integer type, the “invalid” floating-point exception shall be raised and the result of the conversion is unspecified.

[2] When a value of integer type is converted to an interchange or extended floating type, if the value being converted can be represented exactly in the new type, it is unchanged. If the value being converted cannot be represented exactly, the result shall be correctly rounded with exceptions raised as specified in IEC 60559.

### X.4.2 Usual arithmetic conversions

[1] If either operand is of floating type, the common real type is determined as follows:

If one operand has decimal floating type, the other operand shall not have standard floating type, binary floating type, complex type, or imaginary type.

If both operands have floating types and neither of the sets of values of their corresponding real types is a subset of (or equivalent to) the other, the behavior is undefined.

If both operands have the same corresponding real type, no further conversion is needed.

Otherwise, if both operands are floating types and the sets of values of their corresponding real types are equivalent, then the following rules are applied:

If the corresponding real type of either operand is an interchange floating type, the other operand is converted, without change of type domain, to a type whose corresponding real type is that same interchange floating type.

Otherwise, if the corresponding real type of either operand is `long double`, the other operand is converted, without change of type domain, to a type whose corresponding real type is `long double`.

Otherwise, if the corresponding real type of either operand is `double`, the other operand is converted, without change of type domain, to a type whose corresponding real type is `double`.

(All cases where `float` might have the same format as another type are covered above.)

Otherwise, if the corresponding real type of either operand is `_Float128x` or `_Decimal128x`, the other operand is converted, without change of type domain, to a type whose corresponding real type is `_Float128x` or `_Decimal128x`, respectively.



Otherwise, if the corresponding real type of either operand is `_Float64x` or `_Decimal64x`, the other operand is converted, without change of type domain, to a type whose corresponding real type is `_Float64x` or `_Decimal64x`, respectively.

5 Otherwise, if both operands have floating types, the operand, whose set of values of its corresponding real type is a (proper) subset of the set of values of the corresponding real type of the other operand, is converted, without change of type domain, to a type with the corresponding real type of that other operand.

10 Otherwise, if one operand has a floating type, the other operand is converted to the corresponding real type of the operand of floating type.

### X.4.3 Arithmetic and non-arithmetic formats

[1] The operation binding in F.3 for the IEC 60559 `convertFormat` operation applies to IEC 60559 arithmetic and non-arithmetic formats as follows:

- 15 — For conversions between arithmetic formats supported by floating types - casts and implicit conversions.
- For same-radix conversions between non-arithmetic interchange formats - encoding-to-encoding conversion functions (X.11.3.2).
- 20 — For conversions between non-arithmetic interchange formats (same or different radix) – compositions of string-from-encoding functions (X.12.3) (converting exactly) and string-to-encoding functions (X.12.4).
- For same-radix conversions from interchange formats supported by interchange floating types to non-arithmetic interchange formats – compositions of encode functions (X.11.3.1.1, 7.12.16.1, 7.12.16.3) and encoding-to-encoding functions (X.11.3.2).
- 25 — For same radix conversions from non-arithmetic interchange formats to interchange formats supported by interchange floating types – compositions of encoding-to-encoding conversion functions (X.3.2) and decode functions (X.11.3.1.2, 7.12.16.2, 7.12.16.4).
- For conversions from non-arithmetic interchange formats to arithmetic formats supported by floating types (same or different radix) – compositions of string-from-encoding functions (X.12.3) (converting exactly) and numeric conversion functions `strtod`, etc. (7.22.1.5, 7.22.1.6).
- 30 — For conversions from arithmetic formats supported by floating types to non-arithmetic interchange formats (same or different radix) – compositions of numeric conversion functions `strfromd`, etc. (7.22.1.3, 7.22.1.4) (converting exactly) and string-to-encoding functions (X.12.4).

## X.5 Lexical elements

### X.5.1 Keywords

[1] This subclause expands the list of keywords (6.4.1) to also include:

5        \_**Float***N*, where *N* is 16, 32, 64, or  $\geq 128$  and a multiple of 32  
      \_**Float32x**  
      \_**Float64x**  
      \_**Float128x**  
      \_**Decimal***N*, where *N* is 96 or  $> 128$  and a multiple of 32  
10       \_**Decimal64x**  
      \_**Decimal128x**

### X.5.2 Constants

[1] This subclause specifies constants of interchange and extended floating types.

[2] This subclause expands *floating-suffix* (6.4.4.2) to also include:

15        **fN FN fNx FNx dN DN dNx DNx**

[3] A floating suffix **dN**, **DN**, **dNx**, or **DNx** shall not be used in a *hexadecimal-floating-constant*.

[4] A floating suffix shall not designate a type that the implementation does not provide.

[5] If a floating constant is suffixed by **fN** or **FN**, it has type **\_Float***N*. If suffixed by **fNx** or **FNx**, it has type **\_Float***Nx*. If suffixed by **dN** or **DN**, it has type **\_Decimal***N*. If suffixed by **dNx** or **DNx**, it has type **\_Decimal***Nx*.

## X.6 Expressions

[1] This subclause expands the specification of expressions to also cover interchange and extended floating types.

[2] Operators involving operands of interchange or extended floating type are evaluated according to the semantics of IEC 60559, including production of decimal floating-point results with the preferred quantum exponent as specified in IEC 60559 (see 5.2.4.2.3).

[3] For multiplicative operators (6.5.5), additive operators (6.5.6), relational operators (6.5.8), equality operators (6.5.9), and compound assignment operators (6.5.16.2), if either operand has decimal floating type, the other operand shall not have standard floating type, binary floating type, complex type, or imaginary type.

[4] For conditional operators (6.5.15), if either of the second or third operands has decimal floating type, the other operand shall not have standard floating type, binary floating type, complex type, or imaginary type.

[5] The equivalence of expressions noted in F.9.2 apply to expressions of binary floating types, as well as standard floating types.

## X.7 Declarations

[1] This subclause expands the list of type specifiers (6.7.2) to also include:

5        **`_FloatN`**, where  $N$  is 16, 32, 64, or  $\geq 128$  and a multiple of 32  
      **`_Float32x`**  
      **`_Float64x`**  
      **`_Float128x`**  
      **`_DecimalN`**, where  $N$  is 96 or  $> 128$  and a multiple of 32  
      **`_Decimal64x`**  
10       **`_Decimal128x`**

[2] The type specifiers **`_FloatN`** (where  $N$  is 16, 32, 64, or  $\geq 128$  and a multiple of 32), **`_Float32x`**, **`_Float64x`**, **`_Float128x`**, **`_DecimalN`** (where  $N$  is 96 or  $> 128$  and a multiple of 32), **`_Decimal64x`**, and **`_Decimal128x`** shall not be used if the implementation does not support the corresponding types (see 6.10.8.3 and X.2).

15 [3] This subclause also expands the list under Constraints in 6.7.2 to also include:

— **`_FloatN`**, where  $N$  is 16, 32, 64, or  $\geq 128$  and a multiple of 32  
— **`_Float32x`**  
— **`_Float64x`**  
— **`_Float128x`**  
20 — **`_DecimalN`**, where  $N$  is 96 or  $> 128$  and a multiple of 32  
— **`_Decimal64x`**  
— **`_Decimal128x`**  
— **`_FloatN _Complex`**, where  $N$  is 16, 32, 64, or  $\geq 128$  and a multiple of 32  
— **`_Float32x _Complex`**  
25 — **`_Float64x _Complex`**  
— **`_Float128x _Complex`**

## X.8 Identifiers in standard headers

[1] The new identifiers added to C library headers by this annex are defined or declared by their respective headers only if **`__STDC_WANT_IEC_60559_TYPES_EXT__`** is defined as a macro at the point in the source file where the appropriate header is first included. The following subclauses list these identifiers for each applicable standard header.

### X.8.1 `<float.h>`

[1] The following identifiers are defined only if **`__STDC_WANT_IEC_60559_TYPES_EXT__`** is defined as a macro at the point in the source file where `<float.h>` is first included:

for supported types `_FloatN`:

```
    FLTN_MANT_DIG      FLTN_MIN_10_EXP      FLTN_EPSILON
    FLTN_DECIMAL_DIG   FLTN_MAX_EXP        FLTN_MIN
5    FLTN_DIG          FLTN_MAX_10_EXP     FLTN_TRUE_MIN
    FLTN_MIN_EXP       FLTN_MAX
```

for IEC 60559 non-arithmetic binary interchange formats of width *N*:

```
    FLTN_DECIMAL_DIG   FLTN_DIG
```

for supported types `_FloatNx`:

```
10    FLTNx_MANT_DIG      FLTNx_MIN_10_EXP     FLTNx_EPSILON
    FLTNx_DECIMAL_DIG   FLTNx_MAX_EXP        FLTNx_MIN
    FLTNx_DIG          FLTNx_MAX_10_EXP     FLTNx_TRUE_MIN
    FLTNx_MIN_EXP       FLTNx_MAX
```

for supported types `_DecimalN`, where *N* ≠ 32, 64, and 128:

```
15    DECN_MANT_DIG      DECN_MAX             DECN_TRUE_MIN
    DECN_MIN_EXP        DECN_EPSILON
    DECN_MAX_EXP        DECN_MIN
```

for supported types `_DecimalNx`:

```
20    DECNx_MANT_DIG     DECNx_MAX            DECNx_TRUE_MIN
    DECNx_MIN_EXP       DECNx_EPSILON
    DECNx_MAX_EXP       DECNx_MIN
```

## X.8.2 <complex.h>

[1] The following identifiers are declared or defined only if `__STDC_WANT_IEC_60559_TYPES_EXT__` is defined as a macro at the point in the source file where <complex.h> is first included:

25 for supported types `_FloatN`:

```
    cacosfN           catanhfN           csqrtfN
    casinfnN          ccoshfnN           cargfnN
30    catanfN          csinhfnN           cimagfnN
    ccosfnN           ctanhfnN           CMPLXFN
    csinfnN           cexpfnN            conjfnN
    ctanfN           clogfnN            cprojfnN
    cacoshfnN        cabsfnN            crealfnN
    casinhfnN        cpowfnN
```

35 for supported types `_FloatNx`:

```
40    cacosfnNx         catanhfnNx         csqrtfnNx
    casinfnNx         ccoshfnNx         cargfnNx
    catanfNx         csinhfnNx         cimagfnNx
    ccosfnNx         ctanhfnNx         CMPLXFNx
    csinfnNx         cexpfnNx         conjfnNx
    ctanfNx         clogfnNx         cprojfnNx
```

**cacoshfNx**  
**casinhfNx**

**cabsfNx**  
**cpowfNx**

**crealfNx**

### X.8.3 <math.h>

[1] The following identifiers are defined or declared only if

5 **\_\_STDC\_WANT\_IEC\_60559\_TYPES\_EXT\_\_** is defined as a macro at the point in the source file where <math.h> is first included:

**long\_double\_t**

for supported types **\_FloatN**:

10	<b>_FloatN_t</b>	<b>log1pfN</b>	<b>fromfpfN</b>
	<b>HUGE_VAL_FN</b>	<b>log2fN</b>	<b>ufromfpfN</b>
	<b>SNANFN</b>	<b>logbfN</b>	<b>fromfpxfN</b>
	<b>FP_FAST_FMAFN</b>	<b>modffN</b>	<b>ufromfpxfN</b>
	<b>acosfN</b>	<b>scalbnfN</b>	<b>fmodfN</b>
15	<b>asinfN</b>	<b>scalblnfN</b>	<b>remainderfN</b>
	<b>atanfN</b>	<b>cbrtfN</b>	<b>remquofN</b>
	<b>atan2fN</b>	<b>fabsfN</b>	<b>copysignfN</b>
	<b>cosfN</b>	<b>hypotfN</b>	<b>nanfN</b>
	<b>sinfN</b>	<b>powfN</b>	<b>nextafterfN</b>
20	<b>tanfN</b>	<b>sqrtfN</b>	<b>nextupfN</b>
	<b>acoshfN</b>	<b>erffN</b>	<b>nextdownfN</b>
	<b>asinhfN</b>	<b>erfcfN</b>	<b>canonicalizefN</b>
	<b>atanhfN</b>	<b>lgammafN</b>	<b>encodefN</b>
	<b>coshfN</b>	<b>tgammafN</b>	<b>decodefN</b>
25	<b>sinhfN</b>	<b>ceilfN</b>	<b>fdimfN</b>
	<b>tanhfN</b>	<b>floorfN</b>	<b>fmaxfN</b>
	<b>expfN</b>	<b>nearbyintfN</b>	<b>fminfN</b>
	<b>exp2fN</b>	<b>rintfN</b>	<b>fmaxmagfN</b>
	<b>expm1fN</b>	<b>lrintfN</b>	<b>fminmagfN</b>
30	<b>frexpfN</b>	<b>llrintfN</b>	<b>fmafN</b>
	<b>ilogbfN</b>	<b>roundfN</b>	<b>totalorderfN</b>
	<b>ldexpfN</b>	<b>lroundfN</b>	<b>totalordermagfN</b>
	<b>llogbfN</b>	<b>llroundfN</b>	<b>getpayloadfN</b>
	<b>logfN</b>	<b>truncfN</b>	<b>setpayloadfN</b>
35	<b>log10fN</b>	<b>roundevenfN</b>	<b>setpayloadsigfN</b>

for supported types **\_FloatNx**:

	<b>HUGE_VAL_FNX</b>	<b>log1pfNx</b>	<b>roundevenfNx</b>
	<b>SNANFNx</b>	<b>log2fNx</b>	<b>fromfpfNx</b>
40	<b>FP_FAST_FMAFNx</b>	<b>logbfNx</b>	<b>ufromfpfNx</b>
	<b>acosfNx</b>	<b>modffNx</b>	<b>fromfpxfNx</b>
	<b>asinfNx</b>	<b>scalbnfNx</b>	<b>ufromfpxfNx</b>
	<b>atanfNx</b>	<b>scalblnfNx</b>	<b>fmodfNx</b>
	<b>atan2fNx</b>	<b>cbrtfNx</b>	<b>remainderfNx</b>
	<b>cosfNx</b>	<b>fabsfNx</b>	<b>remquofNx</b>
45	<b>sinfNx</b>	<b>hypotfNx</b>	<b>copysignfNx</b>
	<b>tanfNx</b>	<b>powfNx</b>	<b>nanfNx</b>
	<b>acoshfNx</b>	<b>sqrtfNx</b>	<b>nextafterfNx</b>
	<b>asinhfNx</b>	<b>erffNx</b>	<b>nextupfNx</b>

	<code>atanhfNx</code>	<code>erfcfNx</code>	<code>nextdownfNx</code>
	<code>expfNx</code>	<code>lgammafNx</code>	<code>canonicalizefNx</code>
	<code>coshfNx</code>	<code>tgammafNx</code>	<code>fdimfNx</code>
5	<code>sinhfNx</code>	<code>ceilfNx</code>	<code>fmaxfNx</code>
	<code>tanhfNx</code>	<code>floorfNx</code>	<code>fminfNx</code>
	<code>exp2fNx</code>	<code>nearbyintfNx</code>	<code>fmaxmagfNx</code>
	<code>expm1fNx</code>	<code>rintfNx</code>	<code>fminmagfNx</code>
	<code>frexpfNx</code>	<code>lrintfNx</code>	<code>fmafNx</code>
10	<code>ilogbfNx</code>	<code>llrintfNx</code>	<code>totalorderfNx</code>
	<code>llogbfNx</code>	<code>roundfNx</code>	<code>totalordermagfNx</code>
	<code>ldexpfNx</code>	<code>lroundfNx</code>	<code>getpayloadfNx</code>
	<code>logfNx</code>	<code>llroundfNx</code>	<code>setpayloadfNx</code>
	<code>log10fNx</code>	<code>truncfNx</code>	<code>setpayloadsigfNx</code>

for supported types `_FloatM` and `_FloatN` where  $M < N$ :

15	<code>FP_FAST_FMADDFN</code>	<code>FP_FAST_FMFMAFN</code>	<code>fMmulfN</code>
	<code>FP_FAST_FMSUBFN</code>	<code>FP_FAST_FMSQRTFN</code>	<code>fMdivfN</code>
	<code>FP_FAST_FMMULFN</code>	<code>fMaddfN</code>	<code>fMfmfN</code>
	<code>FP_FAST_FMDIVFN</code>	<code>fMsubfN</code>	<code>fMsqrtfN</code>

for supported types `_FloatM` and `_FloatNx` where  $M \leq N$ :

20	<code>FP_FAST_FMADDFNx</code>	<code>FP_FAST_FMFMAFNx</code>	<code>fMmulfNx</code>
	<code>FP_FAST_FMSUBFNx</code>	<code>FP_FAST_FMSQRTFNx</code>	<code>fMdivfNx</code>
	<code>FP_FAST_FMMULFNx</code>	<code>fMaddfNx</code>	<code>fMfmfNx</code>
	<code>FP_FAST_FMDIVFNx</code>	<code>fMsubfNx</code>	<code>fMsqrtfNx</code>

for supported types `_FloatMx` and `_FloatN` where  $M < N$ :

25	<code>FP_FAST_FMXADDFN</code>	<code>FP_FAST_FMXFMAFN</code>	<code>fMxmulfN</code>
	<code>FP_FAST_FMXSUBFN</code>	<code>FP_FAST_FMXSQRTFN</code>	<code>fMxdivfN</code>
	<code>FP_FAST_FMXMULFN</code>	<code>fMxaddfN</code>	<code>fMxfmfN</code>
	<code>FP_FAST_FMXDIVFN</code>	<code>fMxsubfN</code>	<code>fMxsqrtfN</code>

for supported types `_FloatMx` and `_FloatNx` where  $M < N$ :

30	<code>FP_FAST_FMXADDFNx</code>	<code>FP_FAST_FMXFMAFNx</code>	<code>fMxmulfNx</code>
	<code>FP_FAST_FMXSUBFNx</code>	<code>FP_FAST_FMXSQRTFNx</code>	<code>fMxdivfNx</code>
	<code>FP_FAST_FMXMULFNx</code>	<code>fMxaddfNx</code>	<code>fMxfmfNx</code>
	<code>FP_FAST_FMXDIVFNx</code>	<code>fMxsubfNx</code>	<code>fMxsqrtfNx</code>

for supported IEC 60559 arithmetic or non-arithmetic binary interchange formats of widths  $M$  and  $N$ :

`fMencfN`

for supported types `_DecimalN`, where  $N \neq 32, 64, \text{ and } 128$ :

40	<code>_DecimalN_t</code>	<code>logbdN</code>	<code>fmoddN</code>
	<code>HUGE_VAL_DN</code>	<code>modfdN</code>	<code>remainderdN</code>
	<code>SNANDN</code>	<code>scalbndN</code>	<code>copysigndN</code>
	<code>FP_FAST_FMADN</code>	<code>scalblndN</code>	<code>nandN</code>
	<code>acosdN</code>	<code>cbtrdN</code>	<code>nextafterdN</code>
	<code>asindN</code>	<code>fabsdN</code>	<code>nextupdN</code>

	atand <i>N</i>	hypotd <i>N</i>	nextdownd <i>N</i>
	atan2d <i>N</i>	powd <i>N</i>	canonicalized <i>N</i>
	cosd <i>N</i>	sqrtd <i>N</i>	quantized <i>N</i>
5	sind <i>N</i>	erfd <i>N</i>	samequantumd <i>N</i>
	tand <i>N</i>	erfcd <i>N</i>	quantumd <i>N</i>
	acoshd <i>N</i>	lgammad <i>N</i>	llquantexpd <i>N</i>
	asinhd <i>N</i>	tgammad <i>N</i>	encodedecd <i>N</i>
	atanhd <i>N</i>	ceild <i>N</i>	decodedecd <i>N</i>
10	coshd <i>N</i>	floord <i>N</i>	encodebind <i>N</i>
	sinhd <i>N</i>	nearbyintd <i>N</i>	decodebind <i>N</i>
	tanhd <i>N</i>	rintd <i>N</i>	fdimd <i>N</i>
	expd <i>N</i>	lrintd <i>N</i>	fmaxd <i>N</i>
	exp2d <i>N</i>	llrintd <i>N</i>	fmind <i>N</i>
15	expm1d <i>N</i>	roundd <i>N</i>	fmaxmagd <i>N</i>
	frexp <i>N</i>	lroundd <i>N</i>	fminmagd <i>N</i>
	ilogbd <i>N</i>	llroundd <i>N</i>	fmad <i>N</i>
	llogbd <i>N</i>	truncd <i>N</i>	totalorderd <i>N</i>
	ldexpd <i>N</i>	roundevend <i>N</i>	totalordermagd <i>N</i>
	logd <i>N</i>	fromfpd <i>N</i>	getpayloadd <i>N</i>
20	log10d <i>N</i>	ufromfpd <i>N</i>	setpayloadd <i>N</i>
	log1pd <i>N</i>	fromfpd <i>N</i>	setpayloadsigd <i>N</i>
	log2d <i>N</i>	ufromfpd <i>N</i>	

for supported types `_DecimalNx`:

25	HUGE_VAL_D <i>Nx</i>	log2d <i>Nx</i>	ufromfpd <i>Nx</i>
	SNAND <i>Nx</i>	logbd <i>Nx</i>	fromfpd <i>Nx</i>
	FP_FAST_FMAD <i>Nx</i>	modfd <i>Nx</i>	ufromfpd <i>Nx</i>
	acosd <i>Nx</i>	scalbnd <i>Nx</i>	fmodd <i>Nx</i>
	asind <i>Nx</i>	scalblnd <i>Nx</i>	remainderd <i>Nx</i>
30	atand <i>Nx</i>	cbrrd <i>Nx</i>	copysignd <i>Nx</i>
	atan2d <i>Nx</i>	fabsd <i>Nx</i>	nand <i>Nx</i>
	cosd <i>Nx</i>	hypotd <i>Nx</i>	nextafterd <i>Nx</i>
	sind <i>Nx</i>	powd <i>Nx</i>	nextupd <i>Nx</i>
	tand <i>Nx</i>	sqrtd <i>Nx</i>	nextdownd <i>Nx</i>
35	acoshd <i>Nx</i>	erfd <i>Nx</i>	canonicalized <i>Nx</i>
	asinhd <i>Nx</i>	erfcd <i>Nx</i>	quantized <i>Nx</i>
	atanhd <i>Nx</i>	lgammad <i>Nx</i>	samequantumd <i>Nx</i>
	coshd <i>Nx</i>	tgammad <i>Nx</i>	quantumd <i>Nx</i>
	sinhd <i>Nx</i>	ceild <i>Nx</i>	llquantexpd <i>Nx</i>
40	tanhd <i>Nx</i>	floord <i>Nx</i>	fdimd <i>Nx</i>
	expd <i>Nx</i>	nearbyintd <i>Nx</i>	fmaxd <i>Nx</i>
	exp2d <i>Nx</i>	rintd <i>Nx</i>	fmind <i>Nx</i>
	expm1d <i>Nx</i>	lrintd <i>Nx</i>	fmaxmagd <i>Nx</i>
	frexp <i>Nx</i>	llrintd <i>Nx</i>	fminmagd <i>Nx</i>
45	ilogbd <i>Nx</i>	roundd <i>Nx</i>	fmad <i>Nx</i>
	llogbd <i>Nx</i>	lroundd <i>Nx</i>	totalorderd <i>Nx</i>
	ldexpd <i>Nx</i>	llroundd <i>Nx</i>	totalordermagd <i>Nx</i>
	logd <i>Nx</i>	truncd <i>Nx</i>	getpayloadd <i>Nx</i>
	log10d <i>Nx</i>	roundevend <i>Nx</i>	setpayloadd <i>Nx</i>
	log1pd <i>Nx</i>	fromfpd <i>Nx</i>	setpayloadsigd <i>Nx</i>

for supported types `_DecimalM` and `_DecimalN` where  $M < N$  and  $M$  and  $N$  are not both one of 32, 64, and 128:

	<code>FP_FAST_DMADDDN</code>	<code>FP_FAST_DMFMADN</code>	<code>dMmuldN</code>
	<code>FP_FAST_DMSUBDN</code>	<code>FP_FAST_DMSQRTDN</code>	<code>dMdivdN</code>
5	<code>FP_FAST_DMMULDN</code>	<code>dMadddN</code>	<code>dMfmadN</code>
	<code>FP_FAST_DMDIVDN</code>	<code>dMsubdN</code>	<code>dMsqrtdN</code>

for supported types `_DecimalM` and `_DecimalNx` where  $M \leq N$ :

	<code>FP_FAST_DMADDDNx</code>	<code>FP_FAST_DMFMADNx</code>	<code>dMmuldNx</code>
	<code>FP_FAST_DMSUBDNx</code>	<code>FP_FAST_DMSQRTDNx</code>	<code>dMdivdNx</code>
10	<code>FP_FAST_DMMULDNx</code>	<code>dMadddNx</code>	<code>dMfmadNx</code>
	<code>FP_FAST_DMDIVDNx</code>	<code>dMsubdNx</code>	<code>dMsqrtdNx</code>

for supported types `_DecimalMx` and `_DecimalN` where  $M < N$ :

	<code>FP_FAST_DMxADDDN</code>	<code>FP_FAST_DMxFMADN</code>	<code>dMxmuldN</code>
	<code>FP_FAST_DMxSUBDN</code>	<code>FP_FAST_DMxSQRTDN</code>	<code>dMxdivdN</code>
15	<code>FP_FAST_DMxMULDN</code>	<code>dMxadddN</code>	<code>dMxfmadN</code>
	<code>FP_FAST_DMxDIVDN</code>	<code>dMxsubdN</code>	<code>dMxsqrtdN</code>

for supported types `_DecimalMx` and `_DecimalNx` where  $M < N$ :

	<code>FP_FAST_DMxADDDNx</code>	<code>FP_FAST_DMxFMADNx</code>	<code>dMxmuldNx</code>
	<code>FP_FAST_DMxSUBDNx</code>	<code>FP_FAST_DMxSQRTDNx</code>	<code>dMxdivdNx</code>
20	<code>FP_FAST_DMxMULDNx</code>	<code>dMxadddNx</code>	<code>dMxfmadNx</code>
	<code>FP_FAST_DMxDIVDNx</code>	<code>dMxsubdNx</code>	<code>dMxsqrtdNx</code>

for supported IEC 60559 arithmetic and non-arithmetic decimal interchange formats of widths  $M$  and  $N$ :

	<code>dMencdecN</code>	<code>dMencbindN</code>
--	------------------------	-------------------------

## 25 X.8.4 <stdlib.h>

[1] The following identifiers are declared only if `__STDC_WANT_IEC_60559_TYPES_EXT__` is defined as a macro at the point in the source file where `<stdlib.h>` is first included:

for supported types `_FloatN`:

	<code>strfromfN</code>	<code>strtofN</code>
--	------------------------	----------------------

30 for supported types `_FloatNx`:

	<code>strfromfNx</code>	<code>strtofNx</code>
--	-------------------------	-----------------------

for supported types `_DecimalN`, where  $N \neq 32, 64, \text{ and } 128$ :

	<code>strfromdN</code>	<code>strtodN</code>
--	------------------------	----------------------

for supported types `_DecimalNx`:

35	<code>strfromdNx</code>	<code>strtodNx</code>
----	-------------------------	-----------------------



for supported IEC 60559 arithmetic and non-arithmetic binary interchange formats of width  $N$ :

**strfromencfN**                      **strtoencfN**

for supported IEC 60559 arithmetic and non-arithmetic decimal interchange formats of width  $N$ :

**strfromencdecN**                      **strtoencdecN**  
5        **strfromencbindN**                      **strtoencbindN**

## X.9 Complex arithmetic <complex.h>

[1] This subclause specifies complex functions for corresponding real types that are binary floating types.

[2] Each function synopsis in 7.3 specifies a family of functions including a principal function with one or more **double complex** parameters and a **double complex** or **double** return value. This subclause expands the synopsis to also include other functions, with the same name as the principal function but with **fN** and **fNx** suffixes, which are corresponding functions whose parameters and return values have corresponding real types **\_FloatN** and **\_FloatNx**.

[3] The following function prototypes are added to the synopses of the respective subclauses in 7.3. For each binary floating type that the implementation provides, <complex.h> shall declare the associated functions. Conversely, for each such type that the implementation does not provide, <complex.h> shall not declare the associated functions.

### 7.3.5 Trigonometric functions

```
20        _FloatN complex cacosfN(_FloatN complex z);  
      _FloatNx complex cacosfNx(_FloatNx complex z);  
  
      _FloatN complex casinfnN(_FloatN complex z);  
      _FloatNx complex casinfnNx(_FloatNx complex z);  
  
25        _FloatN complex catanfN(_FloatN complex z);  
      _FloatNx complex catanfNx(_FloatNx complex z);  
  
      _FloatN complex ccosfnN(_FloatN complex z);  
      _FloatNx complex ccosfnNx(_FloatNx complex z);  
  
30        _FloatN complex csinfnN(_FloatN complex z);  
      _FloatNx complex csinfnNx(_FloatNx complex z);  
  
      _FloatN complex ctanfN(_FloatN complex z);  
35        _FloatNx complex ctanfNx(_FloatNx complex z);
```

### 7.3.6 Hyperbolic functions

```
40        _FloatN complex cacoshfnN(_FloatN complex z);  
      _FloatNx complex cacoshfnNx(_FloatNx complex z);  
  
      _FloatN complex casinhfnN(_FloatN complex z);  
      _FloatNx complex casinhfnNx(_FloatNx complex z);  
  
      _FloatN complex catanhfnN(_FloatN complex z);  
45        _FloatNx complex catanhfnNx(_FloatNx complex z);
```

```

_FloatN complex ccoshfN(_FloatN complex z);
_FloatNx complex ccoshfNx(_FloatNx complex z);

5  _FloatN complex csinhfN(_FloatN complex z);
   _FloatNx complex csinhfNx(_FloatNx complex z);

   _FloatN complex ctanhfN(_FloatN complex z);
   _FloatNx complex ctanhfNx(_FloatNx complex z);

```

### 7.3.7 Exponential and logarithmic functions

```

_FloatN complex cexpfN(_FloatN complex z);
_FloatNx complex cexpfNx(_FloatNx complex z);

   _FloatN complex clogfN(_FloatN complex z);
   _FloatNx complex clogfNx(_FloatNx complex z);

```

### 7.3.8 Power and absolute value functions

```

_FloatN cabsfN(_FloatN complex z);
_FloatNx cabsfNx(_FloatNx complex z);

   _FloatN complex cpowfN(_FloatN complex x,
   _FloatN complex y);
   _FloatNx complex cpowfNx(_FloatNx complex x,
   _FloatNx complex y);

   _FloatN complex csqrtfN(_FloatN complex z);
   _FloatNx complex csqrtfNx(_FloatNx complex z);

```

### 7.3.9 Manipulation functions

```

_FloatN cargfN(_FloatN complex z);
_FloatNx cargfNx(_FloatNx complex z);

   _FloatN cimagfN(_FloatN complex z);
   _FloatNx cimagfNx(_FloatNx complex z);

   _FloatN complex CMPLXFN(_FloatN x, _FloatN y);
   _FloatNx complex CMPLXFNx(_FloatNx x, _FloatNx y);

   _FloatN complex conjfN(_FloatN complex z);
   _FloatNx complex conjfNx(_FloatNx complex z);

   _FloatN complex cprojfN(_FloatN complex z);
   _FloatNx complex cprojfNx(_FloatNx complex z);

   _FloatN crealfN(_FloatN complex z);
   _FloatNx crealfNx(_FloatNx complex z);

```

[4] For the functions listed in 7.31.1 (Future library directions for `<complex.h>`), the possible suffixes are expanded to also include `fN` and `fNx`

## X.10 Floating-point environment

[1] This subclause broadens the effects of the floating-point environment for standard floating types to also apply to binary floating types.

[2] The rounding control pragma (7.6.2)

5           **#pragma STDC FENV\_ROUND *direction***

affects operations for standard and binary floating types. *direction* shall be one of the names of the supported rounding direction macros for use with **fegetround** and **fesetround** (7.6), or **FE\_DYNAMIC**.

[3] In 7.6.2, the table entitled

### 10           **Functions affected by constant rounding modes – for standard floating types**

applies to binary floating types, as well as standard floating types. Each **<math.h>** function family listed in the table indicates the family of functions of all standard and binary floating types (for example, the **acos** family includes **acosf**, **acosl**, **acosfN**, and **acosfNx** as well as **acos**). The **fMencfN**, **strfromencfN**, and **strtoencfN** functions for binary interchange types are also  
15 affected by constant rounding modes.

[4] In 7.6.3, in the table of functions affected by constant rounding modes for decimal floating types, each **<math.h>** function family indicates the family of functions of all decimal floating types (for example, the **acos** family includes **acosdN** and **acosdNx**). The **dMencbindN**, **dMencdecN**, **strfromencbindN**, **strfromencdecN**, **strtoencbindN**, and **strtoencdecN** functions for  
20 decimal interchange types are also affected by constant rounding modes.

[5] The **fegetround** function (7.6.5.2) gets the current value of the dynamic rounding direction mode for operations for standard and binary floating types. The **fesetround** function (7.6.5.5) sets the dynamic rounding direction mode to the rounding direction represented by its argument **round** for operations for standard and binary floating types.

## 25           **X.11 Mathematics <math.h>**

[1] This subclause specifies functions and macros for interchange and extended floating types, generally corresponding to those specified in 7.12 and F.10.

[2] All classification macros (7.12.3) and comparison macros (7.12.17) naturally extend to handle interchange and extended floating types. For comparison macros, if neither of the sets of values of the  
30 argument formats is a subset of (or equivalent to) the other, the behavior is undefined.

[3] This subclause also specifies encoding conversion functions that are part of support for the non-arithmetic interchange formats in IEC 60559 (see X.2.2).

[4] Most function synopses in 7.12 specify a family of functions including a principal function with one or more **double** parameters, a **double** return value, or both. The synopses are expanded to also  
35 include functions with the same name as the principal function but with **fN**, **fNx**, **dN**, and **dNx** suffixes, which are corresponding functions whose parameters, return values, or both are of types **\_FloatN**, **\_FloatNx**, **\_DecimalN**, and **\_DecimalNx**, respectively.

[5] For each interchange or extended floating type that the implementation provides, **<math.h>** shall define the associated macros and declare the associated functions. Conversely, for each such type that  
40 the implementation does not provide, **<math.h>** shall not define the associated macros or declare the associated functions unless explicitly specified otherwise.

[6] With the types

```
float_t
double_t
```

5 in 7.12 are included the type

```
long_double_t
```

and for each supported type `_FloatN`, the type

```
_FloatN_t
```

10

and for each supported type `_DecimalN`, the type

```
_DecimalN_t
```

These are floating types, such that:

15

- each of the types has at least the range and precision of the corresponding real floating type;
- `long_double_t` has at least the range and precision of `double_t`;
- `_FloatN_t` has at least the range and precision of `_FloatM_t` if  $N > M$ ;
- `_DecimalN_t` has at least the range and precision of `_DecimalM_t` if  $N > M$ .

20 If `FLT_RADIX` is 2 and `FLT_EVAL_METHOD` (X.3) is nonnegative, then each of the types corresponding to a standard or binary floating type is the type whose range and precision are specified by `FLT_EVAL_METHOD` (X.3) to be used for evaluating operations and constants of that standard or binary floating type. If `DEC_EVAL_METHOD` (X.3) is nonnegative, then each of the types corresponding to a decimal floating type is the type whose range and precision are specified by `DEC_EVAL_METHOD` (X.3) to be used for evaluating operations and constants of that decimal floating type.

25 [Jens, I'm preparing an example to go here.]

### X.11.1 Macros

[1] This subclause adds macros in 7.12 as follows.

[2] The macros

```
30 HUGE_VAL_FN
HUGE_VAL_DN
HUGE_VAL_FNX
HUGE_VAL_DNX
```

expand to constant expressions of types `_FloatN`, `_DecimalN`, `_FloatNx`, and `_DecimalNx`, respectively, representing positive infinity.

35 [3] The signaling NaN macros

```
SNANFN
SNANDN
SNANFNX
SNANDNX
```

expand to constant expressions of types `_FloatN`, `_DecimalN`, `_FloatNx`, and `_DecimalNx`, respectively, representing a signaling NaN. If a signaling NaN macro is used for initializing an object of the same type that has static or thread-local storage duration, the object is initialized with a signaling NaN value.

5 [4] The macros

```

FP_FAST_FMAFN
FP_FAST_FMADN
FP_FAST_FMAFNx
FP_FAST_FMADNx

```

10 are, respectively, `_FloatN`, `_DecimalN`, `_FloatNx`, and `_DecimalNx` analogues of `FP_FAST_FMA`.

[5] The macros in the following lists are interchange and extended floating type analogues of `FP_FAST_FADD`, `FP_FAST_FADDL`, `FP_FAST_DADDL`, etc.

[6] For  $M < N$ , the macros

```

15 FP_FAST_FMADDFN
FP_FAST_FMSUBFN
FP_FAST_FMMULFN
FP_FAST_FMDIVFN
FP_FAST_FMFMAFN
20 FP_FAST_FMSQRTFN
FP_FAST_DMADDDN
FP_FAST_DMSUBDN
FP_FAST_DMMULDN
FP_FAST_DMDIVDN
25 FP_FAST_DMFMAFN
FP_FAST_DMSQRFDN

```

characterize the corresponding functions whose arguments are of an interchange floating type of width  $N$  and whose return type is an interchange floating type of width  $M$ .

[7] For  $M \leq N$ , the macros

```

30 FP_FAST_FMADDFNx
FP_FAST_FMSUBFNx
FP_FAST_FMMULFNx
FP_FAST_FMDIVFNx
FP_FAST_FMFMAFNx
35 FP_FAST_FMSQRTFNx
FP_FAST_DMADDDNx
FP_FAST_DMSUBDNx
FP_FAST_DMMULDNx
FP_FAST_DMDIVDNx
40 FP_FAST_DMFMAFNx
FP_FAST_DMSQRFDNx

```

characterize the corresponding functions whose arguments are of an extended floating type that extends a format of width  $N$  and whose return type is an interchange floating type of width  $M$ .

[8] For  $M < N$ , the macros

```
FP_FAST_FMXADDFN
FP_FAST_FMXSUBFN
FP_FAST_FMXMULFN
5 FP_FAST_FMXDIVFN
FP_FAST_FMXFMAFN
FP_FAST_FMXSQRTFN
FP_FAST_DMxADDDN
FP_FAST_DMXSUBDN
10 FP_FAST_DMXMULDN
FP_FAST_DMXDIVDN
FP_FAST_DMXFMDN
FP_FAST_DMXSQRTDN
```

characterize the corresponding functions whose arguments are of an interchange floating type of width  $N$  and whose return type is an extended floating type that extends a format of width  $M$ .

[9] For  $M < N$ , the macros

```
FP_FAST_FMXADDFNX
FP_FAST_FMXSUBFNX
FP_FAST_FMXMULFNX
20 FP_FAST_FMXDIVFNX
FP_FAST_FMXFMAFNX
FP_FAST_FMXSQRTFNX
FP_FAST_DMxADDDNX
FP_FAST_DMXSUBDNX
25 FP_FAST_DMXMULDNX
FP_FAST_DMXDIVDNX
FP_FAST_DMXFMDNX
FP_FAST_DMXSQRTDNX
```

characterize the corresponding functions whose arguments are of an extended floating type that extends a format of width  $N$  and whose return type is an extended floating type that extends a format of width  $M$ .

### X.11.2 Function prototypes

[1] This subclause adds the following function prototypes to the synopses of the respective subclauses in 7.12.

#### 7.12.4 Trigonometric functions

```
_FloatN acosfN(_FloatN x);
_FloatNx acosfNx(_FloatNx x);
_DecimalN acosdN(_DecimalN x);
_DecimalNx acosdNx(_DecimalNx x);
40
_FloatN asinfN(_FloatN x);
_FloatNx asinfNx(_FloatNx x);
_DecimalN asindN(_DecimalN x);
_DecimalNx asindNx(_DecimalNx x);
45
```

```

5  _FloatN atanfN(_FloatN x);
   _FloatNx atanfNx(_FloatNx x);
   _DecimalN atandN(_DecimalN x);
   _DecimalNx atandNx(_DecimalNx x);

   _FloatN atan2fN(_FloatN y, _FloatN x);
   _FloatNx atan2fNx(_FloatNx y, _FloatNx x);
   _DecimalN atan2dN(_DecimalN y, _DecimalN x);
   _DecimalNx atan2dNx(_DecimalNx y, _DecimalNx x);

10  _FloatN cosfN(_FloatN x);
   _FloatNx cosfNx(_FloatNx x);
   _DecimalN cosdN(_DecimalN x);
   _DecimalNx cosdNx(_DecimalNx x);

15  _FloatN sinfN(_FloatN x);
   _FloatNx sinfNx(_FloatNx x);
   _DecimalN sindN(_DecimalN x);
   _DecimalNx sindNx(_DecimalNx x);

20  _FloatN tanfN(_FloatN x);
   _FloatNx tanfNx(_FloatNx x);
   _DecimalN tandN(_DecimalN x);
   _DecimalNx tandNx(_DecimalNx x);

25

```

#### 7.12.5 Hyperbolic functions

```

   _FloatN acoshfN(_FloatN x);
   _FloatNx acoshfNx(_FloatNx x);
   _DecimalN acoshdN(_DecimalN x);
   _DecimalNx acoshdNx(_DecimalNx x);

30  _FloatN asinhfN(_FloatN x);
   _FloatNx asinhfNx(_FloatNx x);
   _DecimalN asinhdN(_DecimalN x);
   _DecimalNx asinhdNx(_DecimalNx x);

35  _FloatN atanhfN(_FloatN x);
   _FloatNx atanhfNx(_FloatNx x);
   _DecimalN atanhdN(_DecimalN x);
   _DecimalNx atanhdNx(_DecimalNx x);

40  _FloatN coshfN(_FloatN x);
   _FloatNx coshfNx(_FloatNx x);
   _DecimalN coshdN(_DecimalN x);
   _DecimalNx coshdNx(_DecimalNx x);

45  _FloatN sinhN(_FloatN x);
   _FloatNx sinhNx(_FloatNx x);
   _DecimalN sinhN(_DecimalN x);
   _DecimalNx sinhNx(_DecimalNx x);

50

```

```

_FloatN tanhfN(_FloatN x);
_FloatNx tanhfNx(_FloatNx x);
_DecimalN tanhdN(_DecimalN x);
_DecimalNx tanhdNx(_DecimalNx x);

```

5

### 7.12.6 Exponential and logarithmic functions

```

_FloatN expfN(_FloatN x);
_FloatNx expfNx(_FloatNx x);
_DecimalN expdN(_DecimalN x);
_DecimalNx expdNx(_DecimalNx x);

```

10

```

_FloatN exp2fN(_FloatN x);
_FloatNx exp2fNx(_FloatNx x);
_DecimalN exp2dN(_DecimalN x);
_DecimalNx exp2dNx(_DecimalNx x);

```

15

```

_FloatN expm1fN(_FloatN x);
_FloatNx expm1fNx(_FloatNx x);
_DecimalN expm1dN(_DecimalN x);
_DecimalNx expm1dNx(_DecimalNx x);

```

20

```

_FloatN frexpfN(_FloatN value, int *exp);
_FloatNx frexpfNx(_FloatNx value, int *exp);
_DecimalN frexpdN(_DecimalN value, int *exp);
_DecimalNx frexpdNx(_DecimalNx value, int *exp);

```

25

```

int ilogbfN(_FloatN x);
int ilogbfNx(_FloatNx x);
int ilogbdN(_DecimalN x);
int ilogbdNx(_DecimalNx x);

```

30

```

_FloatN ldexpfN(_FloatN value, int exp);
_FloatNx ldexpfNx(_FloatNx value, int exp);
_DecimalN ldexpdN(_DecimalN value, int exp);
_DecimalNx ldexpdNx(_DecimalNx value, int exp);

```

35

```

long int llogbfN(_FloatN x);
long int llogbfNx(_FloatNx x);
long int llogbdN(_DecimalN x);
long int llogbdNx(_DecimalNx x);

```

40

```

_FloatN logfN(_FloatN x);
_FloatNx logfNx(_FloatNx x);
_DecimalN logdN(_DecimalN x);
_DecimalNx logdNx(_DecimalNx x);

```

45

```

_FloatN log10fN(_FloatN x);
_FloatNx log10fNx(_FloatNx x);
_DecimalN log10dN(_DecimalN x);
_DecimalNx log10dNx(_DecimalNx x);

```

50



```

5  _FloatN loglpdfN(_FloatN x);
   _FloatNx loglpdfNx(_FloatNx x);
   _DecimalN loglpdfN(_DecimalN x);
   _DecimalNx loglpdfNx(_DecimalNx x);

10  _FloatN log2fN(_FloatN x);
   _FloatNx log2fNx(_FloatNx x);
   _DecimalN log2dN(_DecimalN x);
   _DecimalNx log2dNx(_DecimalNx x);

15  _FloatN logbfN(_FloatN x);
   _FloatNx logbfNx(_FloatNx x);
   _DecimalN logbdN(_DecimalN x);
   _DecimalNx logbdNx(_DecimalNx x);

20  _FloatN modfffN(_FloatN x, _FloatN *iptr);
   _FloatNx modfffNx(_FloatNx x, _FloatNx *iptr);
   _DecimalN modfdN(_DecimalN x, _DecimalN *iptr);
   _DecimalNx modfdNx(_DecimalNx x, _DecimalNx *iptr);

25  _FloatN scalbnfN(_FloatN value, int exp);
   _FloatNx scalbnfNx(_FloatNx value, int exp);
   _DecimalN scalbndN(_DecimalN value, int exp);
   _DecimalNx scalbndNx(_DecimalNx value, int exp);

30  _FloatN scalblnfN(_FloatN value, long int exp);
   _FloatNx scalblnfNx(_FloatNx value, long int exp);
   _DecimalN scalblndN(_DecimalN value, long int exp);
   _DecimalNx scalblndNx(_DecimalNx value, long int exp);

30  7.12.7 Power and absolute-value functions

   _FloatN cbrtfN(_FloatN x);
   _FloatNx cbrtfNx(_FloatNx x);
   _DecimalN cbrtdN(_DecimalN x);
   _DecimalNx cbrtdNx(_DecimalNx x);

35  _FloatN fabsfN(_FloatN x);
   _FloatNx fabsfNx(_FloatNx x);
   _DecimalN fabsdN(_DecimalN x);
   _DecimalNx fabsdNx(_DecimalNx x);

40  _FloatN hypotfN(_FloatN x, _FloatN y);
   _FloatNx hypotfNx(_FloatNx x, _FloatNx y);
   _DecimalN hypotdN(_DecimalN x, _DecimalN y);
   _DecimalNx hypotdNx(_DecimalNx x, _DecimalNx y);

45  _FloatN powfN(_FloatN x, _FloatN y);
   _FloatNx powfNx(_FloatNx x, _FloatNx y);
   _DecimalN powdN(_DecimalN x, _DecimalN y);
   _DecimalNx powdNx(_DecimalNx x, _DecimalNx y);

50  _FloatN sqrtfN(_FloatN x);
   _FloatNx sqrtfNx(_FloatNx x);

```

```
_DecimalN sqrtDN(_DecimalN x);
_DecimalNx sqrtDNx(_DecimalNx x);
```

#### 7.12.8 Error and gamma functions

```
5  _FloatN erffN(_FloatN x);
   _FloatNx erffNx(_FloatNx x);
   _DecimalN erfdN(_DecimalN x);
   _DecimalNx erfdNx(_DecimalNx x);

10 _FloatN erfcfN(_FloatN x);
   _FloatNx erfcfNx(_FloatNx x);
   _DecimalN erfcdN(_DecimalN x);
   _DecimalNx erfcdNx(_DecimalNx x);

15 _FloatN lgammafN(_FloatN x);
   _FloatNx lgammafNx(_FloatNx x);
   _DecimalN lgammadN(_DecimalN x);
   _DecimalNx lgammadNx(_DecimalNx x);

20 _FloatN tgammafN(_FloatN x);
   _FloatNx tgammafNx(_FloatNx x);
   _DecimalN tgammadN(_DecimalN x);
   _DecimalNx tgammadNx(_DecimalNx x);
```

#### 7.12.9 Nearest integer functions

```
25 _FloatN ceilfN(_FloatN x);
   _FloatNx ceilfNx(_FloatNx x);
   _DecimalN ceildN(_DecimalN x);
   _DecimalNx ceildNx(_DecimalNx x);

30 _FloatN floorfN(_FloatN x);
   _FloatNx floorfNx(_FloatNx x);
   _DecimalN floordN(_DecimalN x);
   _DecimalNx floordNx(_DecimalNx x);

35 _FloatN nearbyintfN(_FloatN x);
   _FloatNx nearbyintfNx(_FloatNx x);
   _DecimalN nearbyintdN(_DecimalN x);
   _DecimalNx nearbyintdNx(_DecimalNx x);

40 _FloatN rintfN(_FloatN x);
   _FloatNx rintfNx(_FloatNx x);
   _DecimalN rintdN(_DecimalN x);
   _DecimalNx rintdNx(_DecimalNx x);

45 long int lrintfN(_FloatN x);
   long int lrintfNx(_FloatNx x);
   long int lrintdN(_DecimalN x);
   long int lrintdNx(_DecimalNx x);
```

```

long long int llrintfN(_FloatN x);
long long int llrintfNx(_FloatNx x);
long long int llrintdN(_DecimalN x);
long long int llrintdNx(_DecimalNx x);
5
    _FloatN roundfN(_FloatN x);
    _FloatNx roundfNx(_FloatNx x);
    _DecimalN rounddN(_DecimalN x);
    _DecimalNx rounddNx(_DecimalNx x);
10
long int lroundfN(_FloatN x);
long int lroundfNx(_FloatNx x);
long int lrounddN(_DecimalN x);
long int lrounddNx(_DecimalNx x);
15
long long int llroundfN(_FloatN x);
long long int llroundfNx(_FloatNx x);
long long int llrounddN(_DecimalN x);
long long int llrounddNx(_DecimalNx x);
20
    _FloatN roundevenfN(_FloatN x);
    _FloatNx roundevenfNx(_FloatNx x);
    _DecimalN roundevendN(_DecimalN x);
    _DecimalNx roundevendNx(_DecimalNx x);
25
    _FloatN truncfN(_FloatN x);
    _FloatNx truncfNx(_FloatNx x);
    _DecimalN truncdN(_DecimalN x);
    _DecimalNx truncdNx(_DecimalNx x);
30
intmax_t fromfpfN(_FloatN x, int round, unsigned int width);
intmax_t fromfpfNx(_FloatNx x, int round, unsigned int width);
intmax_t fromfpdN(_DecimalN x, int round, unsigned int width);
intmax_t fromfpdNx(_DecimalNx x, int round,
35
    unsigned int width);
uintmax_t ufromfpfN(_FloatN x, int round, unsigned int width);
uintmax_t ufromfpfNx(_FloatNx x, int round,
    unsigned int width);
uintmax_t ufromfpdN(_DecimalN x, int round,
40
    unsigned int width);
uintmax_t ufromfpdNx(_DecimalNx x, int round,
    unsigned int width);

intmax_t fromfpxfN(_FloatN x, int round, unsigned int width);
45
intmax_t fromfpxfNx(_FloatNx x, int round, unsigned int width);
intmax_t fromfpxdN(_DecimalN x, int round, unsigned int width);
intmax_t fromfpxdNx(_DecimalNx x, int round,
    unsigned int width);
uintmax_t ufromfpxfN(_FloatN x, int round, unsigned int width);
50
uintmax_t ufromfpxfNx(_FloatNx x, int round,
    unsigned int width);
uintmax_t ufromfpxdN(_DecimalN x, int round,
    unsigned int width);
uintmax_t ufromfpxdNx(_DecimalNx x, int round,
55
    unsigned int width);

```

### 7.12.10 Remainder functions

```
5  _FloatN fmodfN(_FloatN x,_FloatN y);
   _FloatNx fmodfNx(_FloatNx x,_FloatNx y);
   _DecimalN fmoddN(_DecimalN x,_DecimalN y);
   _DecimalNx fmoddNx(_DecimalNx x,_DecimalNx y);

   _FloatN remainderfN(_FloatN x,_FloatN y);
   _FloatNx remainderfNx(_FloatNx x,_FloatNx y);
10  _DecimalN remainderdN(_DecimalN x,_DecimalN y);
   _DecimalNx remainderdNx(_DecimalNx x,_DecimalNx y);

   _FloatN remquofN(_FloatN x,_FloatN y, int *quo);
   _FloatNx remquofNx(_FloatNx x,_FloatNx y, int *quo);
```

### 7.12.11 Manipulation functions

```
15  _FloatN copysignfN(_FloatN x,_FloatN y);
   _FloatNx copysignfNx(_FloatNx x,_FloatNx y);
   _DecimalN copysignndN(_DecimalN x,_DecimalN y);
   _DecimalNx copysignndNx(_DecimalNx x,_DecimalNx y);

20  _FloatN nanfN(const char *tagp);
   _FloatNx nanfNx(const char *tagp);
   _DecimalN nandN(const char *tagp);
   _DecimalNx nandNx(const char *tagp);

25  _FloatN nextafterfN(_FloatN x,_FloatN y);
   _FloatNx nextafterfNx(_FloatNx x,_FloatNx y);
   _DecimalN nextafterdN(_DecimalN x,_DecimalN y);
   _DecimalNx nextafterdNx(_DecimalNx x,_DecimalNx y);

30  _FloatN nexttupfN(_FloatN x);
   _FloatNx nexttupfNx(_FloatNx x);
   _DecimalN nexttupdN(_DecimalN x);
   _DecimalNx nexttupdNx(_DecimalNx x);

35  _FloatN nextdownfN(_FloatN x);
   _FloatNx nextdownfNx(_FloatNx x);
   _DecimalN nextdownndN(_DecimalN x);
   _DecimalNx nextdownndNx(_DecimalNx x);

40  int canonicalizefN(_FloatN * cx, const _FloatN * x);
   int canonicalizefNx(_FloatNx * cx, const _FloatNx * x);
   int canonicalizedN(_DecimalN * cx, const _DecimalN * x);
   int canonicalizedNx(_DecimalNx * cx, const _DecimalNx * x);

45  _DecimalN quantizedN(_DecimalN x,_DecimalN y);
   _DecimalNx quantizedNx(_DecimalNx x,_DecimalNx y);

   _Bool samequantumdN(_DecimalN x,_DecimalN y);
   _Bool samequantumdNx(_DecimalNx x,_DecimalNx y);

50  _DecimalN quantumdN(_DecimalN x);
   _DecimalNx quantumdNx(_DecimalNx x);
```

```

long long int llquantexpdN(_DecimalN x);
long long int llquantexpdNx(_DecimalNx x);

5 void encodedecdN(unsigned char * restrict encptr,
  const _DecimalN * restrict xptr);
void decodedecdN(_DecimalN * restrict xptr,
  const unsigned char * restrict encptr);
void encodebindN(unsigned char * restrict encptr,
10 const _DecimalN * restrict xptr);
void decodebindN(_DecimalN * restrict xptr,
  const unsigned char * restrict encptr);

```

#### 7.12.12 Maximum, minimum, and positive difference functions

```

15 _FloatN fdimfN(_FloatN x, _FloatN y);
_FloatNx fdimfNx(_FloatNx x, _FloatNx y);
_DecimalN fdimdN(_DecimalN x, _DecimalN y);
_DecimalNx fdimdNx(_DecimalNx x, _DecimalNx y);

20 _FloatN fmaxfN(_FloatN x, _FloatN y);
_FloatNx fmaxfNx(_FloatNx x, _FloatNx y);
_DecimalN fmaxdN(_DecimalN x, _DecimalN y);
_DecimalNx fmaxdNx(_DecimalNx x, _DecimalNx y);

25 _FloatN fminfN(_FloatN x, _FloatN y);
_FloatNx fminfNx(_FloatNx x, _FloatNx y);
_DecimalN fmindN(_DecimalN x, _DecimalN y);
_DecimalNx fmindNx(_DecimalNx x, _DecimalNx y);

30 _FloatN fmaxmagfN(_FloatN x, _FloatN y);
_FloatNx fmaxmagfNx(_FloatNx x, _FloatNx y);
_DecimalN fmaxmagdN(_DecimalN x, _DecimalN y);
_DecimalNx fmaxmagdNx(_DecimalNx x, _DecimalNx y);

35 _FloatN fminmagfN(_FloatN x, _FloatN y);
_FloatNx fminmagfNx(_FloatNx x, _FloatNx y);
_DecimalN fminmagdN(_DecimalN x, _DecimalN y);
_DecimalNx fminmagdNx(_DecimalNx x, _DecimalNx y);

```

#### 7.12.13 Floating multiply-add

```

40 _FloatN fmafN(_FloatN x, _FloatN y, _FloatN z);
_FloatNx fmafNx(_FloatNx x, _FloatNx y, _FloatNx z);
_DecimalN fmadN(_DecimalN x, _DecimalN y, _DecimalN z);
_DecimalNx fmadNx(_DecimalNx x, _DecimalNx y, _DecimalNx z);

```

#### 7.12.14 Functions that round result to narrower type

```

45 _FloatM fMaddfN(_FloatN x, _FloatN y); // M < N
_FloatM fMaddfNx(_FloatNx x, _FloatNx y); // M <= N
_FloatMx fMxaddfN(_FloatN x, _FloatN y); // M < N
_FloatMx fMxaddfNx(_FloatNx x, _FloatNx y); // M < N

```

```

_DecimalM dMaddN(_DecimalN x, _DecimalN y); // M < N
_DecimalM dMaddNx(_DecimalNx x, _DecimalNx y); // M <= N
_DecimalMx dMxaddN(_DecimalN x, _DecimalN y); // M < N
_DecimalMx dMxaddNx(_DecimalNx x, _DecimalNx y); // M < N
5
_FloatM fMsubfN(_FloatN x, _FloatN y); // M < N
_FloatM fMsubfNx(_FloatNx x, _FloatNx y); // M <= N
_FloatMx fMxsubfN(_FloatN x, _FloatN y); // M < N
_FloatMx fMxsubfNx(_FloatNx x, _FloatNx y); // M < N
10
_DecimalM dMsubdN(_DecimalN x, _DecimalN y); // M < N
_DecimalM dMsubdNx(_DecimalNx x, _DecimalNx y); // M <= N
_DecimalMx dMxsubdN(_DecimalN x, _DecimalN y); // M < N
_DecimalMx dMxsubdNx(_DecimalNx x, _DecimalNx y); // M < N

_FloatM fMmulfN(_FloatN x, _FloatN y); // M < N
_FloatM fMmulfNx(_FloatNx x, _FloatNx y); // M <= N
_FloatMx fMxmulfN(_FloatN x, _FloatN y); // M < N
_FloatMx fMxmulfNx(_FloatNx x, _FloatNx y); // M < N
_DecimalM dMmuldN(_DecimalN x, _DecimalN y); // M < N
_DecimalM dMmuldNx(_DecimalNx x, _DecimalNx y); // M <= N
20
_DecimalMx dMxmuldN(_DecimalN x, _DecimalN y); // M < N
_DecimalMx dMxmuldNx(_DecimalNx x, _DecimalNx y); // M < N

_FloatM fMdivfN(_FloatN x, _FloatN y); // M < N
_FloatM fMdivfNx(_FloatNx x, _FloatNx y); // M <= N
_FloatMx fMxdivfN(_FloatN x, _FloatN y); // M < N
_FloatMx fMxdivfNx(_FloatNx x, _FloatNx y); // M < N
_DecimalM dMdivdN(_DecimalN x, _DecimalN y); // M < N
_DecimalM dMdivdNx(_DecimalNx x, _DecimalNx y); // M <= N
30
_DecimalMx dMxdivdN(_DecimalN x, _DecimalN y); // M < N
_DecimalMx dMxdivdNx(_DecimalNx x, _DecimalNx y); // M < N

_FloatM fMfmfN(_FloatN x, _FloatN y, _FloatN z); // M < N
_FloatM fMfmfNx(_FloatNx x, _FloatNx y,
35
_FloatNx z); // M <= N
_FloatMx fMxfmfN(_FloatN x, _FloatN y, _FloatN z); // M < N
_FloatMx fMxfmfNx(_FloatNx x, _FloatNx y,
_FloatNx z); // M < N
_DecimalM dMfmadN(_DecimalN x, _DecimalN y,
40
_DecimalN z); // M < N
_DecimalM dMfmadNx(_DecimalNx x, _DecimalNx y,
_DecimalNx z); // M <= N
_DecimalMx dMxfmadN(_DecimalN x, _DecimalN y,
_DecimalN z); // M < N
45
_DecimalMx dMxfmadNx(_DecimalNx x, _DecimalNx y,
_DecimalNx z); // M < N

_FloatM fMsqrtfN(_FloatN x); // M < N
_FloatM fMsqrtfNx(_FloatNx x); // M <= N
50
_FloatMx fMxsqrtfN(_FloatN x); // M < N
_FloatMx fMxsqrtfNx(_FloatNx x); // M < N

```

```

    _DecimalM dMsqrtdN(_DecimalN x); // M < N
    _DecimalM dMsqrtdNx(_DecimalNx x); // M <= N
    _DecimalMx dMxsqrtD(_DecimalN x); // M < N
    _DecimalMx dMxsqrtDx(_DecimalNx x); // M < N

```

5

#### F.10.12 Total order functions

```

int totalorderfN(const _FloatN *x, const _FloatN *y);
int totalorderfNx(const _FloatNx *x, const _FloatNx *y);
int totalorderdN(const _DecimalN *x, const _DecimalN *y);
int totalorderdNx(const _DecimalNx *x, const _DecimalNx *y);

```

10

```

int totalordermagfN(const _FloatN *x, const _FloatN *y);
int totalordermagfNx(const _FloatNx *x, const _FloatNx *y);
int totalordermagdN(const _DecimalN *x, const _DecimalN *y);
int totalordermagdNx(const _DecimalNx *x, const _DecimalNx *y);

```

15

#### F.10.13 Payload functions

```

_FloatN getpayloadfN(const _FloatN *x);
_FloatNx getpayloadfNx(const _FloatNx *x);
_DecimalN getpayloaddN(const _DecimalN *x);
_DecimalNx getpayloaddNx(const _DecimalNx *x);

```

20

```

int setpayloadfN(_FloatN *res, _FloatN pl);
int setpayloadfNx(_FloatNx *res, _FloatNx pl);
int setpayloaddN(_DecimalN *res, _DecimalN pl);
int setpayloaddNx(_DecimalNx *res, _DecimalNx pl);

```

25

```

int setpayloadsigfN(_FloatN *res, _FloatN pl);
int setpayloadsigfNx(_FloatNx *res, _FloatNx pl);
int setpayloadsigdN(_DecimalN *res, _DecimalN pl);
int setpayloadsigdNx(_DecimalNx *res, _DecimalNx pl);

```

30

[2] The specification of the **frexp** functions (7.12.6.7) applies the functions for binary floating types like those for standard floating types: the exponent is an integral power of 2 and, when applicable, **value** equals  $x \times 2^{\text{exp}}$ .

[3] The specification of the **ldexp** functions (7.12.6.9) applies to the functions for binary floating types like those for standard floating types: they return  $x \times 2^{\text{exp}}$ .

35

[4] The specification of the **logb** functions (7.12.6.17) applies to binary floating types, with  $b = 2$ .

[5] The specification of the **scalbn** and **scalbln** functions (7.12.6.19) applies to binary floating types, with  $b = 2$ .

### X.11.3 Encoding conversion functions

[1] This subclause introduces `<math.h>` functions that, together with the numerical conversion functions for encodings in X.12, support the non-arithmetic interchange formats specified by IEC 60559.

40

[2] Non-arithmetic interchange formats are not associated with floating types. Arrays of element type **unsigned char** are used as parameters for conversions functions, to represent encodings in interchange formats that might be non-arithmetic formats.

45

### X.11.3.1 Encode and decode functions

[1] This subclause specifies functions to map representations in binary floating types to and from encodings in **unsigned char** arrays.

#### X.11.3.1.1 The **encodefN** functions

##### 5 Synopsis

```
[1] #include <math.h>
    void encodefN(unsigned char * restrict encptr,
        const _FloatN * restrict xptr);
```

##### Description

10 [2] The **encodefN** functions convert **\*xptr** into an IEC 60559 binary $N$  encoding and store the resulting encoding as an  $N/8$  element array, with 8 bits per array element, in the object pointed to by **encptr**. The order of bytes in the array is implementation-defined. These functions preserve the value of **\*xptr** and raise no floating-point exceptions. If **\*xptr** is non-canonical, these functions may or may not produce a canonical encoding.

##### 15 Returns

[3] The **encodefN** functions return no value.

#### X.11.3.1.2 The **decodefN** functions

##### Synopsis

```
[1] #include <math.h>
20 void decodefN (_FloatN * restrict xptr,
    const unsigned char * restrict encptr);
```

##### Description

25 [2] The **decodefN** functions interpret the  $N/8$  element array pointed to by **encptr** as an IEC 60559 binary $N$  encoding, with 8 bits per array element. The order of bytes in the array is implementation-defined. These functions convert the given encoding into a representation in the type **\_FloatN**, and store the result in the object pointed to by **xptr**. These functions preserve the encoded value and raise no floating-point exceptions. If the encoding is non-canonical, these functions may or may not produce a canonical representation.

##### Returns

30 [3] The **decodefN** functions return no value.

#### X.11.3.2 Encoding-to-encoding conversion functions

[1] An implementation shall declare an **fMencfN** function for each  $M$  and  $N$  equal to the width of a supported IEC 60559 arithmetic or non-arithmetic binary interchange format. An implementation shall provide both **dMencdecN** and **dMencbindN** functions for each  $M$  and  $N$  equal to the width of a supported IEC 60559 arithmetic or non-arithmetic decimal interchange format.

35



### X.11.3.2.1 The `fMencfN` functions

#### Synopsis

```
[1] #include <math.h>
    void fMencfN(unsigned char * restrict encMptr,
5         const unsigned char * restrict encNptr);
```

#### Description

[2] These functions convert between IEC 60559 binary interchange formats. These functions interpret the  $N/8$  element array pointed to by `encNptr` as an encoding of width  $N$  bits. They convert the encoding to an encoding of width  $M$  bits and store the resulting encoding as an  $M/8$  element array in the object pointed to by `encMptr`. The conversion rounds and raises floating-point exceptions as specified in IEC 60559. The order of bytes in the arrays is implementation-defined.

#### Returns

[3] These functions return no value.

### X.11.3.2.2 The `dMencdecN` and `dMencbindN` functions

#### 15 Synopsis

```
[1] #include <math.h>
    void dMencdecN(unsigned char * restrict encMptr,
        const unsigned char * restrict encNptr);
    void dMencbindN(unsigned char * restrict encMptr,
20         const unsigned char * restrict encNptr);
```

#### Description

[2] These functions convert between IEC 60559 decimal interchange formats that use the same encoding scheme. The `dMencdecN` functions convert between formats using the encoding scheme based on decimal encoding of the significand. The `dMencbindN` functions convert between formats using the encoding scheme based on binary encoding of the significand. These functions interpret the  $N/8$  element array pointed to by `encNptr` as an encoding of width  $N$  bits. They convert the encoding to an encoding of width  $M$  bits and store the resulting encoding as an  $M/8$  element array in the object pointed to by `encMptr`. The conversion rounds and raises floating-point exceptions as specified in IEC 60559. The order of bytes in the arrays is implementation-defined.

#### 30 Returns

[3] These functions return no value.

## X.12 Numeric conversion functions in `<stdlib.h>`

[1] This clause expands the specification of numeric conversion functions in `<stdlib.h>` (7.22.1) to also include conversions of strings from and to interchange and extended floating types. The conversions from floating are provided by functions analogous to the `strfromd` function. The conversions to floating are provided by functions analogous to the `strtod` function.

[2] This clause also specifies functions to convert strings from and to IEC 60559 interchange format encodings.

[3] For each interchange or extended floating type that the implementation provides, `<stdlib.h>` shall declare the associated functions specified below in X.12.1 and X.12.2. Conversely, for each such type that the implementation does not provide, `<stdlib.h>` shall not declare the associated functions.

5 [4] For each IEC 60559 arithmetic or non-arithmetic format that the implementation supports, `<stdlib.h>` shall declare the associated functions specified below in X.12.3 and X.12.4. Conversely, for each such format that the implementation does not provide, `<stdlib.h>` shall not declare the associated functions.

### X.12.1 String from floating

10 [1] This subclause expands 7.22.1.3 and 7.22.1.4 to also include functions for the interchange and extended floating types. It adds to the synopsis in 7.22.1.3 the prototypes

```
int strfromfN(char * restrict s, size_t n,  
              const char * restrict format, _FloatN fp);  
15 int strfromfNx(char * restrict s, size_t n,  
                const char * restrict format, _FloatNx fp);
```

It encompasses the prototypes in 7.22.1.4 by replacing them with

```
int strfromdN(char * restrict s, size_t n,  
              const char * restrict format, _DecimalN fp);  
20 int strfromdNx(char * restrict s, size_t n,  
                const char * restrict format, _DecimalNx fp);
```

[2] The descriptions and returns for the added functions are analogous to the ones in 7.22.1.3 and 7.22.1.4.

### X.12.2 String to floating

25 [1] This subclause expands 7.22.1.5 and 7.22.1.6 to also include functions for the interchange and extended floating types. It adds to the synopsis in 7.22.1.5 the prototypes

```
_FloatN strttofN(const char * restrict nptr,  
                 char ** restrict endptr);  
_FloatNx strttofNx(const char * restrict nptr,  
                   char ** restrict endptr);
```

30 It encompasses the prototypes in 7.22.1.6 by replacing them with

```
_DecimalN strtodN(const char * restrict nptr,  
                  char ** restrict endptr);  
_DecimalNx strtodNx(const char * restrict nptr,  
                    char ** restrict endptr);
```

35 [2] The descriptions and returns for the added functions are analogous to the ones in 7.22.1.5 and 7.22.1.6.

### X.12.3 String from encoding

[1] An implementation shall declare the `strfromencfN` function for each  $N$  equal to the width of a supported IEC 60559 arithmetic or non-arithmetic binary interchange format. An implementation shall  
40 declare both the `strfromencdecN` and `strfromencbindN` functions for each  $N$  equal to the width of a supported IEC 60559 arithmetic or non-arithmetic decimal interchange format.

### X.12.3.1 The `strfromencfN` functions

#### Synopsis

```
[1] #include <stdlib.h>
    int strfromencfN(char * restrict s, size_t n,
5      const char * restrict format,
      const unsigned char * restrict encptr);
```

#### Description

[2] The `strfromencfN` functions are similar to the `strfromfN` functions, except the input is the value of the  $N/8$  element array pointed to by `encptr`, interpreted as an IEC 60559 binary $N$  encoding.  
10 The order of bytes in the arrays is implementation-defined.

#### Returns

[3] The `strfromencfN` functions return the same values as corresponding `strfromfN` functions.

### X.12.3.2 The `strfromencdecN` and `strfromencbindN` functions

#### Synopsis

```
[1] #include <stdlib.h>
    int strfromencdecN(char * restrict s, size_t n,
      const char * restrict format,
      const unsigned char * restrict encptr);
    int strfromencbindNx(char * restrict s, size_t n,
20     const char * restrict format,
      const unsigned char * restrict encptr);
```

#### Description

[2] The `strfromencdecN` functions are similar to the `strfromdN` functions except the input is the value of the  $N/8$  element array pointed to by `encptr`, interpreted as an IEC 60559 decimal $N$  encoding  
25 in the coding scheme based on decimal encoding of the significand. The `strfromencbindN` functions are similar to the `strfromdN` functions except the input is the value of the  $N/8$  element array pointed to by `encptr`, interpreted as an IEC 60559 decimal $N$  encoding in the coding scheme based on binary encoding of the significand. The order of bytes in the arrays is implementation-defined.

#### Returns

[3] The `strfromencdecN` and `strfromencbindN` functions return the same values as corresponding `strfromdN` functions.

## X.12.4 String to encoding

### X.12.4.1 The `strtoencfN` functions

#### Synopsis

```
[1] #include <stdlib.h>
    void strtoencfN(unsigned char * restrict encptr,
      const char * restrict nptr, char ** restrict endptr);
```

## Description

[2] The `strtoencfN` functions are similar to the `strtofN` functions, except they store an IEC 60559 encoding of the result as an  $N/8$  element array in the object pointed to by `encptr`. The order of bytes in the arrays is implementation-defined.

## 5 Returns

[3] These functions return no value.

### X.12.4.2 The `strtoencdecN` and `strtoencbindN` functions

#### Synopsis

```
[1] #include <stdlib.h>
10 void strtoencdecN(unsigned char * restrict encptr,
    const char * restrict nptr, char ** restrict endptr);
void strtoencbindN(unsigned char * restrict encptr,
    const char * restrict nptr, char ** restrict endptr);
```

## Description

15 [2] The `strtoencdecN` and `strtoencbindN` functions are similar to the `strtodN` functions, except they store an IEC 60559 encoding of the result as an  $N/8$  element array in the object pointed to by `encptr`. The `strtoencdecN` functions produce an encoding in the encoding scheme based on decimal encoding of the significand. The `strtoencbindN` functions produce an encoding in the encoding scheme based on binary encoding of the significand. The order of bytes in the arrays is  
20 implementation-defined.

## Returns

[3] These functions return no value.

### X.13 Type-generic macros `<tgmath.h>`

25 [1] This clause enhances the specification of type-generic macros in `<tgmath.h>` (7.25) to apply to interchange and extended floating types, as well as standard floating types.

[2] If arguments for generic parameters of a type-generic macro are such that some argument has a corresponding real type that is a standard floating type or a floating type of radix 2 and another argument is of decimal floating type, the behavior is undefined.

30 [3] Use of the macro `carg`, `cimag`, `conj`, `cproj`, or `creal` with any argument of standard floating type, floating type of radix 2, or complex type, invokes a complex function. Use of the macro with an argument of a decimal floating type results in undefined behavior.

[4] The functions that round result to a narrower type have type-generic macros whose names are obtained by omitting any suffix from the function names. Thus, the macros with `f` or `d` prefix are (as in 7.25):

```
35      fadd          fmul          ffma
      dadd          dmul          dfma
      fsub          fdiv          fsqrt
      dsub          ddiv          dsqrt
```

and the macros with `fM`, `fMx`, `dM`, or `dMx` prefix are:

	<code>fMadd</code>	<code>fMxm<sub>l</sub></code>	<code>dMfma</code>
	<code>fMsub</code>	<code>fMxd<sub>iv</sub></code>	<code>dMsq<sub>rt</sub></code>
	<code>fMmul</code>	<code>fMxf<sub>ma</sub></code>	<code>dMxadd</code>
	<code>fMdiv</code>	<code>fMxsq<sub>rt</sub></code>	<code>dMxsub</code>
5	<code>fMfma</code>	<code>dMadd</code>	<code>dMxm<sub>l</sub></code>
	<code>fMsq<sub>rt</sub></code>	<code>dMsub</code>	<code>dMxd<sub>iv</sub></code>
	<code>fMxadd</code>	<code>dMmul</code>	<code>dMxf<sub>ma</sub></code>
	<code>fMxsub</code>	<code>dMdiv</code>	<code>dMxsq<sub>rt</sub></code>

10 All arguments are generic. If any argument is not real, use of the macro results in undefined behavior. If the macro prefix is `f` or `d`, use of an argument of interchange or extended floating type results in undefined behavior. If the macro prefix is `fM`, or `fMx`, use of an argument of standard or decimal floating type results in undefined behavior. If the macro prefix is `dM` or `dMx`, use of an argument of standard or binary floating type results in undefined behavior. The function invoked is determined as follows:

- 15 — Arguments that have integer type are regarded as having type `double` if the macro prefix is `f` or `d`, as having type `_Float64` if the macro prefix is `fM` or `fMx`, and as having type `_Decimal64` if the macro prefix is `dM` or `dMx`.
- If the function has exactly one generic parameter, the type determined is the type of the argument.
- 20 — If the function has exactly two generic parameters, the type determined is the type determined by the usual arithmetic conversions (X.4.2) applied to the arguments.
- If the function has three generic parameters, the type determined is the type determined by applying the usual arithmetic conversions (X.4.2) twice, first to the first two arguments, then to that result type and the third argument.
- 25 — If no function with the given prefix has the parameter type determined above, the parameter type is determined from the prefix as follows:

<code>f</code>	<code>double</code>
<code>d</code>	<code>long double</code>
<code>fM</code>	<code>_Float<sub>N</sub></code> for minimum $N > M$ if supported, else <code>_Float<sub>M</sub></code>
<code>fMx</code>	<code>_Float<sub>N</sub>x</code> for minimum $N > M$ if supported, else <code>_Float<sub>N</sub></code> for minimum $N > M$
<code>dM</code>	<code>_Decimal<sub>N</sub></code> for minimum $N > M$ if supported, else <code>_Decimal<sub>M</sub></code>
<code>dMx</code>	<code>_Decimal<sub>N</sub>x</code> for minimum $N > M$ if supported, else <code>_Decimal<sub>N</sub></code> for minimum $N > M$

EXAMPLE With the declarations

```

30 #define __STDC_WANT_IEC_60559_TYPES_EXT__
#include <tgmath.h>
int n;
double d;
long double ld;
double complex dc;
35 _Float32x f32x;
_Float64 f64;
_Float64x f64x;
_Float128 f128;
_Float64x complex f64xc;

```

functions invoked by use of type-generic macros are shown in the following table:

macro use	invokes	
5	<b>cos(f64xc)</b> <b>pow(dc, f128)</b> <b>fmax(f64, d)</b> <b>fmax(d, f32x)</b>	<b>ccosf64x(f64xc)</b> <b>cpowf128(dc, f128)</b> <b>fmaxf64(f64, d)</b> <b>fmax(d, f32x)</b> , the function, if the set of values of <b>_Float32x</b> is a subset of (or equivalent to) the set of values of <b>double</b> , or
10	<b>fmaxf32x(d, f32x)</b> , if the set of values of <b>double</b> is a proper subset of the set of values of <b>_Float32x</b> , or undefined, if neither of the sets of values of <b>double</b> and <b>_Float32x</b> is a subset of the other (and the sets are not equivalent)	<b>fmaxf32x(d, f32x)</b> , if the set of values of <b>double</b> is a proper subset of the set of values of <b>_Float32x</b> , or undefined, if neither of the sets of values of <b>double</b> and <b>_Float32x</b> is a subset of the other (and the sets are not equivalent)
15	<b>pow(f32x, n)</b> <b>fsub(d, ld)</b> <b>f32add(f64x, f64)</b> <b>f32xsqrt(n)</b> <b>f32mul(f128, f32x)</b>	<b>powf32x(f32x, n)</b> <b>fsubl</b> <b>f32addf64x</b> <b>f32xsqrtf64</b> <b>f32mulf128</b> if <b>_Float128</b> is at least as wide as <b>_Float32x</b> , or
20	<b>f32fma(f32x, n, f32x)</b>	<b>f32mulf32x</b> if <b>_Float32x</b> is wider than <b>_Float128</b> <b>f32fmaf64</b> if <b>_Float64</b> is at least as wide as <b>_Float32x</b> , or <b>f32fmaf32x</b> if <b>_Float32x</b> is wider than <b>_Float64</b>
25	<b>ddiv(ld, f128)</b> <b>f32fma(f64, d, f64)</b> <b>fmul(dc, d)</b> <b>f32add(f32, f32)</b> <b>f32xsqrt(f32)</b>	undefined undefined undefined <b>f32addf32x(f32, f32)</b> <b>f32xsqrtf64x(f32)</b>
30	<b>f64div(f32x, f32x)</b>	<b>f64divf128be(f32x, f32x)</b>