

**ISO/IEC JTC 1/SC 22/WG 14 N1789**

Date: yyyy-mm-dd

Reference number of document: **ISO/IEC TS 18661-3**

5 Committee identification: ISO/IEC JTC 1/SC 22/WG 14  
Secretariat: ANSI

**Information Technology — Programming languages, their environments,  
and system software interfaces — Floating-point extensions for C —  
Part 3: Interchange and extended types**

10 *Technologies de l'information — Langages de programmation, leurs environnements et interfaces du logiciel système — Extensions à virgule flottante pour C — Partie 3: Types d'échange et prolongée*

**Warning**

15 This document is not an ISO International Standard. It is distributed for review and comment. It is subject to change without notice and may not be referred to as an International Standard.

Recipients of this draft are invited to submit, with their comments, notification of any relevant patent rights of which they are aware and to provide supporting documentation.

### Copyright notice

5 This ISO document is a working draft or committee draft and is copyright-protected by ISO. While the reproduction of working drafts or committee drafts in any form for use by participants in the ISO standards development process is permitted without prior permission from ISO, neither this document nor any extract from it may be reproduced, stored or transmitted in any form for any other purpose without prior written permission from ISO.

Requests for permission to reproduce this document for the purpose of selling it should be addressed as shown below or to ISO's member body in the country of the requester:

10 *ISO copyright office*  
*Case postale 56 CH-1211 Geneva 20*  
*Tel. +41 22 749 01 11*  
*Fax + 41 22 749 09 47*  
*E-mail [copyright@iso.org](mailto:copyright@iso.org)*  
*Web [www.iso.org](http://www.iso.org)*

15 Reproduction for sales purposes may be subject to royalty payments or a licensing agreement.

Violators may be prosecuted.

## Contents

	Page
Introduction .....	v
Background .....	v
IEC 60559 floating-point standard .....	v
5    C support for IEC 60559.....	vi
Purpose .....	vii
Additional background on formats.....	vii
1  Scope .....	1
2  Conformance .....	1
10 3 Normative references .....	1
4  Terms and definitions.....	1
5  C standard conformance.....	2
5.1  Freestanding implementations.....	2
5.2  Predefined macros.....	2
15  5.3  Standard headers.....	2
6  Types.....	7
7  Characteristics .....	12
8  Conversions .....	17
9  Constants .....	18
20 10 Expressions.....	19
11 Non-arithmetic interchange formats .....	20
12 Mathematics <math.h>.....	21
12.1  Macros .....	21
12.2  Floating-point environment .....	24
25  12.3  Functions.....	26
12.4  Encoding conversion functions .....	36
13 Numeric conversion functions in <stdlib.h> .....	37
14 Complex arithmetic <complex.h>.....	41
15 Type-generic macros <tgmath.h> .....	43
30 Bibliography .....	46

## Foreword

ISO (the International Organization for Standardization) is a worldwide federation of national standards bodies (ISO member bodies). The work of preparing International Standards is normally carried out through ISO technical committees. Each member body interested in a subject for which a technical committee has been established has the right to be represented on that committee. International organizations, governmental and non-governmental, in liaison with ISO, also take part in the work. ISO collaborates closely with the International Electrotechnical Commission (IEC) on all matters of electrotechnical standardization.

International Standards are drafted in accordance with the rules given in the ISO/IEC Directives, Part 2.

The main task of technical committees is to prepare International Standards. Draft International Standards adopted by the technical committees are circulated to the member bodies for voting. Publication as an International Standard requires approval by at least 75 % of the member bodies casting a vote.

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO shall not be held responsible for identifying any or all such patent rights.

ISO/IEC TS 18661 was prepared by Technical Committee ISO JTC 1, *Information Technology*, Subcommittee SC 22, *Programming languages, their environments, and system software interfaces*.

ISO/IEC TS 18661 consists of the following parts, under the general title *Floating-point extensions for C*:

- *Part 1: Binary floating-point arithmetic*
- *Part 2: Decimal floating-point arithmetic*
- *Part 3: Interchange and extended types*
- *Part 4: Supplementary functions*
- *Part 5: Supplementary attributes*

Part 1 updates ISO/IEC 9899:2011 (*Information technology — Programming languages, their environments and system software interfaces — Programming Language C*), Annex F in particular, to support all required features of ISO/IEC/IEEE 60559:2011 (*Information technology — Microprocessor Systems — Floating-point arithmetic*).

Part 2 supersedes ISO/IEC TR 24732:2009 (*Information technology – Programming languages, their environments and system software interfaces – Extension for the programming language C to support decimal floating-point arithmetic*).

Parts 3-5 specify extensions to ISO/IEC 9899:2011 for features recommended in ISO/IEC/IEEE 60559:2011.

## Introduction

### Background

#### IEC 60559 floating-point standard

5 The IEEE 754-1985 standard for binary floating-point arithmetic was motivated by an expanding diversity in floating-point data representation and arithmetic, which made writing robust programs, debugging, and moving programs between systems exceedingly difficult. Now the great majority of systems provide data formats and arithmetic operations according to this standard. The IEC 60559:1989 international standard was equivalent to the IEEE 754-1985 standard. Its stated goals were:

- 10 1 Facilitate movement of existing programs from diverse computers to those that adhere to this standard.
- 2 Enhance the capabilities and safety available to programmers who, though not expert in numerical methods, may well be attempting to produce numerically sophisticated programs. However, we recognize that utility and safety are sometimes antagonists.
- 15 3 Encourage experts to develop and distribute robust and efficient numerical programs that are portable, by way of minor editing and recompilation, onto any computer that conforms to this standard and possesses adequate capacity. When restricted to a declared subset of the standard, these programs should produce identical results on all conforming systems.
- 4 Provide direct support for
  - a. Execution-time diagnosis of anomalies
  - 20 b. Smoother handling of exceptions
  - c. Interval arithmetic at a reasonable cost
- 5 Provide for development of
  - a. Standard elementary functions such as exp and cos
  - b. Very high precision (multiword) arithmetic
  - 25 c. Coupling of numerical and symbolic algebraic computation
- 6 Enable rather than preclude further refinements and extensions.

To these ends, the standard specified a floating-point model comprising:

*formats* – for binary floating-point data, including representations for Not-a-Number (NaN) and signed infinities and zeros

30 *operations* – basic arithmetic operations (addition, multiplication, etc.) on the format data to compose a well-defined, closed arithmetic system; also specified conversions between floating-point formats and decimal character sequences, and a few auxiliary operations

*context* – status flags for detecting exceptional conditions (invalid operation, division by zero, overflow, underflow, and inexact) and controls for choosing different rounding methods

35 The IEC 60559:2011 international standard is equivalent to the IEEE 754-2008 standard for floating-point arithmetic, which is a major revision to IEEE 754-1985.

The revised standard specifies more formats, including decimal as well as binary. It adds a 128-bit binary format to its basic formats. It defines extended formats for all of its basic formats. It specifies data interchange

formats (which may or may not be arithmetic), including a 16-bit binary format and an unbounded tower of wider formats. To conform to the floating-point standard, an implementation must provide at least one of the basic formats, along with the required operations.

5 The revised standard specifies more operations. New requirements include – among others – arithmetic operations that round their result to a narrower format than the operands (with just one rounding), more conversions with integer types, more classifications and comparisons, and more operations for managing flags and modes. New recommendations include an extensive set of mathematical functions and seven reduction functions for sums and scaled products.

10 The revised standard places more emphasis on reproducible results, which is reflected in its standardization of more operations. For the most part, behaviors are completely specified. The standard requires conversions between floating-point formats and decimal character sequences to be correctly rounded for at least three more decimal digits than is required to distinguish all numbers in the widest supported binary format; it fully specifies conversions involving any number of decimal digits. It recommends that transcendental functions be correctly rounded.

15 The revised standard requires a way to specify a constant rounding direction for a static portion of code, with details left to programming language standards. This feature potentially allows rounding control without incurring the overhead of runtime access to a global (or thread) rounding mode.

20 Other features recommended by the revised standard include alternate methods for exception handling, controls for expression evaluation (allowing or disallowing various optimizations), support for fully reproducible results, and support for program debugging.

25 The revised standard, like its predecessor, defines its model of floating-point arithmetic in the abstract. It neither defines the way in which operations are expressed (which might vary depending on the computer language or other interface being used), nor does it define the concrete representation (specific layout in storage, or in a processor's register, for example) of data or context, except that it does define specific encodings that are to be used for data that may be exchanged between different implementations that conform to the specification.

30 IEC 60559 does not include bindings of its floating-point model for particular programming languages. However, the revised standard does include guidance for programming language standards, in recognition of the fact that features of the floating-point standard, even if well supported in the hardware, are not available to users unless the programming language provides a commensurate level of support. The implementation's combination of both hardware and software determines conformance to the floating-point standard.

### **C support for IEC 60559**

35 The C standard specifies floating-point arithmetic using an abstract model. The representation of a floating-point number is specified in an abstract form where the constituent components (sign, exponent, significand) of the representation are defined but not the internals of these components. In particular, the exponent range, significand size, and the base (or radix) are implementation-defined. This allows flexibility for an implementation to take advantage of its underlying hardware architecture. Furthermore, certain behaviors of operations are also implementation-defined, for example in the area of handling of special numbers and in exceptions.

40 The reason for this approach is historical. At the time when C was first standardized, before the floating-point standard was established, there were various hardware implementations of floating-point arithmetic in common use. Specifying the exact details of a representation would have made most of the existing implementations at the time not conforming.

45 Beginning with ISO/IEC 9899:1999 (C99), C has included an optional second level of specification for implementations supporting the floating-point standard. C99, in conditionally normative Annex F, introduced nearly complete support for the IEC 60559:1989 standard for binary floating-point arithmetic. Also, C99's informative Annex G offered a specification of complex arithmetic that is compatible with IEC 60559:1989.

ISO/IEC 9899:2011 (C11) includes refinements to the C99 floating-point specification, though is still based on IEC 60559:1989. C11 upgrades Annex G from “informative” to “conditionally normative”.

5 ISO/IEC Technical Report 24732:2009 introduced partial C support for the decimal floating-point arithmetic in IEC 60559:2011. TR 24732, for which technical content was completed while IEEE 754-2008 was still in the later stages of development, specifies decimal types based on IEC 60559:2011 decimal formats, though it does not include all of the operations required by IEC 60559:2011.

## Purpose

10 The purpose of this Technical Specification is to provide a C language binding for IEC 60559:2011, based on the C11 standard, that delivers the goals of IEC 60559 to users and is feasible to implement. It is organized into five Parts.

Part 1 provides changes to C11 that cover all the requirements, plus some basic recommendations, of IEC 60559:2011 for binary floating-point arithmetic. C implementations intending to support IEC 60559:2011 are expected to conform to conditionally normative Annex F as enhanced by the changes in Part 1.

15 Part 2 enhances TR 24732 to cover all the requirements, plus some basic recommendations, of IEC 60559:2011 for decimal floating-point arithmetic. C implementations intending to provide an extension for decimal floating-point arithmetic supporting IEC 60559:2011 are expected to conform to Part 2.

Part 3 (Interchange and extended types), Part 4 (Supplementary functions), and Part 5 (Supplementary attributes) cover recommended features of IEC 60559:2011. C implementations intending to provide extensions for these features are expected to conform to the corresponding Parts.

## 20 Additional background on formats

The 2011 revision of the ISO/IEC 60559 standard for floating-point arithmetic introduces a variety of new formats, both fixed and extendable. The new fixed formats include

- a 128-bit basic binary format (the 32 and 64 bit basic binary formats are carried over from ISO/IEC 60559:1989)
- 25 — 64 and 128 bit basic decimal formats
- interchange formats, whose precision and range are determined by the width  $k$ , where
  - for binary,  $k = 16, 32, 64$ , and  $k \geq 128$  and a multiple of 32, and
  - for decimal,  $k \geq 32$  and a multiple of 32
- 30 — extended formats, for each basic format, with minimum range and precision specified

Thus IEC 60559 defines five basic formats - binary32, binary64, binary128, decimal64, and decimal128 - and five corresponding extended formats, each with somewhat more precision and range than the basic format it extends. IEC 60559 defines an unlimited number of interchange formats, which include the basic formats.

35 Interchange formats may or may not be supported as arithmetic formats. If not, they may be used for the interchange of floating-point data but not for arithmetic computation. IEC 60559 provides conversions between non-arithmetic interchange formats and arithmetic formats which can be used for computation.

40 Extended formats are intended for intermediate computation, not input or output data. The extra precision often allows the computation of extended results which when converted to a narrower output format differ from the ideal results by little more than a unit in the last place. Also, the extra range often avoids any intermediate overflow or underflow that might occur if the computation were done in the format of the data. The essential property of extended formats is their sufficient extra widths, not their specific widths. Extended formats for any given basic format may vary among implementations.

Extendable formats, which provide user control over range and precision, are not covered in Technical Specification 18661.

45 The 32 and 64 bit binary formats are supported in C by types `float` and `double`. If a C implementation defines the macro `__STDC_IEC_60559_BFP__` (see Part 1 of Technical Specification 18661) signifying that it

supports Annex F of the C Standard, then its `float` and `double` formats must be IEC 60559 binary32 and binary64.

5 Part 2 of Technical Specification 18661 defines types `_Decimal132`, `_Decimal164`, and `_Decimal128` with IEC 60559 formats `decimal32`, `decimal64`, and `decimal128`. Although IEC 60559 does not require arithmetic support (other than conversions) for its `decimal32` interchange format, Part 2 of Technical Specification 18661 has full arithmetic and library support for `_Decimal132`, just like for `_Decimal164` and `_Decimal128`.

10 The C Standard provides just three standard floating types (`float`, `double`, and `long double`) that are required of all implementations. Annex F of the C Standard requires the standard floating types to be binary. The `long double` type must be at least as wide as `double`, but C does not further specify details of its format, even in Annex F.

Part 3 of Technical Specification 18661, this document, provides nomenclatures for types with IEC 60559 arithmetic interchange formats and extended formats. The nomenclatures allow portable use of the formats as envisioned in IEC 60559. This document covers these aspects of the types:

- names
- 15 — characteristics
- conversions
- constants
- function suffixes
- character sequence conversion interfaces

20 This specification includes interchange and extended nomenclatures for `formats` that, in some cases, already have C nomenclatures. For example, `types` with the IEC 60559 `double` format may include `double`, `_Float64` (the type for the binary64 interchange format), and maybe `_Float32x` (the type for the binary32-extended format). This redundancy is intended to support the different programming models appropriate for the types with arithmetic interchange formats and extended formats and C standard floating types.

25 This document also supports the IEC 60559 non-arithmetic interchange formats with functions that convert among encodings and between encodings and character sequences, for all interchange formats.



# Information Technology — Programming languages, their environments, and system software interfaces — Floating-point extensions for C — Part 3: Interchange and extended types

## 1 Scope

- 5 This document, Part 3 of Technical Specification 18661, 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.

## 2 Conformance

An implementation conforms to Part 3 of Technical Specification 18661 if

- 10 a) It meets the requirements for a conforming implementation of C11 with all the changes to C11 as specified in Parts 1-3 of Technical Specification 18661;
- b) It conforms to Part 1 or Part 2 (or both) of Technical Specification 18661; and
- 15 c) It defines `__STDC_IEC_60559_TYPES__` to 201~~ymm~~L.

## 3 Normative references

The following referenced documents are indispensable for the application of this document. Only the editions cited apply.

- 20 ISO/IEC 9899:2011, *Information technology — Programming languages, their environments and system software interfaces — Programming Language C*

ISO/IEC 9899:2011/Cor.1:2012, *Technical Corrigendum 1*

ISO/IEC/IEEE 60559:2011, *Information technology — Microprocessor Systems — Floating-point arithmetic* (with identical content to IEEE 754-2008, *IEEE Standard for Floating-Point Arithmetic*. The Institute of Electrical and Electronic Engineers, Inc., New York, 2008)

- 25 ISO/IEC 18661-1:yyyy, *Information Technology — Programming languages, their environments, and system software interfaces — Floating-point extensions for C — Part 1: Binary floating-point arithmetic*

ISO/IEC 18661-2:yyyy, *Information Technology — Programming languages, their environments, and system software interfaces — Floating-point extensions for C — Part 2: Decimal floating-point arithmetic*

- 30 Changes specified in Part 3 of Technical Specification 18661 are relative to ISO/IEC 9899:2011, including *Technical Corrigendum 1* (ISO/IEC 9899:2011/Cor. 1:2012), together with the changes from Parts 1 and 2 of Technical Specification 18661.

## 4 Terms and definitions

For the purposes of this document, the terms and definitions given in ISO/IEC 9899:2011 and ISO/IEC/IEEE 60559:2011 and the following apply.

#### 4.1 C11

standard ISO/IEC 9899:2011, *Information technology — Programming languages, their environments and system software interfaces — Programming Language C*, including *Technical Corrigendum 1* (ISO/IEC 9899:2011/Cor. 1:2012)

## 5 C standard conformance

### 5.1 Freestanding implementations

The specification in C11 + TS18661-1 + TS18661-2 allows freestanding implementations to conform to this Part of Technical Specification 18661.

### 5.2 Predefined macros

#### Change to C11 + TS18661-1 + TS18661-2:

In 6.10.8.3#1, add:

`__STDC_IEC_60559_TYPES__` The integer constant `2011yymmL`, intended to indicate support of interchange and extended floating types according to IEC 60559.

### 5.3 Standard headers

The new identifiers added to C11 library headers by this Part of Technical Specification 18661 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 changes to C11 + TS18661-1 + TS18661-2 list these identifiers in each applicable library subclause.

#### Changes to C11 + TS18661-1 + TS18661-2:

After 5.2.4.2.2#6b, insert the paragraph:

[6c] 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`:

<code>FLT/MANT_DIG</code>	<code>FLT/MIN_10_EXP</code>	<code>FLT/EPSILON</code>
<code>FLT/DECIMAL_DIG</code>	<code>FLT/MAX_EXP</code>	<code>FLT/MIN</code>
<code>FLT/DIG</code>	<code>FLT/MAX_10_EXP</code>	<code>FLT/TRUE_MIN</code>
<code>FLT/MIN_EXP</code>	<code>FLT/MAX</code>	

for supported types `_FloatNx`:

<code>FLT/X_MANT_DIG</code>	<code>FLT/X_MIN_10_EXP</code>	<code>FLT/X_EPSILON</code>
<code>FLT/X_DECIMAL_DIG</code>	<code>FLT/X_MAX_EXP</code>	<code>FLT/X_MIN</code>
<code>FLT/X_DIG</code>	<code>FLT/X_MAX_10_EXP</code>	<code>FLT/X_TRUE_MIN</code>
<code>FLT/X_MIN_EXP</code>	<code>FLT/X_MAX</code>	

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

<code>DEC/N_MANT_DIG</code>	<code>DEC/N_MAX</code>	<code>DEC/N_TRUE_MIN</code>
<code>DEC/N_MIN_EXP</code>	<code>DEC/N_EPSILON</code>	
<code>DEC/N_MAX_EXP</code>	<code>DEC/N_MIN</code>	

for supported types `_DecimalNx`:

<code>DECNx_MANT_DIG</code>	<code>DECNx_MAX</code>	<code>DECNx_TRUE_MIN</code>
<code>DECNx_MIN_EXP</code>	<code>DECNx_EPSILON</code>	
<code>DECNx_MAX_EXP</code>	<code>DECNx_MIN</code>	

5 After 7.3#2, insert the paragraph:

[2a] 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:

for supported types `_FloatN`:

10

<code>cacosfN</code>	<code>catanhfN</code>	<code>csqrtfN</code>
<code>casinfN</code>	<code>ccoshfN</code>	<code>cargfN</code>
<code>catanfN</code>	<code>csinhfN</code>	<code>cimagfN</code>
<code>ccosfN</code>	<code>ctanhfN</code>	<code>CMPLXfN</code>
15 <code>csinfN</code>	<code>cexpfN</code>	<code>conjfN</code>
<code>ctanfN</code>	<code>clogfN</code>	<code>cprojfN</code>
<code>cacoshfN</code>	<code>cabsfN</code>	<code>crealfN</code>
<code>casinhfN</code>	<code>cpowfN</code>	

for supported types `_FloatNx`:

20

25

<code>cacosfNx</code>	<code>catanhfNx</code>	<code>csqrtfNx</code>
<code>casinfNx</code>	<code>ccoshfNx</code>	<code>cargfNx</code>
<code>catanfNx</code>	<code>csinhfNx</code>	<code>cimagfNx</code>
<code>ccosfNx</code>	<code>ctanhfNx</code>	<code>CMPLXfNx</code>
<code>csinfNx</code>	<code>cexpfNx</code>	<code>conjfNx</code>
<code>ctanfNx</code>	<code>clogfNx</code>	<code>cprojfNx</code>
<code>cacoshfNx</code>	<code>cabsfNx</code>	<code>crealfNx</code>
<code>casinhfNx</code>	<code>cpowfNx</code>	

After 7.12#1c, insert the paragraph:

30

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

for supported types `_FloatN`:

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

for supported types `_FloatNx`:

30	<code>HUGE_VAL_FNx</code>	<code>logbfNx</code>	<code>fromfpNx</code>
	<code>SNANFNx</code>	<code>modffNx</code>	<code>ufromfpNx</code>
	<code>FP_FAST_FMAFNx</code>	<code>scalbnfNx</code>	<code>fromfpNx</code>
	<code>acosfNx</code>	<code>scalblnfNx</code>	<code>ufromfpNx</code>
	<code>asinfNx</code>	<code>cbrtfNx</code>	<code>fmodfNx</code>
35	<code>atanfNx</code>	<code>fabsfNx</code>	<code>remainderfNx</code>
	<code>atan2fNx</code>	<code>hypotfNx</code>	<code>remquoNx</code>
	<code>cosfNx</code>	<code>powfNx</code>	<code>copysignfNx</code>
	<code>sinfNx</code>	<code>sqrtfNx</code>	<code>nanfNx</code>
	<code>tanfNx</code>	<code>erffNx</code>	<code>nextafterfNx</code>
	<code>acoshfNx</code>	<code>erfcfNx</code>	<code>nexttupfNx</code>
40	<code>asinhfNx</code>	<code>lgammafNx</code>	<code>nextdownfNx</code>
	<code>atanhfNx</code>	<code>tgammafNx</code>	<code>canonicalizefNx</code>
	<code>expfNx</code>	<code>ceilfNx</code>	<code>fdimfNx</code>
	<code>exp2fNx</code>	<code>floorfNx</code>	<code>fmaxfNx</code>
	<code>expm1fNx</code>	<code>nearbyintfNx</code>	<code>fminfNx</code>
45	<code>frexpfNx</code>	<code>rintfNx</code>	<code>fmaxmagfNx</code>
	<code>ilogbfNx</code>	<code>lrintfNx</code>	<code>fminmagfNx</code>
	<code>llogbfNx</code>	<code>llrintfNx</code>	<code>fmafNx</code>
	<code>ldexpfNx</code>	<code>roundfNx</code>	<code>totalorderfNx</code>
	<code>logfNx</code>	<code>lroundfNx</code>	<code>totalordermagfNx</code>
50	<code>log10fNx</code>	<code>llroundfNx</code>	<code>getpayloadfNx</code>
	<code>log1pfNx</code>	<code>truncfNx</code>	<code>setpayloadfNx</code>
	<code>log2fNx</code>	<code>roundevenfNx</code>	<code>setpayloadsigfNx</code>

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

<code>fMaddfN</code>	<code>fMmulfN</code>	<code>fMsqrtfN</code>
<code>fMsubfN</code>	<code>fMdivfN</code>	<code>fMfmafN</code>

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

5	<code>fMaddfNx</code>	<code>fMmulfNx</code>	<code>fMsqrtfNx</code>
	<code>fMsubfNx</code>	<code>fMdivfNx</code>	<code>fMfmafNx</code>

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

<code>fMxaddfN</code>	<code>fMxmulfN</code>	<code>fMxsqrtfN</code>
<code>fMxsubfN</code>	<code>fMxdivfN</code>	<code>fMxfmafN</code>

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

<code>fMxaddfNx</code>	<code>fMxmulfNx</code>	<code>fMxsqrtfNx</code>
<code>fMxsubfNx</code>	<code>fMxdivfNx</code>	<code>fMxfmafNx</code>

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

`fMencfN`

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

	<code>HUGE_VAL_DN</code>	<code>scalblndN</code>	<code>remquodN</code>
	<code>SNANDN</code>	<code>cbtrtdN</code>	<code>copysigndN</code>
	<code>FP_FAST_FMADN</code>	<code>fabsdN</code>	<code>nandN</code>
20	<code>acosdN</code>	<code>hypotdN</code>	<code>nextafterdN</code>
	<code>asindN</code>	<code>powdN</code>	<code>nextupdN</code>
	<code>atandN</code>	<code>sqrtdN</code>	<code>nextdowndN</code>
	<code>atan2dN</code>	<code>erfdN</code>	<code>canonicalizedN</code>
	<code>cosdN</code>	<code>erfcdN</code>	<code>quantizedN</code>
25	<code>sindN</code>	<code>lgammadN</code>	<code>samequantumdN</code>
	<code>tandN</code>	<code>tgammadN</code>	<code>quantumdN</code>
	<code>acoshdN</code>	<code>ceildN</code>	<code>llquantexpdN</code>
	<code>asinhdN</code>	<code>floordN</code>	<code>encodedecdN</code>
	<code>atanhdN</code>	<code>nearbyintdN</code>	<code>decodedecdN</code>
30	<code>expdN</code>	<code>rintdN</code>	<code>encodebindN</code>
	<code>exp2dN</code>	<code>lrintdN</code>	<code>decodebindN</code>
	<code>expm1dN</code>	<code>llrintdN</code>	<code>fdimdN</code>
	<code>frexpdpN</code>	<code>rounddN</code>	<code>fmaxdN</code>
	<code>ilogbdN</code>	<code>lrounddN</code>	<code>fmindN</code>
35	<code>llogbdN</code>	<code>llrounddN</code>	<code>fmaxmagdN</code>
	<code>ldexpdN</code>	<code>truncdN</code>	<code>fminmagdN</code>
	<code>logdN</code>	<code>roundevendN</code>	<code>fmadN</code>
	<code>log10dN</code>	<code>fromfpdN</code>	<code>totalorderdN</code>
	<code>log1pdN</code>	<code>ufromfpdN</code>	<code>totalordermagdN</code>
40	<code>log2dN</code>	<code>fromfpdN</code>	<code>getpayloaddN</code>
	<code>logbdN</code>	<code>ufromfpdN</code>	<code>setpayloaddN</code>
	<code>modfdN</code>	<code>fmoddN</code>	<code>setpayloadsigdN</code>
	<code>scalbndN</code>	<code>remainderdN</code>	

for supported types `_DecimalNx`:

	<code>HUGE_VAL_DNx</code>	<code>scalbndNx</code>	<code>fmoddNx</code>
	<code>SNANDNx</code>	<code>scalblndNx</code>	<code>remainderdNx</code>
5	<code>FP_FAST_FMADNx</code>	<code>cbrtdNx</code>	<code>remquoddNx</code>
	<code>acosdNx</code>	<code>fabsdNx</code>	<code>copysigndNx</code>
	<code>asindNx</code>	<code>hypotdNx</code>	<code>nandNx</code>
	<code>atandNx</code>	<code>powdNx</code>	<code>nextafterdNx</code>
	<code>atan2dNx</code>	<code>sqrtdNx</code>	<code>nexttupdNx</code>
10	<code>cosdNx</code>	<code>erfdNx</code>	<code>nextdowndNx</code>
	<code>sindNx</code>	<code>erfcdNx</code>	<code>canonicalizedNx</code>
	<code>tandNx</code>	<code>lgammadNx</code>	<code>quantizedNx</code>
	<code>acoshdNx</code>	<code>tgammadNx</code>	<code>samequantumdNx</code>
	<code>asinhdNx</code>	<code>ceildNx</code>	<code>quantumdNx</code>
15	<code>atanhdNx</code>	<code>floordNx</code>	<code>llquantexpdNx</code>
	<code>expdNx</code>	<code>nearbyintdNx</code>	<code>fdimdNx</code>
	<code>exp2dNx</code>	<code>rintdNx</code>	<code>fmaxdNx</code>
	<code>expm1dNx</code>	<code>lrintdNx</code>	<code>fmindNx</code>
	<code>frexpdx</code>	<code>llrintdNx</code>	<code>fmaxmagdNx</code>
20	<code>ilogbdNx</code>	<code>rounddNx</code>	<code>fminmagdNx</code>
	<code>llogbdNx</code>	<code>lrounddNx</code>	<code>fmadNx</code>
	<code>ldexpdNx</code>	<code>llrounddNx</code>	<code>totalorderdNx</code>
	<code>logdNx</code>	<code>truncdNx</code>	<code>totalordermagdNx</code>
	<code>log10dNx</code>	<code>roundevendNx</code>	<code>getpayloaddNx</code>
25	<code>log1pdNx</code>	<code>fromfpdNx</code>	<code>setpayloaddNx</code>
	<code>log2dNx</code>	<code>ufromfpdNx</code>	<code>setpayloadsigdNx</code>
	<code>logbdNx</code>	<code>fromfpxdNx</code>	
	<code>modfdNx</code>	<code>ufromfpxdNx</code>	

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

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

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

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

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

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

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

	<code>FP_FAST_DMxADDDNx</code>	<code>FP_FAST_DMxSQRTDNx</code>	<code>dMxmuldNx</code>
	<code>FP_FAST_DMxSUBDNx</code>	<code>FP_FAST_DMxFMADNx</code>	<code>dMxdivdNx</code>
	<code>FP_FAST_DMxMULDNx</code>	<code>dMxadddNx</code>	<code>dMxsqrdNx</code>
5	<code>FP_FAST_DMxDIVDNx</code>	<code>dMxsubdNx</code>	<code>dMxfmadNx</code>

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

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

After 7.22#1b, insert the paragraph:

10 [1c] 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>
------------------------	----------------------

for supported types `_FloatNx`:

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

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

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

for supported types `_DecimalNx`:

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

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

<code>strfromencfN</code>	<code>strtoencfN</code>
---------------------------	-------------------------

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

<code>strfromencdecN</code>	<code>strtoencdecN</code>
<code>strfromencbindN</code>	<code>strtoencbindN</code>

## 25 6 Types

This clause specifies changes to C11 + TS18661-1 + TS18661-2 to include types that support the arithmetic interchange formats and extended formats specified in IEC 60559. The encoding conversion functions (11.3) and numeric conversion functions for encodings (12) support the non-arithmetic interchange formats specified in IEC 60559.

### 30 Changes to C11 + TS18661-1 + TS18661-2:

Replace 6.2.5#10a-10b:

[10a] There are three *decimal floating types*, designated as `_Decimal32`, `_Decimal64`, and `_Decimal128`. Respectively, they have the IEC 60559 formats: `decimal32`, `decimal64`, and `decimal128`. Decimal floating types are real floating types.

[10b] Together, the standard floating types and the decimal floating types comprise the *real floating types*.

with:

[10a] 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.

**Table 1 – 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  $emax = 262143$ .

**Table 2 – 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
$w$ , exponent 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  $emax = 1572864$ .



[10b] Types designated

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

and types designated

`_DecimalN`, where  $N \geq 32$  and a multiple of 32

5 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 `_FloatN`, 10 for `_DecimalN`). Interchange floating types are not compatible with any other types.

10 [10c] 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 `_FloatN` as an interchange floating type with the same representation and alignment requirements as `long double`. The implementation may provide other binary interchange floating types.

15 [10d] 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.

20 [10e] Note that 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 (7.12.11.7c), string-to-encoding functions (7.22.1.3c), and string-from-encoding functions (7.22.1.3d). An implementation that defines `__STDC_IEC_60559_TYPES__` shall support the IEC 60559 binary16 format, at least as a non-arithmetic interchange format.

25 [10f] 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:

**Table 3 – Extended format parameters for floating-point numbers**

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

30 [10g] Types designated `_Float32x`, `_Float64x`, `_Float128x`, `_Decimal64x`, and `_Decimal128x` support the corresponding IEC 60559 extended formats and are collectively called the *extended floating types*. Extended floating types are not compatible with any other types. 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

35 `__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`.

[10h] The standard floating types, interchange floating types, and extended floating types are collectively called the *real floating types*.

40 [10i] 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*. Thus the binary floating types and the decimal floating types are real floating types.

Replace 6.2.5#11:

5 [11] There are three *complex types*, designated as `float _Complex`, `double _Complex`, and `long double _Complex`.<sup>43)</sup> (Complex types are a conditional feature that implementations need not support; see 6.10.8.3.) The real floating and complex types are collectively called the *floating types*.

with:

10 [11] For the standard real types `float`, `double`, and `long double`, the interchange floating types `_FloatN`, and the extended floating types `_FloatNx`, there are *complex types* designated respectively as `float _Complex`, `double _Complex`, `long double _Complex`, `_FloatN _Complex`, and `_FloatNx _Complex`.<sup>43)</sup> (Complex types are a conditional feature that implementations need not support; see 6.10.8.3.) The real floating and complex types are collectively  
15 called the *floating types*.

In the list of keywords in 6.4.1, replace:

```
_Decimal32
_Decimal64
_Decimal128
```

20 with:

```
_FloatN, where N is 16, 32, 64, or ≥ 128 and a multiple of 32
_Float32x
_Float64x
_Float128x
_DecimalN, where N ≥ 32 and a multiple of 32
_Decimal64x
_Decimal128x
```

25

In the list of type specifiers in 6.7.2, replace:

```
_Decimal32
_Decimal64
_Decimal128
```

30

with:

```
_FloatN, where N is 16, 32, 64, or ≥ 128 and a multiple of 32
_Float32x
_Float64x
_Float128x
_DecimalN, where N ≥ 32 and a multiple of 32
_Decimal64x
_Decimal128x
```

35

40

In the list of constraints in 6.7.2#2, replace:

- `_Decimal32`
- `_Decimal64`
- `_Decimal128`

5 with:

- `_FloatN`, where  $N$  is 16, 32, 64, or  $\geq 128$  and a multiple of 32
- `_Float32x`
- `_Float64x`
- `_Float128x`

10 — `_DecimalN`, where  $N \geq 32$  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`

15 — `_Float64x _Complex`

- `_Float128x _Complex`

Replace 6.7.2#3a:

[3a] The type specifiers `_Decimal32`, `_Decimal64`, and `_Decimal128` shall not be used if the implementation does not support decimal floating types (see 6.10.8.3).

20 with:

[3a] The type specifiers `_FloatN` (where  $N$  is 16, 32, 64, or  $\geq 128$  and a multiple of 32), `_Float32x`, `_Float64x`, `_Float128x`, `_DecimalN` (where  $N \geq 32$  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).

25 Replace 6.5#8a:

[8a] Operators involving decimal floating types are evaluated according to the semantics of IEC 60559, including production of results with the preferred quantum exponent as specified in IEC 60559.

with:

30 [8a] 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 5.2.4.2.2b).

Replace G.2#2:

[2] There are three *imaginary types*, designated as `float _Imaginary`, `double _Imaginary`, and `long double _Imaginary`. The imaginary types (along with the real floating and complex types) are floating types.

5 with:

[2] For the standard floating types `float`, `double`, and `long double`, the interchange floating types `_FloatN`, and the extended floating types `_FloatNx`, there are *imaginary types* designated respectively as `float _Imaginary`, `double _Imaginary`, `long double _Imaginary`, `_FloatN _Imaginary`, and `_FloatNx _Imaginary`. The imaginary types (along with the real floating and complex types) are floating types.

## 7 Characteristics

This clause specifies new `<float.h>` macros, analogous to the macros for standard floating types, that characterize the interchange and extended floating types. Some specification for decimal floating types introduced in Part 2 of Technical Specification 18661 is subsumed under the general specification for interchange floating types.

### Changes to C11 + TS18661-1 + TS18661-2:

Re-number and rename 5.2.4.2.2a:

#### 5.2.4.2.2a Characteristics of decimal floating types in `<float.h>`

to:

#### 5.2.4.2.2b Alternate model for decimal floating-point numbers

and remove paragraphs 1-3:

[1] This subclause specifies macros in `<float.h>` that provide characteristics of decimal floating types in terms of the model presented in 5.2.4.2.2. The prefixes `DEC32_`, `DEC64_`, and `DEC128_` denote the types `_Decimal32`, `_Decimal64`, and `_Decimal128` respectively.

[2] `DEC_EVAL_METHOD` is the decimal floating-point analogue of `FLT_EVAL_METHOD` (5.2.4.2.2). Its implementation-defined value characterizes the use of evaluation formats for decimal floating types:

- 1 indeterminate;
- 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] 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



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 lists.

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

$N + 1$ , where `_FloatNx` is a supported extended floating type  
evaluate operations and constants, whose semantic type has at most the range and precision of the `_FloatNx` type, to the range and precision of the `_FloatNx` 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` characterizes the use of evaluation formats (see analogous `FLT_EVAL_METHOD` in 5.2.4.2.2) 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 `_DecimalN` is a supported interchange floating type  
evaluate operations and constants, whose semantic type has at most the range and precision of the `_DecimalN` type, to the range and precision of the `_DecimalN` type; evaluate all other operations and constants to the range and precision of the semantic type;

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

5 [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)

10 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,  
 since the radix is an inherent property of the types.

— 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,

15 `FLT_N_DECIMAL_DIG`  
`FLT_NX_DECIMAL_DIG`

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

20 `FLT_N_MANT_DIG`  
`FLT_NX_MANT_DIG`

— number of digits in the coefficient,  $p$

25 `DEC_N_MANT_DIG`  
`DEC_NX_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,  
 $\lceil 1 + p \log_{10} 2 \rceil$

30 `FLT_N_DECIMAL_DIG`  
`FLT_NX_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,  $\lceil (p - 1) \log_{10} 2 \rceil$

35 `FLT_N_DIG`  
`FLT_NX_DIG`

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

40 `FLT_N_MIN_EXP`  
`FLT_NX_MIN_EXP`  
`DEC_N_MIN_EXP`  
`DEC_NX_MIN_EXP`

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

45 `FLT_N_MIN_10_EXP`  
`FLT_NX_MIN_10_EXP`

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

5  
**FLT/N\_MAX\_EXP**  
**FLT/X\_MAX\_EXP**  
**DEC/N\_MAX\_EXP**  
**DEC/X\_MAX\_EXP**

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

10  
**FLT/N\_MAX\_10\_EXP**  
**FLT/X\_MAX\_10\_EXP**

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

15  
**FLT/N\_MAX**  
**FLT/X\_MAX**  
**DEC/N\_MAX**  
**DEC/X\_MAX**

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

20  
**FLT/N\_EPSILON**  
**FLT/X\_EPSILON**  
**DEC/N\_EPSILON**  
**DEC/X\_EPSILON**

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

25  
**FLT/N\_MIN**  
**FLT/X\_MIN**  
**DEC/N\_MIN**  
**DEC/X\_MIN**

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

30  
**FLT/N\_TRUE\_MIN**  
**FLT/X\_TRUE\_MIN**  
**DEC/N\_TRUE\_MIN**  
**DEC/X\_TRUE\_MIN**

With the following change, **DECIMAL\_DIG** characterizes conversions of supported IEC 60559 encodings, which may be wider than supported floating types.

#### Change to C11 + TS18661-1 + TS18661-2:

40 In 5.2.4.2.2#11, change the bullet defining **DECIMAL\_DIG** from:

- number of decimal digits,  $n$ , such that any floating-point number in the widest supported floating type with ...



to:

- number of decimal digits,  $n$ , such that any floating-point number in the widest of the supported floating types and the supported IEC 60559 encodings with ...

## 8 Conversions

- 5 The following change to C11 + TS18661-1 + TS18661-2 enhances the usual arithmetic conversions to handle interchange and extended floating types. IEC 60559 recommends 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. The following change supports this restriction.

### Changes to C11 + TS18661-1 + TS18661-2:

- 10 Replace 6.3.1.4#1a:

[1a] When a finite value of decimal 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.

- 15 with:

[1a] 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.

- 20 Replace 6.3.1.4#2a:

[2a] When a value of integer type is converted to a decimal 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.

- 25 with:

[2a] 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.

- 30 In 6.3.1.8#1, replace the following items after “This pattern is called the *usual arithmetic conversions*”:

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

First, if the type of either operand is `_Decimal128`, the other operand is converted to `_Decimal128`.

- 35 Otherwise, if the type of either operand is `_Decimal164`, the other operand is converted to `_Decimal164`.

Otherwise, if the type of either operand is `_Decimal132`, the other operand is converted to `_Decimal132`.

If there are no decimal floating types in the operands:

First, 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`.

5 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`.

Otherwise, if the corresponding real type of either operand is `float`, the other operand is converted, without change of type domain, to a type whose corresponding real type is `float.62`)

with:

10 If one operand has decimal floating type, the other operand shall not have standard floating, complex, or imaginary type, nor shall it have a floating type of radix 2.

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.

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

Otherwise, 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.

20 Otherwise, if the corresponding real type of either operand is a standard floating type, the other operand is converted, without change of type domain, to a type whose corresponding real type is that same standard floating type.

25 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.

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

## 9 Constants

30 The following changes to C11 + TS18661-1 + TS18661-2 provide suffixes that designate constants of interchange and extended floating types.

### Changes to C11 + TS18661-1 + TS18661-2:

Change *floating-suffix* in 6.4.4.2 from:

35 *floating-suffix*: one of  
f l F L df dd dL DF DD DL

to:

*floating-suffix*: one of  
f l F L df dd dL DF DD DL fN FN fNx FNx dN DN dNx DNx

Replace 6.4.4.2#2a:

[2a] A *floating-suffix* `df`, `dd`, `d1`, `DF`, `DD`, or `DL` shall not be used in a *hexadecimal-floating-constant*.

with:

5 [2a] A *floating-suffix* `df`, `dd`, `d1`, `DF`, `DD`, `DL`, `dN`, `DN`, `dNx`, or `DNx` shall not be used in a *hexadecimal-floating-constant*.

[2b] A *floating-suffix* shall not designate a type that the implementation does not provide.

Replace 6.4.4.2#4a:

[4a] If a floating constant is suffixed by `df` or `DF`, it has type `_Decimal32`. If suffixed by `dd` or `DD`, it has type `_Decimal64`. If suffixed by `d1` or `DL`, it has type `_Decimal128`.

10 with:

[4a] If a floating constant is suffixed by `fN` or `FN`, it has type `_FloatN`. If suffixed by `fNx` or `FNx`, it has type `_FloatNx`. If suffixed by `df` or `DF`, it has type `_Decimal32`. If suffixed by `dd` or `DD`, it has type `_Decimal64`. If suffixed by `d1` or `DL`, it has type `_Decimal128`. If suffixed by `dN` or `DN`, it has type `_DecimalN`. If suffixed by `dNx` or `DNx`, it has type `_DecimalNx`.

15 Replace the second sentence of 6.4.4.2#5a:

The quantum exponent is specified to be the same as for the corresponding `strtod32`, `strtod64`, or `strtod128` function for the same numeric string.

with:

20 The quantum exponent is specified to be the same as for the corresponding `strtodN` or `strtodNx` function for the same numeric string.

## 10 Expressions

The following changes to C11 + TS18661-1 + TS18661-2 specify operator constraints for interchange and extended floating types.

### Changes to C11 + TS18661-1 + TS18661-2:

25 Replace 6.5.5#2a:

[2a] If either operand has decimal floating type, the other operand shall not have standard floating type, complex type, or imaginary type.

with:

30 [2a] If either operand has decimal floating type, the other operand shall not have standard floating type, binary floating type, complex type, or imaginary type.

Replace 6.5.6#3a:

[3a] If either operand has decimal floating type, the other operand shall not have standard floating type, complex type, or imaginary type.

with:

[3a] If either operand has decimal floating type, the other operand shall not have standard floating type, binary floating type, complex type, or imaginary type.

Replace 6.5.8#2a:

[2a] If either operand has decimal floating type, the other operand shall not have standard floating type.

with:

[2a] If either operand has decimal floating type, the other operand shall not have standard floating type or binary floating type.

Replace 6.5.9#2a:

[2a] If either operand has decimal floating type, the other operand shall not have standard floating type, complex type, or imaginary type.

with:

[2a] If either operand has decimal floating type, the other operand shall not have standard floating type, binary floating type, complex type, or imaginary type.

In 6.5.15#3, replace the bullet:

- one operand has decimal floating type, and the other has arithmetic type other than standard floating type, complex type, and imaginary type;

with:

- one operand has decimal floating type, and the other has arithmetic type other than standard floating type, binary floating types, complex type, and imaginary type;

Replace 6.5.16#2a:

[2a] If either operand has decimal floating type, the other operand shall not have standard floating type, complex type, or imaginary type.

with:

[2a] If either operand has decimal floating type, the other operand shall not have standard floating type, binary floating type, complex type, or imaginary type.

In F.9.2#1, replace the first sentence:

[1] The equivalences noted below apply to expressions of standard floating types.

with:

[1] The equivalences noted below apply to expressions of standard floating types and binary floating types.

## 11 Non-arithmetic interchange formats

An implementation supports IEC 60559 arithmetic interchange formats by providing the corresponding interchange floating types. An implementation supports IEC 60559 non-arithmetic formats by providing the encoding-to-encoding conversion functions in `<math.h>` and the string-to-encoding and string-from-encoding

functions in `<stdlib.h>`. See 6.2.5. These functions, together with functions required for interchange floating types, provide conversions between any two of the supported IEC 60559 arithmetic and non-arithmetic interchange formats and between character sequences and any supported IEC 60559 arithmetic or non-arithmetic format.

## 5 **12 Mathematics** `<math.h>`

This clause specifies changes to C11 + TS18661-1 + TS18661-2 to include functions and macros for interchange and extended floating types. The binary types are supported by functions and macros corresponding to those specified for standard floating types (`float`, `double`, and `long double`) in C11 + TS18661-1, including Annex F. The decimal types are supported by functions and macros corresponding to those specified for decimal floating types in TS18661-2.

All classification (7.12.3) and comparison (7.12.14) macros specified in C11 + TS18661-1 + TS18661-2 naturally extend to handle interchange and extended floating types.

This clause also specifies encoding conversion functions that are part of support for the non-arithmetic interchange formats in IEC 60559 (see 6.2.5).

### 15 **Changes to C11 + TS18661-1 + TS18661-2:**

In 7.12#1, change the second sentence from:

Most synopses specify a family of functions consisting of a principal function with one or more `double` parameters, a `double` return value, or both; and other functions with the same name but with `f` and `l` suffixes, which are corresponding functions with `float` and `long double` parameters, return values, or both.

to:

Most synopses specify a family of functions consisting of:

a principal function with one or more `double` parameters, a `double` return value, or both; and,  
other functions with the same name but with `f`, `l`, `fN`, `fNx`, `dN`, and `dNx` suffixes, which are corresponding functions whose parameters, return values, or both are of types `float`, `long double`, `_FloatN`, `_FloatNx`, `_DecimalN`, and `_DecimalNx`, respectively.

Add after 7.12#1d:

[1e] 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 the implementation does not provide, `<math.h>` shall not define the associated macros or declare the associated functions unless explicitly specified otherwise.

### **12.1 Macros**

#### **Changes to C11 + TS18661-1 + TS18661-2:**

Replace 7.12#3a:

[3a] The macro

```
HUGE_VAL_D32
```

expands to a constant expression of type `_Decimal64` representing positive infinity. The macros

**HUGE\_VAL\_D64**  
**HUGE\_VAL\_D128**

are respectively `_Decimal64` and `_Decimal128` analogues of `HUGE_VAL_D32`.

5 with:

[3a] The macros

**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.

Replace 7.12#5b:

[5b] The decimal signaling NaN macros

**SNAND32**  
**SNAND64**  
**SNAND128**

each expands to a constant expression of the respective decimal floating type 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.

with:

[5b] 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.

Replace 7.12#7b:

[7b] The macros

**FP\_FAST\_FMAD32**  
**FP\_FAST\_FMAD64**  
**FP\_FAST\_FMAD128**

are, respectively, `_Decimal32`, `_Decimal64`, and `_Decimal128` analogues of `FP_FAST_FMA`.

40 with:

[7b] The macros

**FP\_FAST\_FMAFN**

```

FP_FAST_FMADN
FP_FAST_FMAFN
FP_FAST_FMADNX

```

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

Replace 7.12#7c:

[7c] The macros

```

10 FP_FAST_D32ADDD64
FP_FAST_D32ADDD128
FP_FAST_D64ADDD128
FP_FAST_D32SUBD64
FP_FAST_D32SUBD128
FP_FAST_D64SUBD128
15 FP_FAST_D32MULD64
FP_FAST_D32MULD128
FP_FAST_D64MULD128
FP_FAST_D32DIVD64
FP_FAST_D32DIVD128
20 FP_FAST_D64DIVD128
FP_FAST_D32FMAD64
FP_FAST_D32FMAD128
FP_FAST_D64FMAD128
FP_FAST_D32SQRTD64
25 FP_FAST_D32SQRTD128
FP_FAST_D64SQRTD128

```

are decimal analogues of `FP_FAST_FADD`, `FP_FAST_FADDL`, `FP_FAST_DADDL`, etc.

with:

30 [7c] The macros in the following lists are interchange and extended floating type analogues of `FP_FAST_FADD`, `FP_FAST_FADDL`, `FP_FAST_DADDL`, etc.

[7d] For  $M < N$ , the macros

```

35 FP_FAST_FMADDFN
FP_FAST_FMSUBFN
FP_FAST_FMMULFN
FP_FAST_FMDIVFN
FP_FAST_FMFMAFN
FP_FAST_FMSQRTFN
40 FP_FAST_DMADDDN
FP_FAST_DMSUBDN
FP_FAST_DMMULDN
FP_FAST_DMDIVDN
FP_FAST_DMFMADN
45 FP_FAST_DMSQRTDN

```

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  $M$ .

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

```

FP_FAST_FMADDFNX
FP_FAST_FMSUBFNX
FP_FAST_FMMULFNX
5 FP_FAST_FMDIVFNX
FP_FAST_FMFMFNX
FP_FAST_FMSQRTFNX
FP_FAST_DMADDDNX
FP_FAST_DMSUBDNX
10 FP_FAST_DMMULDNX
FP_FAST_DMDIVDNX
FP_FAST_DMFMADNX
FP_FAST_DMSQRTDNX

```

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  $M$ .

[7f] For  $M < N$ , the macros

```

FP_FAST_FMXADDFN
FP_FAST_FMXSUBFN
FP_FAST_FMXMULFN
20 FP_FAST_FMXDIVFN
FP_FAST_FMXFMFN
FP_FAST_FMXSQRTFN
FP_FAST_DMxADDDN
FP_FAST_DMXSUBDN
25 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  $M$ .

[7g] For  $M < N$ , the macros

```

FP_FAST_FMXADDFNX
FP_FAST_FMXSUBFNX
FP_FAST_FMXMULFNX
35 FP_FAST_FMXDIVFNX
FP_FAST_FMXFMFNX
FP_FAST_FMXSQRTFNX
FP_FAST_DMxADDDNX
FP_FAST_DMXSUBDNX
40 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  $M$ .

## 12.2 Floating-point environment

Changes to C11 + TS18661-1 + TS18661-2:



In 7.6.1a#2, change the first sentence from:

The `FENV_ROUND` pragma provides a means to specify a constant rounding direction for floating-point operations for standard floating types within a translation unit or compound statement.

to:

- 5 The `FENV_ROUND` pragma provides a means to specify a constant rounding direction for floating-point operations for standard and binary floating types within a translation unit or compound statement.

In 7.6.1a#3, change the first sentence from:

*direction* shall be one of the names of the supported rounding direction macros for operations for standard floating types (7.6), or `FE_DYNAMIC`.

10 to:

*direction* shall be one of the names of the supported rounding direction macros for use with `fegetround` and `fesetround` (7.6), or `FE_DYNAMIC`.

In 7.6.1a#4, change the first sentence from:

- 15 The `FENV_ROUND` directive affects operations for standard floating types. Within the scope of an `FENV_ROUND` directive establishing a mode other than `FE_DYNAMIC`, floating-point operators, ...

to:

The `FENV_ROUND` directive affects operations for standard and binary floating types. Within the scope of an `FENV_ROUND` directive establishing a mode other than `FE_DYNAMIC`, floating-point operators, ...

- 20 In 7.6.1a#4, change the table title from:

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

to:

**Functions affected by constant rounding modes – for standard and binary floating types**

In 7.6.1a#4, replace the sentence following the table:

- 25 Each `<math.h>` function listed in the table above indicates the family of functions of all standard floating types (for example, `acosf` and `acosl` as well as `acos`).

with:

Each `<math.h>` function listed in the table above indicates the family of functions of all standard and binary floating types (for example, `acosf`, `acosl`, `acosfN`, and `acosfNx` as well as `acos`).

- 30 After 7.6.1a#4, add:

[4a] The `fMencfN`, `strfromencfN`, and `strtoencfN` functions for binary interchange types are also affected by constant rounding modes.

In 7.6.1b#2 after the table, add:

- 35 Each `<math.h>` function listed in the table above indicates the family of functions of all decimal floating types (for example, `acosdNx`, as well as `acosdN`).

After 7.6.1b#2, add:

[3] The `dMencbindN`, `dMencdecN`, `strfromencbindN`, `strfromencdecN`, `strtoencbindN`, and `strtoencdecN` functions for decimal interchange types are also affected by constant rounding modes.

5 Change 7.6.3 from:

The `fegetround` and `fesetround` functions provide control of rounding direction modes.

to:

The functions in this subclause provide control of rounding direction modes.

Change 7.6.3.1#2 from:

The `fegetround` function gets the current rounding direction.

to:

The `fegetround` function gets the current rounding direction for operations for standard and binary floating types.

In 7.6.3.2#2, change the first sentence from:

The `fesetround` function establishes the rounding direction represented by its argument `round`.

to:

The `fesetround` function establishes the rounding direction represented by its argument `round` for operations for standard and binary floating types.

## 12.3 Functions

20 **Changes to C11 + TS18661-1 + TS18661-2:**

Add the following list of function prototypes to the synopsis of the respective subclauses:

### 7.12.4 Trigonometric functions

```

25  _FloatN acosfN(_FloatN x);
    _FloatNx acosfNx(_FloatNx x);
    _DecimalN acosdN(_DecimalN x);
    _DecimalNx acosdNx(_DecimalNx x);

30  _FloatN asinfN(_FloatN x);
    _FloatNx asinfNx(_FloatNx x);
    _DecimalN asindN(_DecimalN x);
    _DecimalNx asindNx(_DecimalNx x);

35  _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);

```

```

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

```

```

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

```

```

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

```

#### 7.12.5 Hyperbolic functions

```

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

```

```

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

```

```

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

```

```

35   _FloatN coshfN(_FloatN x);
    _FloatNx coshfNx(_FloatNx x);
    _DecimalN coshdN(_DecimalN x);
    _DecimalNx scoshdNx(_DecimalNx x);

```

```

40   _FloatN sinhN(_FloatN x);
    _FloatNx sinhNx(_FloatNx x);
    _DecimalN sinhdN(_DecimalN x);
    _DecimalNx sinhdNx(_DecimalNx x);

```

```

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

```

#### 50 7.12.6 Exponential and logarithmic functions

```

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

```

```

    _DecimalNx expdNx(_DecimalNx x);

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

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

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

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

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

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

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

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

    _FloatN log1pfN(_FloatN x);
    _FloatNx log1pfNx(_FloatNx x);
45  _DecimalN log1pdN(_DecimalN x);
    _DecimalNx log1pdNx(_DecimalNx x);

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

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

```

```

_FloatN modffN(_FloatN x, _FloatN *iptr);
_FloatNx modffNx(_FloatNx x, _FloatNx *iptr);
_DecimalN modfdN(_DecimalN x, _DecimalN *iptr);
_DecimalNx modfdNx(_DecimalNx x, _DecimalNx *iptr);

```

5

```

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

```

10

```

_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);

```

## 15 7.12.7 Power and absolute-value functions

```

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

```

20

```

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

```

25

```

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

```

30

```

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

```

35

```

_FloatN sqrtfN(_FloatN x);
_FloatNx sqrtfNx(_FloatNx x);
_DecimalN sqrtdN(_DecimalN x);
_DecimalNx sqrtdNx(_DecimalNx x);

```

## 40 7.12.8 Error and gamma functions

```

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

```

45

```

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

```

50

```

_FloatN lgammafN(_FloatN x);
_FloatNx lgammafNx(_FloatNx x);

```

```

    _DecimalN lgammadN(_DecimalN x);
    _DecimalNx lgammadNx(_DecimalNx x);

    _FloatN tgammafN(_FloatN x);
    _FloatNx tgammafNx(_FloatNx x);
    _DecimalN tgammaadN(_DecimalN x);
    _DecimalNx tgammaadNx(_DecimalNx x);

```

## 7.12.9 Nearest integer functions

```

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

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

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

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

    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);

    _FloatN roundfN(_FloatN x);
    _FloatNx roundfNx(_FloatNx x);
    _DecimalN rounddN(_DecimalN x);
    _DecimalNx rounddNx(_DecimalNx x);

    long int lroundfN(_FloatN x);
    long int lroundfNx(_FloatNx x);
    long int lrounddN(_DecimalN x);
    long int lrounddNx(_DecimalNx x);

    long long int llroundfN(_FloatN x);
    long long int llroundfNx(_FloatNx x);
    long long int llrounddN(_DecimalN x);
    long long int llrounddNx(_DecimalNx x);

    _FloatN roundevenfN(_FloatN x);

```

```

    _FloatNx roundevenfNx(_FloatNx x);
    _DecimalN roundevendN(_DecimalN x);
    _DecimalNx roundevendNx(_DecimalNx x);

5    _FloatN truncfN(_FloatN x);
    _FloatNx truncfNx(_FloatNx x);
    _DecimalN truncdN(_DecimalN x);
    _DecimalNx truncdNx(_DecimalNx x);

10   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, unsigned int width);
    uintmax_t ufromfpfN(_FloatN x, int round, unsigned int width);
15   uintmax_t ufromfpfNx(_FloatNx x, int round, unsigned int width);
    uintmax_t ufromfpdN(_DecimalN x, int round, unsigned int width);
    uintmax_t ufromfpdNx(_DecimalNx x, int round, unsigned int width);

    intmax_t fromfpxfN(_FloatN x, int round, unsigned int width);
20   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);
    uintmax_t ufromfpxfNx(_FloatNx x, int round, unsigned int width);
25   uintmax_t ufromfpxdN(_DecimalN x, int round, unsigned int width);
    uintmax_t ufromfpxdNx(_DecimalNx x, int round, unsigned int width);

```

#### 7.12.10 Remainder functions

```

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

    _FloatN remainderfN(_FloatN x, _FloatN y);
    _FloatNx remainderfNx(_FloatNx x, _FloatNx y);
35   _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);
40

```

#### 7.12.11 Manipulation functions

```

    _FloatN copysignfN(_FloatN x, _FloatN y);
    _FloatNx copysignfNx(_FloatNx x, _FloatNx y);
    _DecimalN copysigndN(_DecimalN x, _DecimalN y);
    _DecimalNx copysigndNx(_DecimalNx x, _DecimalNx y);
45

    _FloatN nanfN(const char *tagp);
    _FloatNx nanfNx(const char *tagp);
    _DecimalN nandN(const char *tagp);
    _DecimalNx nandNx(const char *tagp);
50

    _FloatN nextafterfN(_FloatN x, _FloatN y);
    _FloatNx nextafterfNx(_FloatNx x, _FloatNx y);

```

```

    _DecimalN nextafterdN(_DecimalN x, _DecimalN y);
    _DecimalNx nextafterdNx(_DecimalNx x, _DecimalNx y);

    _FloatN nextupfN(_FloatN x);
    _FloatNx nextupfNx(_FloatNx x);
5   _DecimalN nextupdN(_DecimalN x);
    _DecimalNx nextupdNx(_DecimalNx x);

    _FloatN nextdownfN(_FloatN x);
    _FloatNx nextdownfNx(_FloatNx x);
10  _DecimalN nextdowndN(_DecimalN x);
    _DecimalNx nextdowndNx(_DecimalNx x);

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

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

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

25  _DecimalN quantumdN(_DecimalN x);
    _DecimalNx quantumdNx(_DecimalNx x);

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

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

```

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

```

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

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

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

55  _FloatN fmaxmagfN(_FloatN x, _FloatN y);

```



```

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

```

```

5  _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

```

10  _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 format

```

15  _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
20  _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

```

```

25  _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
30  _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

```

```

35  _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
40  _DecimalM dMmuldN(_DecimalN x, _DecimalN y); // M < N
     _DecimalM dMmuldNx(_DecimalNx x, _DecimalNx y); // M <= N
     _DecimalMx dMxmuldN(_DecimalN x, _DecimalN y); // M < N
     _DecimalMx dMxmuldNx(_DecimalNx x, _DecimalNx y); // M < N

```

```

45  _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
50  _DecimalM dMdivdN(_DecimalN x, _DecimalN y); // M < N
     _DecimalM dMdivdNx(_DecimalNx x, _DecimalNx y); // M <= N
     _DecimalMx dMxdivdN(_DecimalN x, _DecimalN y); // M < N
     _DecimalMx dMxdivdNx(_DecimalNx x, _DecimalNx y); // M < N

```

```

_FloatM fMsqrtfN(_FloatN x); // M < N
_FloatM fMsqrtfNx(_FloatNx x); // M <= N

```

```

_FloatMx fMxsqrtfN(_FloatN x); // M < N
_FloatMx fMxsqrtfNx(_FloatNx x); // M < N
_DecimalM dMsqrtfN(_DecimalN x); // M < N
_DecimalM dMsqrtfNx(_DecimalNx x); // M <= N
5  _DecimalMx dMxsqrtfN(_DecimalN x); // M < N
   _DecimalMx dMxsqrtfNx(_DecimalNx x); // M < N

_FloatM fMfmafN(_FloatN x, _FloatN y, _FloatN z); // M < N
_FloatM fMfmafNx(_FloatNx x, _FloatNx y, _FloatNx z); // M <= N
10  _FloatMx fMxfmafN(_FloatN x, _FloatN y, _FloatN z); // M < N
   _FloatMx fMxfmafNx(_FloatNx x, _FloatNx y, _FloatNx z); // M < N
   _DecimalM dMfmadN(_DecimalN x, _DecimalN y, _DecimalN z); // M < N
   _DecimalM dMdfmadNx(_DecimalNx x, _DecimalNx y, _DecimalNx z);
   // M <= N
15  _DecimalMx dMxfmadN(_DecimalN x, _DecimalN y, _DecimalN z);
   // M < N
   _DecimalMx dMxfmadNx(_DecimalNx x, _DecimalNx y, _DecimalNx z);
   // M < N

```

## F.10.12 Total order functions

```

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

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

```

## F.10.13 Payload functions

```

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

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

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

```

In 7.12.6.4#2, change the third sentence from:

If the type of the function is a standard floating type, the exponent is an integral power of 2.

to:

If the type of the function is a standard or binary floating type, the exponent is an integral power of 2.

In 7.12.6.4#3, change the second sentence from:

Otherwise, the `frexp` functions return the value `x`, such that: `x` has a magnitude in the interval  $[1/2, 1)$  or zero, and `value` equals  $x \times 2^{\text{exp}}$ , when the type of the function is a standard floating type; ...

to:

5 Otherwise, the `frexp` functions return the value `x`, such that: `x` has a magnitude in the interval  $[1/2, 1)$  or zero, and `value` equals  $x \times 2^{\text{exp}}$ , when the type of the function is a standard or binary floating type; ...

In 7.12.6.6#2, change the first sentence from:

10 The `ldexp` functions multiply a floating-point number by an integral power of 2 when the type of the function is a standard floating type, or by an integral power of 10 when the type of the function is a decimal floating type.

to:

15 The `ldexp` functions multiply a floating-point number by an integral power of 2 when the type of the function is a standard or binary floating type, or by an integral power of 10 when the type of the function is a decimal floating type.

Change 7.12.6.6#3 from:

[3] The `ldexp` functions return  $x \times 2^{\text{exp}}$  when the type of the function is a standard floating type, or return  $x \times 10^{\text{exp}}$  when the type of the function is a decimal floating type.

to:

20 [3] The `ldexp` functions return  $x \times 2^{\text{exp}}$  when the type of the function is a standard or binary floating type, or return  $x \times 10^{\text{exp}}$  when the type of the function is a decimal floating type.

In 7.12.6.11#2, change the second sentence from:

If `x` is subnormal it is treated as though it were normalized; thus, for positive finite `x`,

$$1 \leq x \times b^{-\text{logb}(x)} < b$$

25 where  $b = \text{FLT\_RADIX}$  if the type of the function is a standard floating type, or  $b = 10$  if the type of the function is a decimal floating type.

to:

If `x` is subnormal it is treated as though it were normalized; thus, for positive finite `x`,

$$1 \leq x \times b^{-\text{logb}(x)} < b$$

30 where  $b = \text{FLT\_RADIX}$  if the type of the function is a standard floating type,  $b = 2$  if the type of the function is a binary floating type, or  $b = 10$  if the type of the function is a decimal floating type.

In 7.12.6.13#2, change the first sentence from:

The `scalbn` and `scalbln` functions compute  $x \times b^n$ , where  $b = \text{FLT\_RADIX}$  if the type of the function is a standard floating type, or  $b = 10$  if the type of the function is a decimal floating type.

to:

The `scalbn` and `scalbln` functions compute  $x \times b^n$ , where  $b = \text{FLT\_RADIX}$  if the type of the function is a standard floating type,  $b = 2$  if the type of the function is a binary floating type, or  $b = 10$  if the type of the function is a decimal floating type.

## 5 12.4 Encoding conversion functions

The functions in this subclause, together with the numerical conversion functions for encodings in clause 13, support the non-arithmetic interchange formats specified by IEC 60559.

### Change to C11 + TS18661-1 + TS18661-2:

After 7.12.11.7, add:

#### 10 7.12.11.7a The `encodefN` functions

##### Synopsis

```
[1] #define __STDC_WANT_IEC_60559_TYPES_EXT__
#include <math.h>
void encodefN(unsigned char * restrict encptr, const _FloatN *
15 restrict xptr);
```

##### Description

[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.

##### Returns

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

#### 25 7.12.11.7b The `decodefN` functions

##### Synopsis

```
[1] #define __STDC_WANT_IEC_60559_TYPES_EXT__
#include <math.h>
void decodefN (_FloatN * restrict xptr, const unsigned char *
30 restrict encptr);
```

##### Description

[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

[3] The `decodefN` functions return no value.

### 7.12.11.7c Encoding-to-encoding conversion functions

[1] An implementation shall declare a  $fMencfN$  function for each M and N equal 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 the width of a supported IEC 60559 arithmetic or non-arithmetic decimal interchange format.

#### 7.12.11.7c.1 The $fMencfN$ functions

##### Synopsis

```
[1] #define __STDC_WANT_IEC_60559_TYPES_EXT__
#include <math.h>
void fMencfN(unsigned char * restrict encMptr, 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.

#### 7.12.11.7c.2 The $dMencdecN$ and $dMencbindN$ functions

##### Synopsis

```
[1] #define __STDC_WANT_IEC_60559_TYPES_EXT__
#include <math.h>
void dMencdecN(unsigned char * restrict encMptr, const unsigned char
* restrict encNptr);
void dMencbindN(unsigned char * restrict encMptr, 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.

##### Returns

[3] These functions return no value.

## 13 Numeric conversion functions in `<stdlib.h>`

This clause specifies functions to convert between character sequences and the interchange and extended floating types. Conversions from character sequences are provided by functions analogous to the `strtod`

function in `<stdlib.h>`. Conversions to character sequences are provided by functions analogous to the `strfromd` function in `<stdlib.h>`.

This clause also specifies functions to convert between character sequences and IEC 60559 interchange format encodings.

## 5 Changes to C11 + TS18661-1 + TS18661-2:

After 7.22.1#1, insert

[3a] For each interchange or extended floating type that the implementation provides, `<stdlib.h>` shall declare the associated functions. Conversely, for each such type that the implementation does not provide, `<stdlib.h>` shall not declare the associated functions unless specified otherwise.

10 After 7.22.1.2b, insert:

### 7.22.1.2c The `strfromfN`, `strfromfNx`, `strfromdN`, and `strfromdNx` functions

#### Synopsis

```
[1] #define __STDC_WANT_IEC_60559_TYPES_EXT__
#include <stdlib.h>
15 int strfromfN(char * restrict s, size_t n, const char * restrict
    format, _FloatN fp);
    int strfromfNx(char * restrict s, size_t n, const char * restrict
        format, _FloatNx fp);
    int strfromdN(char * restrict s, size_t n, const char * restrict
20     format, _DecimalN fp);
    int strfromdNx(char * restrict s, size_t n, const char * restrict
        format, _DecimalNx fp);
```

#### Description

25 [2] The `strfromfN` and `strfromfNx` functions are similar to the `strfromd` function, except they convert to the types `_FloatN` and `_FloatNx`, respectively. The `strfromdN` and `strfromdNx` functions are similar to the `strfromd64` function, except they convert from the types `_DecimalN` and `_DecimalNx`, respectively.

#### Returns

30 [3] The `strfromfN` and `strfromfNx` functions return values similar to the `strfromd` function. The `strfromdN` and `strfromdNx` functions return values similar to the `strfromd64` function.

After 7.22.1.3a, insert:

### 7.22.1.3b The `strttofN`, `strttofNx`, `strtodN`, and `strtodNx` functions

#### Synopsis

```
[1] #define __STDC_WANT_IEC_60559_TYPES_EXT__
#include <stdlib.h>
35 _FloatN strttofN(const char * restrict nptr, char ** restrict
    endptr);
    _FloatNx strttofNx(const char * restrict nptr, char ** restrict
40     endptr);
    _DecimalN strtodN(const char * restrict nptr, char ** restrict
        endptr);
```

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

### Description

5 [2] The `strtofN` and `strtofNx` functions are similar to the `strtod` function, except they convert to the types `_FloatN` and `_FloatNx`, respectively. The `strtodN` and `strtodNx` functions are similar to the `strtod64` function, except they convert to the types `_DecimalN` and `_DecimalNx`, respectively.

### Returns

10 [3] The `strtofN` and `strtofNx` functions return values similar to the `strtod` function, except in the types `_FloatN` and `_FloatNx`, respectively. The `strtodN` and `strtodNx` functions return values similar to the `strtod64` function, except in the types `_DecimalN` and `_DecimalNx`, respectively.

## 7.22.1.3c String-to-encoding functions

15 [1] An implementation shall declare the `strtoencfN` function for each  $N$  equal the width of a supported IEC 60559 arithmetic or non-arithmetic binary interchange format. An implementation shall declare both the `strtoencdecN` and `strtoenbindN` functions for each  $N$  equal the width of a supported IEC 60559 arithmetic or non-arithmetic decimal interchange format.

### 7.22.1.3c.1 The `strtoencfN` functions

#### Synopsis

```
[1] #define __STDC_WANT_IEC_60559_TYPES_EXT__
    #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.

#### Returns

[3] These functions return no value.

### 7.22.1.3c.2 The `strtoencdecN` and `strtoenbindN` functions

#### Synopsis

```
[1] #define __STDC_WANT_IEC_60559_TYPES_EXT__
    #include <stdlib.h>
    void strtoencdecN(unsigned char * restrict encptr, const char *
        restrict nptr, char ** restrict endptr);
    void strtoenbindN(unsigned char * restrict encptr, const char *
        restrict nptr, char ** restrict endptr);
```

#### Description

[2] The `strtoencdecN` and `strtoenbindN` 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 implementation-defined.

## Returns

[3] These functions return no value.

### 7.22.1.3d String-from-encoding functions

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

#### 7.22.1.3d.1 The `strfromencfN` functions

##### Synopsis

```
[1] #define __STDC_WANT_IEC_60559_TYPES_EXT__
#include <stdlib.h>
int strfromencfN(char * restrict s, size_t n, 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. The order of bytes in the arrays is implementation-defined.

##### Returns

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

#### 7.22.1.3d.2 The `strfromencdecN` and `strfromencbindN` functions

##### Synopsis

```
[1] #define __STDC_WANT_IEC_60559_TYPES_EXT__
#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, 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 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 `strfromendecdN` and `strfromendcbindN` functions return the same values as corresponding `strfromdN` functions.

## 14 Complex arithmetic <complex.h>

- 5 This clause specifies complex functions for corresponding real types that are interchange and extended floating types.

### Changes to C11 + TS18661-1 + TS18661-2:

Change 7.3.1#3 from:

10 [3] Each synopsis specifies a family of functions consisting of a principal function with one or more `double complex` parameters and a `double complex` or `double` return value; and other functions with the same name but with `f` and `l` suffixes which are corresponding functions with `float` and `long double` parameters and return values.

to:

[3] Each synopsis specifies a family of functions consisting of:

15 a principal function with one or more `double complex` parameters and a `double complex` or `double` return value; and,

other functions with the same name but with `f`, `l`, `fN`, and `fNx` suffixes which are corresponding functions whose parameters and return values have corresponding real types `float`, `long double`, `_FloatN`, and `_FloatNx`.

20 Add after 7.3.1#3:

[3a] For each interchange or extended 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.

Add the following list of function prototypes to the synopsis of the respective subclauses:

25 7.3.5 Trigonometric functions

```
_FloatN complex cacosfN(_FloatN complex z);
_FloatNx complex cacosfNx(_FloatNx complex z);
```

30 `_FloatN complex casinfN(_FloatN complex z);`  
`_FloatNx complex casinfNx(_FloatNx complex z);`

```
_FloatN complex catanfN(_FloatN complex z);
_FloatNx complex catanfNx(_FloatNx complex z);
```

35 `_FloatN complex ccosfN(_FloatN complex z);`  
`_FloatNx complex ccosfNx(_FloatNx complex z);`

```
_FloatN complex csinfN(_FloatN complex z);
_FloatNx complex csinfNx(_FloatNx complex z);
```

40 `_FloatN complex ctanfN(_FloatN complex z);`  
`_FloatNx complex ctanfNx(_FloatNx complex z);`

## 7.3.6 Hyperbolic functions

```

5  _FloatN complex cacoshfN(_FloatN complex z);
   _FloatNx complex cacoshfNx(_FloatNx complex z);

   _FloatN complex casinhfN(_FloatN complex z);
   _FloatNx complex casinhfNx(_FloatNx complex z);

10  _FloatN complex catanhfN(_FloatN complex z);
   _FloatNx complex catanhfNx(_FloatNx complex z);

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

15  _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

```

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

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

```

## 7.3.8 Power and absolute value functions

```

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

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

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

```

## 7.3.9 Manipulation functions

```

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

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

   _FloatN complex CmplxfN(_FloatN x, _FloatN y);
   _FloatNx complex CmplxfNx(_FloatNx x, _FloatNx y);

45  _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);

```

## 5 15 Type-generic macros <tgmath.h>

The following changes to C11 + TS18661-1 + TS18661-2 enhance the specification of type-generic macros in <tgmath.h> to apply to interchange and extended floating types, as well as standard floating types.

### Changes to C11 + TS18661-1 + TS18661-2:

In 7.25, replace paragraphs [3b]:

- 10 [3b] If arguments for generic parameters of a type-generic macro are such that some argument has a corresponding real type that is of standard floating type and another argument is of decimal floating type, the behavior is undefined.

with:

- 15 [3b] 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.

In 7.25#3c, replace the bullets:

- First, if any argument for generic parameters has type `_Decimal128`, the type determined is `_Decimal128`.
- 20 — Otherwise, if any argument for generic parameters has type `_Decimal64`, or if any argument for generic parameters is of integer type and another argument for generic parameters has type `_Decimal32`, the type determined is `_Decimal64`.
- Otherwise, if any argument for generic parameters has type `_Decimal32`, the type determined is `_Decimal32`.
- 25 — Otherwise, if the corresponding real type of any argument for generic parameters is `long double`, the type determined is `long double`.
- Otherwise, if the corresponding real type of any argument for generic parameters is `double` or is of integer type, the type determined is `double`.
- 30 — Otherwise, if any argument for generic parameters is of integer type, the type determined is `double`.
- Otherwise, the type determined is `float`.

with:

- If two arguments 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.
- 35 — If any arguments for generic parameters have type `_DecimalM` where  $M \geq 64$  or `_DecimalNx` where  $N \geq 32$ , the type determined is the widest of the types of these arguments. If `_DecimalM` and `_DecimalNx` are both widest types (with equivalent sets of values) of these arguments, the type determined is `_DecimalM`.

- Otherwise, if any argument for generic parameters is of integer type and another argument for generic parameters has type `_Decimal132`, the type determined is `_Decimal164`.
- Otherwise, if any argument for generic parameters has type `_Decimal132`, the type determined is `_Decimal132`.
- 5 — Otherwise, if the corresponding real type of any argument for generic parameters has type `long double`, `_FloatM` where  $M \geq 128$ , or `_FloatNx` where  $N \geq 64$ , the type determined is the widest of the corresponding real types of these arguments. If `_FloatM` and either `long double` or `_FloatNx` are both widest corresponding real types (with equivalent sets of values) of these arguments, the type determined is `_FloatM`. Otherwise, if `long double` and `_FloatNx` are both widest corresponding real types (with equivalent sets of values) of these arguments, the type determined is `long double`.
- 10 — Otherwise, if the corresponding real type of any argument for generic parameters has type `double`, `_Float64`, or `_Float32x`, the type determined is the widest of the corresponding real types of these arguments. If `_Float64` and either `double` or `_Float32x` are both widest corresponding real types (with equivalent sets of values) of these arguments, the type determined is `_Float64`. Otherwise, if `double` and `_Float32x` are both widest corresponding real types (with equivalent sets of values) of these arguments, the type determined is `double`.
- 15 — Otherwise, if any argument for generic parameters is of integer type, the type determined is `double`.
- 20 — Otherwise, if the corresponding real type of any argument for generic parameters has type `_Float32`, the type determined is `_Float32`.
- Otherwise, the type determined is `float`.

In the second bullet 7.25#3c, attach a footnote to the wording:

the type determined is the widest

25 where the footnote is:

\*) The term widest here refers to a type whose set of values is a superset of (or equivalent to) the sets of values of the other types.

In 7.25#6, replace:

Use of the macro with any argument of standard floating or complex type invokes a complex function.  
 30 Use of the macro with an argument of decimal floating type results in undefined behavior.

with:

Use of the macro 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.

35 After 7.25#6c, add the paragraph:

[6d] For an implementation that provides the following real floating types:

type	IEC 60559 format
<code>float</code>	binary32
<code>double</code>	binary64
<code>long double</code>	binary128

```

    _Float32      binary32
    _Float64      binary64
    _Float128     binary128
    _Float32x     binary64
5   _Float64x     binary128

```

a type-generic macro `cbrt` that conforms to the specification in this clause and that is affected by constant rounding modes could be implemented as follows:

```

10   #if defined(__STDC_WANT_IEC_60559_TYPES_EXT__)
        #define cbrt(X)  _Generic((X),
                                \
                                _Float128: cbrtf128(X),
                                \
                                _Float64: cbrtf64(X),
                                \
                                _Float32: cbrtf32(X),
                                \
15   _Float64x: cbrtf64x(X),
                                \
                                _Float32x: cbrtf32x(X),
                                \
                                long double: cbrtl(X),
                                \
                                default: _Roundwise_cbrt(X),
                                \
                                float: cbrtf(X)
                                \
                                )
20   #else
        #define cbrt(X)  _Generic((X),
                                \
                                long double: cbrtl(X),
                                \
                                default: _Roundwise_cbrt(X),
                                \
25   float: cbrtf(X)
                                \
                                )
        #endif

```

where `_Roundwise_cbrt()` is equivalent to `cbrt()` invoked without macro-replacement suppression.

30 In 7.25#7, insert at the beginning of the example:

```
#define __STDC_WANT_IEC_60559_TYPES_EXT__
```

In 7.25#7, append to the declarations:

```

35   #if __STDC_IEC_60559_TYPES__ >= 201ymmL
        _Float32x f32x;
        _Float64 f64;
        _Float128 f128;
        _Float64x complex f64xc;
        #endif

```

40 In 7.25#7, append to the table:

	<code>cos(f64xc)</code>	<code>ccosf64x(f64xc)</code>
	<code>pow(dc, f128)</code>	<code>cpowf128(dc, f128)</code>
	<code>fmax(f64, d)</code>	<code>fmaxf64(f64, d)</code>
45	<code>fmax(d, f32x)</code>	<code>fmax(d, f32x)</code> , the function, if the set of values of <code>_Float32x</code> is a subset of (or equivalent to) the set of values of <code>double</code> , or <code>fmaxf32x(d, f32x)</code> , if the set of values of <code>double</code> is a proper subset of the set of values of <code>_Float32x</code> , or undefined, if neither of the sets of values of <code>double</code> and <code>_Float32x</code> is a subset of the other (and the sets are not equivalent)
50	<code>pow(f32x, n)</code>	<code>powf32x(f32x, n)</code>

## Bibliography

- [1] ISO/IEC 9899:2011, *Information technology — Programming languages, their environments and system software interfaces — Programming Language C*
- [2] ISO/IEC 9899:2011/Cor.1:2012, *Technical Corrigendum 1*
- 5 [3] ISO/IEC/IEEE 60559:2011, *Information technology — Microprocessor Systems — Floating-point arithmetic*
- [4] ISO/IEC TR 24732:2009, *Information technology – Programming languages, their environments and system software interfaces – Extension for the programming language C to support decimal floating-point arithmetic*
- 10 [5] IEC 60559:1989, *Binary floating-point arithmetic for microprocessor systems, second edition*
- [6] IEEE 754-2008, *IEEE Standard for Floating-Point Arithmetic*
- [7] IEEE 754-1985, *IEEE Standard for Binary Floating-Point Arithmetic*
- [8] IEEE 854-1987, *IEEE Standard for Radix-Independent Floating-Point Arithmetic*